

# Project 1: Multiple Sequence Alignment

Name: 吴曾宇

Student ID: 519021910981

GitHub: <https://github.com/DDuanCang/SJTU-AICourse-Proj1---MSA>

## 目录

Project 1: Multiple Sequence Alignment.....	1
Implementation .....	2
DP .....	2
A*.....	3
GA .....	3
Alignment and cost.....	5
DP: .....	5
A*: .....	6
GA: .....	8
Running time.....	9
DP: .....	9
A*: .....	9
GA: .....	9
Time complexity.....	9

## Implementation

### DP

DP 解决序列比对问题，我采用了序列比对动态规划矩阵来进行搜索。

	<i>i</i>	C	A	T	G	T
<i>j</i>	0	-1	-2	-3	-4	-5
A	-1	-1	1	0	-1	-2
C	-2	1	0	0	-1	-2
G	-3	0 <sup>③</sup>	0	-1	2	1
C	-4	-1	-1	-1	1	1
T	-5	-2	-2	1	0	3 <sup>①</sup>
G	-6	-3	-3	0	3	2

首先生成 score matrix，并填写出第一行第一列。再进行从第二行到最后一行，行中从左到右的 cost value 填写。填写规则为：将 blank 的左、左上、上方 blank 中的 cost value 加上这一次移动代表的 cost，再进行 cost value 比对，取其中的最小 cost 为结果，并记录下移动方向。其中右移代表横轴进行一个 GAP，下移代表纵轴进行一个 GAP，斜移代表进行一次对应位置字符的比对 MATCH or MIXMATCH。根据 proj1 任务书中给出的 Cost matrix 可以有，一次 GAP 增加 3cost，一次 MATCH 增加 0cost，一次 MISMATCH 增加 4cost。重复此法填写 score matrix 直至填写完整。

Score matrix 填写完毕后，从右下角向左上角进行回溯，获得 least cost 结果。结果中包含 least cost value 和对应的其中一种比对结果（score matrix 中的回溯结果可以有多种，即可能有多种有 least cost value 的比对方法，只是这里根据 proj1 的要求只取其中一种为比对结果）。

重复上述流程即可完成一次 query sequence 与 database sequence 的序列比对。记录下每次的比对 least cost value 与比对结果，遍历比对 database 中的所有 sequence，并进行比较，获得 cost value 最低的比对序列组。

多序列对比则在普通序列比对情况基础上增加一次，即对一个 database 中的 sequence 用两个 query sequence 同时进行比对，同样取最低 cost。

## A\*

A\*的大致流程与 DP 相似，不过在填写 score matrix 中进行 cost 计算时增加  $h(n)$ ，这里的  $h(n)$  取当前位置到最左上角的曼哈顿距离的三倍。取曼哈顿距离是因为这里的所在位置到左上角位置可以近似表示为回溯流程中当前位置到最终位置的距离；取三倍是因为曼哈顿距离的计算中取水平距离与垂直距离的和，这在表格填写过程中表示一次右移 GAP 或左移 GAP，因此在 cost 计算过程中以 GAP value 计算。

多序列对比则在普通序列比对情况基础上增加一次，即对一个 database 中的 sequence 用两个 query sequence 同时进行比对，同样取最低 cost。

## GA

（在经过各种测试后，由于运行时间问题，将迭代次数定为 100）

此次 GA 解决序列比对问题我使用了 python3.9 版本中的 scikit-opt 库进行辅助。该库的文档链接为：<https://scikit-opt.github.io/scikit-opt/#/zh/README>。具体使用函数为 `from sko.GA import GA`。该函数的参数如下：

入参	默认值	意义
func	-	目标函数
n_dim	-	目标函数维度
size_pop	50	种群规模
max_iter	200	最大迭代次数
probab_mut	0.001	变异概率
lb	-1	每个自变量的最小值
ub	1	每个自变量的最大值
constraint_eq	empty array	等式约束
constraint_uneq	empty array	不等式约束
precision	1e-7	精准度, int/float 或者它们组成的列表

本次序列比对中，使用 GA 的思路为：

eg:

```

# seq1: KJXXJAJKP
# seq2: KJOBL
# =>
# seq1: [K, J, X, X, J, A, J, K, P, -, -, -, -, -]
# seq2: [-, -, -, -, -, -, -, -, K, J, O, B, L]
# =>
# seq1: [0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1]
# seq2: [1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0]
# gap_num1 = len(test_seq2)
# gap_num2 = len(test_seq1)
# =>
# p: [0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0]

```

按照上述方法构成 p 为种群个体，fitness 函数（即引用库中的 func 参）为比对的 cost 结果，constraint\_eq 参数为保证每个种群个体中的 gap 数一定，lb 为 0（表示非 gap），ub 为 1（表示 gap），precaution 为 1（表示每个位置只存在 gap 与非 gap 情况）。

