Advanced Programming

Edit distance

Ramaseshan Ramachandran

Minimum Edit Distance Minimum Edit Distance Applications Operations Examples
Alignment
Properties of Edit Distance
Dynamic Programming
Application of Rules

MINIMUM EDIT DISTANCE

In computational linguistics, edit distance is a way to quantify how similar/dissimilar two strings are. This (dis)similarity is found by counting the minimum number of operations required to transform one string into the other. The count is known as Levenshtein minimum edit distance.

APPLICATIONS

- Retrieving similar words
- Spelling correction
- Speech recognition
- OCR error corrections
- Computational Biology alignment of DNA sequence
- ► Identification of duplicate invoices

MINIMUM EDIT DISTANCE - OPERATIONS

- ► Two strings S1 and S2 are given
- ▶ S1 is the input word with wrong spelling and S2 is the Correct word
- ▶ Operations allowed to transform S1 to S2 insert, delete and substitute/replace

S1	W	Н	I	Z	Z	Τ	Ν	G
S2	М	U	Z	Z	Z	1	Ν	G
Operations	r	r	r	n	n	n	n	n
Cost of operation	1	1	1	0	0	0	0	0

Cost of operation = 3

S1	W	Н	I	Z	Z	1	Ν	G
S2	Q	U	1	Z	Z	Е	R	S
Operations	r	r	n	n	n	r	r	r
Cost of operation	1	1	0	0	0	1	1	1

Cost of operation = 5

S1	С	Α	R	N	Α	П	Ι	С	
S2	K	Α	R	N	Α	Т	Α	K	Α
Operations	r	r	n	n	n	r	r	r	i
Cost of operation	1	0	0	0	0	0	1	1	1

Cost of operation = 4

FIND SIMILARITY BETWEEN DNA SEQUENCE

Α	Α	G	Т	С	Т	Т	Α	Т	Α	С	Α	G	G	С
Α	Т	G	Α	С	Т	Α	Т	Α	G	G	G	С	Α	

Α	Α	G	Т	C	Т	Т	Α	Т	Α	C	Α	G	G	C	-
	Χ		Χ								Χ				
Α	Т	G	Α	C	Т		Α	Т	Α	G	G		G	C	Α
	S		S			D					S				1

X between the pair	Transformation is necessary
— in the first string	Insertion operation is required
— in the second string	Deletion operation is needed

How about these two strings?
AGGCTATCACCTGACCTCCAGGCCGATGCCC
TAGCTATCACGACCGCGGTCGATTTGCCCGAC

ANOTHER EXAMPLE

S1	Е	Χ	ı	Т	Е	М	Е	N	Т	
S2	Е	Χ	С	1	Т	Е	М	Ε	Ν	Т
Operations										
Cost of operation										

Cost of operation =

IMPORTANCE OF ALIGNMENT

S1	E	Х		T	E	*	M	Ε	N	T
S2	Е	Χ	C	1	Т	Ε	М	Е	Ν	Т
Operations										
Cost of operation										

Cost of operation =

IMPORTANCE OF ALIGNMENT

S1	Е	Χ	ı	*	Т	Е	М	Е	N	Т
S2	Ε	Χ	С		Т	Е	М	Ε	N	Т
Operations										
Cost of operation										
	\overline{C}	ct of		rati	2 12					

Cost of operation =

S1	E	X	*	1	T	Е	М	Е	N	Т
S2	Е	Χ	С	1	Т	Е	М	Е	N	Т
Operations										
Cost of operation										

Cost of operation =

SEARCH SPACE

Searching for sequence of edit operation (path in a graph) from the S_1 to the S_2 :

- ▶ Start State: $S_1 \rightarrow S_2$
- Operators: Insert, delete and substitute
- Final State: S₂
- Cost: Number of edit operations
- Search space of all edit sequences is huge
- Several distinct path to reach the final state
- Do we need to remember all of them?
- Can we just find the shortest path?

PROPERTIES OF EDIT DISTANCE

- If two strings are equal, then Levenshtein distance is zero or $d(x,x) = \forall x \in M$, where M is the metric space
- The simplest operation of the replacement.
- ▶ If two strings differ in k^{th} position, then $d_e = D(x_k, y_k)$
- Penalty for multiple replacements = $\sum_{k=1}^{r} D(x_k, y_k)$
- If the deletion/insertion of an entry $x_{\hat{k}}$ gives the required string, then the penalty is $D(x_{\hat{v}},y_{\hat{v}})$
- Penalty for multiple deletions $=\sum_{k=1}^{r}D(x_{\hat{k}},y_{\hat{k}})$
- ► The cost of the operation =

$$\rho_r \sum_{k=1}^r D(x_k, y_k) + \rho_d \sum_{k=1}^d D(x_{\hat{k}d}, y_{\hat{k}d}) + \rho_i \sum_{k=1}^i D(x_{\hat{k}i}, y_{\hat{k}i}) \text{ where } \rho_s, \rho_d, \rho_i \text{ are the penalties for replacement, deletion and substitution operations, respectively}$$

