

BIOINFORMATICS ASSIGNMENT

Syed Hussain Ahmed, NED University

16/11/2024

Transitions and Transversions - Alignment

For DNA strings s_1 and s_2 having the same length, their *transition/transversion ratio* $R(s_1, s_2)$ is the ratio of the total number of transitions to the total number of transversions, where symbol substitutions are inferred from mismatched corresponding symbols as when calculating Hamming distance (see “Counting Point Mutations”).

Given

Two DNA strings s_1 and s_2 of equal length (at most 1 kbp).

Return

The transition/transversion ratio $R(s_1, s_2)$.

Sample Dataset

```
>Rosalind_0209
GCACGCCAGAGAAACCTTATGGGAAGTGAATTATTTCTGGTATCGTTGTAGTTATTGGA
AGTAGGGCAGTACACCCAGTT
>Rosalind_2200
TTATCTGCAAAGAAAGGCCTGACCGGTGGATATTCCGCGATCGTCTCGGTGTTTACTGGC
GTTTACGAGTTCTCTTGGTGGGT
```

Sample Output

1.21428571429

```
1 from readFASTA import readFASTA
2
3 def ti_tv(seqs: list[str]) -> float:
4     """Returns the Transitio/Transversion ratio of the sequences"""
5
6     purines = ["A", "G"]
7     pyrimidines = ["C", "T"]
8
9     ti_count = 0
10    tv_count = 0
11
12    for i in range(len(seqs[1])):
13        if seqs[0][i] != seqs[1][i]:
14            if (seqs[0][i] in purines and seqs[1][i] in purines) or (
15                seqs[0][i] in pyrimidines and seqs[1][i] in pyrimidines
16            ):
17                ti_count += 1
18            else:
19                tv_count += 1
20
21    return ti_count / tv_count
```

```

17         ti_count += 1
18     elif (seqs[0][i] in purines and seqs[1][i] in pyrimidines) or ↵
19         (
20             seqs[0][i] in pyrimidines and seqs[1][i] in purines
21         ):
22         tv_count += 1
23     return ti_count / tv_count
24 if __name__ == "__main__":
25     filepath = <path_to_file>
26     _, seqs = readFASTA(filepath)
27     ratio = ti_tv(seqs)
28     print(ratio)

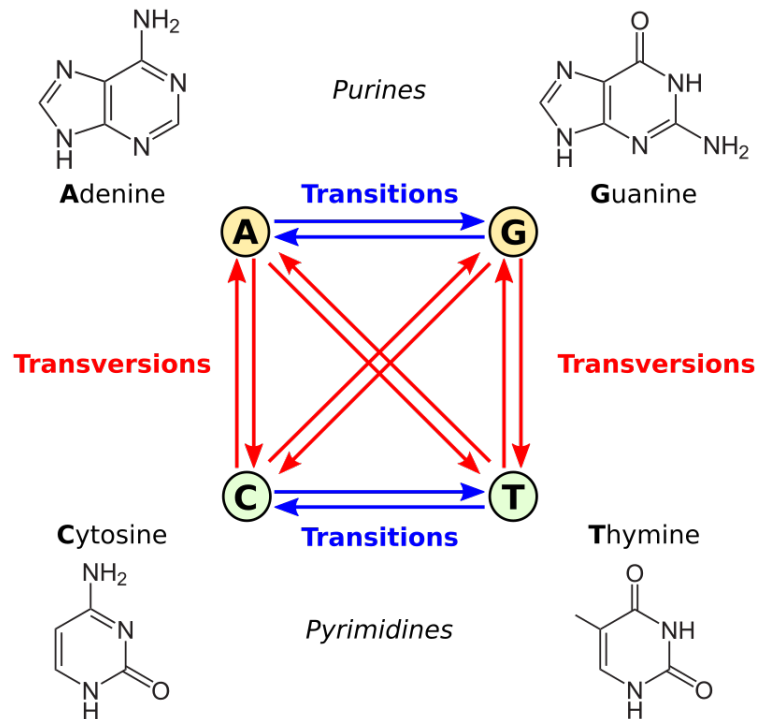
```

1.9895833333333333

Explanation

The function `ti_tv` implements the logic to identify transitions and transversions by checking each nucleotide pair in the sequences. If both belong to the same category, it increments the transition count; otherwise, it increments the transversion count. Finally, it returns the ratio of transitions to transversions.