此结构的本质为通过 GA 找出两序列中能形成最低 cost 的 gap 位置。

多序列比对则是：

```

eg:
# seq1: KJXXJAJKP
# seq2: KJOBL
# seq3: KJP
# =>
# seq1: [K, J, X, X, J, A, J, K, P, -, -, -, -, -, -, -]
# seq2: [-, -, -, -, -, -, -, -, K, J, O, B, L, -, -, -]
# seq3: [-, -, -, -, -, -, -, -, -, -, -, -, K, J, P]
# =>
# seq1: [0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1]
# seq2: [1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 1, 1, 1]
# seq3: [1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0]
# gap_num1 = len(test_seq2) + len(test_seq3)

```

```
# gap_num2 = len(test_seq1) + len(test_seq3)
# gap_num3 = len(test_seq1) + len(test_seq2)
# =>
#
p[0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1,
, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0]
```

剩余部分与普通序列比对步骤相同。

## Alignment and cost

DP:

```
#####
KJXXJAJKPX_KJ__JXJKPX_K__JXXJAJK___PXKJJXJK_P_XKJXXJAJKPXKJXXJAJKHXX_JXXJAJKPXK
_JXXJAJKH__XKJXX
XHAPXJAJ_XXXJAJXDJ_AJXXXJAPXJAHXXXJAJXDJAJXXXJAJXXXJPP_X__JAJXX_X__HAPXJAJXX_XJ
AJXXXJA_JXXXJA__
cost: 161.0

#####
ILOTGJ_JLABWTSTGGONXJ MUTUXSJH_KWJHC__TOQHWGAGIWLZHWPK_ZUL_JTZWAKBWHXMIKLZJ_GL__
_X__BPAHO__HVOLZWOS_JJLPO
IPOT_WJJLAB__KSGZG_WJJ__KSPPOPHTSJEWJHCT__OOTHK__AXBKLBHWPKZULJPZKAKVKHXUIKLZ
JGLXBTG_HOHBJLZPP_SJJPJO_
cost: 148.0

#####
IHKKKKKKKKKKXG_WG_KKKP__KSKKKKBKKKPKHKKKXKBSKKPKWKKLKSRRKKWXPKBKKPKTSKHKKKK
LA__DKKYPKKKOP_HKKBWWLPPWKK
IKBSKKKK__WKKKKKKKKWWKKKKPGKKKKXXGGKRKWWK_WKKPKKKKKXKKRWK__MKKPKK_W_P__
KKKKKPGKKK__LBKWWK__KJ__
cost: 163.0

#####
```



#####

IPOTWJ,LABKSGZGWJJKSPPOPHTSJ\_EWJHCTOOTHRAXBKLBHWPKZULJPZKAKVKHXUIKLZJGLXBTGHOHB  
JLZPPSJJPJO

#####

IKBSKKKKWKKKKKKKKKWWKKK\_KPGKKKK\_KXXGGKRKWWK\_WKKPK\_KK\_K\_KXKK\_\_KRWK\_MKKPK\_\_KWPKKK\_  
 \_KK\_PGKK\_\_KLBKW\_\_WKKJ

#####

\_\_OPJPXJPPMJPPMX\_PP\_MJ\_PPXJPPOXPPXJJPJXXPXJPPOJPPMXPPOGP\_\_PXXP\_P\_\_\_\_OM\_\_PPXXPP  
OXPPXJPQXPPBJPPPPPPX

#####

ITPVKWSKXKUAXPVHXVOMMKHYBPABLLOBGKOLLJGXZGXL SOLAMOGKIGXBATBXMPJTCVMTAXVMPWWAW  
OMOUPHHZBITKKXLK

#####

#####

IPZJJMLTKJULOST\_KTJGLKJOBBLTXGKTP\_LUWWKOMOYJ\_BGALJUKLGLSVHWPBG\_WSLUKOBSOLOOK  
UKSARPPJ





CKTWLPGJTJJG-LG-JLS--P--PA---XAJPJ-P--J-PJP--J-P-JP-JPJP-JP-J-K-JP-J-XS-T-LXH--  
J-JOSAKPBHVBAGCPJ

cost: 276

#####

----I-P-P--V-KBKXWXKHS-AP--HV--X--XVO-J-M-RAKKP-J---VLL---J-B-WKOL--L-J--K--X-  
H-GXL

ITOJ-L-LLULST--KS-P---P-MOG-KJLWXOGOOKWL-L-SK--P-KWSOB-VOQGRBS---PPP-BXUA-

LLPKPO--M

cost: 256

#####

#####

run time: 1116 seconds

## Running time

**DP:**

2: 8 seconds

3: 3 seconds

**A\*:**

2: 11 seconds

3: 4 seconds

**GA:**

2: 1116 seconds

3: 由于我使用的方法对计算量要求过高，因此在迭代次数与种群规模较小时计算不出结果  
(代码是可以运行的能获得部分解，但仍存在问题，且由于时间问题，没有时间进行大量  
计算得出结果了)

## Time complexity

DP:  $O(n^2)$

A\*:  $O(n^2)$

GA: 与迭代次数(max\_iter)有关