



# مقدمه‌ای بر بیوانفورماتیک

نیم‌سال اول ۱۴۰۰-۱۳۹۹

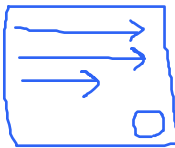
مدرس: دکتر شریفی زارچی، دکتر کوهی

میافانترم

راحت

## مسئله‌ی ۱.

دو رشته‌ی ACGTCATCA و TAGTGTCA را در اختیار داریم. جدول alignment این دو رشته را تشکیل دهید و تغییرات لازم در رشته‌ی اول برای این‌که به رشته‌ی دوم تبدیل شود را ذکر کنید.



## مسئله‌ی ۲.

در مسئله‌ی alignment فرض کنید که هرچه قدر نواحی مشابه طولانی‌تر باشند برای ما مطلوب‌تر است. برای مثال در align کردن دو رشته‌ی CAGCAGTGG و CAGCACTTGGATTCTCGG ، رشته‌ی تطبیق‌یافته‌ی CAGCAGTGG - - - - - را به رشته‌ی تطبیق‌یافته‌ی CAGCAGTGG - - - - - GG - T - - - - - ترجیح می‌دهیم. با این‌که تعداد match در رشته‌ی دوم بیشتر است. راه حلی ارائه دهید که بتواند این مسئله را حل کند.

Recurrence Relation:

- $H[i][j] = H[i-1][j]$
- Assign a score based on the length of the current contiguous match. For example, let the bonus for a contiguous match be  $b$ , and the base match score be  $m$ .
- $dp[i][j] = dp[i-1][j-1] + m + b \times (\text{current match length})$
- You can keep track of the current match length using an auxiliary table.
- Otherwise:
- $dp[i][j] = \max(dp[i-1][j] - \text{gap penalty}, dp[i][j-1] - \text{gap penalty})$

```

# Fill DP table
for i in range(1, n + 1):
    for j in range(1, m + 1):
        if S1[i-1] == S2[j-1]:
            match_length[i][j] = match_length[i-1][j-1] + 1
            dp[i][j] = dp[i-1][j-1] + match_score + bonus * match_length[i][j]
        else:
            match_length[i][j] = 0
            dp[i][j] = max(dp[i-1][j] - gap_penalty, dp[i][j-1] - gap_penalty)
    
```

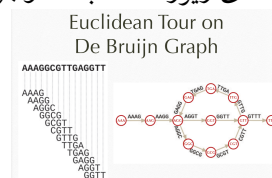
## مسئله‌ی ۳.

با استفاده از گراف *de Bruijn* اویلری کوتاه‌ترین رشته‌ای که شامل تمامی زیررشته‌های زیر باشد را بیابید.

{ACGT, AGAC, CGTG, GACG, GGTA, GTAC, GTAG, GTGT, TACG, TAGA, TGTA}

سپس توضیح دهید در صورتی که در هنگام بازسازی یک ژنوم از  $k$ -merهای آن، بیش از یک رشته با کوتاه‌ترین طول که شامل همه‌ی زیررشته‌ها باشد وجود داشته باشد، چه گونه می‌توان رشته‌ی مربوط به ژنوم را از سایر رشته‌ها تشخیص داد؟

<https://chatgpt.com/c/674e0d73-b488-8005-b71d-9ec4b6034ff4>



## مسئله‌ی ۴.

یک رشته پس از اعمال تبدیل BWT<sup>۱</sup> به فرم ATTAT\$G درآمده‌است، ابتدا رشته‌ی اولیه را بدست آورید سپس توضیح دهید که چگونه می‌توان با انجام پیش‌پردازش بر روی تبدیل BWT این رشته، تعداد تکرارهای زیررشته‌ی TA در رشته‌ی اصلی را محاسبه کرد.

