

Spell Checker Design

Dr. G. Bharadwaja Kumar

Edit Distance

- The **Edit Distance** (or **Levenshtein distance**) is a metric for measuring the amount of difference between two string.
 - The Edit Distance is defined as the minimum number of edits needed to transform one string into the other.
- It has many applications, such as spell checkers, natural language translation, and bioinformatics.
 - An example of one application in Bioinformatics, measuring the amount of difference between two DNA sequences.

Applications

- Spell correction
 - The user typed "graffe"
 Which is closest?
 - graf
 - graft
 - grail
 - giraffe

- Computational Biology
 - Align two sequences of nucleotides

AGGCTATCACCTGACCTCCAGGCCGATGCCC
TAGCTATCACGACCGCGGTCGATTTGCCCGAC

Resulting alignment:

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

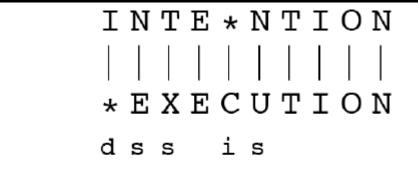
Also for Machine Translation, Information Extraction, Speech Recognition

Edit Distance

 The problem of finding an edit distance between two strings is as follows:

 Given an initial string x, and a target string y, what is the minimum number of changes that have to be applied to x to turn it into y?

- The list of valid changes are:
 - 1) Inserting a character
 - 2) Deleting a character
 - 3) Substituting a character to another character.

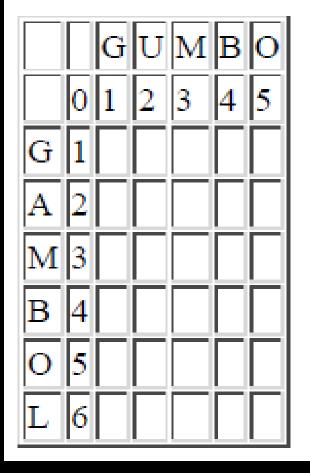


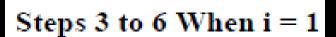
- If each operation has cost of 1
 - Distance between these is 5
- If substitutions cost 2 (Levenshtein)
 - Distance between them is 8

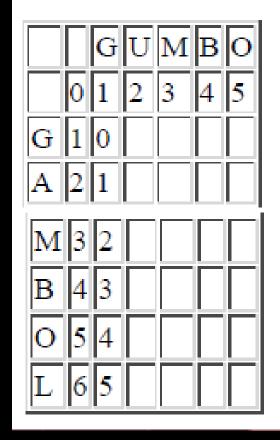
| Step | Description |
|------|--|
| 1 | Set n to be the length of s. |
| | Set m to be the length of t. |
| | If $n = 0$, return m and exit. |
| | If $m = 0$, return n and exit. |
| | Construct a matrix containing 0m rows and 0n columns. |
| 2 | Initialize the first row to 0n. |
| | Initialize the first column to 0m. |
| 3 | Examine each character of s (i from 1 to n). |
| 4 | Examine each character of t (j from 1 to m). |
| 5 | If s[i] equals t[j], the cost is 0. |
| | If s[i] doesn't equal t[j], the cost is 1. |
| 6 | Set cell d[i,j] of the matrix equal to the minimum of: |
| | a. The cell immediately above plus 1: d[i-1,j] + 1. |
| | b. The cell immediately to the left plus 1: d[i,j-1] + 1. |
| | c. The cell diagonally above and to the left plus the cost: d[i-1,j-1] + cost. |
| 7 | After the iteration steps (3, 4, 5, 6) are complete, the distance is found in cell d[n,m]. |

See how the Levenshtein distance is computed when the source string is "GUMBO" and the target string is "GAMBOL".

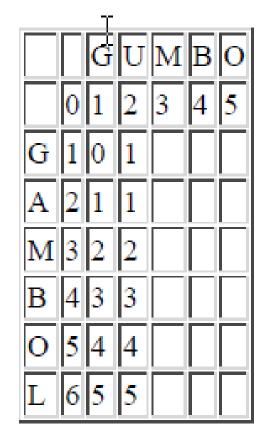
Steps 1 and 2







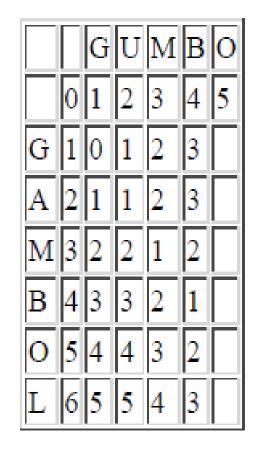
Steps 3 to 6 When i = 2



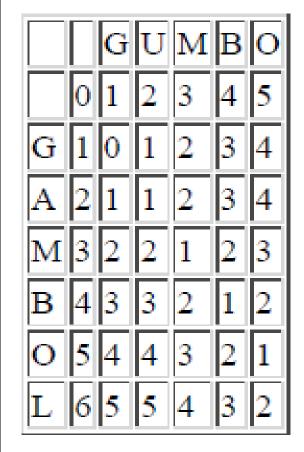
Steps 3 to 6 When i = 3

| | | G | U | M | В | o |
|---|---|---|---|---|---|---|
| | 0 | 1 | 2 | 3 | 4 | 5 |
| G | 1 | 0 | 1 | 2 | | |
| A | 2 | 1 | 1 | 2 | | |
| M | 3 | 2 | 2 | 1 | | |
| В | 4 | 3 | 3 | 2 | | |
| О | 5 | 4 | 4 | 3 | | |
| L | 6 | 5 | 5 | 4 | | |

L Steps 3 to 6 When i = 4



Steps 3 to 6 When i = 5



Initialization

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

Recurrence Relation:

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

For each $j = 1...N$

$$D(i,j) = \min \begin{cases} D(i-1,j) + 1 \\ D(i,j-1) + 1 \\ D(i-1,j-1) + 2; \text{ if } X(i) \neq Y(j) \\ 0; \text{ if } X(i) = Y(j) \end{cases}$$
Termination:

Termination:

D(N,M) is distance

Finally

 The distance is in the lower right hand corner of the matrix, i.e. 2. This corresponds to our intuitive realization that "GUMBO" can be transformed into "GAMBOL" by substituting "A" for "U" and adding "L" (one substitution and 1 insertion = 2 changes).