



## 最小编辑距离

**编辑距离**是针对二个字符串的差异程度的量化量测，量测方式是看至少需要多少次的处理才能将一个字符串变成另一个字符串。编辑操作中，许可的操作一般包括将一个字符**替换**成另一个字符，**插入**一个字符，**删除**一个字符。一般来说，编辑距离越小，两个串的相似度越大。由一个字符串转换成另外一个字符串所需的最少编辑次数，也被称之为**最小编辑距离**。

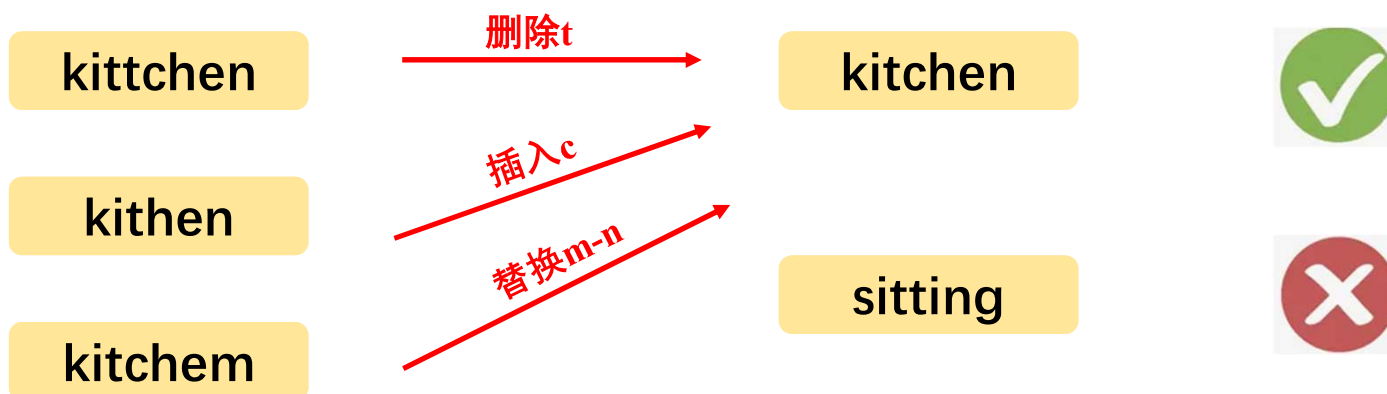
编辑距离由俄罗斯科学家 Vladimir Levenshtein 在1965年提出，也因此而得名 Levenshtein Distance。可以用在信息论、语言学和计算机科学领域，用于DNA分析、拼字检测等等，核心是比较两者的相似性。



# 最小编辑距离

问题背景：

输入法自动更正



问题：如何衡量序列的相似程度



# 最小编辑距离

kittchen

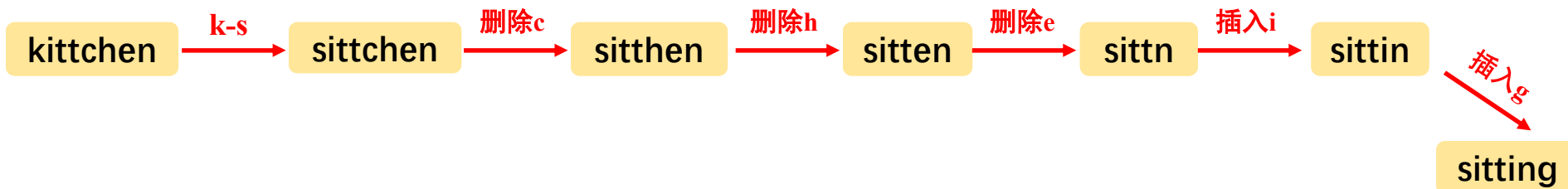
编辑操作

插入、删除、替换

sitting

方案1

编辑操作：6次



方案2

编辑操作：5次



问题：如何求解最少的编辑操作次数（最小编辑距离）



# 最小编辑距离

形式化表示：

输入：长度为 $n$ 的字符串 $X: \{a_1, a_2, \dots, a_n\}$ ，长度为 $m$ 的字符串 $Y: \{b_1, b_2, \dots, b_m\}$