Biological Connection



Transitions mean when there is substitution between purine and purine or pyrimidine and pyrimidine. Transversions is substitution between purine and pyrimidine. Normally the ti/tv ratio is approximately two but near coding regions it can get higher than 3. We can calculate the ti/tv ratio to identify the coding regions

Longest Increasing Subsequence - Dynamic Programming

A **subsequence** of a permutation is a collection of elements of the permutation in the order that they appear. For example, (5, 3, 4) is a subsequence of (5, 1, 3, 4, 2).

A subsequence is **increasing** if the elements of the subsequence increase, and **decreasing** if the elements decrease. For example, given the permutation (8, 2, 1, 6, 5, 7, 4, 3, 9), an increasing subsequence is (2, 6, 7, 9), and a decreasing subsequence is (8, 6, 5, 4, 3). You may verify that these two subsequences are as long as possible.

Given: A positive integer $n \leq 10,000$ followed by a permutation π of length n .

Return: A longest increasing subsequence of π , followed by a longest decreasing subsequence of π .

Sample Dataset

```
5
5 1 4 2 3
```

Sample Output

```
1 2 3
5 4 2
```

```
1 def read_file(filename):
2     with open(filename) as f:
3         n = int(f.readline().strip())
4         perm = list(map(int, f.read().strip().split()))
5         return n, perm
6
7
8 def print_seq(seq):
9     for i in range(len(seq)):
10         print(seq[i], end=" ")
11     print()
12
13
14 def longest_increasing_subsequence(X, N):
15     P = [-1] * N # Predecessor array
16     M = [-1] * (N + 1) # Index array for LIS of each length
17     L = 0 # Length of the LIS
18
19     for i in range(N):
```

```

20     # Binary search for the smallest positive l      L
21     # such that X[M[l]] >= X[i]
22     lo = 1
23     hi = L + 1
24     while lo < hi:
25         mid = lo + (hi - lo) // 2
26         if X[M[mid]] >= X[i]:
27             hi = mid
28         else:
29             lo = mid + 1
30
31     # After searching, lo == hi
32     newL = lo
33
34     # Update the predecessor of X[i]
35     P[i] = M[newL - 1]
36     M[newL] = i
37
38     # If newL is longer than the current LIS length, update L
39     if newL > L:
40         L = newL
41
42     # Reconstruct the LIS
43     S = [0] * L
44     k = M[L]
45     for j in range(L - 1, -1, -1):
46         S[j] = X[k]
47         k = P[k]
48
49     return S
50
51
52 def longest_decreasing_subsequence(X, N):
53     P = [-1] * N
54     M = [-1] * (N + 1)
55     L = 0
56
57     for i in range(N):
58         # Binary search for the smallest positive l      L
59         # such that X[M[l]] <= X[i]
60         lo = 1
61         hi = L + 1
62         while lo < hi:
63             mid = lo + (hi - lo) // 2
64             if X[M[mid]] <= X[i]:
65                 hi = mid
66             else:

```

```

67         lo = mid + 1
68
69     newL = lo
70     P[i] = M[newL - 1]
71     M[newL] = i
72     if newL > L:
73         L = newL
74
75     S = [0] * L
76     k = M[L]
77     for j in range(L - 1, -1, -1):
78         S[j] = X[k]
79         k = P[k]
80
81     return S
82
83
84 if __name__ == "__main__":
85     filepath = <path_to_file>
86     n, seq = read_file(filepath)
87
88     lis = longest_increasing_subsequence(seq, n)
89     print_seq(lis)
90
91     lds = longest_decreasing_subsequence(seq, n)
92     print_seq(lds)

