

NLP - Lecture 01 - Handout

Minimum Edit Distance (MED)

Given String X and Y , three operations allowed for editing string X or Y :

- Insertion (Ins)
- Deletion (Del)
- Substitution (Sub)

Example $X = \text{INTENTION}$
 $Y = \text{EXECUTION}$

Let $O = [O_1, O_2, \dots, O_k]$ be k operations that we can take X to Y , where each $O_k \in \{ \text{Ins}, \text{Del}, \text{Sub} \}$.

$X \xrightarrow{\text{Del: I}} \text{INTENTION} \xrightarrow{\text{Sub: N} \rightarrow \text{E}} \text{ETENTION}$

$\xrightarrow{\text{Sub: T} \rightarrow \text{X}} \text{EXENTION} \xrightarrow{\text{Ins: C}} \text{EXECUTION}$

$\xrightarrow{\text{Sub: U} \rightarrow \text{U}} \text{EXECUTION}$

$O_1: \text{Del}, O_2: \text{Sub}, O_3: \text{Sub}, O_4: \text{Ins}, O_5: \text{Sub}. k=5.$

If Ins, Del cost 1 and Sub costs 2.
then. this alignment costs 8.

We call $\underline{D} = [D_1, D_2, \dots, D_k]$ that takes
 X to Y an alignment.

Q: How to find the optimal alignment?

(By optimal, we mean the cost of \underline{D} is minimum). Let us define:

$X = x_1, x_2, \dots, x_n$, i indexing i th char of X

$Y = y_1, y_2, \dots, y_m$, j indexing j -th char of Y .

Naive method: See slide P77.

An alternative way, consider a smaller
problem (subproblem).

$X_i = [x_1, x_2, \dots, x_i]$, $Y_j = [y_1, y_2, \dots, y_j]$.

Let $X_i = [x_1, x_2, \dots, x_i]$, $Y_j = [y_1, y_2, \dots, y_j]$

Define $D[i, j]$: MED from $X_i \rightarrow Y_j$.

Then, we know:

- $D[i-1, j]$: MED from $X_{i-1} \rightarrow Y_j$
- $D[i, j-1]$: MED from $X_i \rightarrow Y_{j-1}$
- $D[i-1, j-1]$: MED from $X_{i-1} \rightarrow Y_{j-1}$.

Let $D = D_1 D_2 \dots D_k$ be the optimal operations occurred in the MED process.

Example: $I N T E * N T I O N$
 $\downarrow \downarrow \downarrow \vdots \downarrow \downarrow \vdots \vdots \vdots \vdots$
 $* E X E C U T I O N$

$$X_i = X_i^0 \xrightarrow{D_1} X_i^1 \xrightarrow{D_2} X_i^2 \rightarrow \dots \xrightarrow{D_{k-1}} X_i^{k-1} \xrightarrow{D_k} X_i^k = Y_j$$

where $D_k: a \rightarrow b$

a : could be either $*$ or $x \in X_i$.

b : could be either $*$ or $y \in Y_j$.

• optimal substructure:

$$X_i = X_i^0 \xrightarrow{O_1} X_i^1 \xrightarrow{O_2} X_i^2 \rightarrow \dots \xrightarrow{O_{k-1}} X_i^{k-1} \xrightarrow{O_k} X_i^k = Y_j$$

$X_i \rightarrow X_i^{k-1}$: O_1, O_2, \dots, O_{k-1} is still optimal operations from X_i to X_i^{k-1} . (Note

X_i^{k-1} could be either Y_j or Y_{j-1} .

Why? (Hint: Make a contradiction).

$O_k = \begin{cases} \text{case 1. } x_i \text{ is deleted from } X_{i-1} \\ \text{case 2. } y_i \text{ is inserted into } Y_{i-1} \\ \text{case 3. } x_i \text{ is substituted by } y_i \end{cases}$

For case 1. : $D[i-1, j] + \text{Del}(x_i)$

$$D[i, j] = D[i-1, j] + \text{Del}(x_i)$$

For case 2:

$$D[i, j] = D[i, j-1] + \text{Ins}(y_i)$$

For case 3:

$$D[i, j] = D[i-1, j-1] + \text{Sub}(x_i, y_i)$$

Clearly, $D[i, j]$ must be one of them, i.e.,

$$D[i, j] = \min \begin{cases} D[i-1, j] + Del(x_i) \\ D[i, j-1] + Ins(y_i) \\ D[i-1, j-1] + Sub(x_i, y_i) \end{cases}$$

$i = 1, 2, \dots, n$
 $j = 1, 2, \dots, m$

Note For $i = 0$, $D[i, j] = j$.

For $j = 0$, $D[i, j] = i$.

We call the above method a dynamic programming method.

Remarks and Questions:

- ①. Prove or disprove the uniqueness of operation,
- ②. Find out different type of costs (edit distance metrics).
- ③. If n, m are very large, can you find approximate solutions?

(Hint: find "approximate string matching").