

# Multiple Sequence Alignment

Xin Xu

## Dynamic Programming

Dynamic programming: Break up a problem into a series of overlapping problems, and build up solutions to larger and larger sub-problems. It's a little like greedy algorithm, but contrast to greedy, it can override past optimal local solution because it records all local solutions so that dynamic programming can always find the global optimal solution while greedy may result in local optimal solution.

Dynamic programming has three important parts: state, state transfer equation and boundary conditions, and I will explain my implementation in these three aspects.

## Pairwise Alignment

### State

The procession to go through each sequence is the state. In pairwise alignment, I use a 2D matrix to denote state. The axis of  $x$  means it has go through  $x$  chars of the  $sx$  sequence and it's the same with  $y$ .

### State Transfer Equation

The state transfer equation is:

```
1 for line in range(1, lenx):
2     cost_matrix[0][line] = line * GAP
3     for row in range(1, leny):
4         cost_matrix[row][line] = min([PairCost(sx[line], sy[row]) +
            cost_matrix[row - 1][line - 1], GAP + cost_matrix[row - 1][line], GAP +
            cost_matrix[row][line - 1]))
```

There are three conditions:

- $x_i$  and  $y_i$  matches. The cost is consist of  $x_i - y_i$  mismatch and past least cost.
- Leaving  $x_i$  unmatched. The cost is consist of GAP and past least cost.
- Leaving  $y_i$  unmatched. The cost is consist of GAP and past least cost.

The `Pairwise` function is:

```
1 def PairCost(sx, sy):
2     if sx == sy:
3         return MATCH
4     elif sx == "-":
5         return GAP
6     elif sy == "-":
7         return GAP
8     else:
9         return MISMATCH
```

## Boundary Conditions

The boundary conditions happen when all *sx* sequence chars are matched with GAP or all *sy* sequence chars are matched with GAP:

```
1 for row in range(leny):
2     cost_matrix[row][0] = row * GAP
3 for line in range(1, lenx):
4     cost_matrix[0][line] = line * GAP
```

## Result

The result and running are as follows:

[illegible]

The least cost of pairwise alignment by dynamic programming is: 107  
The base sequences is: IL0TGJLABWTSTGGONXJMTUXSXJHKWJHCTOQHGWAGIWLZHWPKZULJTZWAKBWHXMIKLZJGLXBPAHOHVOLZWOSJLJLP  
The best match is: IP0TWJLABKSGZGWJZKSPPOPHIT-SJ-EWJHCTOOTH-RAXBKLBIHWPCKZULJPZKAKVHXUJLKLZJGLXBTHGHOHBJLZPPSJJPJO  
IL0TGJLABWTSTG-GONXJ- --- -MTUXSXJHKWJHCT-OQHGWAGIWLZHWPKZULJTZWAKBWHXMIKLZJGLXBPAHOHVOLZWOSJLJLP-O  
IP0TWJLAB-KS-GZG-WJJKSPPOPHIT--SJ-EWJHCTOOTH-RAXBKLBIHWPCKZULJPZKAKVHXUJLKLZJGLXBTHGHOHBJLZPPSJJP-JJO  
The running time of pairwise alignment by dynamic programming is: 0.0009175

