

CSE 566 Spring 2023

Sequence Alignment

Instructor: Mingfu Shao

non-crossing Sequence Alignment

- Optimization problem: given two string S and T, find an alignment A such that cost(A) is minimized.

S = A C C G A C T A C T C G T T G A C

T = A C C G G A C C T C G T T G T C

S = A C C - G A C T A C T C G T T G A C

| | | | | | | | | | | | |

T = A C C G G A C - - C T C G T T G T C

A C C - G A C - T A C - T C G T T G A C

| | | | | | | | | | | |

A C C G G A C C T C G T T G - T - - - C

Cost Functions

- Cost of mismatch(substitution)

- Unit cost: each cost 1

- Arbitrary cost for different x-y $A-T: 3$, $G-A: 4$

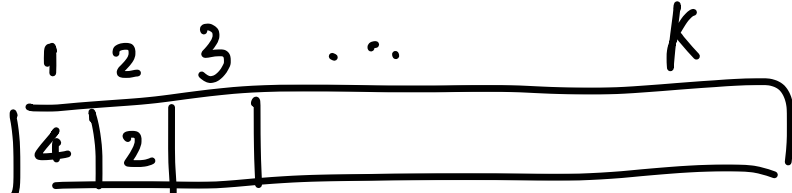
- Cost of gap (function of length of a gap)

- Unit cost: $\text{gap}(k) = k$

- Affine gap cost: $\text{gap}(k) = \overset{20}{\text{gap-open}} + k * \overset{3}{\text{gap-extend}}$

- Convex gap cost: for any $a < b < c < d$,
 $\text{gap}(d-a) + \text{gap}(c-b) \geq \text{gap}(c-a) + \text{gap}(d-b)$

- General cost: arbitrary function of k



The diagram shows a horizontal bar representing a gap of total length k . The bar is divided into segments by vertical lines. Above the segments are labels for their lengths: 1, 2, 3, ..., k . Below the first three segments are labels for their costs: 4, 2, and an empty box. The remaining segments are not labeled with costs.

Alignment with Unit Cost

- Edit distance problem: the smallest number of ^{mismatch} substitutions, insertions and deletions that can transform S into T.

S = A C C - G A C T A C T C G T T G A C
 | | | | | | | | | | | | | |
 T = A C C G G A C - - C T C G T T G T C

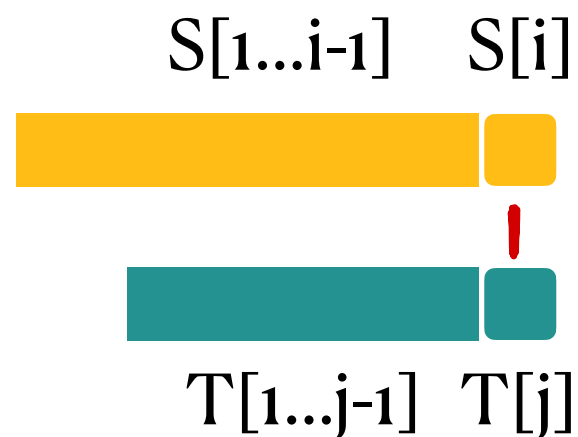
cost = 4

A C C - G A C - T A C - T C G T T G A C
 | | | | | | | | | | | | |
 A C C G G A C C T C G T T G - T - - - C

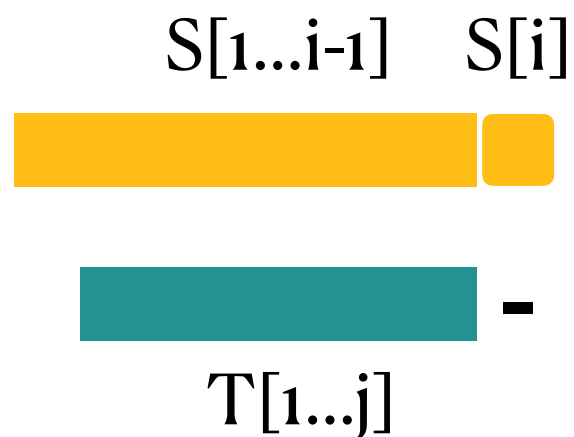
Dynamic Programming Algorithm

$m \times n$ subproblems. $m = |S|$, $n = |T|$ $1 \leq i \leq m, 1 \leq j \leq n$

