

# Выравнивание последователей

Беляков Дмитрий

December 29, 2021

## 1 Описание задачи

Задача: Найти и исследовать алгоритм кратчайшего выравнивания последовательностей  $a$  и  $b$  (Нахождения наименьшей по длине последовательности в которой  $a$  и  $b$  являются подпоследовательностями)

## 2 Описание решения

Используем следующий алгоритм(алгоритм Ниддлмана-Вунша): Возьмем две последовательности:  $a$  длины  $n$  и  $b$  длины  $m$

- Составим таблицу  $n+1$  на  $m+1$
- Первый столбец заполним числами начиная с 0 до  $-n$
- Первую строку заполним числами начиная с 0 до  $-m$
- С оставшимися клетками таблицами будем действовать следующим образом:

$$F(i, j) = \max \begin{cases} F(i-1, j-1) + g(i, j), \\ F(i-1, j) - 1, \\ F(i, j-1) - 1. \end{cases}$$
$$\text{где } g(i, j) = \begin{cases} 1, & \text{if } s1[i] = s2[j], \\ -1, & \text{if } s1[i] \neq s2[j] \end{cases}$$

Затем выполним обратный обход для нахождения ответа

### 3 Описание данных

Рассмотрим несколько наборов данных:

### 3.1 1

Строки: aaaaaaabbbbbbb aaacaaaabbbbcbbbb Результат: aaa-aaaabbbb-bbb  
aaacaaaabbbbcbbbb

### 3.2 2

Строки: ааааabb bбааааа Результат -ааааabb bбааааа-

### 3.3 3

Строки aaabcbaaa bbaaaaccc Результат: aaabcb-aaa— —b-baaaaccc

## 3.4 4

Рассмотрим геномы человека (более 15 тысяч символов):

### 3.5 Первая строка

aagaggactccaagcgccatggccgctgccgcccagcccgggtcgcgctacttgcttaggcaactgcaacggcgagcgtgccagtgcccaactcattctcatac

### 3.6 Вторая строка

gaagaggactccaagcgccatggccgctgccgcccagcccgggtcgcgctacttgctgaggcaactgcaacgcgcagcgtgccagtgcccaactcattctcata

### 3.7 Решение

-aagaggactccaagcgccatggccgctgccgccgagccgggtcgctacttgctt-aggcaactgcaacgcgcagcgtgccagtgccc-  
 a-a-ctc-a-tt-ct-c-a-t-ac-tt-a-c-tc-c-c-a-agatggctgtttcaaatattagatatggagcagcttacaaggaagtagga  
 atcacttcgagagccatcaaaccacactgggactgattgatcctctgcacacctccacatgcctgccgagtggtcgcaacagtggtttgatgtgctttgc  
 gtcagaaatcccgatgatcttgaagcaaggtctcatatgcacttggcaagtgttttctggcatcggtttggaaatgctggtgttcatctgtgcatggaatgt  
 gccagatgtttccagagcgacacctggagatggcagaaatac-tgggagccgacaccgcactgccaggatccaagatgcagggctggtgttggcagacacg  
 gtccccagtcagaagaggatctggctgctctgtttgaagcttcaatgaaactgtattaattgtcattttaactgaaagaattaccgctggccattgtagtgtgag  
 ttcatgtagcaaatgagcaaaaactcagtatctatcaaaagtgtaaattatatttctatgcctagtaattcacttcatgtctaaaaatttatctgatagaacac  
 tgaggcaactgcaacgcgcagcgtgccagtgcccaactcattctcatacttactcccaagcccctggactttcaccttctgggaaaaacaacagattatgcctttga  
 tatcacttcgagagccatcaaaccacactgggactgattgatcctctgcacacctccacatgcctgccgagtggtcgcaacagtggtttgatgtgctttgc  
 cgtcagaaatcccgatgatcttgaagcaaggtctcatatgcacttggcaagtgttttctggcatcggtttggaaatgctggtgttcatctgtgcatggaatg  
 cgcccgatgtttccagagcgacacctggagatggcagaaata-ttgggagccgacaccgcactgccaggatccaagatgcagggctggtgttggcagaca  
 ccctgtccccagtcagaagaggatctggctgctctgtttgaagcttcaatgaaactgtattaattgtcattttaactgaaagaattaccgctggccattgtagtgc  
 ttcatgtagcaaatgagcaaaaactcagtatctatcaaaagtgtaaattatatttctatgcctagtaattcacttcatgtctaaaaatttatctgatagaaca  
 -a-c-a-ga--ag--c-at--g-g--