[illegible]

```
The least cost of pairwise alignment by dynamic programming is: 148
The base sequences is: MPPPJXPXGPGPPXXPPXJJPJPPXKPPSSPJJJJPXXPPPPPJPJPPXPPXCPJMMXXPKPSVGULMHHZPAWHTHKAHHUPAONAPJSMWPJGA
The best match is: OPJPXJPMJPMXMPMXMPPXJJPJPPXKXJJPJXXKJPPQJPPMXKPOGGPPXKPOGPXKPOGXPPBJJQXQPPBJJPPXPPX
MPPPJPX-PGP-JPPXPPPPJJPJPPXPPPPSPJJJJPXXK-PPPPJPJPPXPIJMMXXPKPSVGULMHZPAWHTHKAHHUPAONAP-JSMWP-JGA
--OPJXJP-PMJPMX-PP-MJ-PPXJPOGXPPXJJPJXXKJJPJOJPPMXKPOGP--PXXK-P----QM---PPXQOXKXJJP-QXQPBJ-PPXPPX
The running time of pairwise alignment by dynamic programming is: 0.572254s
```

```
The least cost of pairwise alignment by dynamic programming is: 131
The base sequences is: IPPVKBKXKKXKHSAPHVXXVQJMRAPKPVLLJBWKOLLJKXHXGLLCPAJOBKPGXBATGXMPOMCVZTAXVPAGKXGQJQOLJGWGKXLQ
The best match is: ITPVKMKKKKKXUAPVHVXVMMKHYBPALLBOBKOLLJGXZGKSLAMOGKIGXBATBXMPJTCVMTAXVMPWAWOMOPHHZBITKXK
IPPVKBKXKKXKHSAPHVXXVQJMRAPKPVLLJBWKOLLJKXHXGLLCPAJOBKPGXBATGXMPOMCVZTAXVP--AGKXGQJQO---LJGWGKXLQ
ITPVKMKKKKKXUAPVHVXVMMKHYBPALLBOBKOLLJGXZGKSL--AMOGKIGXBATBXMPJTCVMTAXVMPWAW--WOM--OUPHHZBITKXK
The running time of pairwise alignment by dynamic programming is: 0.574570s
```

## Time Complexity

Suppose the longest length of sequences  $sx$  and  $sy$  is  $l$ . The time to construct a cost matrix is  $O(l^2)$  and the time to trace the match path is  $O(l)$ . So, the total time complexity is  $O(l^2)$ .

## Three-Sequence Alignment

The implementation is almost the same with pairwise alignment.

## State

The procession to go through each sequence is the state. In three-sequence alignment, I use a 3D matrix to denote state. The axis of  $x$  means it has go through  $x$  chars of the  $sx$  sequence and it's the same with  $y$  and  $z$ .

## State Transfer Equation

The state transfer equation is:

```
1 for i in range(1, leni):
2     for j in range(1, lenj):
3         for k in range(1, lenk):
4             cost_matrix[i][j][k] = min([cost_matrix[i - 1][j - 1][k - 1] +
5                                         ThreeSequenceCost(si[i], sj[j], sk[k]), \
6                                         cost_matrix[i - 1][j][k - 1] + ThreeSequenceCost(si[i],
7                                         "-", sk[k]), \
8                                         cost_matrix[i][j - 1][k - 1] + ThreeSequenceCost("-",
9                                         sj[j], sk[k]), \
10                                         cost_matrix[i - 1][j - 1][k] +
11                                         ThreeSequenceCost(si[i], sj[j], "-"), \
12                                         cost_matrix[i][j][k - 1] +
13                                         ThreeSequenceCost("-", "-", sk[k]), \
14                                         cost_matrix[i][j - 1][k] +
15                                         ThreeSequenceCost("-", sj[j], "-"), \
16                                         cost_matrix[i - 1][j][k] +
17                                         ThreeSequenceCost(si[i], "-", "-"))])
```

There are three conditions:

- $x_i, y_i$  and  $z_i$  matches. The cost is consist of  $x_i, y_i$  and  $z_i$  mismatch and past least cost.
- one of  $x_i, y_i$  and  $z_i$  unmatched. The cost is consist of one block '-' and two chars mismatch and past least cost.
- two of  $x_i, y_i$  and  $z_i$  unmatched. The cost is consist of two blocks '-' and one chars mismatch and past least cost.

