

Sequence Alignment Problem

Dynamic Programming

A DNA strand consists of a string of molecules called bases:

Adenine A

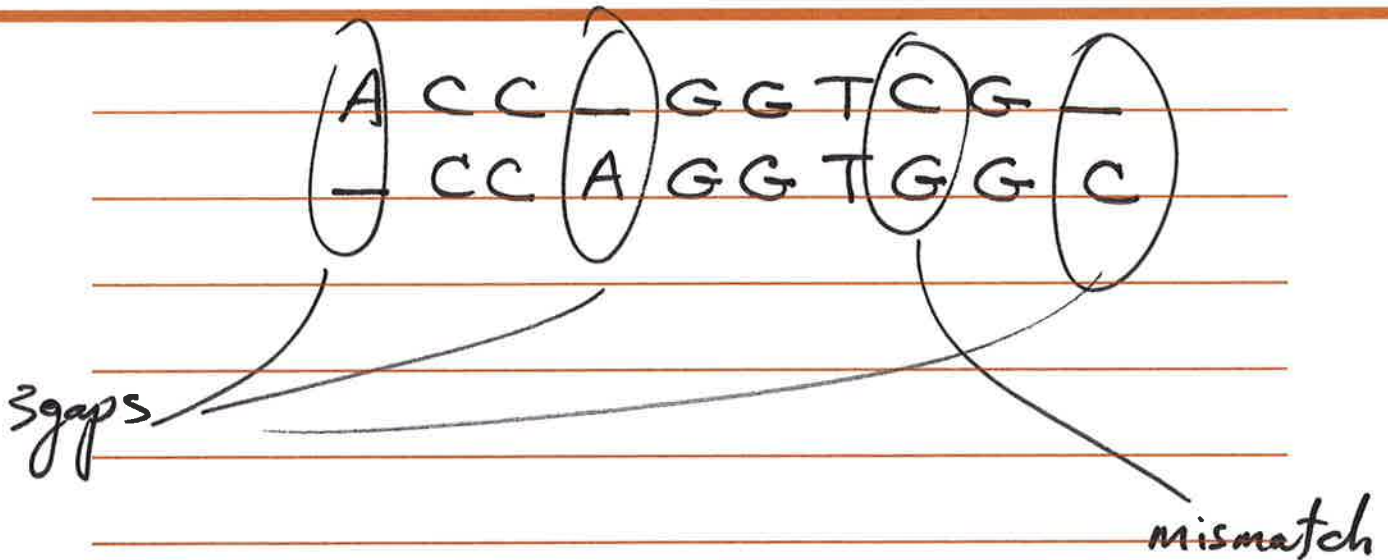
Cytosine C

Guanine G

Thymine T

$S_1 =$ ACCGGTCG

$S_2 =$ CCAGGTGGC



Suppose we have 2 strings X & Y

$$X = \{x_1, x_2, \dots, x_m\}$$

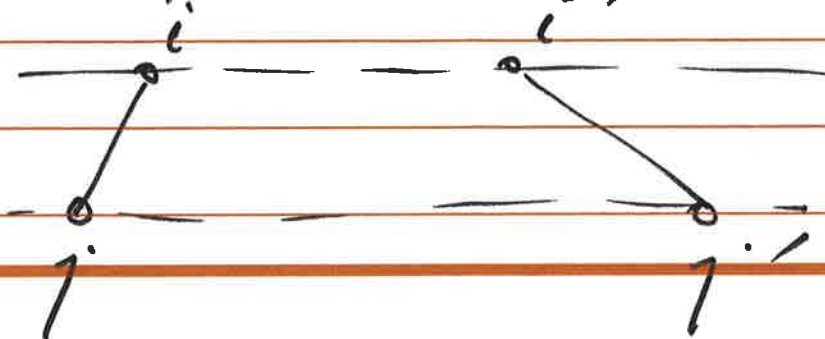
$$Y = \{y_1, y_2, \dots, y_n\}$$

Def. A matching is a set of ordered pairs w/ property that each item occurs at most once.

Car
ing
X
racing

Def. A matching is an alignment if there are no crossing pairs

$$(i, j), (i', j') \in M \text{ \& } i < i' \Rightarrow j < j'$$



For a given alignment M between X & Y

1- We incur a "gap penalty"
of 8 (cost of 8)


2- For each mismatch (of letters g & p)
we incur a mismatch cost
 d_{pg}

#	A	C	G	T
A	0	x	x	x
C		0	x	x
G			0	x
T				0

Similarity between strings X and Y

is the Min. Cost of an alignment
between X and Y

$$X = \{x_1 \dots x_m\}$$

$$Y = \{y_1 \dots y_n\}$$


either $(x_m, y_n) \in M$ or $(x_m, y_n) \notin M$

Define $OPT(i, j)$ as the min. cost of an alignment between $x_1 \dots x_i$ & $y_1 \dots y_j$

