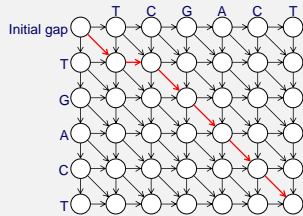


1

BTRY 4840/6840, CS 4775

Computational Genetics and Genomics

TCGACT
T-GACT

September 11, 2018

2

Announcements

- Problem set 1 due tonight at midnight
 - Grace period until 2am
- Problem set 2 out today, due in two weeks
 - Dynamic programming
 - Approximating distributions
 - Sequence alignment
- Make sure you are:
 - Signed up for Piazza
 - On CMS (for submitting problem sets)
- Collaboration:
 - Allowed, but don't study each other's code
 - Can discuss concepts: otherwise you may not learn as much

3

Today's lecture

- Sequence alignment
 - Global alignment via Needleman-Wunsch algorithm
 - Local alignment via Smith-Waterman algorithm
 - Affine gaps
 - Optimizations

4

Recall: traveling salesman

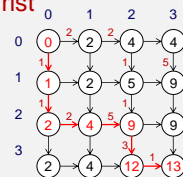
- Problem: given a set of places (cities) determine the cycle that visits all places with minimum distance using some metric
- Why fundamentally different than Manhattan Tourist?
 - May be best to choose path between some cities that is less optimal locally but leads to global optimum
- Can't solve sub-problem of few cities to build path
 - Breaking into sub-problems essential to dynamic programming



5

Solving sub-problem works for
Manhattan Tourist

- What value is stored in each cell in this problem?
 - Maximum weight to arrive at that position
- Why is solving the sub-problem
 - the optimal path to each node – sufficient for finding the global optimum in this problem?
 - Choices made about the path to a given node have no effect on optimal path thereafter, so best to choose maximum
 - When reasoning about later positions, only need the weight of earlier nodes: later positions independent of earlier path
- Affine gaps (later) somewhat more involved

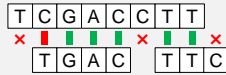


6

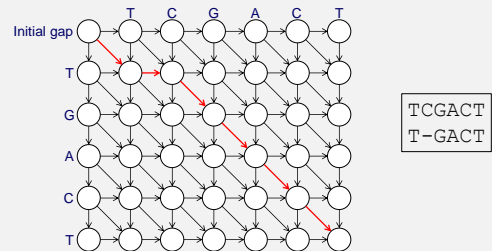
Global sequence alignment:
Needleman-Wunsch algorithm

Sequence alignment problem

- Goal: infer mutations, insertion/deletions, and identical positions between two sequences
 - Can decide whether the alignment is biologically sensible using score of alignment
 - Mutations are also called substitutions

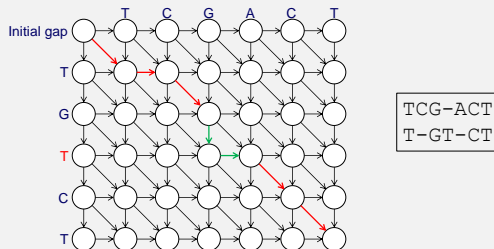


Can encode alignment using graph



Diagonal \Rightarrow bases on row and column aligned
 Horizontal \Rightarrow gap in left side sequence
 Vertical \Rightarrow gap in top sequence

Can encode alignment using graph, example 2



Diagonal \Rightarrow bases on row and column aligned
 Horizontal \Rightarrow gap in left side sequence
 Vertical \Rightarrow gap in top sequence

Defining scoring metric

- Need way to score potential alignments
 - Want to define score rigorously
 - Convenient for it to be additive
- Given two sequences x, y , both length n , assume each base (nucleotide) evolves independently
 - Independence not true for all sequences, but OK most of the time

Scoring: probabilistic foundations

- Assuming x and y are homologous – i.e., descend from shared ancestral sequence – can consider

$$P(x, y | \text{homologous}) = \prod_{i=1}^n P(x_i, y_i) = \prod_{i=1}^n \phi_{x_i, y_i}$$

- If x and y are non-homologous, bases are independent, so we have

$$P(x, y | \text{non-homologous}) = \prod_{i=1}^n P(x_i)P(y_i) = \prod_{i=1}^n \pi_{x_i} \pi_{y_i}$$

- Can estimate above using empirical training data
- Compute likelihood ratio (odds) of homology as

$$\frac{P(x, y | \text{homologous})}{P(x, y | \text{non-homologous})} = \prod_{i=1}^n \frac{\phi_{x_i, y_i}}{\pi_{x_i} \pi_{y_i}}$$