The `ThreeSequenceCost` function is:

```
1 def ThreeSequenceCost(c1, c2, c3):
2     return PairCost(c1, c2) + PairCost(c1, c3) + PairCost(c2, c3)
```

## Boundary Conditions

Suppose the search space of three-sequence alignment in dynamic programming is a cube. The boundary conditions happen in the face of the cube: it reduces to the situation in 2D. So, the boundary condition is:

```
1 for i in range(leni):
2     cost_matrix[i][0][0] = i * GAP * 2
3 for j in range(lenj):
4     cost_matrix[0][j][0] = j * GAP * 2
5 for k in range(lenk):
6     cost_matrix[0][0][k] = k * GAP * 2
7
8 for i in range(1, leni):
9     for j in range(1, lenj):
10         cost_matrix[i][j][0] = min(ThreeSequenceCost(si[i], sj[j], "-") +
11                                     cost_matrix[i - 1][j - 1][0], ThreeSequenceCost(si[i], "-", "-") +
12                                     cost_matrix[i - 1][j][0], ThreeSequenceCost("-", sj[j], "-") +
13                                     cost_matrix[i][j - 1][0])
14
15 for i in range(1, leni):
16     for k in range(1, lenk):
```

## Result

```
The least cost of three-sequence alignment by dynamic programming is: 456  
The base sequences are: ["TPTZJLMLTXIULOSTKTOJGLKIOBLTXGKPLUWKKOMYQJBAGLUKLGLOSVHWBPGLSUKOB SOPLOOKUSARPP"], "TWBTGTGTGTGBTPKHAXHAGTCSJJPTJPAPJHHOHHDHJDHJDHJKPSTJUMXHP  
HGALKLPJTJPGVXPLBHJHKPKDPDJSJ"]  
  
The best match is: TWBTGTGTGTGTBTPKHAXHAGTCDJCKPJTJHHOHHDHJDHJDHJHKUTJUWXGHGHGALKLPJTJPGVXPLBJHH  
I---PZJ-J-LMLT-K-----JULOS-T-P-----J-OLGLK-OBLTXGK-TPLUWKKOMYQJBAGLUKLGLOSVHWBPGL-SLUKOB SOPLOOKUSARPP--J--  
TWBTGTGTGTGTBTPKHAXHAGT--JSJJPTJPAPJHHOHHDHJDHJDHJKPSTJUW-W-XGPHGAL-K-L-PJTJPGVXPL----B--JHHDPK-----WPPDJSG  
TWBTGTGTGTGTBTPKHAXHAGT--JDXCKPJTJHHOHHDHJDHJDHJHKUTJUW-W-XGHGHGAL-K-L-PJTJPGVXPL----B--JHH-----  
The running time of three-sequence alignment by dynamic programming is: 401.121597s
```

# A-star Algorithm

## Pairwise Alignment

```
1 def Astar_Pairwise(self):
2     frontier = PriorityQueue()
3     frontier.put(self.start, 0)
4     self.came from[self.start] = None
```

## Result

The least cost of pairwise alignment by Astar Algorithm is: 131  
The base sequences is: ITPVKWKXKXIXUAXP-VHXVO-M-KMHHYPABLLBGKOLLJGXGXLSOL--AMOGKIGXBATXMPJTCVMTAXVMPMWA--WOM--OUPHHZBITKXKXK  
The best match is: ITPVKWKXKXIXUAXP-VHXVO-M-KMHHYPABLLBGKOLLJGXGXLSOL--AMOGKIGXBATXMPJTCVMTAXVMPMWA--WOM--OUPHHZBITKXKXK  
IPPVKBKXKXKXHS-A-PHVXVOJMRKAK--KPJVLJJBWIKOLLJJKXHGXL--LCPAJOBKPGXBATGXMPOMCVZTAXV-P--AGKXGOMJQO---LJGWGXKLQ  
ITPVKWKXKXIXUAXP-VHXVO-M-KMHHYPABLLBGKOLLJGXGXLSOL--AMOGKIGXBATXMPJTCVMTAXVMPMWA--WOM--OUPHHZBITKXKXK  
The running time of pairwise alignment by A-star algorithm is: 9.351278s

The implementation is almost the same with pairwise alignment.



reproduce the next generation until end condition happens. It's random to get an optimal solution.

In MSA, I use the position of '-' in a sequence to denote a match. For example, (0, 3, 5, 8) means  $-CT - C - CT - AG$ .