```

```

4 69 169 232 249 351 426 542 707 742 756 772 801 842 844 932 1044 1073 1087 1275 1288 1320
1343 1413 1416 1437 1473 1479 1481 1504 1608 1610 1612 1736 1769 1814 1845 1851 1887 2040
2042 2170 2199 2216 2254 2287 2380 2396 2403 2428 2470 2477 2528 2538 2546 2547 2628 2642
2689 2701 2732 2796 2829 2838 2921 2937 3166 3184 3197 3198 3221 3304 3318 3548 3616 3654
3660 3665 3745 3763 4005 4087 4119 4140 4141 4176 4254 4287 4344 4468 4485 4488 4489 4525
4562 4674 4710 4735 4801 4903 4970 5003 5008 5038 5105 5186 5228 5290 5301 5322 5355 5403
5629 5632 5642 5719 5839 5848 5873 6127 6187 6188 6238 6250 6289 6305 6323 6344 6411 6420
6436 6578 6631 6636 6805 6817 6836 6879 6929 6951 7003 7008 7012 7134 7161 7218 7224 7236
7255 7353 7358 7371 7573 7619 7755 7775 7780 7821 7854 7887 7968 8053 8176 8331 8445 8450
8465 8558 8569 8575 8639 8684 8759 8777 8813
8810 8789 8735 8732 8731 8686 8673 8648 8592 8590 8588 8486 8419 8417 8339 8307 8295 8292
8270 8236 8229 8189 8186 8183 8179 8169 8138 8130 8095 8059 7992 7980 7951 7910 7849 7791
7790 7746 7736 7685 7643 7627 7624 7604 7567 7508 7467 7419 7393 7373 7275 7190 7157 7081
6942 6924 6873 6778 6774 6669 6639 6632 6571 6556 6554 6537 6484 6475 6449 6448 6429 6409
6216 6196 6167 6153 6144 6109 6032 6015 5881 5800 5790 5778 5723 5657 5655 5645 5592 5583
5501 5488 5480 5445 5403 5276 5270 5163 5159 5118 5099 4990 4941 4938 4922 4858 4850 4803
4727 4664 4630 4627 4622 4544 4526 4517 4456 4414 4315 4274 4223 4218 4078 4048 4034 3974
3908 3805 3790 3787 3773 3760 3749 3743 3606 3600 3590 3512 3402 3315 3305 3226 3210 3070
2943 2804 2656 2560 2495 2355 2059 2037 2011 1922 1909 1877 1846 1611 1576 1569 1517 1498
1448 1377 1360 1199 1188 1137 1050 977 878 813 811 782 721 632 614 610 509 432 391 378 370
140 130 |

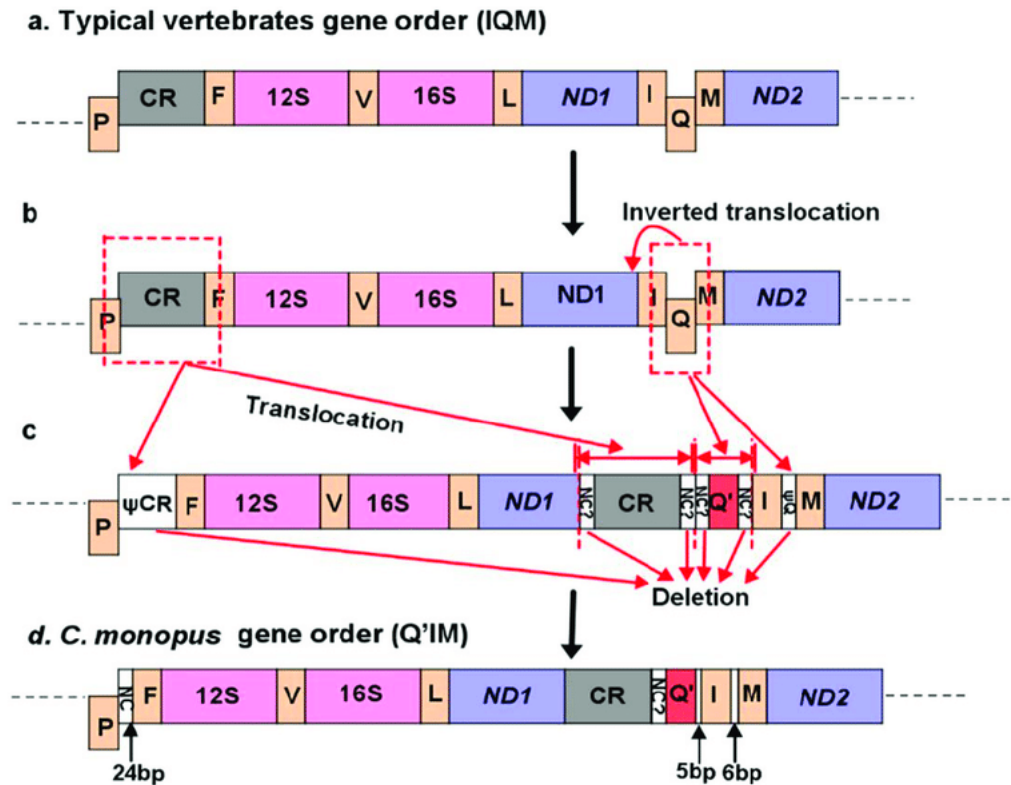
```