输出：求出一组编辑操作 $D = \langle d_1, d_2, \dots, d_l \rangle$

令

$$\min |D|$$

优化目标

*s.t.* 字符串经过了 $D$ 次的操作之后满足

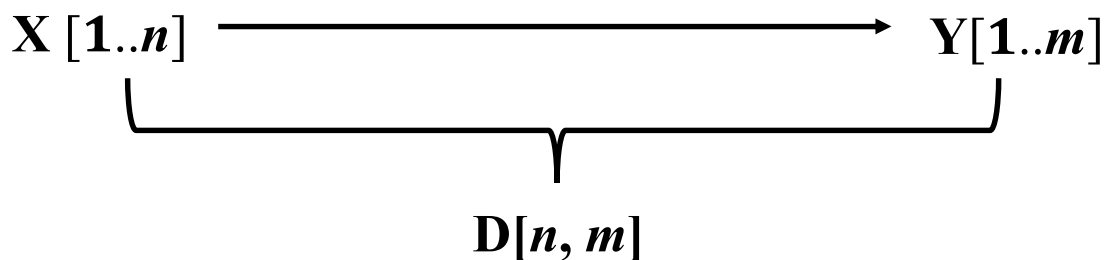
$$\{a_1, a_2, \dots, a_n\} = \{b_1, b_2, \dots, b_m\}$$

约束条件



## 递归结构

数组  $D[i, j]$ : 表示字符串  $X[1..i]$  变成字符串  $Y[1..j]$  的最少操作次数



$X = \{a_1, a_2, \dots, a_n\} \quad (1 \leq i_1 < i_2, \dots, i_l \leq n)$

$Y = \{b_1, b_2, \dots, b_m\} \quad (1 \leq j_1 < j_2, \dots, j_l \leq m)$



$X=Y$



对于任意两条序列 $X[1..m]$ 和 $Y[1..n]$ 的最长公共子序列的长度

The diagram illustrates the dynamic programming recurrence for the Longest Common Subsequence (LCS) problem. It shows three cases for calculating  $C[m, n]$  based on the last elements of sequences  $X$  and  $Y$ .

On the left, a vertical bracket labeled  $\max|Z|$  groups the three cases. Each case consists of a sequence  $X$  and a sequence  $Y$ , followed by a yellow arrow pointing to the recurrence formula for  $C[m, n]$ .

- Case 1:**  $X = \{x_1, x_2, \dots, x_m\}$  and  $Y = \{y_1, y_2, \dots, y_n\}$ . The recurrence is  $C[m, n] = C[m-1, n-1] + 1$ , with the condition  $x_m = y_n$ .
- Case 2:**  $X = \{x_1, x_2, \dots, x_m\}$  and  $Y = \{y_1, y_2, \dots, y_{n-1}, y_n\}$ . The recurrence is  $C[m, n] = C[m, n-1] + 0$ , with the condition  $x_m \neq y_n$ .
- Case 3:**  $X = \{x_1, x_2, \dots, x_{m-1}, x_m\}$  and  $Y = \{y_1, y_2, \dots, y_n\}$ . The recurrence is  $C[m, n] = C[m-1, n] + 0$ , with the condition  $x_m \neq y_n$ .

The recurrence formulas are shown in red dashed boxes, and the conditions are shown below them. Red arrows point to the conditions in the third case.