I have used randomly initialization and greedy initialization in start. But in practice, there isn't much difference in performance between them.

The fitness of each solution is the same with dynamic programming.

## Pairwise Alignment

### The Process of Reproduce

```
1 def Reproduce(self):
2     for i in range(ceil(self.population * self.cross_rate)):
3         # select two individuals at random
4         father = randint(0, len(self.group) - 1)
5         mother = randint(0, len(self.group) - 1)
6         while father == mother:
7             mother = randint(0, len(self.group) - 1)
8
9         x = self.group[father] # individual
10        y = self.group[mother]
11        # crossover these two individuals
12        self.Crossover(x, y)
13
14        # mutation
15        self.Mutation(x)
16        self.Mutation(y)
17
18        # add new children to the mate pool
19        self.group.append(x)
20        self.group.append(y)
21
22        self.group_fitness.append(self.Individual_Fitness(x))
23        self.group_fitness.append(self.Individual_Fitness(y))
```

### Crossover

```
1 def Crossover(self, x, y):
2     for i in range(2):
3         point = self.CrossoverRange(x[i], y[i]) # get the range to crossover
4         pos = 0
5         while x[i][pos] < point: # crossover between x and y
6             tmp = x[i][pos]
7             x[i][pos] = y[i][pos]
8             y[i][pos] = tmp
9             pos += 1
10            if pos == len(x[i]):
11                break
12            x[i].sort() # sort the position of '-'
13            y[i].sort()
```



And the *CrossoverRange* is to find the range to perform crossover. Because the length of expanded sequence(expanded by '-') is fixed, I request the number of chars to perform crossover between  $x$  and  $y$  is the same. The function is as follows:

```

1 def CrossoverRange(self, father_block, mother_block):
2     point = 1
3
4     # find the point in the left of which there are the same number of block
5     '-' between father and mother
6     count_left_father = sum(i < point for i in father_block)
7     count_left_mother = sum(i < point for i in mother_block)
8     while count_left_father != count_left_mother or count_left_father == 0:
9         point +=1
10        count_left_father = sum(i < point for i in father_block)
11        count_left_mother = sum(i < point for i in mother_block)
12
13    return point

```

## Mutation

```

1 def Mutation(self, individual):
2     for seq in range(2):
3         for i in range(len(individual[seq])): # for every block
4             if random() < self.mutation_rate: # mutate in a probability of
mutation_rate
5                 new_block = randint(0, self.seq_len - 1)
6                 while new_block in individual[seq]:
7                     new_block = randint(0, self.seq_len - 1)
8                 individual[seq][i] = new_block
9
10    individual[seq].sort()

```

## Result

The result and running time are as follows:

[illegible]

The least cost of pairwise alignment by Genetic Algorithm is: 177  
The base sequences is: ILOTGJ3LABWTSXGGONXJUTUXS3JHKWJHJCTOQHWGAGIWLZHWPKZULJZTWAKBWHXMIKLZJGLXBPAPHOHVOLZWOSJ3LP  
The best match is: IPOTWJ3LABKSGZGWJJKSPPOPTH3JEWHJCTOOTHRAKBLBHPKZULJPZAKVXKHJUIKLZJGLXBTGHOHBJLZPPSJJ3PJO  
ILOTGJ-3L-AB-W-TSXTGGONXJ-MUTU-XS3HKW--JHCTOQ-HWGAGIWLZHW--PKZULJ-ZTW-AKBWHX-MIKLZJGLXBPAPHOH-WO-LZ-W-OSJ3LP-  
IPOTWJ3LABKS--GZGWJ-JKSPPOP-HTSJEWJH-CTOOTH-AR--BKLBPMPKZUL--LJP-ZK---AKVXKHUIKLZJGLX-BT--GHOHBJL-ZPPSJJ-PJO  
The running time of pairwise alignment by genetic algorithm is: 152.76932s

[illegible]

