

# Implementação do algoritmo Needleman-Wunsch

O algoritmo Needleman-Wunsch é utilizado para alinhamento global de sequências de DNA ou aminoácidos. Esse algoritmo retorna o melhor alinhamento possível, porém, o resultado obtido não necessariamente tem uma significância biológica, nesse sentido, é muito importante a escolha de um sistema de pontuação e/ou matriz de substituição adequados ao problema.

Implementou-se o algoritmo, em Python 3.8, para o alinhamento de duas sequências. O código fonte está disponível no *Github* (<https://github.com/glenjasper/needleman-wunsch.git>). O algoritmo aceita um arquivo fasta contendo as duas sequências biológicas, para o qual é requerido escolher a matriz de substituição (BLOSUM ou PAM) e o *gap penalty* para sequências proteicas, entretanto para sequências de DNA podem-se configurar o *match* e o *mismatch*. Além do alinhamento feito, o algoritmo retorna a matriz de pontuação em um arquivo TXT. Também, é possível visualizar o alinhamento através do *framework* Dash (<https://dash.plotly.com>), que oferece uma visualização simples e amigável em HTML.

## Pré-requisitos

pip install dash

pip install dash-bio

## Uso básico

```
$ python needleman_wunsch.py --help
```

```
usage: needleman_wunsch.py [-h] -t {nt,aa} [-sm {BLOSUM45, BLOSUM50, BLOSUM62, BLOSUM80,
BLOSUM90, PAM30, PAM70, PAM250}] -f FILE [-m MATCH] [-mi MISMATCH_PENALTY] [-gap
GAP_PENALTY] [-o FOLDER] [--version]
```

Implementation of the Needleman-Wunsch algorithm

optional arguments:

-h, --help show this help message and exit

-t {nt,aa}, --type {nt,aa}

nt: Nucleotide sequence | aa: Amino acid sequence

-sm {BLOSUM45, BLOSUM50, BLOSUM62, BLOSUM80, BLOSUM90, PAM30, PAM70, PAM250}, --  
substitution\_matrix {BLOSUM45, BLOSUM50, BLOSUM62, BLOSUM80, BLOSUM90, PAM30, PAM70,  
PAM250}

Substitution Matrix type (Only for amino acid sequence) [default: BLOSUM62].

-f FILE, --fasta FILE

Fasta file

-m MATCH, --match MATCH

Match value (Only for nucleotide sequence) [default: 1].

-mi MISMATCH\_PENALTY, --mismatch\_penalty MISMATCH\_PENALTY

Mismatch penalty value (Only for nucleotide sequence) [default: 0].

-gap GAP\_PENALTY, --gap\_penalty GAP\_PENALTY

Gap penalty value [default: 0].

-o FOLDER, --output FOLDER

Output folder

--version show program's version number and exit

## Parâmetros

Parâmetro	Descrição	Possíveis valores	Default
<b>-t   --type</b>	Tipo de sequências a alinhar. Podem ser de aminoácidos ou nucleotídeos.	nt, aa	
<b>-sm   --substitution_matrix</b>	Matriz de substituição BLOSUM ou PAM, usado quando as sequencias forem proteicas.	BLOSUM45, BLOSUM50, BLOSUM62, BLOSUM80, BLOSUM90, PAM30, PAM70, PAM250	BLOSUM62
<b>-f   --fasta</b>	Arquivo fasta que contem as sequências biológicas a serem alinhadas.		
<b>-m   --match</b>	Valor de <i>match</i> . Usado quando as sequencias forem de DNA.		1
<b>-mi   --mismatch_penalty</b>	Valor de <i>mismatch</i> . Usado quando as sequencias forem de DNA.		0
<b>-gap   --gap_penalty</b>	Valor de <i>gap penalty</i> .		0
<b>-o   --output</b>	Pasta de saída.		

