

## 2.10 编辑距离问题(EDP)

编辑距离问题(EDP)在[16,125页]中被称作“能用DP解决的不精确匹配问题”。它在[22, 284 - 286页]中也被称作“字串编辑问题”。[10,364-367页]还描述了这个问题的变种。

令 $\Sigma$ 为一个有限的字母表。给定两个串 $x \in \Sigma^m$ 和 $y \in \Sigma^n$ , 假定 $x = x_1 \dots x_m$ 且 $y = y_1 \dots y_n$ 。我们的任务就是用最少的编辑操作使 $x$ 变成 $y$ 。这里说的编辑操作有以下三种:

- 删除操作D:花费 $c(D)$ 从串中删掉一个字母。
- 插入操作I:花费 $c(I)$ 从 $\Sigma$ 往串中插入一个字母。
- 复位操作(或者说是替换操作)R: 花费 $c(R)$ 用 $\Sigma$ 中的一个字母把串中的一个字母换掉。

我们的目标是寻找一个编辑序列, 让 $x$ 以最小的花费变成 $y$ 。这里我们定义编辑序列的费用为各操作的费用之和。

下面的递归式能够计算出 $x$ 的前缀 $X_i$ 变成 $y$ 的前缀 $Y_j$ 的最小花费 $f(i, j)$ 。

$$f(i, j) = \begin{cases} jD & \text{if } i = 0 \\ iD & \text{if } j = 0 \\ \min \{f(i-1, j) + c(D), \\ f(i, j-1) + c(I), f(i-1, j-1) + c(R)\} & \text{if } i > 0 \text{ and } j > 0 \end{cases}$$

这里费用函数 $c$ 定义如下:

$$\begin{aligned} c(D) &= c_D && \text{在任意位置删掉任意字母} \\ c(I) &= c_I && \text{在任意位置插入任意字母} \\ c(R) &= \begin{cases} 0 & \text{if } x_i = y_j \text{ (字母相等)} \\ c_R & \text{if } x_i \neq y_j \text{ (一个可行的复位)} \end{cases} \end{aligned}$$

状态转移函数定义如下:

$$\begin{aligned} t(X_i, Y_j, D) &= (X_{i-1}, Y_j) \\ t(X_i, Y_j, I) &= (X_i, Y_{j-1}) \\ t(X_i, Y_j, R) &= (X_{i-1}, Y_{j-1}) \end{aligned}$$

我们的目标是计算 $f(x, y)$ , 即编辑序列的最小花费。

考虑[16,223页]中的一个例子,  $x = \text{"CAN"}$ 且 $y = \text{"ANN"}$ , 插入花费 $c_I = 1$ , 删除花费 $c_D = 1$ , 复位花费 $c_R = 1$ 。有以下几种花费 $f(x, y) = 2$  (花费最小)的编辑序列:

- $\text{CAN} \vdash_R \text{CNN} \vdash_R \text{ANN}$ ,
- $\text{CAN} \vdash_D \text{AN} \vdash_I \text{ANN}$ ,
- $\text{CAN} \vdash_I \text{CANN} \vdash_D \text{ANN}$ ,

这个问题和2.23节点最长公共子序列问题(LCS)有很密切的联系。