The least cost of pairwise alignment by Genetic Algorithm is: 196  
The base sequences is: MPPPJPXPGPJXPPXPPPJJPJPPXHPKPPSPJJJPXPPPPPPJPJPPXPPXIJMMXXPKPSVGULMHHPZPAWHTHKAAHHUAPONAPJSPWPJGA  
The best match is: OPJXPJPMJXPMXPMJXPPXJPPXQXPPXJJPJXJXPPJPOJPMXPKPOGPPXPKPOMPXXPKOXPPXJQXQXPBJJPPXPPX  
MPPPJPX-PXG-PJPPXPPPPJJPJPP--XPPPPS-PP-JJJPPXXP-PP-PPJP-PPXPPXI-PJMMXX-PKPSVG-ULMHHPZPAWH-T-HKAA-HHU-PAO-NA-P-JSW-PPJGA  
-OP-JP--XJPPM-J-PMXPMX-JPP-XJPOXQXPPXJJ-P-JXX-PXJ-PP-O--J-PJPM-XP-PQG-PPX-XP-PO-MPPX-X-PPOX-PP--X--JQXQ-PPB-JPPXPPX-  
The running time of pairwise alignment by genetic algorithm is: 173.687169s



The least cost of pairwise alignment by Genetic Algorithm is: 188  
The base sequences is: IPPVKIKXKXKHSAPHVXVVOIMRAKKPJVLVLJBWKOLLJGXHGXLCPAJOBKPGXBATGXPMPOMCVZTAXVPAGIXGOMQJOLJGWGXQLQ  
The best match is: ITPVKWKSIXKXUAXPVXVOMMKHYBPABLLOBGKOLLJGXZGXL SOLAMOGKIGXBATBXPMPJTCVMITAXVPWPAWOMOUHPHZBITKXKXK  
IPPV-K-KXKX--WKXKSA-PHVXV---VOJM-RAKKPJVLVLJBWKOL-L-JXKHGX-LLCPA-JO-BKPGXBATGXPMPOMC--VZT-AXVPAGIX-GOMQJOLJ-GWGXQL-Q-ITPV--KWKSKIXKXUAXPVXV-V-OMMK-HYBPABLLOB--GKO-L-LJGXZGXL-SO--LAMOG-K-IGXBATBX-MPJTCVM-TA-XVMPWPAWOMOU-PHZBITKXKXK  
The running time of pairwise alignment by genetic algorithm is: 168.689891s

We can compare these results to dynamic programming and A-star algorithm, it shows that although genetic algorithm cannot get the best match, it can find the best-matched sequence.

## Time Complexity

Suppose the population of one generation is  $p$ , the longest length of sequence is  $l$ , the expansion rate is  $r_{sp}$ , the generation times is  $g$ , the crossover rate is  $r_c$ , the mutation rate is  $r_m$ . The time complexity in initialization is  $O(pl)$ . The time complexity in evolution is  $O(g * (p * l + p \log p))$ . So, the total time complexity is  $O(g * (p * l + p \log p))$ .

## Three-Sequence Alignment

The implementation of three-sequence alignment is almost the same with pairwise alignment. So, I'm not going to explain it in here.

## Result

The result and running time are as follows:

```
The least cost of three-sequence alignment by genetic algorithm is: 666
The base sequences are: ['ITPZJJMLTKJULOSTKTOGGLKJ0BLTGXKTP-LW-WKK-OMOY--J-B-GA-LJUK-LG-LOS-VHN--B-PG--WSLUKO-BSOPLOO-KUK-SARPP'] , ['ITWBTGTGTGTGKBTPKHAXHAGJCSJJPPAPJHJHDHJHDHJHDHJPKSTJUMWHP
HGALKPTTJPJPGVPLBHPKHPKAPDJSJG']
The best match is: ['ITWBTGTGTGTGKBTPKHAXHAGJCSJ]KPTTJHJHDHJHDHJHDHJHKTUTJUMWGHGHGALKLPJTPJPGVPLB]H
ITPZJJ--MLTKJUL--OSTKTOGGLKJ0B-LTX-GKTP-LW-WKK-OMOY--J-B-GA-LJUK-LG-LOS-VHN--B-PG--WSLUKO-BSOPLOO-KUK-SARPP]--
I-WBTGTGTGT--WGBTPKHAXHAGJCSJJPPAPJHJHDHJHDH--HJHDHJPKST-JJ-U-WKHGPH---G-AL-KLPJ-TPJPGVPLBHPKHPKAPDJSJG
-IWBTGTGTGTGKBTPK-H-AX-H-AGJCSJ--KPTTJPHJ-HJHDH--HD-HJH-KUTJJUMWGH-GH-GL-LK-LP--JTPJPG-V-X-PLB--J-JH
The running time of the three-sequence alignment by genetic algorithm is: 348.6437488
```