<https://chatgpt.com/c/674e0d73-b488-8005-b71d-9ec4b6034ff4>

last to first  
count array

$s = \text{TAGTTA\$}$

## مسئله‌ی ۵.

آ. الگوریتم additive phylogeny را مرحله به مرحله روی ماتریس فاصله‌ی زیر اجرا کنید.

Introducing Count Array

$i$	FirstColumn	LastColumn	LastToFirst(i)	COUNT
0	$\$1$	$s1$	13	0 0 0 0 0 0 0
1	$a1$	$m1$	8	0 0 0 0 0 0 1
2	$a2$	$n1$	9	0 0 0 1 0 0 1
3	$a3$	$p1$	12	0 0 0 1 1 0 1
4	$a4$	$b1$	7	0 0 0 1 1 1 1
5	$a5$	$n2$	10	0 0 1 1 1 1 1
6	$a6$	$n3$	11	0 0 1 1 2 1 1
7	$b1$	$a1$	1	0 0 1 1 3 1 1
8	$m1$	$a2$	2	0 1 1 1 3 1 1
9	$n1$	$a3$	3	0 2 1 1 3 1 1
10	$n2$	$a4$	4	0 3 1 1 3 1 1
11	$n3$	$a5$	5	0 4 1 1 3 1 1
12	$p1$	$\$1$	0	0 5 1 1 3 1 1
13	$s1$	$a6$	6	1 5 1 1 3 1 1
				1 6 1 1 3 1 1

	$v_1$	$v_2$	$v_3$	$v_4$
$v_1$	0	9	10	9
$v_2$	9	0	9	8
$v_3$	10	9	0	7
$v_4$	9	8	7	0

Burrows-Wheeler Transform<sup>۱</sup>

Count<sub>symbol(i, LastColumn)</sub>:

#occurrences of symbol in the first  $i$  positions of LastColumn

ب. وزن یال‌ها در درخت فیلوژنی می‌تواند نشان‌دهنده‌ی چه چیزی باشد؟

## مسئله‌ی ۶.

آ. الگوریتم UPGMA را مرحله به مرحله روی ماتریس فاصله‌ی زیر اجرا کنید. سوال 6 تمرین 2

not ultrametric

سرعت تکامل گونه‌ها یکی نباشد

	$v_1$	$v_2$	$v_3$	$v_4$	$v_5$
$v_1$	۰	۲	۵	۴	۳
$v_2$	۲	۰	۴	۳	۵
$v_3$	۵	۴	۰	۲	۳
$v_4$	۴	۳	۲	۰	۱
$v_5$	۳	۵	۳	۱	۰



ب. الگوریتم UPGMA در چه صورتی ممکن است درختی اشتباه را به عنوان خروجی تولید کند؟

The UPGMA algorithm (Unweighted Pair Group Method with Arithmetic Mean) can produce an incorrect phylogenetic tree due to its underlying assumptions and limitations. These issues arise primarily because UPGMA assumes a molecular clock hypothesis, which means it expects all lineages to evolve at the same constant rate. If this assumption is violated or there are inconsistencies in the input data, UPGMA may construct a tree that does not accurately represent the evolutionary relationships.

UPGMA assumes that evolutionary distances are ultrametric, meaning the distance from the root to any leaf is the same for all leaves. If some lineages (دودمان) evolve faster than others (rate heterogeneity), the ultrametric property is violated, and UPGMA fails to correctly infer relationships.

For example: Consider three taxa: A, B, and C. If A and B have evolved slowly and C has evolved rapidly, UPGMA might group A and B correctly but place C incorrectly because it interprets the greater distance to C as a later divergence, rather than faster evolution.

UPGMA assumes that pairwise distances are additive (i.e., the distance between two taxa equals the sum of the branch lengths connecting them). If distances are not strictly additive (e.g., due to homoplasy or back mutations), UPGMA may produce an inaccurate tree.

UPGMA clusters taxa step-by-step, and once two taxa or groups are merged, the decision cannot be undone. If an incorrect merge occurs early, subsequent steps propagate the error, producing a globally incorrect tree.