



دانشگاه صنعتی شریف

دانشکده‌ی مهندسی کامپیووتر

مقدمه‌ای بر بیوانفورماتیک

پاییز ۱۴۰۳

استاد: علی شریفی زارچی

مسئول تمرین: ساجده فدائی

مهلت ارسال نهایی: ۷ آذر

تمرین دوم

مهلت ارسال بدون تاخیر: ۳ آذر

- مهلت ارسال پاسخ تا ساعت ۲۳:۵۹ روزهای مشخص شده است.

• در طول ترم، برای هر تمرین می‌توانید تا ۴ روز تأخیر داشته باشید و در مجموع حداقل ۸ روز تأخیر مجاز خواهد داشت. توجه داشته باشید که تأخیر در تمرین‌های عملی و تئوری به صورت مشترک محاسبه می‌شود. پس از اتمام تأخیرهای مجاز، می‌توانید با تاخیری ساعتی ۱ درصد تمرین خود را ارسال کنید.

• حتماً تمرین‌ها را بر اساس موارد ذکر شده در صورت سوالات حل کنید. در صورت وجود هرگونه ابهام، آن را در صفحه تمرین در سایت کوئرا مطرح کنید و به پاسخ‌هایی که از سوی دستیار آموزشی مربوطه ارائه می‌شود، توجه کنید.

• در صورت هم‌فکری و یا استفاده از هر منابع خارج درسی، نام هم‌فکران و آدرس منابع مورد استفاده برای حل سوال مورد نظر را ذکر کنید.

• فایل پاسخ‌های سوالات نظری را در قالب یک فایل pdf به فرمت $HW2_{[STD_ID].pdf}$ آپلود کنید.

• گردآورندگان تمرین: ساجده فدائی، نیکان واسعی، محمد مولوی، امیرحسین علیشاھی

سوالات نظری (۱۰۰ نمره)

۱. (۱۵ نمره) در تکنولوژی‌های توالی‌یابی نسل جدید، (NGS) تعیین توالی بخش‌های تکراری ژنوم چالش‌های خاصی به همراه دارد. به سوالات زیر پاسخ دهید:

الف) چرا توالی‌یابی بخش‌های تکراری ژنوم دشوار است و چگونه می‌توان این چالش‌ها را کاهش داد؟

ب) با توجه به تکنیک‌های موجود، چگونه از استراتژی‌های جفت-خوانشی (paired-end) برای بهبود دقیق توالی‌یابی در این بخش‌های تکراری استفاده می‌شود؟

۲. (۱۵ نمره) دانشمندان در حین تحقیقات خود در شناسایی ژن‌های عجیب دو گونه، توانستند read‌های مربوط به آن‌ها را بدست آورند و در فایل‌های جداگانه ذخیره کنند. اما متاسفانه به دلیل اشتباه یکی از دانشمندان، این فایل‌ها با یک‌دیگر ادغام شدند و درون یک فایل قرار گرفتند و بازسازی ژن‌ها را سخت‌تر کردند. با توجه به read‌های مخلوط شده، سعی کنید دو تا ژن اولیه را بازسازی کنید.

AAGT - TCTT - AGTA - CCAA - GTAG - ATCT - TAGG - ACCA - AGGA -
GGAT - GATT - ATTT - ACTC - CTCG - TCGC - TTTG - GCAT - CATC -
CTTA - CGCA - TTAC - TACC

۳. (۳۰ نمره) فرض کنید یک متاداده جدید برای alignment Multiple که قطعه کد آن در زیر معرفی شده است به شما داده شده، با این قطعه کد و جدول امتیازات زیر چهار رشته داده شده را با هم ترازو کنید. (توجه داشته باشید که حتماً جداول همترازوی را به صورت کامل رسم کنید و صرفاً نمایش خروجی همترازو شده کامل نمی‌باشد و نمره‌ای به آن تعلق نمی‌گیرد)

	-	A	T	C	G
-	0	-1	-1	-1	-1
A	-1	2	1	0	0
T	-1	1	2	-1	1
C	-1	0	-1	3	2
G	-1	0	1	2	3

Multiple Sequence Alignment Algorithm

- Require:** A set S of sequences
- Ensure:** A multiple alignment of M with sum of pair distances at most twice that of the optimal alignment of S
- 1 Find $D(S_i, S_j)$ for all i, j .
 - 2 Find the center sequence S_c which minimizes $\sum_{i=1}^k D(S_c, S_i)$.
 - 3 **for** each $S_i \in S - \{S_c\}$
 - 4 Choose an optimal alignment between S_c and S_i .
 - 5 Introduce spaces into S_c so that the multiple alignment M satisfies the alignments found in Step 3.

رشته‌های داده شده به صورت زیر هستند:

ACCCCTGAACC
ACTCGGAGC
CTGGAATCT
GCTAGGACC

۴. (۱۰ نمره) امروزه روش همترازی ساختاری (structural alignment) برای همترازی ماکرومولکول‌هایی مانند پروتئین‌ها به کار می‌رودند.

- الف) در مورد این روش تحقیق کنید و مراحل اصلی و کلیدی را به ترتیب و بطور کامل توضیح دهید.
- ب) دو مورد از الگوریتم‌های کاربردی در این روش را نام برد و در خصوص آنها توضیح دهید.
- ج) در خصوص تفاوت دو روش همترازی مبتنی بر توالی و همترازی ساختاری توضیح دهید.
- د) چگونه روش‌های تراز ساختاری می‌توانند همترازی توالی‌های چندگانه مبتنی بر توالی سنتی را برای استباط روابط تکاملی و عملکردی در پروتئین‌ها تکمیل کنند، در مورد چالش‌ها و محدودیت‌های ادغام داده‌های ساختاری در MSA توضیح دهید.

۵. (۱۵ نمره) علاوه بر الگوریتم‌های سنتی ساخت درخت فیلوژنی، از روش‌های دیگری هم مانند Maximum Likelihood و یا استنتاج بیزی (Bayesian Inference) در ساخت این درختها استفاده می‌شوند. به پرسش‌های زیر در رابطه با هر یک از این الگوریتمها پاسخ دهید:

- الف) نحوه کار این الگوریتم چگونه است؟
- ب) آن را با الگوریتم Neighbor-Joining مقایسه کنید و بررسی کنید که در چه موقعی بهتر است از این الگوریتم استفاده کنیم.

- ج) دو الگوریتم استنتاج بیزی و Maximum-Likelihood را با هم مقایسه کنید و برتری هر یک را نسبت به دیگری بیان کنید.
۶. (۱۵ نمره) الگوریتم UPGMA را برای ماتریس فاصله زیر پیاده سازی کنید.

	A	B	C	D	E	F	G
A							
B	5						
C	9	10					
D	9	10	8				
E	8	7	7	3			
F	7	4	10	9	6		
G	12	11	6	5	4	9	

شکل ۱ : ماتریس فاصله

سوالات عملی (۱۰۰ نمره)

۱. (۱۰۰ نمره) برای سوالات عملی به *quera* مراجعه کنید..

(۱) (۷۰ نمره) برای تمرین عملی اول و دوم به *quera* مراجعه کنید. برای این تمرین می‌توانید از هر یک از زیان‌های *python*, *Java*, *C*, *C++* و *C#* استفاده کنید.

(۲) (۳۰ نمره) برای تمرین عملی دوم یک فایل در قالب *jupyter notebook* در اختیار شما قرار گرفته است که بایستی آن را دانلود و تمامی بخش‌های خواسته شده را به صورت کامل و بدون خطا اجرا نموده و آن را در محل تعیین شده آپلود نمایید.