## Time Complexity

Use the same symbols in pairwise alignment. The time complexity in initialization is  $O(pl)$ . The time complexity in evolution is  $O(g * (p * l + p \log p))$ . So, the total time complexity is  $O(g * (p * l + p \log p))$ .

## Reflection and Problems

- In A-star algorithm for three-sequence alignment, I use a more accurate heuristic function at first: the sum of pairwise actual cost. But after it runs for four hours, there isn't any signal to show its end. So, I pick up the old method used in pairwise alignment. It's more rough but simpler and faster. In this condition, although more accurate function results in less search space, the cost of precomputing in accurate function exceeds the save of less search space. So, it's inefficient. The first heuristic function of three-sequence is as follows:

```
1 def Heuristic(self, current, goal):
2     self.HeuristicMatrix()
3     return self.cost_matrix_xy[current[1]][current[0]] +
self.cost_matrix_yz[current[2]][current[1]] + self.cost_matrix_zx[current[0]]
[current[2]]
```

`HeuristicMatrix()` is to precompute the dynamic cost matrix for each two sequences. It's the same with dynamic programming.

```
1 def HeuristicMatrix(self):
2     self.cost_matrix_xy = DP_MinCost_Matrix(self.sx, self.sy)
3     total_cost_xy = self.cost_matrix_xy[self.leny - 1][self.lenx - 1]
```

```

4
5     # the cost to reach goal - the cost to reach current point = needed
least cost = h(n)
6     for row in range(len(self.cost_matrix_xy)):
7         for line in range(len(self.cost_matrix_xy[row])):
8             self.cost_matrix_xy[row][line] = total_cost_xy -
self.cost_matrix_xy[row][line]
9
10    self.cost_matrix_yz = DP_MinCost_Matrix(self.sy, self.sz)
11    total_cost_yz = self.cost_matrix_yz[self.lenz - 1][self.leny - 1]
12
13    # the cost to reach goal - the cost to reach current point = needed
least cost = h(n)
14    for row in range(len(self.cost_matrix_yz)):
15        for line in range(len(self.cost_matrix_yz[row])):
16            self.cost_matrix_yz[row][line] = total_cost_yz -
self.cost_matrix_yz[row][line]
17
18    self.cost_matrix_zx = DP_MinCost_Matrix(self.sz, self.sx)
19    total_cost_zx = self.cost_matrix_zx[self.lenx - 1][self.lenz - 1]
20
21    # the cost to reach goal - the cost to reach current point = needed
least cost = h(n)
22    for row in range(len(self.cost_matrix_zx)):
23        for line in range(len(self.cost_matrix_zx[row])):
24            self.cost_matrix_zx[row][line] = total_cost_zx -
self.cost_matrix_zx[row][line]

```

- Although genetic algorithm cannot get the best match, it can find the best-matched sequence.
- I don't understand why my A-star algorithm runs longest. I think it would be the shortest. I suppose that the reason is STL of python. I use STL only in A-star algorithm, it makes codes clean and simple. But there may be some extra and inexplicit time for python to maintain the state of STL.
- I don't know how to configure the value of population, generation, crossover rate and mutation rate in genetic algorithm to get a better performance.