

# Выравнивание двух строк по алгоритму Нидлмана-Вунша. Отчет по задаче.

Владимирова Элина

23 декабря 2020 г.

## 1 Резюме

Для решения задачи об оптимальном выравнивании двух последовательностей был реализован классический алгоритм Нидлмана-Вунша.

Входные данные - длины двух строк. Преобразованиям подвергаются построенные по ним путем случайной генерации символов из заданного заранее алфавита буквенные последовательности. Результат работы программы - две строки, полученные из исходных вставками пропусков для достижения наибольшей посимвольной схожести последовательностей при их наименьшей длине.

## 2 Постановка задачи

Необходимо найти оптимальное выравнивание 2 заданных последовательностей, состоящих из символов алфавита  $\{A, G, C, T, \}$  т.е. выявить наименее трудозатратную комбинацию элементарных преобразований строк, позволяющую получить последовательности с максимальным совпадением символов на одинаковых местах, имеющие при этом минимальную длину.

## 3 Допущения

Строки генерируются внутри программы, их длины задаются вручную. Алгоритм, однако, не привязан к содержимому строк, что позволяет регулировать как алфавит, так и непосредственно строки.

## 4 Описание решения

Был реализован классический алгоритм Нидлмана-Вунша, предусматривающий работу с матрицей размера  $(n + 2) * (m + 2)$ , где  $n, m$  - длины строк.

## 5 Результаты

Правильность решения была проверена на строках малой длины (4, 15 и 40 символов), в том числе на строках с сравнительно большим относительно длин модулем разности длин (строки

длин 2 и 4, 3 и 15, 5 и 40 символов). Были проведены тесты как с генерируемыми случайно и потому значительно отличающимися друг от друга, так и со схожими, отличающимися 1-3 символами строками, а также со строкой и некоторым ее фрагментом:

			40	10
			20	3
AGA	ACGT	abracaabra	GCTGCCGGGATGCCACAGAGGCCCATAGCTGATGCCTGG	CAATCATCTA
AGA	TGCA	abracadabra	ATGGCTATGCAAAGTAAGTA	CAT
AGA	ACG_T	abraca_abra	GCTGCCGGGATGCCACAGAGGCCCATAGCTGATGCCTGG	CAATCATCTA
AGA	_TGCA	abracadabra	_ATG__GCTATGC_A_A_AG_____TAAG_T_A_____	CA_T_____

120

120

C66T6CTAGACT6TAT6C6TCACCTTCCAACCTC6AGTCTAATATGTACCGAGTTGACGTAGATCGTGATAGAGTACCGCACACGTAGGTGATCGATAAATTGGAATGCA6TTTGTCAATAG  
CATGCCCTGGC6TC66CC6TCCTTGTCAAACCTTCCAATTGTAAAGCATCGCTTCGAATTACTGATGGTCGCTCAGAGGCTCCAGAGCCTGGACCTTTGTTCACACGGACCGGACGACG

C66T6C\_\_TAGACT6TAT6\_C\_6TCACCTT\_\_CCAA\_\_CT\_CGAGTCT\_\_AATA\_\_T\_G\_TACCGAGTTGACGTAGATCG\_T\_GAT\_AGAG\_\_TACCGCACAC\_C\_\_GTAGGT\_\_GAT\_CGA\_TAAATTGGAATGCA\_GTTTGTCAATAG  
C\_ATGCCCT\_66C\_6T\_C66CC6TC\_CTTGT\_C\_AAACTTCCAAT\_TGTAA\_AGCATCGCTTC\_GAATT\_AC\_T\_GAT\_GGTCGCTCAGAGGCT\_CC\_\_AGAGCCTGGACCTTTG\_TTC\_AACA\_\_C6GACCGGACG\_\_\_\_\_CA\_\_G

Программа успешно обрабатывает последовательности длиной до 10 тыс. знаков - по преодолении этого порога возникают проблемы с памятью, причиной которым - невозможность хранения выстраиваемой в ходе решения матрицы и произведения внутри нее предусмотренных алгоритмом операций.