## 最优子结构性质

## 1. 考察末尾元素

## 2. 操作步骤:

插入  
删除  
替换

## 递归式？

在最小编辑距离问题中，  
只操作其中的一条序列，另一  
条序列固定不变。



## 递归结构

固定Y序列，只操作X序列

$X = \{A, B, C, B, D, A, \mathbf{B}\}$

$Y = \{B, D, C, A, B, \mathbf{A}\}$

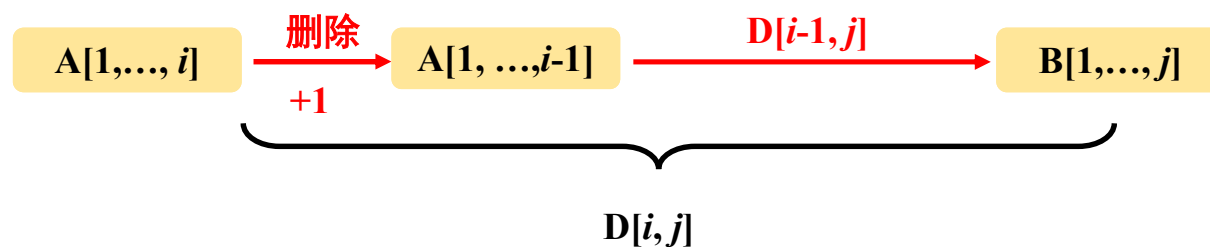
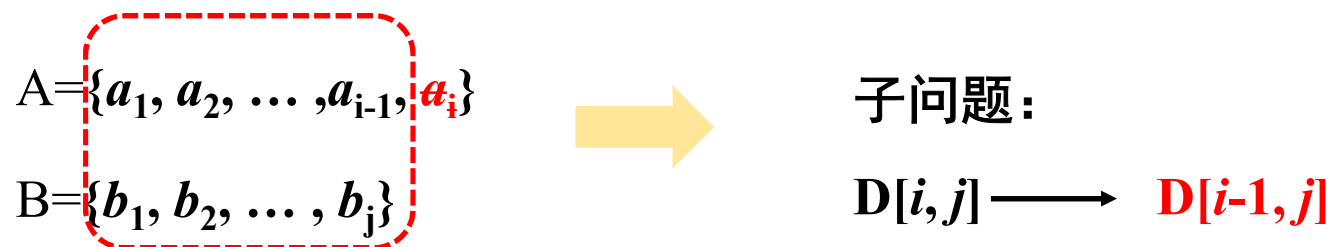
min	{	$A = \{A, B, C, B, D, \mathbf{A}, \mathbf{B}\}$	删除B
		$B = \{B, D, C, A, B, \mathbf{A}\}$	
		$A = \{A, B, C, B, D, A, B, \mathbf{A}\}$	插入A
		$B = \{B, D, C, A, B, \mathbf{A}\}$	
		$A = \{A, B, C, B, D, A, \mathbf{A}\}$	替换B-A
		$B = \{B, D, C, A, B, \mathbf{A}\}$	



## 递归结构

对于任意两条序列  $A = \{a_1, a_2, \dots, a_i\}$  和  $B = \{b_1, b_2, \dots, b_j\}$

$D[i, j]$



在删除操作中:  $D[i, j] = D[i-1, j] + 1$

最优子结构性质





# 递归结构

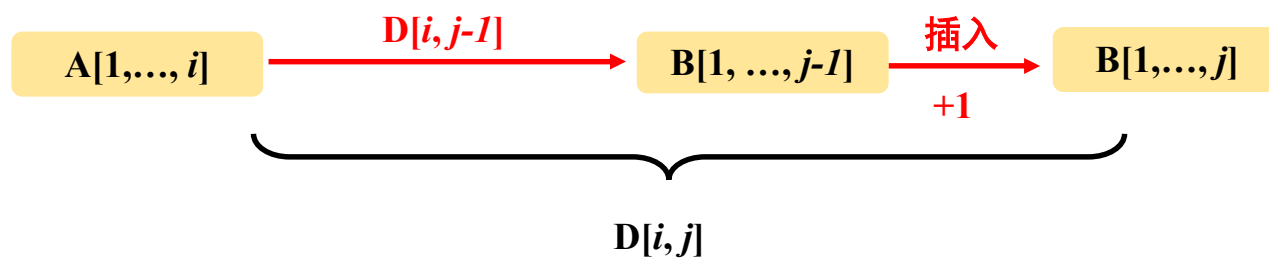
$D[i, j]$

$$A = \{a_1, a_2, \dots, a_{i-1}, a_i, \mathbf{b_j}\}$$
$$B = \{b_1, b_2, \dots, b_{j-1}, b_j\}$$



子问题:

$$D[i, j] \longrightarrow \mathbf{D[i, j-1]}$$



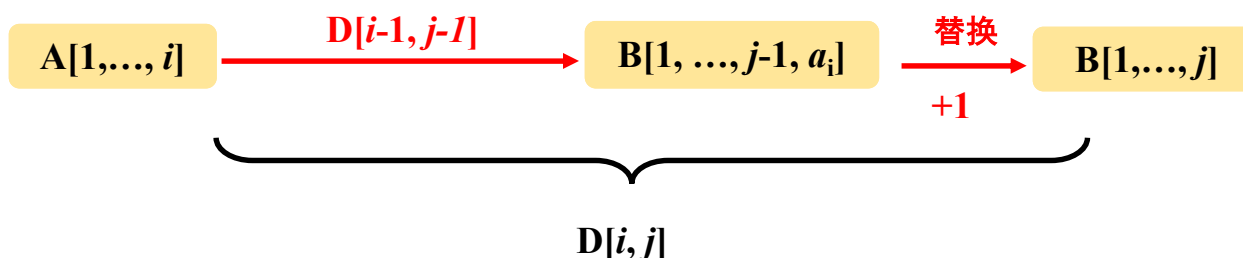
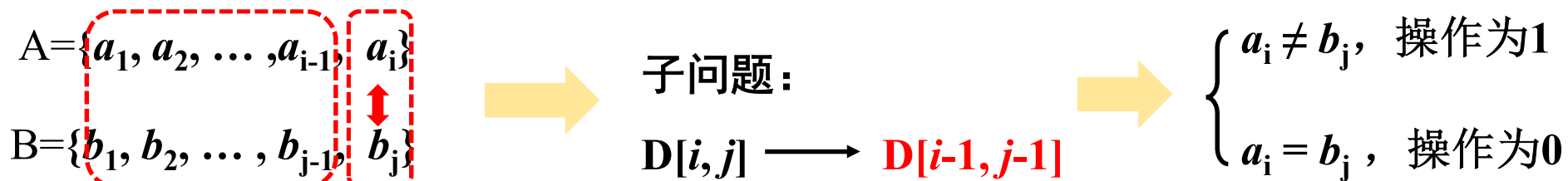
在插入操作中  $\mathbf{D[i, j] = D[i, j-1] + 1}$

最优子结构性质



# 递归结构

$D[i, j]$



在替换操作中  $D[i, j] = D[i-1, j-1] + \begin{cases} 1, & \text{if } a_i \neq b_j, \\ 0, & \text{if } a_i = b_j \end{cases}$

最优子结构性质

如果  $a_i = b_j$  时，则不需要替换，操作步骤为0，否则则需要替换，操作步骤为1



# 递归结构

综上所述：最小编辑距离的递归式

$$D[i, j] = \min \begin{cases} D[i-1, j] + 1 & \text{删除} \\ D[i, j-1] + 1 & \text{插入} \\ D[i-1, j-1] + \begin{cases} 1, & \text{if } a_i \neq b_j, \\ 0, & \text{if } a_i = b_j \end{cases} & \text{替换} \end{cases}$$

X \ Y	$j = 0$	$j = 1$	$j = 2$	...	$j = m$
$i = 0$	0	1	2		m
$i = 1$	1				
$i = 2$	2				
...					
$i = n$	n				

← 插入

- m次操作将空串变成长度为j的字符串
- n次操作将长度为i的字符串变成空串

↑ 删除



## 实例

**X**={A, G, G, C, A, T, A, G, C, T}

**Y**={C, A, G, T, A, T, C, C, T}

**$m = \text{length}(X) = 10$**

**$n = \text{length}(Y) = 9$**

**求解将序列X变成序列Y的最小编辑距离？**



# 实例

## 构建D[i,j]数组

$$D[1,1] = \min \begin{cases} D[0,1] + 1 = 2 \\ D[1,0] + 1 = 2 \\ D[0,0] + 1 = 1 \end{cases}$$

$A \neq C, x_1 \neq y_1$

$$D[1,2] = \min \begin{cases} D[0,2] + 1 = 3 \\ D[1,1] + 1 = 2 \\ D[0,1] + 0 = 1 \end{cases}$$

$A = A, x_1 = y_2$

$$D[1,3] = \min \begin{cases} D[0,3] + 1 = 4 \\ D[1,2] + 1 = 2 \\ D[0,2] + 1 = 3 \end{cases}$$

$A \neq G, x_1 \neq y_3$

	Y	0	C	A	G	T	A	T	C	C	T
X	D[i,j]	j=0	j=1	j=2	j=3	j=4	j=5	j=6	j=7	j=8	j=9
0	i=0	0	1	2	3	4	5	6	7	8	9
A	i=1	1	C[1,1]	C[1,2]	C[1,3]						
G	i=2	2									
G	i=3	3									
C	i=4	4									
A	i=5	5									
T	i=6	6									
A	i=7	7									
G	i=8	8									
C	i=9	9									
T	i=10	10									

以此类推……

$$D[i,j] = \min \begin{cases} D[i-1,j] + 1 & \text{删除} \\ D[i,j-1] + 1 & \text{插入} \\ D[i-1,j-1] + \begin{cases} 1, & \text{if } a_i \neq b_j \\ 0, & \text{if } a_i = b_j \end{cases} & \text{替换} \end{cases}$$