Explanation

The functions `longest_increasing_subsequence` and `longest_decreasing_subsequence` both compute subsequences efficiently using an $O(n \log n)$ algorithm described here .

The `longest_increasing_subsequence` function finds the Longest Increasing Subsequence (LIS) in a given sequence by maintaining a predecessor array `P` to track the indices of elements forming the LIS and an index array `M` to store the smallest end elements of LIS for various lengths. The function iterates through the sequence, using binary search to determine where the current element can extend or replace an existing LIS, updating `P`, `M`, and the LIS length (`L`) as needed. Finally, it reconstructs the LIS by backtracking through `P`. Similarly, the `longest_decreasing_subsequence` function finds the Longest Decreasing Subsequence (LDS) with the same structure, but the binary search logic is adjusted to ensure the subsequence elements decrease. Both functions return the reconstructed subsequences as lists, demonstrating their modularity and adaptability for different subsequence criteria.

Biological Connection



Biological sequences are prone to mutations. Despite these mutations, closely related organisms often have similar regions in their DNA. These conserved regions are called syntenic blocks. Although these syntenic blocks may be rearranged within the genome, they typically remain largely conserved. Due to genomic rearrangements, the order of genes within these blocks—the specific locations where a gene exists—may differ. One simple way to compare genes between two chromosomes is to search for the largest collection of genes that are found in the same order on both chromosomes. To identify the largest set of genes appearing in the same order, we need only to find the largest collection of increasing elements in the permutation.

Consensus and Profile - String Algorithms

A matrix is a rectangular table of values divided into rows and columns. An $m \times n$ matrix has m rows and n columns. Given a matrix A , we write A_{ij} to indicate the value found at the intersection of row i and column j .

Say that we have a collection of DNA strings, all having the same length n . Their profile matrix is a $4 \times n$ matrix P in which P_{ij} represents the number of times that 'A' occurs in the j th position of one of the strings, P_{ij} represents the number of times that 'C' occurs in the j th position, and so on (see below).

A *consensus string* c is a string of length n formed from our collection by taking the most common symbol at each position; the j th symbol of c therefore corresponds to the symbol having the maximum value in the j th column of the profile matrix. Of course, there may be more than one most common symbol, leading to multiple possible consensus strings.

DNA Strings

A	T	C	C	A	G	C	T
G	G	G	C	A	A	C	T
A	T	G	G	A	T	C	T
A	A	G	C	A	A	C	C
T	T	G	G	A	A	C	T
A	T	G	C	C	A	T	T
A	T	G	C	A	A	C	T

Profile

A	5	1	0	0	5	5	0	0
C	0	0	1	4	2	0	6	1
G	1	1	6	3	0	1	0	0
T	1	5	0	0	1	1	1	6

Consensus

A T G C A A C T

Given: A collection of at most 10 DNA strings of equal length (at most 1 kbp) in FASTA format.