## Exemplos

1. Alinhar as sequências proteicas DRQTAQAAGTTTIT e DRNTAQLLGTDIT (contidas no arquivo **file.fa**), com a matriz de substituição BLOSUM80 e *gap penalty* -1.

```
$ python needleman_wunsch.py -t aa -f file.fa -gap -1 -o out_align1
```

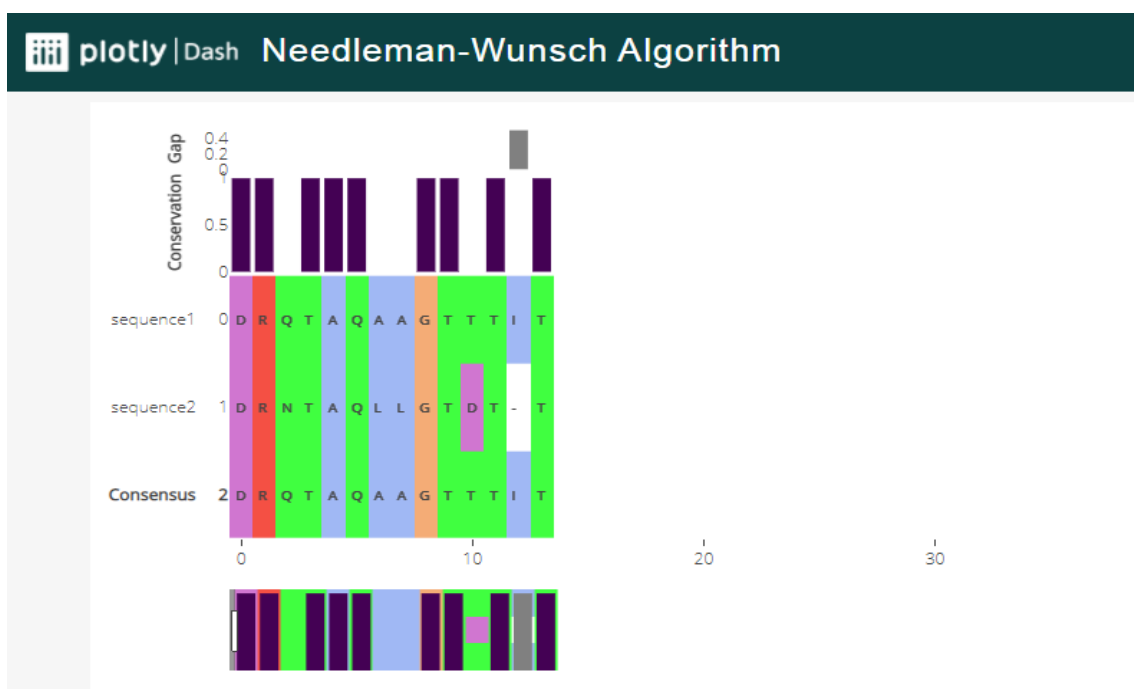
Output:

```
2021-02-05 00:14:21 #####
2021-02-05 00:14:21 ##### Needleman-Wunsch Algorithm #####
2021-02-05 00:14:21 #####
2021-02-05 00:14:21 Input:
2021-02-05 00:14:21   Fasta file: C:\Users\Glen\Dropbox\UFMG\Disciplinas\Bioinformática\Atividades\TP1\script\needleman-wunsch\file.fa
2021-02-05 00:14:21
2021-02-05 00:14:21 Parameters:
2021-02-05 00:14:21   Matrix: BLOSUM62
2021-02-05 00:14:21   Gap penalty: -1
2021-02-05 00:14:21
2021-02-05 00:14:21 Alignment:
2021-02-05 00:14:21   Score: 42
2021-02-05 00:14:21
2021-02-05 00:14:21 sequence1  DRQTAQAAGTTTIT
2021-02-05 00:14:21             || ||| - || | |
2021-02-05 00:14:21 sequence2  DRNTAQLLGTDI-T
2021-02-05 00:14:21
2021-02-05 00:14:21 Matrix file: C:\Users\Glen\Dropbox\UFMG\Disciplinas\Bioinformática\Atividades\TP1\script\needleman-wunsch\out_align1\alignment_matrix.txt
2021-02-05 00:14:21
2021-02-05 00:14:21 Running Dash's alignment viewer
2021-02-05 00:14:21 Run the address http://127.0.0.1:8050 in your browser
Dash is running on http://127.0.0.1:8050/
* Serving Flask app "needleman_wunsch" (lazy loading)
* Environment: production
  WARNING: This is a development server. Do not use it in a production deployment.
  Use a production WSGI server instead.
* Debug mode: on
```

Matriz de substituição (ver arquivo **alignment\_matrix.txt**):