Scoring: probabilistic foundations, continued

- Move to log space

$$\log \left(\frac{P(x, y | \text{homologous})}{P(x, y | \text{non-homologous})} \right) = \sum_{i=1}^n \log \frac{\phi_{x_i, y_i}}{\pi_{x_i} \pi_{y_i}}$$

- Let $s_{a,b} = \log \frac{\phi_{a,b}}{\pi_a \pi_b}$, $S(x, y) = \sum_{i=1}^n s_{x_i, y_i}$
- $S(x, y)$: relative likelihood of x, y homologous vs not
- The $s_{a,b}$ scores:
 - Are additive across bases
 - Form a score matrix (substitution matrix) of size $|\mathcal{A}|^2$, where \mathcal{A} is the alphabet (4 nucleotides / 20 amino acids)
- More on generating score matrices soon

Gap penalties needed

13

- Gaps: insertions and deletions – indels
 - In practice don't know whether insertion or deletion occurred since we don't have ancestral sequence
- We allow gaps in alignment with a penalty
- If gap penalty is low, will end up with many gaps
 AC-TT-C-A-GG
 -CCT-ACG-C--
- If gap penalty is high, many substitutions
 ACTTCAGG
 CCTACGC-
- Alignment is about deciding where to place gaps

Gap penalty functions

14

- Let g be the length of a gap
- Have two common gap penalty functions
 - Linear: fixed cost for each gapped base
 - Each base in gap is equally penalized by $-d$
 - $\gamma(g) = -gd$
 - Affine: more elaborate
 - First base in gap penalized by $-d$
 - Subsequent bases penalized by (lesser) extension penalty $-e$
 - $\gamma(g) = -d - (g - 1)e$
- More complex / empirical models exist, but these are computationally expensive; less common

Edge weights for Needleman-Wunsch

15

- What is diagonal edge weight from $(i-1, j-1)$ to (i, j) ?
 Want log-odds of matching
 – Answer: $s_{x_i y_j}$
- Horizontal edge weight from $(i, j-1)$ to (i, j) ?
 – Answer: $-d$
- Vertical edge weight from $(i-1, j)$ to (i, j) ?
 – Answer: $-d$
- How should we define score of cell (i, j) ?
 - Recursively: maximum over optimal score to neighboring cells
 - Defines and solves sub-problem: what is the sub-problem?
 - Optimal path to given cell

Needleman-Wunsch: recurrence

16

- Define $F(i, j)$ as the optimal score at position (i, j) , calculated as:

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

- Base cases:

$$F(0, 0) = 0$$

$$F(i, 0) = F(i-1, 0) - d$$

$$F(0, j) = F(0, j-1) - d$$

Example amino acid alignment

17

	y									
	H	E	A	G	A	W	G	H	E	E
0	0	-8	-16	-24	-32	-40	-48	-56	-64	-80
P	-8	-2	-9	-17	-25	-33	-42	-49	-57	-73
A	-16	-10	-3	-4	-12	-20	-28	-36	-44	-60
W	-24	-18	-11	-6	-7	-15	-5	-13	-21	-37
H	-32	-14	-18	-13	-8	-9	-13	-7	-3	-11
E	-40	-22	-8	-16	-16	-9	-12	-15	-7	-5
A	-48	-30	-16	-3	-11	-11	-12	-12	-15	-5
E	-56	-38	-24	-11	-6	-12	-14	-15	-12	-9

HEAGAWGHE-E
 --P-AW-HEAE

Durbin et al. (2006)

Pseudo code including trace back pointers

18

alignNW(x, y, d)

Input: sequences $x = (x_1, \dots, x_n)$ and $y = (y_1, \dots, y_m)$, linear gap penalty d

Output: Best alignment: $\text{argmax}_a S(x, y, a)$

$F(0, 0) = 0$, $T(0, 0) = \text{null}$

for $i = 1$ to n

$F(i, 0) = F(i-1, 0) - d$, $T(i, 0) = u_i$

for $j = 1$ to m

$F(0, j) = F(0, j-1) - d$, $T(0, j) = l_j$

for $i = 1$ to n

for $j = 1$ to m

$v_g = F(i-1, j-1) + s_{x_i y_j}$

$v_u = F(i-1, j) - d$

$v_l = F(i, j-1) - d$

$F(i, j) = \max_{z \in \{g, u, l\}} v_z$

$T(i, j) = \arg \max_{z \in \{g, u, l\}} v_z$

$a = \text{traceback}(x, y, T)$

How does traceback function work?

