

Dynamic Programming

- used in DNA sequence alignment

A C G T C T A G

A C - T C T A G



not aligned

- how to align sequences?

- Come up with numbers indicating
m m m

- sequence alignment:
based (partially) on
Levenshtein distance
used in Unix diff
program.

2 diff prog. cpp old. cpp

ready input for versioning
(SVN depends on it)

- for example

WHAT	W	H	A	T
WH-Y				
WH-Y				
	+	-	-	-
cost				

Use scores for

identity (match) + 1

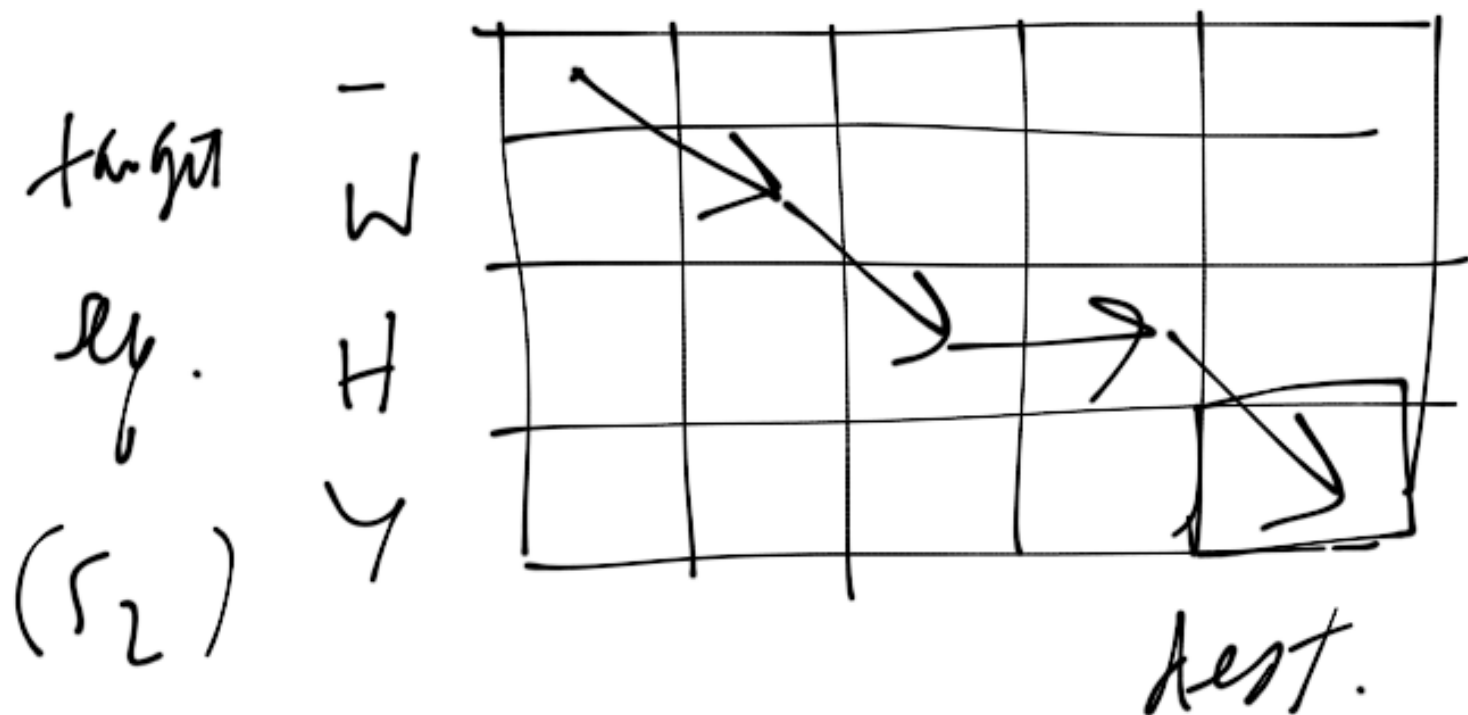
substitution (mismatch) - μ

insertion/deletion
(indel) - δ

-ve sign ensure low scores for mismatches

- Rewrite algorithm as a
matrix: search space (s_1)

- W H A T



- alg. starts at upper left,
proceeds right, down, &
diag. depending on slo't
(stored in matrix)

	-	W	H	A	T
0					
W		1			
H			2	1	
Y					2

- one matrix is filled in,
 work backward, from
 distance to get
 alignment

- say we have three lists:

identity : +1

sub : -1

in del : -2

0 1 2

0	-	-	W	H
1	W			
2	H			

(2, 2)

- to get to (2,2):

i) from (2,1) (better left)

WHAT

WHY

W-HY

$$S_{2,2} = S_{2,1} - 2$$

(ideal)

ii) from (1,2) (shove)

WHAT

WHY

$$S_{2,2} = S_{1,2} - 2$$

(ideal)

PICK

iii) from (1,1)

WHAT
WHY

$$S_{-1,-1} = S_{1,1} + 1$$

- alg.:

init first row, first col

with $(g) \neq i$

(-2)

includ gap cost

matrix C:

	-	W	H	A	T
-	0	-L	-4	-6	-8
W	-L				
H	-4				
A	-6				
T	-8				

- calculate scores for
all cells starting at
top-left:

for ($i=1, i \leq \text{rows}; i++$)
for ($j=1; j \leq \text{cols}; j++$)

$$c_0 = C[i-1][j-1] +$$

$$S(S_1[i-1], S_2[j-1])$$

$$c_1 = C[i-1][j] + g$$

$$c_2 = C[i][j-1] + g$$

$$C[i][j] = \max(c_0, c_1, c_2)$$

$S(s_1[i-1], s_2[j-1]) :$

$S(a, b) :$

if $(a == b)$

return
(max 1 cost)

else

return -1

(substitute cost)

Needleman - Word ab.