

Редакционное расстояние

Расстояние Левенштейна

Задача о выравнивании

SUNNY SNOWY

1. $\begin{matrix} \text{SUNNY} \\ \text{SNOWY} \end{matrix}$ 3 $\begin{matrix} \text{D O X D b} \\ \text{D P O X b} \end{matrix}$ 3

2. $\begin{matrix} \text{SUN-NY} \\ \text{S-NOWY} \end{matrix}$ 3 $\begin{matrix} \text{D - O X D b} \\ \text{D P O X - b} \end{matrix}$ 2

Задача о редакционном расстоянии

Три типа операций:

1. замена $\begin{matrix} a \\ b \end{matrix}$
2. удаление $\begin{matrix} a \\ - \end{matrix}$
3. вставка $\begin{matrix} - \\ b \end{matrix}$

Вопрос: сколько нужно редакций для преобразования одного слова в другое?

$S_1 = \text{POLYNOMIAL}$

$S_2 = \text{EXPONENTIAL}$

$E(S_1, S_2) = ?$

Динамическое программирование:

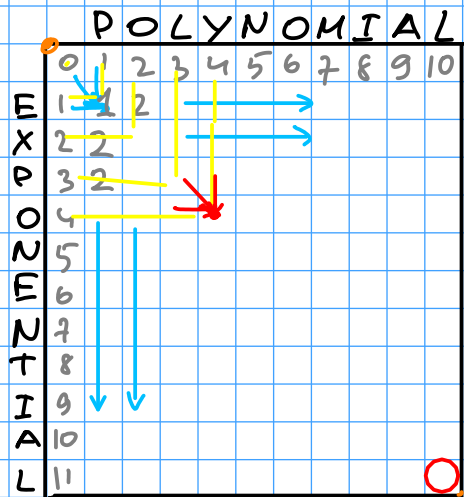
1. $E[i, j]$ - выравнивание $S_1[1, i]$ и $S_2[1, j]$

2. $E[0, j] = j$
 $E[i, 0] = i$

$$E[i, j] = \min \left\{ \begin{array}{l} E[i-1, j-1] + [S_1[i] \neq S_2[j]], \\ \text{вставка} \quad E[i, j-1] + 1, \\ \text{удаление} \quad E[i-1, j] + 1 \end{array} \right\}$$

← стоимость операции

3. Порядок вычисления



По строкам /
по столбцам

Время работы: $O(|S_1| \cdot |S_2|)$

Память: $O(|S_1| \cdot |S_2|)$

Можно восстановить путь, запомнив все рёбра.

Пусть $|S_1| = |S_2| = 10^9$

Время: $|S_1| \cdot |S_2| = 10^{18}$

Умеем: $3 \cdot 10^9$ операций в секунду

День: $3600 \cdot 24 = 10^5$

Потребуется $\sim 3 \cdot 10^3$ дней.

Память: проблема.

Можно хранить только две строки \Rightarrow

Память $O(|S_1| + |S_2|)$

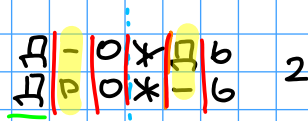
Как восстановить порядок последовательности?

Алгоритм Хурмидера

1. $E(S_{11}, S_2) = E(\text{rev}(S_1), \text{rev}(S_2))$

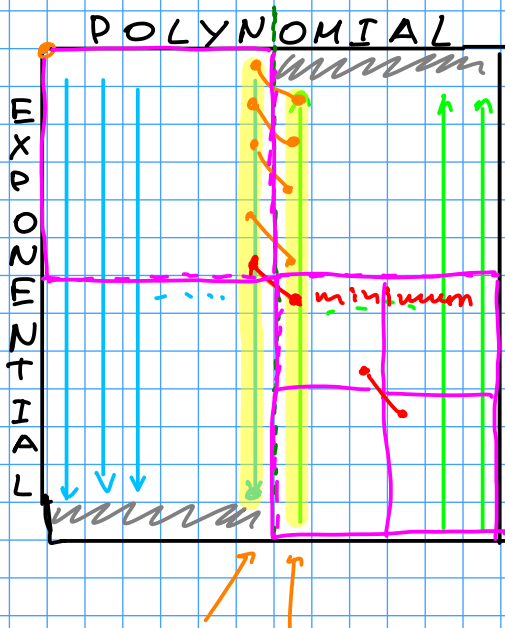
2. $\forall S_{11} \circ S_{12} = S_1 \Rightarrow \exists S_{21} \circ S_{22} :$
 $E(S_{11}, S_{21}) + E(S_{12}, S_{22}) = E(S_1, S_2)$

$S_1 = \text{POLYNOMIAL}$
 $S_2 = \text{EXPONENTIAL}$



$S_{11} = \text{A O X} \quad S_{12} = \text{A B}$

$S_{21} = \text{A P O X} \quad S_{22} = \text{B}$



Первое итерация:

1. Разбиваем $S_1 = S_{11} \circ S_{12}$

$|S_{11}| = |S_{12}| = |S_1| / 2$

2. Вычисляем:

- $E(S_{11}, S_2)$
- $E(\text{rev}(S_{12}), \text{rev}(S_2))$

$E(S_{11}, S_2[0, j])$

$E(\text{rev}(S_{12}), \text{rev}(S_2[j, n]))$

Повторяем рекурсивно для левого верхнего и правого нижнего прямоугольника.

Время: $|S_1| \cdot |S_2| + \frac{|S_1| \cdot |S_2|}{2} + \frac{|S_1| \cdot |S_2|}{4} \dots =$

$= 2|S_1| \cdot |S_2| = O(|S_1| \cdot |S_2|)$

Память: $O(|S_1| + |S_2|)$

NB: Meet-in-the-middle