-	-	D	R	Q	T	A	Q	A	A	G	T	T	T	I	T
-	0	-1	-2	-3	-4	-5	-6	-7	-8	-9	-10	-11	-12	-13	-14
D	-1	6	d	5	l	4	l	3	l	2	l	1	l	0	l
R	-2	5	u	11	d	10	l	9	l	8	l	7	l	6	l
N	-3	4	u	10	u	11	d	10	d	9	l	8	d	7	l
T	-4	3	u	9	u	10	u	16	d	15	l	14	l	13	l
A	-5	2	u	8	u	9	u	15	u	20	d	19	l	18	d
Q	-6	1	u	7	u	13	d	14	u	19	u	25	d	24	l
L	-7	0	u	6	u	12	u	13	u	18	u	24	u	24	d
L	-8	-1	u	5	u	11	u	12	u	17	u	23	u	23	d
G	-9	-2	u	4	u	10	u	11	u	16	u	22	u	23	d
T	-10	-3	u	3	u	9	u	15	d	15	u	21	u	22	d
D	-11	-4	d	2	u	8	u	14	u	14	u	20	u	21	u
T	-12	-5	u	1	u	7	u	13	d	14	d	19	u	20	d
T	-13	-6	u	0	u	6	u	12	d	13	d	18	u	19	d

Visualização do alinhamento com o *framework* Dash (HTML+CSS), através do endereço <http://127.0.0.1:8050>:



2. Alinhar as glicoproteínas Spike sp|P11223|SPIKE\_IBVB e sp|P12651|SPIKE\_IBVM (contidas no arquivo **file2.fa**), com a matriz de substituição BLOSUM62 e *gap penalty* -2.

```
$ python needleman_wunsch.py -t aa -f file2.fa -gap -2 -o out_align2
```

