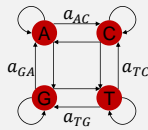


1

## BTRY 4840/6840, CS 4775 Computational Genetics and Genomics



September 13, 2018

2

## Announcements

- Problem set 2 out
  - Skeleton code that deals with input/output forthcoming

3

## Today's lecture

- Finish sequence alignment
  - Affine gaps
  - Optimizations
- Hidden Markov models (HMMs)
  - Example
  - Definition of Markov chain, HMM
  - Viterbi algorithm

4

## Affine gaps

5

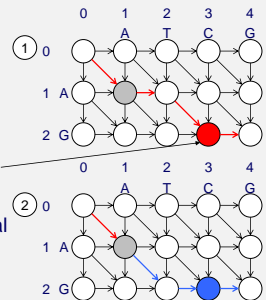
## Affine gap penalties

- A bit more difficult than linear  $\gamma(g) = -d - (g-1)e$
- Need to know whether a prior site is a gap or not
  - If previous alignment is match, gap opening penalty:  $-d$
  - If previous alignment is gap in  $x$ , gap extension penalty:  $-e$
  - If previous alignment is gap in  $y$ , gap extension penalty:  $-e$
- Proposal: inspect trace back pointer of prior site and use  $e$  if so,  $d$  if not **Does not work**
  - Only have gap at prior site if gap length  $\geq 1$  better than match
  - Need to know best possible path that includes gap at prior site in order to decide on whether to extend or start new gap
- Solution: track best path with gap at given site
  - Instead of  $F$ , have  $M$ ,  $I_x$ ,  $I_y$  – latter two always gapped

6

## Example: one matrix insufficient for affine gaps

- Consider two paths
  - Score for ① is  $F(1,1) - d + s_{G,C} - d$
  - Score for ② is  $F(1,1) + s_{G,T} - d - e$
  - If  $s_{G,C} > s_{G,T}$ , then  $F(2,3) = F(1,1) - d + s_{G,C}$
  - But if  $s_{G,C} - d < s_{G,T} - e$ , this assignment not optimal
- Does not “look forward” at  $F(2,3)$ . Need  $I_x$ ,  $I_y$



### Affine gap recurrence relations

(previous)

$M$	$I_x$	$I_y$	
$b\ b$	$b\ b$	$-b$	
$b\ b$	$-b$	$b\ b$	

$$M(i, j) = \max \begin{cases} M(i-1, j-1) + s_{x_i, y_j} \\ I_x(i-1, j-1) + s_{x_i, y_j} \\ I_y(i-1, j-1) + s_{x_i, y_j} \end{cases}$$

(current)

$I_x$	$b\ b$	$b\ b$	$\times$
	$b\ -$	$-$	
	$b\ b$	$-$	

$$I_x(i, j) = \max \begin{cases} M(i-1, j) - d \\ I_x(i-1, j) - e \end{cases}$$

$I_y$	$b\ -$	$\times$	$-$
	$b\ b$	$-$	$b\ b$

$$I_y(i, j) = \max \begin{cases} M(i, j-1) - d \\ I_y(i, j-1) - e \end{cases}$$

- Above is for Needleman-Wunsch
- Runtime complexity?  $O(nm)$ , but  $\sim 3\times$  as long

### Affine gap base cases

- Initialization:  $M(0, 0) = I_x(0, 0) = I_y(0, 0) = 0$
- Boundaries:
 
$$I_x(i, 0) = I_x(i-1, 0) - e, \quad 1 \leq i \leq n$$

$$I_y(0, j) = I_y(0, j-1) - e, \quad 1 \leq j \leq m$$

$$I_x(0, j) = -\infty, \quad 1 \leq j \leq m$$

$$I_y(i, 0) = -\infty, \quad 1 \leq i \leq n$$

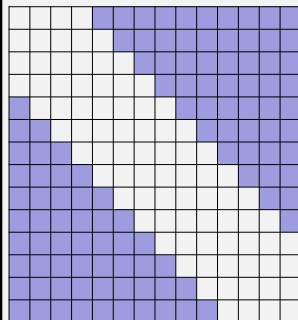
$$M(0, j) = M(i, 0) = -\infty$$
- Final score:
 
$$\max_a S(x, y, a) = \max(M(n, m), I_x(n, m), I_y(n, m))$$
 a an alignment

### Optimizations

### Efficiency in runtime, space usage

- Runtime of these alignment algorithms:
  - $O(nm)$
- Space usage:
  - $O(nm)$
- Can we do better?