# 实例

构建D[i,j]数组

	Y	0	C	A	G	T	A	T	C	C	T
X	D[i,j]	j=0	j=1	j=2	j=3	j=4	j=5	j=6	j=7	j=8	j=9
0	i=0	0	1	2	3	4	5	6	7	8	9
A	i=1	1	1	1	2	3	4	5	6	7	8
G	i=2	2	2	2	1	2	3	4	5	6	7
G	i=3	3	3	3	2	2	3	4	5	6	7
C	i=4	4	3	4	3	3	3	4	4	5	6
A	i=5	5	4	3	4	4	3	4	5	5	6
T	i=6	6	5	4	4	4	4	3	4	5	5
A	i=7	7	6	5	5	5	4	4	4	5	6
G	i=8	8	7	6	5	6	5	5	5	5	6
C	i=9	9	8	7	6	6	6	6	5	5	6
T	i=10	10	9	8	7	6	7	6	6	6	5



## 实例

回溯最小编辑路径  $D[i, j] = 5$ , 即: X和Y序列的最小编辑距离为5

	Y	0	C	A	G	T	A	T	C	C	T
X	D[i, j]	j = 0	j = 1	j = 2	j = 3	j = 4	j = 5	j = 6	j = 7	j = 8	j = 9
0	i = 0	0	1	2	3	4	5	6	7	8	9
A	i = 1	1	1	1	2	3	4	5	6	7	8
G	i = 2	2	2	2	1	2	3	4	5	6	7
G	i = 3	3	3	3	2	2	3	4	5	6	7
C	i = 4	4	3	4	3	3	3	4	4	5	6
A	i = 5	5	4	3	4	4	3	4	5	5	6
T	i = 6	6	5	4	4	4	4	3	4	5	5
A	i = 7	7	6	5	5	5	4	4	4	5	6
G	i = 8	8	7	6	5	6	5	5	5	5	6
C	i = 9	9	8	7	6	6	6	6	5	5	6
T	i = 10	10	9	8	7	6	7	6	6	6	5

↖ : 替换;  
← : 插入;  
↑ : 删除;



# 实例 回溯最小编辑路径

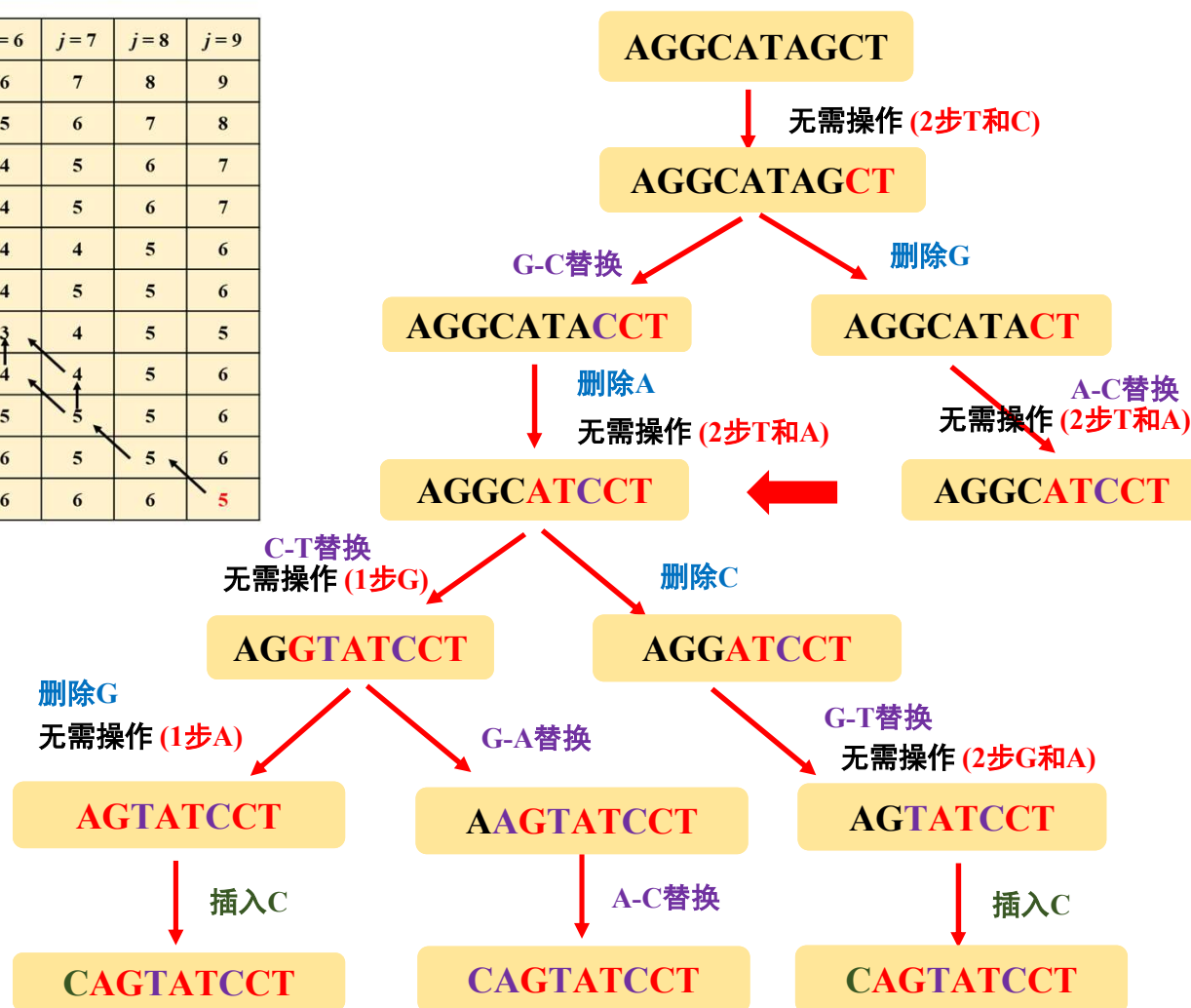
Y	0	C	A	G	T	A	T	C	C	T
---	---	---	---	---	---	---	---	---	---	---