Output:

```

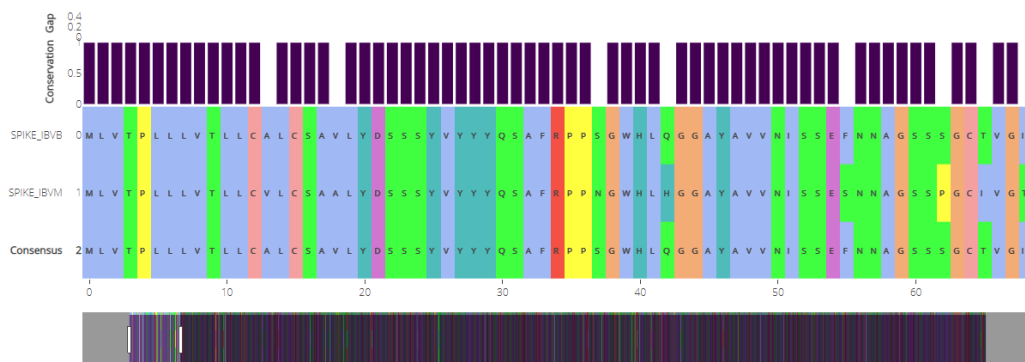
2021-02-05 00:25:56 #####
2021-02-05 00:25:56 ##### Needleman-Wunsch Algorithm #####
2021-02-05 00:25:56 #####
2021-02-05 00:25:56 Input:
2021-02-05 00:25:56 Fasta file: C:\Users\Glen\Dropbox\UFMG\Disciplinas\Bioinformática\Atividades\TP1\script\needleman-wunsch\file2.fa
2021-02-05 00:25:56
2021-02-05 00:25:56 Parameters:
2021-02-05 00:25:56 Matrix: BLOSUM62
2021-02-05 00:25:56 Gap penalty: -2
2021-02-05 00:25:56
2021-02-05 00:25:56 Alignment:
2021-02-05 00:26:13 Score: 5869
2021-02-05 00:26:13
2021-02-05 00:26:13 sp|P11223|SPIKE_IBVB MLVTPLLLVTLLCALCSAVLYDSSSYVYVYQSAFRPPSGMHLGGGAYAVWNISSEFNAG
2021-02-05 00:26:13 |||||-----|||-----|||-----|||-----|||-----|||
2021-02-05 00:26:13 sp|P12651|SPIKE_IBVM MLVTPLLLVTLLCVLCSAALYDSSSYVYVYQSAFRPPNGMHLGGGAYAVWNISSEFNAG
2021-02-05 00:26:13 |||||-----|||-----|||-----|||-----|||-----|||
2021-02-05 00:26:13 sp|P11223|SPIKE_IBVB SSSGCTVGIIHGGRVWASSIAMTAPSSGMAWSSSQFCTAHCNFSDTTVFVTHCYKHGGC
2021-02-05 00:26:13 ||-||-||-|||-----|||-----|||-----|||-----||-||
2021-02-05 00:26:13 sp|P12651|SPIKE_IBVM SSPGCTVGTIHGGRVWASSIAMTAPSSGMAWSSSQFCTAHCNFSDTTVFVTHCYKYDGC
2021-02-05 00:26:13 ||-----||-||-||-|||-----|||-----|||-----|||-----|||
2021-02-05 00:26:13 sp|P11223|SPIKE_IBVB PLTGMQLQNLIRVSAMKNGQLFYNLTVSAKYPTFRSFQCVNLTSSVYNGDLVYTSNET
2021-02-05 00:26:13 ||-----||-||-||-|||-----|||-----|||-----|||-----|||
2021-02-05 00:26:13 sp|P12651|SPIKE_IBVM PETGMQLQNFLRVSAKNGQLFYNLTVSAKYPTFKSFQCVNLTSSVYNGDLVYTSNET
2021-02-05 00:26:13 ||-----||-||-||-|||-----|||-----|||-----|||-----|||
2021-02-05 00:26:13 sp|P11223|SPIKE_IBVB IDVTSAGVYFKAGGPITYKVMREVKALAYFVNGTAQDVLCDGSPRGLLACQYNTGNFSD
2021-02-05 00:26:13 -||-----|||-----|||-----|||-----|||-----|||-----|||
2021-02-05 00:26:13 sp|P12651|SPIKE_IBVM TDVTSAGVYFKAGGPITYKVMRKVKALAYFVNGTAQDVLCDGSPRGLLACQYNTGNFSD
2021-02-05 00:26:13 ||-----|||-----|||-----|||-----|||-----|||-----|||
2021-02-05 00:26:13 sp|P11223|SPIKE_IBVB GFYPFTNSSLWKQKFIYVRENSVNTTCTLHNFIFHNETGANPNPSGVQNIQTYYQTKAQ5
2021-02-05 00:26:13 |||||-----|||-----|||-----|||-----|||-----|||
2021-02-05 00:26:13 sp|P12651|SPIKE_IBVM GFYPFTNSSLWKQKFIYVRENSVNTTFTLHNFIFHNETGANPNPSGVQNIQTYYQTKAQ5
2021-02-05 00:26:13 |||||-----|||-----|||-----|||-----|||-----|||
2021-02-05 00:26:13 sp|P11223|SPIKE_IBVB GYVNFNFSFLSSFYVYKSNFMGYSHPSCKFRLETINGLWFNLSVSIAYGPLQGGCKQ
2021-02-05 00:26:13 |||||-----|||-----|||-----|||-----|||-----|||
2021-02-05 00:26:13 sp|P12651|SPIKE_IBVM GYVNFNFSFLSSFYVYKSNFMGYSHPSCKFRLETINGLWFNLSVSIAYGPLQGGCKQ
2021-02-05 00:26:13 |||||-----|||-----|||-----|||-----|||-----|||
2021-02-05 00:26:13 sp|P11223|SPIKE_IBVB SVFKGRATCCYAYSVGGPSLCKGVYSGELDNFECGLLVVYTKSGGSRIQTATEPPVITQ
2021-02-05 00:26:13 |||-----|||-----|||-----|||-----|||-----|||-----|||
2021-02-05 00:26:13 sp|P12651|SPIKE_IBVM SVFSGRATCCYAYSVGGPSLCKGVYSGELDNFECGLLVVYTKSGGSRIQTATEPPVITR
2021-02-05 00:26:13 |||-----|||-----|||-----|||-----|||-----|||-----|||
2021-02-05 00:26:13 sp|P11223|SPIKE_IBVB HWYNNITLNTCDVYNIYGRTGQGFITNWDTSVSYWYADAGLAIDLTSGSIDIFWQGE
2021-02-05 00:26:13 -||-----|||-----|||-----|||-----|||-----|||-----|||
2021-02-05 00:26:13 sp|P12651|SPIKE_IBVM HWYNNITLNTCDVYNIYGRTGQGFITNWDTSVSYWYADAGLAIDLTSGSIDIFWQGE
2021-02-05 00:26:13 ||-----|||-----|||-----|||-----|||-----|||-----|||
2021-02-05 00:26:13 sp|P11223|SPIKE_IBVB YGLNYYKVNPCEDVMQQFVWSGGKLVGILTSRNETGSQLENQFYIKITNGTRFRFRSIT
2021-02-05 00:26:13 |||-----|||-----|||-----|||-----|||-----|||-----|||
2021-02-05 00:26:13 sp|P12651|SPIKE_IBVM YGLTYKVNPCEDVMQQFVWSGGKLVGILTSRNETGSQLENQFYIKITNGTRFRFRSIT
2021-02-05 00:26:13 |||-----|||-----|||-----|||-----|||-----|||-----|||