### Bounded dynamic programming



- Key idea:  $F(i, j)$  defined only for bounded region
- Initialization:
  - $F(i, 0), F(j, 0)$  undefined outside bounded region
- Iteration for  $|i - j| < k$ :
 
$$F(i, j) = \max \begin{cases} F(i-1, j-1) + s_{x_i, y_j} \\ F(i-1, j) - d \text{ if } j-(i-1) < k \\ F(i, j-1) - d \text{ if } i-(j-1) < k \end{cases}$$
- Runtime, space requirements?
  - $O(nk)$

### Markov chain and Hidden Markov models

### Example: dishonest casino

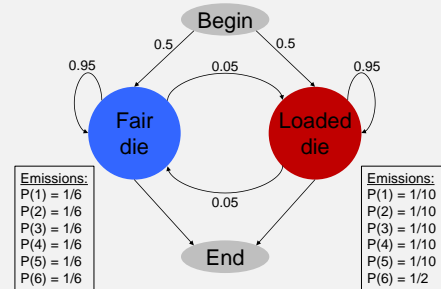
13

- Casino uses two dice:
  - Fair:  $P(1) = P(2) = P(3) = P(4) = P(5) = P(6) = 1/6$
  - Loaded:  $P(1) = P(2) = P(3) = P(4) = P(5) = 1/10$   
 $P(6) = 1/2$
- Casino switches dice on average every 20 rolls



### Markov model of dice

14



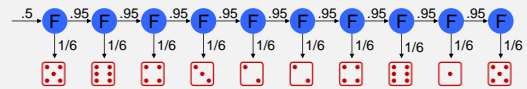
### Probability of fair die for series of die rolls

15

- Suppose we observe 10 die rolls:
- Problem: assuming all rolls are from the same die, was it fair or loaded?
  - Will evaluate likelihood ratio

### Calculating probabilities of hidden states

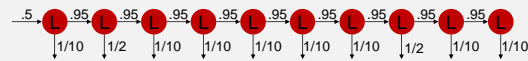
16



- What is the probability of  $z = (F, F, F, F, F, F, F, F, F, F)$  and  $x = (5, 6, 4, 3, 2, 2, 4, 1, 5)$
- $P(x, z) = P(z_1 = F)P(x_1 = 5|z_1 = F)P(z_2 = F|z_1 = F) \dots$   
 $= .5 \times (1/6)^{10} \times (.95)^9$   
 $= 5.2 \times 10^{-9}$
- Why is this probability so small? Any particular sequence is rare
- What is the probability I have a dollar bill with serial number C 00273441 G?
- Likelihood ratio comparing two models still very informative

### Calculating probabilities of hidden states

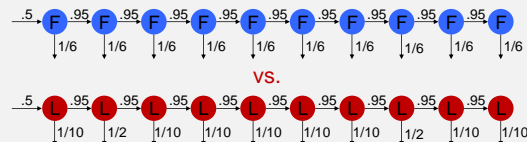
17



- What is the probability of  $z = (L, L, L, L, L, L, L, L, L, L)$  and  $x = (5, 6, 4, 3, 2, 2, 4, 1, 5)$
- $P(x, z) = P(z_1 = L)P(x_1 = 5|z_1 = L)P(z_2 = L|z_1 = L) \dots$   
 $= .5 \times (1/10)^8 \times (1/2)^2 \times (.95)^9$   
 $= 7.9 \times 10^{-10}$

### Compare the two possibilities

18



- Have  $P(x, \text{all fair}) = 5.2 \times 10^{-9}$   
 $P(x, \text{all loaded}) = 7.9 \times 10^{-10}$
- Likelihood ratio:  $\frac{P(x, \text{all fair})}{P(x, \text{all loaded})} = 6.58$
- Nearly 7x more likely that the die was fair!