19

- Where do we begin? $(i, j) = ?$
 $(i, j) = (n, m)$
 - Next step?
 - Inspect $T(i, j)$
 - If $T(i, j) = g$:
 - Alignment contains $x[i], y[j]$;
 - Update $i = i - 1, j = j - 1$;
 - If $T(i, j) = l$:
 - Alignment contains $gap -, y[j]$;
 - Update $j = j - 1$
 - If $T(i, j) = u$:
 - Alignment contains $x[i], gap -$
 - Update $i = i - 1$
- Loop until $i = 0, j = 0$

Local sequence alignment: Smith-Waterman algorithm

20

Local alignments

21

- Local alignment is substring in two sequences that is very similar, shows evidence of homology
- Why do we need local alignments?
 - Two sequences may be highly diverged but contain a conserved region we wish to identify
 - May want to look for short sequence in large chromosome

Smith-Waterman algorithm

22

- Expected score of a random match is negative, since it is non-homologous
- Add new option to recurrence – start of local match:

$$F(i, j) = \max \begin{cases} 0 \\ F(i-1, j-1) + s_{x_i y_j} \\ F(i-1, j) - d \\ F(i, j-1) - d \end{cases}$$

$$F(i, j) = 0, \quad i = 0 \text{ or } j = 0$$

- Can have start / end positions throughout matrix
 - For best alignment, start traceback from $\max_{i,j} F(i, j)$
 - For all local alignments greater than threshold t , do trace back from all i, j where $F(i, j) > t$

Smith-Waterman example

23

	H	E	A	G	A	W	G	H	E	E
P	0	0	0	0	0	0	0	0	0	0
A	0	0	0	5	0	5	0	0	0	0
W	0	0	0	0	2	0	20	12	4	0
H	0	10	2	0	0	0	12	18	22	14
E	0	2	16	8	0	0	4	10	18	28
A	0	0	8	21	13	5	0	4	10	20
E	0	0	6	13	18	12	4	0	4	16

AWGHE
AW-HE

Durbin et al. (2006)

Summary: global and local sequence alignment

24

- Global alignment recurrence

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

- Trace back from $F(n, m)$ – gives alignment of full sequence

- Local alignment recurrence

$$F(i, j) = \max \begin{cases} 0 \\ F(i-1, j-1) + s_{x_i y_j} \\ F(i-1, j) - d \\ F(i, j-1) - d \end{cases}$$

- Trace back from either $\max_{i,j} F(i, j)$ or all i, j where $F(i, j) > t$

Trace back end points

25

Global alignment

	H	E	A	G	A	W	G	H	E	E
P	0	-4	-8	-12	-16	-20	-24	-28	-32	-36
A	-4	0	-2	-6	-10	-14	-18	-22	-26	-30
W	-8	-2	0	-10	-14	0	-18	-22	-26	-30
H	-12	-6	-4	-14	-18	-10	0	-22	-26	-30
E	-16	-10	-8	-18	-22	-14	-2	0	-26	-30
E	-20	-14	-12	-22	-26	-18	-6	-4	0	-30
E	-24	-18	-16	-26	-30	-22	-10	-8	-2	0
E	-28	-22	-20	-30	-34	-26	-14	-12	-6	-2
E	-32	-26	-24	-34	-38	-30	-18	-16	-10	-6
E	-36	-30	-28	-38	-42	-34	-22	-20	-14	-10

Where should trace back terminate in each case?

- Global: at position (0,0) – aligns full sequence
- Local: at first i, j encountered where $F(i, j) = 0$

Durbin *et al.* (2006)

Local alignment

	H	E	A	G	A	W	G	H	E	E
P	0	0	0	0	0	0	0	0	0	0
A	0	0	0	0	0	0	0	0	0	0
W	0	0	0	0	0	0	0	0	0	0
H	0	0	0	0	0	0	0	0	0	0
E	0	0	0	0	0	0	0	0	0	0
E	0	0	0	0	0	0	0	0	0	0
E	0	0	0	0	0	0	0	0	0	0
E	0	0	0	0	0	0	0	0	0	0
E	0	0	0	0	0	0	0	0	0	0
E	0	0	0	0	0	0	0	0	0	0

Affine gaps

26

Affine gap penalties

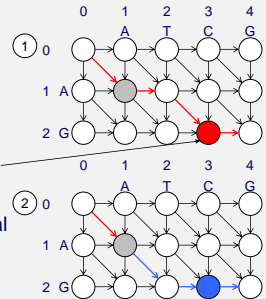
27

- 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 ≥ 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

Example: one matrix insufficient for affine gaps

28

- 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



Final notes

29

- Summary – sequence alignment:
 - Several dynamic programming algorithms for optimal global, local alignments
 - Affine gaps possible, more work than linear
 - Optimize runtime via bounded dynamic programming
- Problem set 2 out today
- Readings:
 - Durbin *et al.* chap. 2: sequence alignment
 - Durbin *et al.* chap. 3: hidden Markov models