In an opt. alignment M , at least one of the following is true:

1 - $(x_m, y_n) \in M \Rightarrow$

$$OPT(m, n) = OPT(m-1, n-1) + \alpha_{x_m y_n}$$

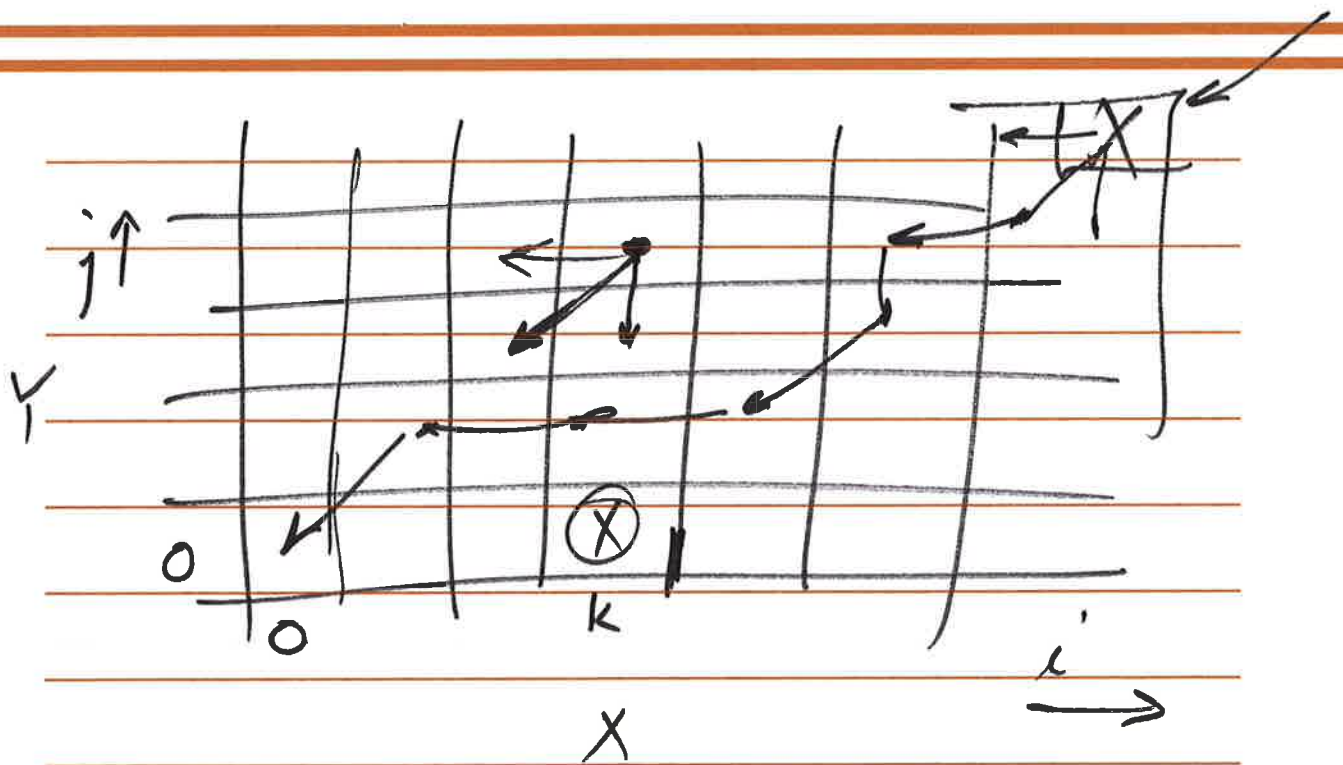
2 - x_m is not matched \Rightarrow

$$OPT(m, n) = OPT(m-1, n) + \delta$$

3 - y_n is not matched \Rightarrow

$$OPT(m, n) = OPT(m, n-1) + \delta$$

$$OPT(i, j) = \min \left[\alpha_{x_i y_j} + \delta OPT(i-1, j-1), \right. \\ \left. \delta + OPT(i-1, j), \right. \\ \left. \delta + OPT(i, j-1) \right]$$



$$10 \overset{7}{\times} 10 \overset{7}{\times} 4 = 4 \overset{14}{\times} 10$$

Alignment (x, y)

Initialize $A[i, 0] = i\delta$ for each i

" $A[0, j] = j\delta$ " " j

for $j = 1$ to n

for $i = 1$ to m

$A[i, j] = \text{Min} (\alpha_{x_i y_j} + A[i-1, j-1],$

$\delta + A[i-1, j],$

$\delta + A[i, j-1]]$

endfor

endfor

return $A[m, n]$

takes $O(mn)$