### Comparison for different sequence

19

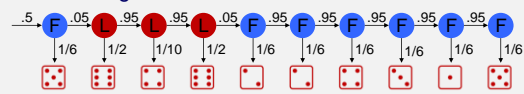


- For different sequence:  $x = (6, 2, 3, 6, 6, 5, 6, 4, 6, 1)$
- Then  $P(x, \text{all fair}) = 5.2 \times 10^{-9}$  (same as before)  
 $P(x, \text{all loaded}) = .5 \times (1/10)^5 \times (1/2)^5 \times (.95)^9$   
 $= 9.8 \times 10^{-8}$
- Likelihood ratio:  $\frac{P(x, \text{all loaded})}{P(x, \text{all fair})} = 18.8$ 
  - Very likely to be loaded, unlikely to be fair

### Likelihood of die being switched

20

- Now using different  $z$ :



- Probability of  $z = (F, L, L, L, F, F, F, F, F, F)$
- $P(x, z) = P(z_1 = F)P(x_1 = 5|z_1 = F)P(z_2 = L|z_1 = F) \dots$   
 $= .5 \times (1/6)^7 \times (1/10)^3 \times (1/2)^2 \times (.95)^7 \times (.05)^2$   
 $= 7.8 \times 10^{-11}$
- This state path unlikely, but
  - What about others state paths  $z$ ?
  - Other observed outcomes  $x$ ?

Want general way  
to analyze data  
given HMM  
specification

### Desired uses of HMMs

21

- Evaluation:**
  - Given:** observed  $x$  and HMM specification
  - Question:** what is the joint probability of  $x$  and a given  $z$ ?
  - Question:** what is the likelihood of  $x$  based on the HMM?
- Decoding:**
  - Given:** observed  $x$  and HMM
  - Question:** what sequence of hidden states produced  $x$ ?
  - Viterbi decoding:** most likely hidden state sequence
  - Posterior probability of hidden states:** probability of each state  $z_i$  producing each  $x_i$ 
    - Technically not a decoding: not path of states, but probabilities
- Learning:**
  - Given:** observed  $x$  and HMM without complete probabilities
  - Question:** what emission, transition probabilities produced  $x$ ?

### Markov chain (model) definition

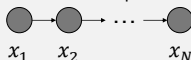
22

- Markov chain:** sequence of states at given times
  - In genomics, we often think of "time" as *position* on sequence
- Formally, Markov chain defined by**
  - Set of states  $S$ : individual state at given time  $i$  denoted  $x_i$
  - Matrix  $A$  of state transitions probabilities  
 Element  $A[k, l] = a_{kl} = P(x_i = l | x_{i-1} = k)$
  - Initial state probabilities  $a_{0k} = P(x_1 = k)$
  - (Formulations vary: states sometimes states denoted  $\pi_i$ )
- Key feature:** all states *observed* (not so in HMMs)
- Markov property:** state  $x_i$  depends only on  $x_{i-1}$ , so
  - $P(x) = P(x_1, x_2, \dots, x_N)$   
 $= P(x_N | x_{N-1})P(x_{N-1} | x_{N-2}) \dots P(x_2 | x_1)P(x_1)$

### Graphical model representation of Markov chain

23

- Using graphical model notation:
  - Each node is a random variable
  - Shaded nodes are observed
  - White (unfilled) nodes are hidden
  - Arrows represent conditional dependence



- Above implies  
 $P(x_1, x_2, \dots, x_N) = P(x_1)P(x_2|x_1) \dots P(x_N|x_{N-1})$   
 $= P(x_1) \prod_{i=2}^N P(x_i|x_{i-1})$

- This is 1<sup>st</sup> order Markov model:  $x_i$  depends on  $x_{i-1}$ 
  - For  $K^{\text{th}}$  order Markov model,  $x_i$  depends on  $x_{i-K}, \dots, x_{i-1}$

### HMM notation (somewhat different in Durbin)

24

- States:**  $z = (z_1, z_2, \dots, z_L)$ , for  $L$  observations
  - Examples:  $z = (F, F, F, F, F, F, F, F, F, F)$  – rolls from fair die  
 $z = (L, L, L, L, L, L, L, L, L, L)$  – rolls from loaded die
- Observations:**  $x = (x_1, x_2, \dots, x_L)$ 
  - Example:  $x = (5, 6, 4, 3, 2, 2, 4, 1, 5)$
- Transition probabilities**  $a_{kl} = P(z_i = l | z_{i-1} = k)$ 
  - Probability of moving to state  $l$  from previous state  $k$
  - Example:  $a_{FL} = P(z_i = L | z_{i-1} = F) = 0.05$
- Emission probabilities**  $e_k(b) = P(x_i = b | z_i = k)$ 
  - Probability of observing  $b$  given current state  $k$
  - Example:  $e_L(6) = P(x_i = 6 | z_i = L) = 1/2$
- Initial probabilities**  $a_{0k} = P(z_1 = k)$