Return: A consensus string and profile matrix for the collection. (If several possible consensus strings exist, then you may return any one of them.)

Sample Dataset

```
>Rosalind_1
ATCCAGCT
```

```
>Rosalind_2
GGGCAACT
>Rosalind_3
ATGGATCT
>Rosalind_4
AAGCAACC
>Rosalind_5
TTGGAACT
>Rosalind_6
ATGCCATT
>Rosalind_7
ATGGCACT
```

Sample Output

```
ATGCAACT
A: 5 1 0 0 5 5 0 0
C: 0 0 1 4 2 0 6 1
G: 1 1 6 3 0 1 0 0
T: 1 5 0 0 0 1 1 6
```

```
1
2 from readFASTA import readFASTA
3
4 def print_matrix(matrix: list[list[int]]) -> None:
5     residues = ["A", "C", "G", "T"]
6     for row in range(len(matrix)):
7         print(residues[row] + ": " + " ".join(map(str, matrix[row])))
8
9 def get_consensus_string(profile: list) -> str:
10     con_str = ""
11     residue_dict = {
12         0: "A",
13         1: "C",
14         2: "G",
15         3: "T",
16     }
17     for col in range(len(profile_mat[0])):
18         max_value = -1
19         for row in range(len(profile_mat)):
20             if profile_mat[row][col] > max_value:
21                 max_value = profile_mat[row][col]
22                 residue_index = row # 0=A; 1=C; 2=G; 3=T
```

```

23         con_str += residue_dict[residue_index]
24     return con_str
25
26
27 def get_profile_matrix(seqs: list[str]) -> list:
28     res_count = {
29         "A": [0] * len(seqs[0]),
30         "C": [0] * len(seqs[0]),
31         "G": [0] * len(seqs[0]),
32         "T": [0] * len(seqs[0]),
33     }
34     for col in range(len(seqs[0])):
35         for row in range(len(seqs)):
36             current_residue = seqs[row][col]
37             res_count.get(current_residue)[col] += 1
38     return list(res_count.values())
39
40
41 if __name__ == "__main__":
42     filename = <path_to_file>
43     _, seqs = readFASTA(filename)
44
45     profile_mat = get_profile_matrix(seqs)
46     print(get_consensus_string(profile_mat))
47     print_matrix(profile_mat)