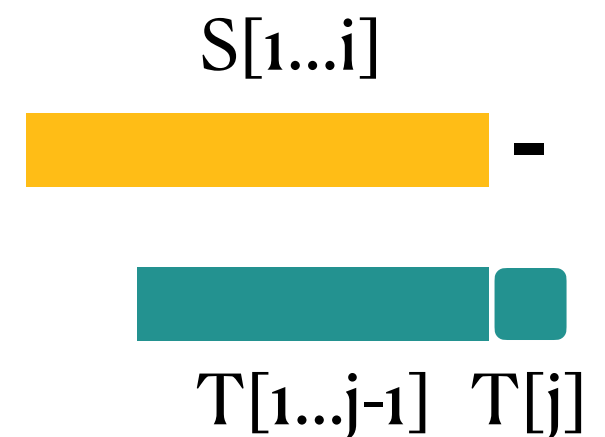
- Subproblem: $OPT(i, j)$ as the edit distance of $S[1..i]$ and $T[1..j]$.
- Consider how $S[i]$ and $T[j]$ are aligned in the optimal alignment:



$$OPT(i, j) = OPT(i-1, j-1) + \underline{cost(S[i], T[j])}$$



$$OPT(i, j) = OPT(i-1, j) + 1$$



$$OPT(i, j) = OPT(i, j-1) + 1$$

Recurrence

$$OPT(i, j) = \min \begin{cases} OPT(i-1, j-1) + cost(S[i], T[j]) \\ OPT(i-1, j) + 1 \\ OPT(i, j-1) + 1 \end{cases}$$

$$cost(S[i], T[j]) = \begin{cases} 1 & \text{if } S[i] \neq T[j] \\ 0 & \text{if } S[i] = T[j] \end{cases}$$

The Algorithm

- Step 1: Initialization
- Step 2: fill up the table following the recurrence
- Step 3: OPT[m,n] gives the edit distance of S and T
- Step 4: backtrace to find the optimal alignment

$T = \text{A C C G T}$

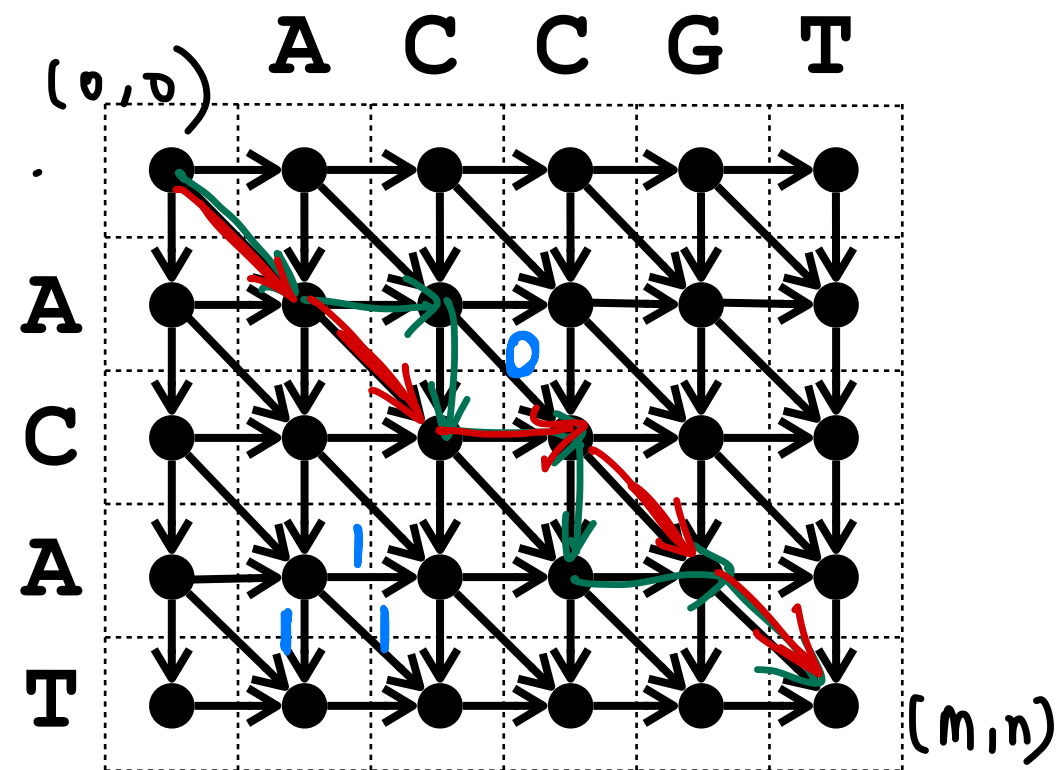
j

S	0	1	2	3	4	5
A	1					
C	2					
A	3					
T	4					

i

$\rightarrow \text{OPT}(m,n)$

Connection with Shortest Path



- ~~is~~ ^(DAG) directed acyclic graph, with mn vertices and $O(mn)$ edges
- One-to-one correspondence between alignments and $(0,0)$ - (m,n) paths
- $\text{cost}(\text{alignment}) = \text{length}(\text{path})$
- Algorithm to find the $(0,0)$ - (m,n) shortest path: $O(E) = O(mn)$.