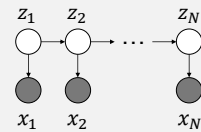
## Hidden Markov model definition

25

- Formally, Hidden Markov model defined by
  - Set of states  $S$ : state at time/position  $i$  denoted  $z_i$
  - Set of possible observations  $O$ : observation  $i$  denoted  $x_i$
  - Matrix  $A$  of state transitions probabilities  
Element  $A[k, l] = a_{kl} = P(z_i = l | z_{i-1} = k)$
  - Initial state probabilities  $a_{0k} = P(z_1 = k)$
  - Emission probabilities  $e_k(b) = P(z_i = b | z_i = k)$
- Key features:
  - States hidden (unobserved)
  - Markov property holds: state  $z_i$  depends only on  $z_{i-1}$ 
    - (For first order HMM: most common type)
  - Observation  $x_i$  depends only on current state  $z_i$

## Graphical model representation of HMM

26



So:

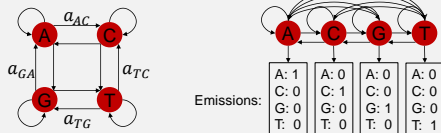
$$P(\mathbf{x}, \mathbf{z}) = P(x_1 | z_1) P(z_1) \prod_{i=2}^N P(x_i | z_i) P(z_i | z_{i-1})$$

Note: in practice will want to evaluate  
 $P(\mathbf{x}) = \sum_{\mathbf{z}} P(\mathbf{x}, \mathbf{z})$  and/or  $\max_{\mathbf{z}} P(\mathbf{x}, \mathbf{z})$   
 (i.e., where  $\mathbf{z}$  unknown)

## Genomic Markov chain, HMM

27

- Can model sequence using Markov chain, HMM



## Hidden Markov models CpG island example

28

## Biology background: CpG dinucleotide

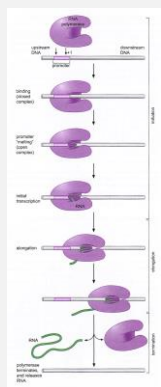
29

- Dinucleotide sequence CG is typically written CpG
  - p for phosphate (between bases in DNA backbone): emphasizes this is dinucleotide, not base pairing
  - Can also talk about GC content: % G or C nucleotides in region, not dinucleotides
- CpG dinucleotides:
  - Cytosine in CpGs are often methylated
  - When methylated, have high rate of C→T mutations
    - Methylation: addition of methyl group (in this case to cytosine)
    - In mammals, 70-80% of CpG cytosines are methylated
  - Consequently CpGs are rarer in genome than expected

## Gene promoter

30

- Gene promoter is:
  - Sequence where transcription is initiated
    - May or may not be transcribed, but near transcription start site (TSS)
  - Between ~100-1000 bp long
  - Subsequence bound by transcription factors:
    - A protein (i.e., product of a gene) that binds a specific DNA sequence
    - Recruits RNA polymerase, and thus controls the rate of transcription
      - In eukaryotes, often works in tandem with other elements (activators, repressors, others)



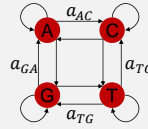
### Biology background: CpG islands

31

- Methylation is suppressed in promoters, other regions of the genome
- Such regions have high rate of: CpG, GC content
  - Called CpG islands
- Problems:
  - Given short sequence, is it from a CpG island?
  - Given (long) genome sequence, locate CpG islands (regions with likely biological importance)
- How would you address the first question?

### Can train Markov chains for CpG island / not

32



- Suppose we were given data labeled as *CpG island* or *not CpG island*
- Could use this to train a model
  - That is, set the parameters in the model based on the data

### Final notes

33

- Summary – sequence alignment
  - Affine gaps possible, more work than linear
  - Optimize runtime via bounded dynamic programming
- Summary – hidden Markov models (HMMs):
  - Flexible tool to evaluate wide range of data (including genetic)
  - Can use to evaluate, decode, and learn parameters of HMMs
  - More detailed HMM example: CpG islands