

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

• در طول ترم، برای هر تمرین می‌توانید تا ۵ روز تأخیر داشته باشید (به غیر از تمرین چهارم که تا ۳ روز تأخیر قابل قبول می‌باشد) و در مجموع ۱۲ روز تأخیر مجاز خواهد داشت. پس از گذشت ۵ روز از مهلت ارسال امکان ارسال پاسخ به اتمام میرسد و به ازای هر میزان از تأخیر غیر مجاز به صورت خطی از نمره آن تمرین کسر خواهد شد.

- در صورت وجود هرگونه ابهام میتوانید سوالات خود را از طریق صفحه آن تمرین در کوئرا مطرح کنید.
- لطفاً پاسخ سوالات را در قالب یک فایل pdf به فرمت *STID\_HW<sup>۳</sup>.pdf* در کوئرا آپلود کنید.

### سوال ۱.

#### غنى‌شده‌گی موتیف در نواحی پروموتر: آزمون دقیق جایگشت (Permutation Test)

در یک مطالعه‌ی RNA-seq (شرایط A در برابر کنترل)، شما ۳ ژن با بیان افزایشی (up-regulated) از میان ۶ ژن بیان شده شناسایی کرده‌اید. فرض شما این است که یک فاکتور رونویسی (TF-X) این پاسخ را از طریق موتیف شناخته‌شده‌ی خود، یعنی موتیف ۶-تایی ACGTGC، تنظیم می‌کند. توالی‌های ۳۰ جفت‌پایه‌ای از نواحی پروموتر برای این شش ژن در جدول زیر آمده است.

بیان افزایشی؟	توالی ۳۰ جفت‌پایه‌ای پروموتر	ژن
خیر	TTAAGGACGTGCTTCCGATGGAATTGACA	A
بله	TTGGAATGCCATTGTTGGAATCCATTGGA	B
بله	GGTCCGTAACGTGCGGATTAACCTGGAAT	C
خیر	GCTTCTGGAATTGCAATGGTTAACCAATT	D
بله	CCATGGAAACGTGCTTAGGCTAACGATGTT	E
خیر	CGGTTGAATGCCATTGTTGGAACCATTG	F

آیا فراوانی وقوع ACGTGC در پرموتی‌ژن‌های با بیان افزایشی بیش از حد انتظار است؟ فرض صفر و فرض جایگزین را دقیقاً تعریف کرده و از یک آزمون دقیق جایگشت استفاده کنید و *p-value* یک‌طرفه را برای غنى‌شده‌گی محاسبه کنید. سطح معناداری را  $\alpha = 0.05$  در نظر بگیرید.

### سوال ۲.

#### تحلیل بیان ژن و محاسبه p-value تعديل شده

در یک مطالعه‌ی تحلیل بیان ژن، شش ژن برای تفاوت بیان میان دو شرایط با استفاده از t-test دو نمونه‌ای بررسی شده‌اند. جدول زیر مقدار log<sup>2</sup> Fold Change (تفاوت میانگین لگاریتمی بیان، حالت تیمارشده در برابر کنترل) و مقدار p-value خام هر آزمون را نشان می‌دهد.

ژن	log2FC	p-value خام
A	+1.2	0.004
B	+0.9	0.011
C	+0.7	0.021
D	+0.5	0.039
E	+0.3	0.081
F	-0.4	0.250

تعداد آزمون‌ها  $m = 6$  است. تمام محاسبات را به صورت دستی انجام دهید و مراحل میانی را نشان دهید.

۱. توضیح دهید که چرا نیاز به استفاده از  $p$ -value های تعديل شده داریم و چرا مقادیر خام گمراه کننده هستند.

۲. با استفاده از روش Bonferroni مقادیر تعديل شده را محاسبه کنید و مشخص کنید کدام ژن‌ها در سطح معناداری  $\alpha = 0.05$  معنی‌دار هستند؟

۳. با استفاده از روش Benjamini-Hochberg مقادیر تعديل شده را محاسبه کنید و مشخص کنید کدام ژن‌ها در سطح معناداری  $\alpha = 0.05$  معنی‌دار هستند؟

۴. توضیح دهید چرا روش BH در این مثال تعداد ژن‌های بیشتری را معنی‌دار تشخیص می‌دهد نسبت به بونفرونی. همچنین توضیح دهید هر روش چه خطای را کنترل می‌کند.

### سوال ۳

فرض کنید  $k$  و  $\ell$  اعداد صحیح مثبت ( $1 \leq k, \ell \leq l$ ) هستند. در هر یک از موارد زیر، رشته‌ی  $S$  با الگوی مشخصی داده شده است. با فرض اینکه یک کاراکتر پایان رشته  $\$$  (که از نظر ترتیب الفبایی از تمام حروف دیگر کوچکتر است) به انتهای  $S$  اضافه می‌شود (یعنی ورودی الگوریتم رشته‌ی  $S\$$  است)، تبدیل Burrows-Wheeler Transform یا همان رشته‌ی BWT خروجی را به صورت دقیق بر حسب  $k$  و  $\ell$  محاسبه کنید.

الف. رشته‌ی  $S$  به صورت  $(AT)^k A$  است (برای مثال اگر  $k = 2$  باشد، رشته برابر  $ATATA$  خواهد بود).

ب. رشته‌ی  $S$  به صورت  $T^k A^\ell$  است (برای مثال اگر  $k = 3$  و  $\ell = 4$  باشد، رشته برابر  $TTTAAAAA$  خواهد بود).

ج. رشته‌ی  $S$  به صورت  $A^k T A^k$  است (برای مثال اگر  $k = 3$  باشد، رشته برابر  $AAATAAAA$  خواهد بود).

### سوال ۴

در یک مطالعه‌ی Microarray، اثر یک داروی جدید بر روی ۱۰ ژن کاندیدا بررسی شده است. نرمافزار آماری، مقادیر P-value خام حاصل از آزمون t را برای این ۱۰ ژن گزارش کرده است. داده‌های خام (P-values) در جدول زیر آمده است:

۰.۰۶	۰.۰۰۹	۰.۰۳۵	۰.۶۰	۰.۰۰۳	۰.۱۵	۰.۰۱۴	۰.۸۲	۰.۰۰۰۲	۰.۰۴۵
------	-------	-------	------	-------	------	-------	------	--------	-------

با در نظر گرفتن سطح معنی‌داری کلی  $\alpha = 0.05$ ، به سوالات پاسخ دهید:

الف. روش Bonferroni: با توجه به تعداد آزمون‌ها ( $m = 10$ )، ابتدا مقدار آستانه جدید را محاسبه کنید. بر اساس این معیار، کدامیک از مقادیر فوق همچنان از نظر آماری معنی‌دار محسوب می‌شوند؟ (مقادیر آنها را ذکر کنید).

ب. روش FDR / Benjamini-Hochberg: فرمول محاسبه مقدار بحرانی را برای هر رتبه  $k$  بنویسید. جدول محاسباتی را تشکیل داده و بزرگترین رتبه  $k$  را که شرط الگوریتم در آن صدق می‌کند، بیابید. لیست نهایی ژن‌های پذیرفته شده در این روش کدامند؟

ج. به طور خاص، وضعیت ژن با مقدار  $P\text{-value} = 0.014$  را در دو روش بالا بررسی کنید. چرا سرنوشت این ژن در دو روش متفاوت است؟ (توضیح دهید که هر روش بر کنترل چه نوع خطای تمرکز دارد).

موفق باشید.