S = A - C - A - T
 T = A C - C - G T

S = A C - A T
 T = A C C G T

Better Algorithms & Bounds

- Four-Russians algorithm: $O(n^2/\log n)$ $\frac{m \approx n}{n = \max\{|S|, |T|\}}$
- Meyer's algorithm (aka the wavefront algorithm): $O(nd)$, where d is the edit distance
- The edit distance cannot be computed in $O(n^{2-\epsilon})$ time, for any $\epsilon > 0$, unless the strong exponential time hypothesis is false.

Alignment with Affine Gap

A	A	A	G	A	A	T	T	C	G	A
A	-	A	-	A	-	T	-	C	T	A

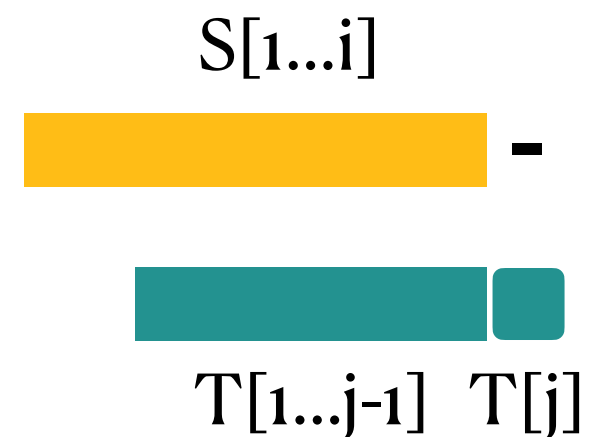
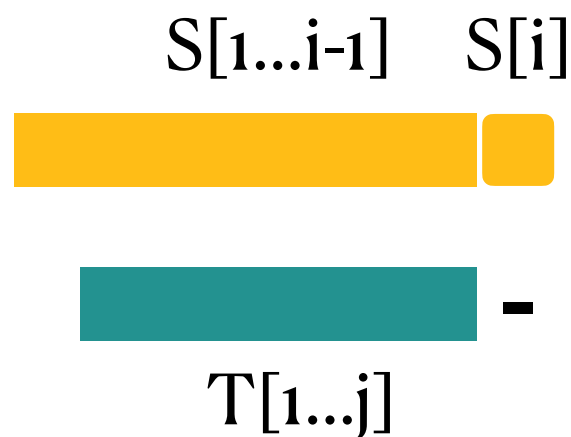
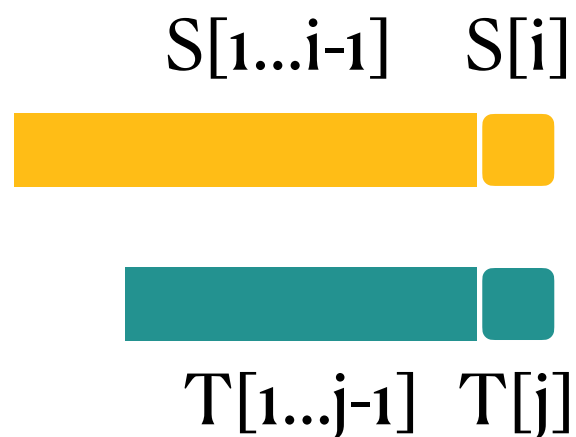
A	A	A	G	A	A	T	T	C	G	A
A	A	A	-	-	-	-	-	T	C	T

- They have the same cost under unit cost, but the 2nd alignment is more plausible.
- Affine gap cost: gap-open + k * gap-extension

$$20 \quad k * 3$$

Does the Previous DP work?

- Define OPT(i, j) as the minimized cost of S[1...i] and T[1...j].
- Consider how S[i] and T[j] are aligned in the optimal alignment:



$$\text{OPT}(i, j) =$$

$$\text{OPT}(i-1, j-1)$$

$$+ \text{cost}(S[i], T[j])$$

$$\text{OPT}(i, j) =$$

$$\text{OPT}(i-1, j)$$