```

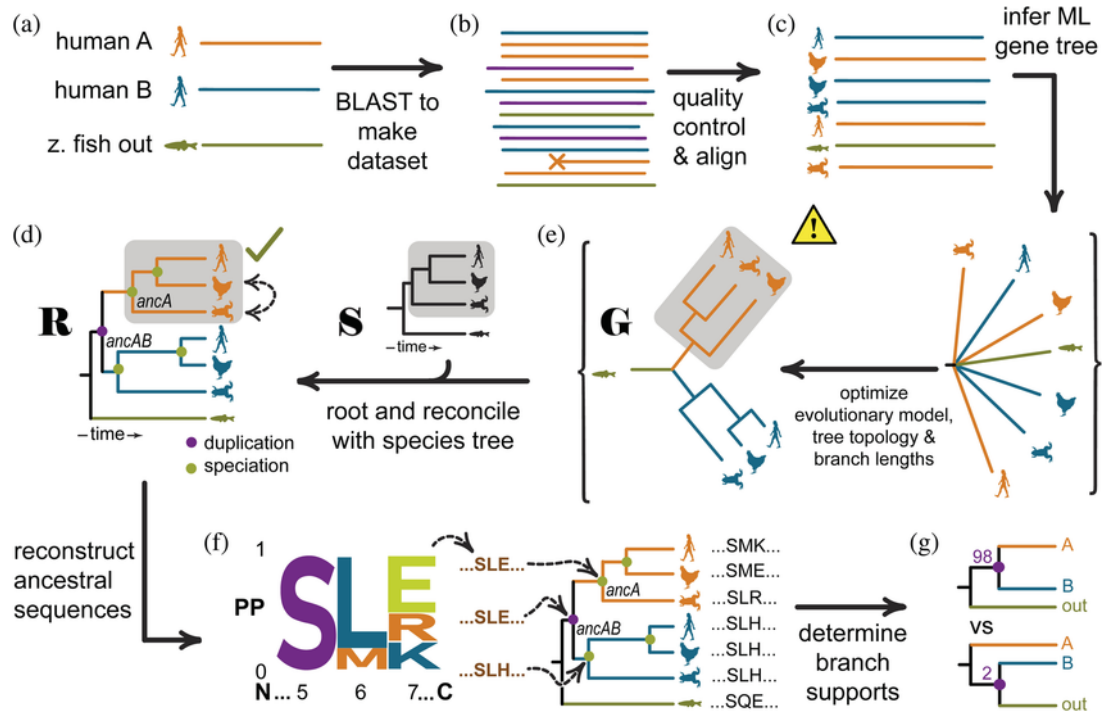
```
TCCTGGAGACTTCATGCGATACACCATAAACACATTACCCGACATCAGAGTGCCGACGTCGTCGGGGCC
ATCGGTGATAATCTTAAGTGAAGAAAGACTTACTGACAAAACAGCCTCCGGCCACGTTGTCCTGCGTCGAGA
GGGACTACCTCGGCAATGAGCTCAAAAATGCCAGCAATAGGATCAATGTAACGATAACCGTGTCCGCGGGC
GTAGCGATGAAGAGGGAATTGAGACAGAGATATAGGCGGGAAGTACGGCCATATCGGGGTATACCCATACCG
CAATCAGAACGAAAAATCCCTCGAAGACAAAGGCGGTACGACGACATCTGGAACACACAAGGCCAGCAGT
ACCTCGCTAATGGATAAACAGCGGTTGACACAGACTCGATTCTGGGAGACAAAGATGAGCATCCAGTATGC
AACCTCTGAACATCAAACTCAGCAAGACTACGTTGATGCGAGAGCGCGGCCAAGAAGCGGGAAGTCCG
CGCCGTAAACCCAAATCGCGGTGTACAACCTAAACACCTTGAAGATGACACAGTGCCAACTGCCA
CAAGAAAGCAAGTCAAACTAGTTACAACCAAGAACTATACGCCCTAAAGACTGCTATAACCTGAATCCCCA
GAAATACCCGAAGGTTAACCGCTCATGTGCAACTATAAAACAAACACGCTATGCTTAAGTGAAGAA
TGTTATTGGGGTAAACCGGAGATGCCAAGTGAATAATAGGCCCACTTGAAGGAGACGGAAGCGGTGA
CGCTGAGTTCATACACTACACAATAATCACTTTCCAAATGCACGAGCATACCGCGTACAACAAAGCGCA
GAGGTCACGAAACGAGAATTAACAAGGCGCATGTGAACTGCAACGAGTCACTAA
A:33222333421234320232435314255054333113222714223331312
12252303222113032024331222307321033433214454242224212524
3431422220222226232222323123223130026224233220442232
13364434122125223435124225313345014143222111022213211361
12410341301333112434351324243422332533342013132432334124
24133313221246113344223344322212313523355441211036324224
153211234626333212352243244212223401052353133332122441
433022343021112151525424124623432260242204522232553161
23344132520463430432212433213231323233113323302232144232
11111223133433334523221023324422422355235212243324323
2330422201632221243332464204302133333141114266313541532
333423211143624121125235033233011323233424201640210460
23212615130234301333453143325031211522312433121445253212
2311221202342222324320262442133454241410310322331532233
41124121232151212242221324253223042452344334241102333311
23125321445031221321341243332232323330412442342540355344
024220342233421233231420242244
C:245120222421422252304414511246215431324352253163232136
5242423413565152343216222222215212152134220313531323271
3315323405342343232221442052422333226103114524225402122
52402323324552241210211333331321442322532323454214142233
42313322331013231321413231322123412232124324431233141114
3135433331431252242422132224532442122513012341131511421
4125221233050432314422530235423132312123323343042232243
2334323123344331213232222211061036123121442363202413321
54323153332122242142120211411223024421543313314132232324
342532122334452121042232412343321311233405321133134232
4344012255235124344113202261333043114323331301344133024
12115034421213462343222345122323544312234244552132131204
40304232105242362322412413260526254212522334032213113221
22232021422244213233325323421322313223321445163213022131
14220153422026133213224212412503535340323241512343232332
43322144210543453124315121141431322245142424221231001403
2353253022231332142443223323233
G:121134351333121325131112231212320113021325123302435241
34322403610104442333207423330232312211641224413212235011
3132333233324221242231112224013225342523231212354131534
3113222104220033222451121114013451132134252310360625004
```

```
45236225246433224042034242131334254124212561333003515511
12312321335212303512333323302232152463131315425601241243
22112552122212224411041142111325312427411213224145223013
2422313440332055404123242421250223124144340121313011312
2114232212422423313034224343533462345332241352436523311
34351232322032230323334041112032433312220233512341142
23322413321131342201334312523241612331234022322342314331
2312053322304222531113130252323324353211201244104503233
2434122042342233213232331211421431233032411347231401243
2211556532224344341052302234242221112246032321153223435
2344241431341143151422431122313331201322111410223223114
323253220242403423221432341433524411321303142532322330
6412432142422233523424231212201
T:412633213144224331351231226102032235633211020422112421
12112170143242113321410243251342553314121212142143340304
123203225222123131443632450344332124222523141311135222
142121262223331324332443135243120351322242432432312512
1115122232331543313312313414231112221432212214442120361
43311143223221442012132311262134203002111342133342034222
3242523312131231432322234203432252161414311323211443413
20242522044136123422302223441313521144124222312543133316
1230114313430120522374424214312330111122233142310211243
32223334432301023415133236343313241243203133344222311413
2104435322130242332233113212223432233241264322111112223
42533222353231014124252132413163110122331423113234266213
2325304345221211424201322331413221414324403261022134344
4365131234222033121241322301341312534133323304413333311
3232442214352233414144126322251221305222421255433322353
12331113253112402523240423223112221422343213122016532033
210322252412322421201223633343]
```

