

CS 598 AGB, Pairwise alignment

Tandy Warnow

February 15, 2016

CS 598 AGB, Pairwise alignment

Tandy Warnow

True Pairwise Alignment

- ▶ Suppose X and Y are two sequences, and X evolves into sequence Y via insertions, deletions, and substitutions.
- ▶ The true pairwise alignment of X and Y represents this true history.
- ▶ Examples:
 - ▶ AAT evolves into ACCAT by the insertion of CC
 - ▶ ATGA evolves into ATTAG by changing G to T, and then adding G
 - ▶ CTAA evolves into CTTAA by inserting a T.

Questions:

1. What are the pairwise alignments?
2. How can we guess at these evolutionary histories (and so pairwise alignment)?

Edit distance

Suppose each event (insertion, deletion, and substitution) costs 1.
Can we compute the minimum cost edit transformation between two sequences?

Computing the edit distance

Input: sequences s_1 and s_2

Output: minimum number of indels and substitutions needed to transform s_1 to s_2 .

A two-dimensional matrix, $m[0..|s_1|, 0..|s_2|]$ is used to hold the edit distance values:

$$m[i, j] = d(s_1[1..i], s_2[1..j])$$

$$m[0, 0] = 0$$

$$m[i, 0] = i, i = 1..|s_1|$$

$$m[0, j] = j, j = 1..|s_2|$$

$$\text{For } i, j \geq 1, m[i, j] = \min\{\begin{array}{l} m[i-1, j-1] + \text{if } s_1[i]=s_2[j] \text{ then } 0 \text{ else } 1, \\ m[i-1, j] + 1, \\ m[i, j-1] + 1 \end{array}\}$$

Dynamic programming algorithm for minimum edit distance

Compute the matrix from the bottom up!

Running time is $O(KL)$ where the first sequence has length K and the second sequence has length L .

Example: *ACAT* and *CCGT*

Finding the actual transformation - backtracing

To find the actual transformation, use backtracing.

Example: *ACAT* and *CCGT*

Extensions

Minimum cost approaches:

- ▶ How would you modify this algorithm if indels cost C and substitutions cost C' ?
- ▶ How would you modify the edit distance algorithm if a indel of length p has cost $C + C'p$?

Maximize similarity approaches:

- ▶ How would you modify the algorithm to maximize score, so matches have value 1 and mismatches and indels each have negative value -1 (i.e., they cost)?
- ▶ How would you modify the algorithm to maximize score, where matches and mismatches have scores (possibly negative) that depend on the pair of letters, and indels all have negative scores?
- ▶ How would you modify the algorithm if you want a *local alignment* (maximize cut off some prefix and suffix)

More generally,

- ▶ How would you align two alignments?

Extending to multiple alignment

Now that we have a way of defining the “cost” of a pairwise alignment, we can extend to a set of three or more sequences in at least two ways:

- ▶ Sum-of-pairs score (find MSA to minimize the sum of costs of all induced pairwise scores)
- ▶ Treelength (find tree and sequences at all nodes of the tree to minimize the sum of costs on the edges of the tree)

Both problems are NP-hard.

Treelength

Finding the sequences at internal nodes of a fixed tree to find the minimum cost is itself NP-hard.

Methods for Treelength include POY and BeeTLe.

Very computationally intensive (worse than MP) ...and also controversial!

Progressive Alignment

- ▶ Given sequences S , find rooted tree (somehow)
- ▶ Align sequences from the bottom up. Note this requires the ability to align two alignments.
- ▶ Return the alignment defined at the root.