

①

## String Similarity

How similar are two strings:

ocurrance

occurrence

Alignment 1

o c u r r a n c e  
o c c u r r e n c e

6 mismatches

1 gap

Alignment 2

o c - u r r a n c e  
o c c u r r e n c e

1 mismatch

1 gap

Alignment 3

o c - u r r a - n c e  
o c c u r r - e n c e

0 mismatches

3 gaps

Applications:

- spelling correction
- speech recognition
- computational biology

(2)

Edit distance (Levenshtein 1966,  
Needleman/Wunsch 1970, ...)

- gap penalty  $\delta > 0$
- mismatch penalty  $\alpha_{p,q}$  ( $p, q$  characters)

Levenshtein:  $\delta = 1$ ,  $\alpha_{pq} = 1$  if  $p \neq q$ ,  
0 otherwise

- cost of alignment: sum of gap and mismatch penalties

Problem: Given  $x = [x_1, \dots, x_m]$ ,  $y = [y_1, \dots, y_n]$ ,  
find (alignment of) minimal cost

Alignment (Formalization): sequence of ordered pairs  $(u_k, v_k)$

- $u_k \in \{x_1, \dots, x_m, -\}$ ,  $v_k \in \{y_1, \dots, y_n, -\}$
- each  $x_i$  and  $y_j$  occurs exactly once in a pair
- order of the  $x_i$  in alignment is the same as in  $x$  (same for  $y_j$  and  $y$ )

(3)

We consider the Levenshtein distance:

Cost of aligning  $x[1..i]$  and  $y[1..j]$ .

The best alignment can be one of three:

Case 1:

$$\begin{array}{cccc}
 | & | & | & \dots & | \\
 & & & & x_i \\
 & & & & | \\
 & & & & y_j
 \end{array}
 \Rightarrow \text{cost} \geq \text{cost}(x[1..i-1], y[1..j-1]) + \alpha_{x_i, y_j}$$

Case 2:

$$\begin{array}{cccc}
 | & | & \dots & | & \overset{\text{something}}{\underbrace{0 \quad 0 \quad \dots \quad 0}} & | \\
 & & & & | & x_i \\
 & & & & y_j & | \\
 & & & & | & | \\
 & & & & | & | \\
 & & & & | & | \\
 & & & & | & |
 \end{array}
 \Rightarrow \text{cost} \geq \text{cost}(x[1..i-1], y[1..j]) + \delta$$

Case 3: symmetric

$$\text{cost} \geq \text{cost}(x[1..i], y[1..j-1]) + \delta$$

Conversely, consider the best alignment for  $x[1..i-1]$ ,  $y[1..j-1]$  and add  $(x_i, y_j)$ :

$$\begin{array}{cccc}
 | & | & \dots & | \\
 & & & + \\
 & & & x_i \\
 & & & | \\
 & & & y_j
 \end{array}
 \Rightarrow \text{cost}(x[1..i-1], y[1..j-1]) + \alpha_{x_i, y_j} \geq \text{cost}$$

(4)

Consider best alignment for  $x[1..i-1]$  and  $y[1..j]$ :

$x[1..i-1]$

$$\left| \begin{array}{c} | \dots | \\ \vdots \\ | \dots | \end{array} \right| + \frac{x_i}{-} \Rightarrow \text{cost}(x[1..i-1], y[1..j]) + \delta$$
$$y[1..j] \geq \text{cost}$$

Also:  $\text{cost}(x[1..i], y[1..j-1]) + \delta \geq \text{cost}$ .

We conclude: Let  $c[i,j] := \text{cost}(x[1..i], y[1..j])$ .

Recurrence for  $c[i,j]$ :

$$c[0,j] = \delta \cdot j$$

$$c[i,0] = \delta \cdot i$$

$$c[i,j] = \min \left( \begin{array}{l} c[i-1, j-1] + \alpha_{x_i, y_j} \\ c[i-1, j] + \delta, \\ c[i, j-1] + \delta \end{array} \right)$$

Compute  $c[i,j]$  with dynamic programming

5

Example: HOUSE vs HOST

		H	O	U	S	E
	0	1	2	3	4	5
H	1	0	1	2	3	4
O	2	1	0	1	2	3
S	3	2	1	1	0	2
T	4	3	2	2	2	1

Trace : where did the minimum come from?

Alignment

H O U S E  
H O - S T