Explanation

The script features three key functions: `get_profile_matrix`, which generates a profile matrix by counting the occurrences of each nucleotide (A, C, G, T) at every position across multiple DNA sequences, storing the counts in lists grouped by nucleotide; `get_consensus_string`, which uses the profile matrix to construct the consensus string by identifying the most frequent nucleotide at each column and appending it to the string; and `print_matrix`, which formats and prints the profile matrix with row labels (A, C, G, T) for clarity, making it easy to interpret nucleotide counts at each position. These functions collectively help analyze DNA sequence alignments and identify conserved regions.

Biological Connection



With the passage of time, mutations are bound to occur, leading to the emergence of many new sequences. These sequences share a common ancestry and are said to be homologous. Given these homologous sequences, we can attempt to reconstruct the most likely ancestral sequence. However, it should be noted that the problem described above is an extreme oversimplification of the actual challenges involved.

Reference

Below is the implementation of readFASTA function. For reading the FASTA file I created a readFasta function as I did not know about the BioPython Module Bio.SeqIO.parse() method. I decided to not update it as it does the same thing.

```
1 def readFASTA(filepath) -> list[str]:
2     header_indices = []
3     seq_start_indices = []
4     seqs = []
5     headers = []
6
7     with open(filepath) as rf:
8         data = list(map(str.strip, rf.readlines()))
9
10    for index in range(len(data)):
11        if data[index].startswith(">"):
12            header_indices.append(index)
13            seq_start_indices.append(index + 1)
14
15    for i in range(len(header_indices)):
16        try:
17            seq = data[seq_start_indices[i] : header_indices[i + 1]]
18        except IndexError:
19            seq = data[seq_start_indices[i] :]
20        finally:
21            seq = "".join(seq)
22            seqs.append(seq)
23
24    for header_index in header_indices:
25        headers.append(data[header_index][1:])
26    return headers, seqs
```