DYNAMIC PROGRAMMING

Breaking Down Complexity

► Efficient mechanism to find optimal solutions which contains a lot of overlapping sub-problems

Optimal Substructure

Each sub-problem has an optimal solution that contributes to the global optimal solution

Memoization

Sub-problems results are stored to avoid redundant calculations

► Bottom-Up Approach

Recombines solutions to sub-problems to form the final solution

SUBPROBLEMS

```
MED({A,B,C}, {D,B,C})

/

MED({B,C}, {D,B,C})

MED({A,B,C}, {B,C})

MED({A,B,C}, {B,C})

MED({C}, {D,B,C})

MED({B,C}, {B,C})

MED({A,B,C}, {C})

MED({A,B,C}, {C})
```

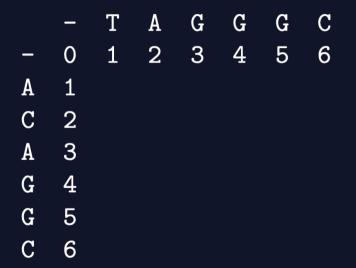
EDIT DISTANCE COMPUTATION

A tabular computation of D(m,n) is achieved by computing minimum edit distance, D(i,j) for small i,j i.e., fill up element D(i,j) for 0 < i < m, 0 < j < n using two rules

Rule 1: if
$$s_j = t_j$$
 then $D(i,j) = D(i-1,j-1)$ (1)

$$\begin{aligned} & \text{\textbf{Rule 2:}} \text{if } s_j \neq t_j \text{ then } D(i,j) = \min \begin{cases} D(i-1,j)+1 & \text{insert} \\ D(i,j-1)+1 & \text{delete} \\ D(i-1,j-1)+1 & \text{replace} \end{cases} \end{aligned} \tag{2}$$

DP - EDIT DISTANCE



	_	T	Α	G	G	G	C	Α
-	0	1	2	3	4	5	6	7
Α	1	1	1	2	3	4	5	6
C	2	2	2	2	3	4	4	5
Α	3	3	2	3	3	4	5	4
G	4	4	3	2	3	3	4	5
G	5	5	4	3	2	3	4	5
С	6	6	5	4	3	4	3	4

DYNAMIC PROGRAMING - APPLICATION OF RULES

Application of Delete operation

		P	L	Α	Y
	0	1	2	3	4
Р	1	0	1	2	3
L	2	1	0	1	2
Υ	3	2	1	1	1

Application of Insert operation

ſ			Р	L	Υ
		0	1	2	3
Γ	Р	1	0	1	2
Γ	L	2	1	0	1
	Α	3	2	1	1
	Υ	4	3	2	1

Application of replace operation

		Р	L	0	Υ
	0	1	2	3	4
Р	1	0	1	2	3
L	2	1	0	1	2
Α	3	2	1	1	2
Υ	4 Adv	3	2	2	1

DYNAMIC PROGRAMMING -EXERCISE - CORRECT THIS TABLE

		Е	Х	Т	Т	Е	М	Е	N	Т
	0	1	2	3	4	5	6	7	8	9
Е	1	0	1	2	3	4	5	6	7	8
Х	2	1	0	1	2	3	4	5	6	7
С	3	2	1	1	2	3	4	5	6	7
- 1	4	3	2	1	2	3	4	5	6	7
Т	5	4	3	2	1	2	3	4	5	6
Е	6	5	4	3	2	1	2	3	4	5
М	7	6	5	4	3	2	1	2	3	4
Е	8	7	6	5	4	3	2	1	2	3
N	9	8	7	6	5	4	3	2	1	2
Т	10	9	8	7	6	5	4	3	2	1

What is the minimum edit distance? what is the minimum distance for the strings below: How about these two strings?

AGGCTATCACCTGACCTCCAGGCCGATGCCC
TAGCTATCACGACCGCGGTCGATTTGCCCGAC

- ▶ Base condition D(i,0) = i and D(0,j) = j
- ightharpoonup Termination D(m,n)

$$\begin{split} \text{for } i = 1 \dots m \\ \text{for } j = 1 \dots n \\ Di, j = min \begin{cases} D(i-1,j) + 1 \text{deletion} \\ D(i,j-1) + 1 \text{insertion} \\ D(i-1,j-1) + 1; \text{ if } X(i) \neq Y(j) \\ D(i-1,j-1) + 0; \text{ if } X(i) = Y(j) \end{cases} \end{split}$$