X	D[i,j]	j=0	j=1	j=2	j=3	j=4	j=5	j=6	j=7	j=8	j=9
0	i=0	0	1	2	3	4	5	6	7	8	9
A	i=1	1	1	1	2	3	4	5	6	7	8
G	i=2	2	2	2	1	2	3	4	5	6	7
G	i=3	3	3	3	2	2	3	4	5	6	7
C	i=4	4	3	4	3	3	3	4	4	5	6
A	i=5	5	4	3	4	4	3	4	5	5	6
T	i=6	6	5	4	4	4	4	3	4	5	5
A	i=7	7	6	5	5	5	4	4	4	5	6
G	i=8	8	7	6	5	6	5	5	5	5	6
C	i=9	9	8	7	6	6	6	5	5	5	6
T	i=10	10	9	8	7	6	7	6	6	6	5

编辑操作：5次

$X = \{A, G, G, C, A, T, A, G, C, T\}$

$Y = \{C, A, G, T, A, T, C, C, T\}$







## 拓展—带权最小编辑距离

最小编辑距离只求解了需要操作的次数，如果对删除、插入和替换三种不同的操作，赋予不同的代价权重呢？

分别用  $w_1$ ,  $w_2$  和  $w_3$  表示进行一次删除、插入和替换的代价

$$D[i, j] = \min \begin{cases} D[i-1, j] + 1 & \text{删除} \\ D[i, j-1] + 1 & \text{插入} \\ D[i-1, j-1] + \begin{cases} 1, & \text{if } a_i \neq b_j, \\ 0, & \text{if } a_i = b_j \end{cases} & \text{替换} \end{cases}$$

递归式

$$s[0, 0] = 0$$

$$s[i, 0] = w_1 * i, 1 \leq i \leq n$$

$$s[0, j] = w_2 * j, 1 \leq j \leq m$$

$$s[i, j] = \min \begin{cases} s[i-1, j] + w_1 \\ s[i, j-1] + w_2 \\ s[i-1, j-1] + w_3/0 \end{cases}$$



## 总结—动态规划算法设计要点

- 引入参数来界定子问题的边界，注意子问题的重叠程度。
- 给出带界参数的优化函数定义与优化函数的递推关系，找到递推关系的初值。
- 判断该优化问题是否满足优化原则。
- 考虑是否需要标记函数。
- 采用自向上的实现技术，从最小的子问题开始迭代计算，计算中用备忘录保留优化函数和标记函数的值。
- 动态规划算法的时间复杂度是对所有子问题(备忘录)的计算工作量求和(可能需要追踪解的工作量)。
- 动态规划算法一般使用较多的存储空间，这往往成为限制动态规划算法使用的瓶颈因素。