```

Matriz de substituição (ver arquivo `alignment_matrix.txt`). Apenas uma parte é mostrada porque a matriz é muito extensa:

-	-	M	L	V	T	P	L	L	L	V	T	L	L	C	A	L	C	S
-	0	-2	-4	-6	-8	-10	-12	-14	-16	-18	-20	-22	-24	-26	-28	-30	-32	-34
M	-2	5	d	3	l	1	l	1	l	-1	l	-3	l	-15	l	-19	l	-27
L	-4	3	u	9	d	7	l	5	l	3	l	1	d	-9	d	-13	l	-21
V	-6	1	u	7	u	13	d	11	l	9	l	7	l	-5	l	-7	l	-13
T	-8	-1	u	5	u	11	u	18	d	16	l	14	l	12	l	10	l	-8
P	-10	-3	u	3	u	9	u	16	u	25	d	23	l	21	l	19	l	1
L	-12	-5	u	1	d	7	u	14	u	23	u	29	d	27	d	25	d	7
L	-14	-7	u	-1	d	5	u	12	u	21	u	27	d	33	d	31	d	1
L	-16	-9	u	-3	d	3	u	10	u	19	u	25	d	31	d	37	d	1
V	-18	-11	u	-5	u	1	d	8	u	17	u	23	u	29	u	35	u	1
T	-20	-13	u	-7	u	-1	d	6	u	15	u	21	u	27	u	33	u	1
L	-22	-15	u	-9	d	-3	u	4	u	13	u	19	d	25	d	31	d	1
L	-24	-17	u	-11	d	-5	u	2	u	11	u	17	d	23	d	29	d	1
C	-26	-19	u	-13	u	-7	u	0	u	9	u	15	u	21	u	27	u	1
V	-28	-21	u	-15	u	-9	d	-2	u	7	u	13	u	19	u	25	u	1
L	-30	-23	u	-17	d	-11	u	-4	u	5	u	11	d	17	d	23	d	1
C	-32	-25	u	-19	u	-13	u	-6	u	3	u	9	u	15	u	21	u	1
S	-34	-27	u	-21	u	-15	u	-8	u	1	u	7	u	13	u	19	u	1
A	-36	-29	u	-23	u	-17	u	-10	u	-1	u	5	u	11	u	17	u	1
A	-38	-31	u	-25	u	-19	u	-12	u	-3	u	3	u	9	u	15	u	1
L	-40	-33	u	-27	d	-21	u	-14	u	-5	u	1	d	7	d	13	d	1
Y	-42	-35	u	-29	u	-23	u	-16	u	-7	u	-1	u	5	u	11	u	1
D	-44	-37	u	-31	u	-25	u	-18	u	-9	u	-3	u	3	u	9	u	1
S	-46	-39	u	-33	u	-27	u	-20	u	-11	u	-5	u	1	u	7	u	1
S	-48	-41	u	-35	u	-29	u	-22	u	-13	u	-7	u	-1	u	5	u	1
S	-50	-43	u	-37	u	-31	u	-24	u	-15	u	-9	u	-3	u	3	u	1

Visualização do alinhamento com o *framework* Dash (HTML+CSS), através do endereço <http://127.0.0.1:8050>:



3. Alinhar as seqüências de DNA AATTTACGCGGCATTATAGATACAATCGTGTCT e GCAATTGGCCGGAATTTAATTGATACAGCGC (contidas no arquivo **file3.fa**), com valores de *match* +2, *mismatch* -1 e *gap penalty* -2.

```
$ python needleman_wunsch.py -t nt -f file3.fa -m 2 -mi -1 -gap -2 -o out_align3
```

Output:

Visualização do alinhamento com o *framework* Dash (HTML+CSS), através do endereço <http://127.0.0.1:8050>:

