Dynamic programming: Edit Distance

Edit Distance

- Given two strings s and t, the edit distance between s and t is the minimum number of editing operations needed to turn s into t
- Editing operations
 - Insertion
 - Deletion
 - Substitution

Example

Distance: 5 (assuming unit cost for each operation)

Application: Computational Biology

Given a sequence of bases

AGGCTATCACCTGACCTCCAGGCCGATGCCC
TAGCTATCACGACCGCGGTCGATTTGCCCGAC

An alignment:

```
-AGGCTATCACCTGACCTCCAGGCCGA--TGCCC---
TAG-CTATCAC--GACCGC--GGTCGATTTGCCCGAC
```

Application: NLP

Evaluating Machine Translation and speech recognition

```
R Spokesman confirms senior government adviser was shot

H Spokesman said the senior adviser was shot dead

S I D
```

Optimal substructure?

• Q: if this is an optimal alignment

will this be optimal for sure?

$$s = s_1 s_2 ... s_m; t = t_1 t_2 ... t_n$$

 We can create sub-problems by considering the prefixes of s and t

D(i,j) = the edit distance between $s_1s_2...s_i$ and $t_1t_2...t_j$

Q: To figure out D(i,j), what are the possible choices we can make regarding the last positions i and j?

A: Three possibilities:

Align i with j

$$s_1 s_2 ... s_{i-1} s_i$$

 $t_1 t_2 ... t_{j-1} t_j$

Align i with *

Align j with *

$$s_1 s_2 ... s_{i-1} s_i *$$
 $t_1 t_2 ... t_{j-1} t_j$

D(i,j): the minimum of the three possible choices:

Align i with j

$$s_1 s_2 ... s_{i-1} s_i$$

 $t_1 t_2 ... t_{j-1} t_j$

Align i with *

$$s_1 s_2 \dots s_{i-1} s_i$$

 $t_1 t_2 \dots t_{j-1} t_j *$

Align j with *

Choice 1:

If
$$s_i = t_j$$
:
 $D(i, j) = D(i-1, j-1)$

Otherwise:

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

Choice 2:

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

Choice 2:

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

Recurrence relation for D(i,j)

For i,
$$j \ge 1$$

$$D(i,j) = \min \begin{array}{c} D(i-1,j) + 1 & \text{deletion} \\ D(i,j-1) + 1 & \text{insertion} \\ D(i-1,j-1) + \text{diff}(s_i,t_j) & \text{align i with } j \end{array}$$

$$0 \quad \text{if } a = b$$

$$1 \quad \text{if } a \ne b$$

Base case? any others?

$$D(0,0) = 0$$

Edit Distance

```
For i=0 to m: D(i,0) = i

For j=1 to n: D(0,j) = j

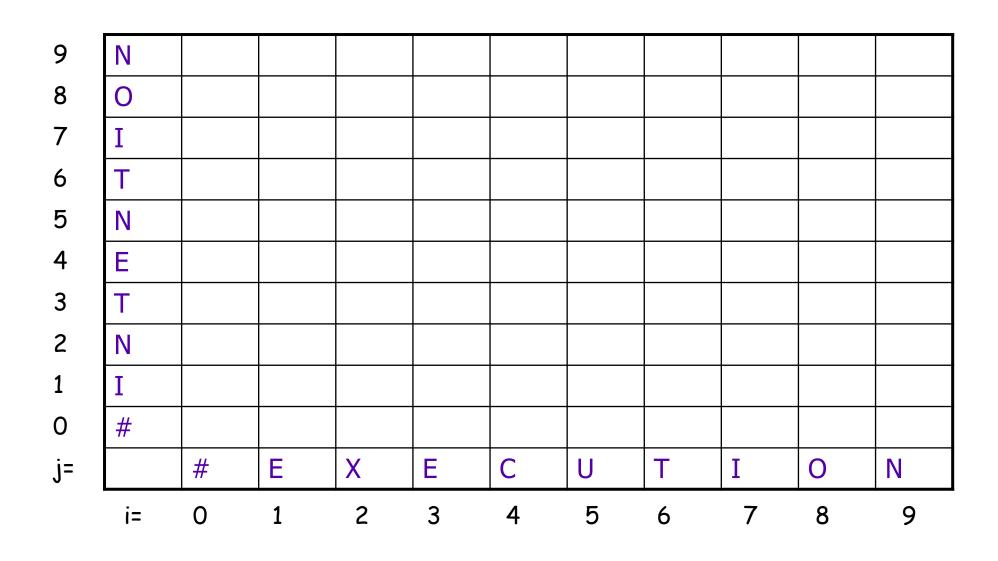
For each i = 1...m

For each j = 1...n

D(i,j) = min \begin{cases} D(i-1,j) + 1 \\ D(i,j-1) + 1 \end{cases}
D(i-1,j-1) + diff(s_i, t_j)
```

Return D(m,n)

The Edit Distance Table



Computing alignments

- · Getting the edit distance isn't sufficient
 - We often need to align each character of the two strings to each other
- We do this by keeping a "backtrace"
- Every time we enter a cell, remember where we came from
- When we reach the end,
 - Trace back the path from the upper right corner to read off the alignment

Adding Backtrace to Minimum Edit Distance

Base conditions:

$$D(i,0) = i$$

$$D(i,0) = i D(0,j) = j$$

Termination:

return D(m,n)

Recurrence Relation:

For each
$$i = 1...M$$

For each $j = 1...N$

$$D(i,j) = \min \begin{cases} D(i-1,j) + 1 & \text{deletion} \\ D(i,j-1) + 1 & \text{insertion} \\ D(i-1,j-1) + \text{diff}(s_i, t_j) & \text{align} \end{cases}$$

$$ptr(i,j) = \begin{cases} LEFT & insertion \\ DOWN & deletion \\ DIAG & align \end{cases}$$

Result of Backtrace

Two strings and their alignment:

Performance

· Time:

O(nm)

• Space:

O(nm)

Backtrace

O(n+m)