



دانشگاه صنعتی شریف

دانشکده مهندسی کامپیوتر

مقدمه‌ای بر بیوانفورماتیک

پاییز ۱۴۰۳

استاد: علی شریفی زارچی

مسئول تمرین: آراد ملکی

مهلت ارسال نهایی: ۸ دی

پاسخ نامه تمرین سوم

مهلت ارسال بدون تاخیر: ۵ دی

- 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 حداکثر ۸ روز تأخیر مجاز خواهید 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 `[STD_ID].pdf` `HW3` آپلود کنید.
- گردآوردندگان تمرین: آراد ملکی، علی حاجی صادقیان، آرزو پاک‌سرشت، سینا نمازی

سوالات نظری (۱۰۰ نمره)

۱. (۱۰ نمره) به سوالات زیر در رابطه با Transcription پاسخ دهید:

- الف) تفاوت‌های اساسی بین فرآیند رونویسی در پروکاریوت‌ها و یوکاریوت‌ها چیست و این تفاوت‌ها چگونه بر تنظیم ژن و پیچیدگی آن تأثیر می‌گذارند؟
- ب) چگونه اپی ژنتیک، از جمله متیلاسیون DNA و تغییرات هیستونی، می‌تواند بر الگوی رونویسی ژن‌ها تأثیر بگذارد و در نتیجه در تنظیم بیان ژن دخیل باشد؟
- پ) توضیح دهید که چگونه فاکتورهای رونویسی اختصاصی می‌توانند بر شروع رونویسی در سلول‌های یوکاریوتی تأثیر بگذارند و نقش آن‌ها در تنظیم بیان ژن چیست؟
- ت) چگونه مهارکننده‌های رونویسی می‌توانند به عنوان داروهای ضد سرطان عمل کنند؟ یک مثال از چنین مهارکننده‌ای را توضیح دهید و مکانیزم عمل آن را بیان کنید.
- ث) مفهوم Alternative Splicing چیست و چگونه می‌تواند به تنوع پروتئینی منجر شود؟
- حل.

الف) در پروکاریوت‌ها رونویسی و ترجمه هم‌زمان در سیتوپلاسم انجام می‌شود و معمولاً mRNA بدون تغییرات پس رونویسی (نظیر اسپلیسینگ) مستقیماً استفاده می‌گردد. در مقابل، یوکاریوت‌ها دارای هسته اند؛ رونویسی در هسته انجام و سپس mRNA با اضافه شدن کلاهک ۵'، دنباله پلی آ و حذف اینترون‌ها پردازش می‌شود. این فرایندهای اضافی و حضور عناصر تنظیمی پیچیده (نظیر پروموتورهای پیشرفته و اینهنسرها) باعث افزایش دقت، تنظیم چندسطحی و تنوع بیان ژن در یوکاریوت‌ها می‌شود.

ب) متیلاسیون DNA و تغییرات شیمیایی در هیستون‌ها ساختار کروماتین را تغییر داده و دسترسی فاکتورهای رونویسی را به ژن‌ها تنظیم می‌کنند. در نتیجه با روشن یا خاموش کردن ژن‌ها، بیان ژنی را کنترل می‌کنند.

پ) فاکتورهای رونویسی اختصاصی در سلول های یوکاریوتی به توالی های خاصی در DNA متصل می شوند (مانند پروموتورها و Enhancer ها) و با جذب یا دفع اجزای رونویسی، شروع رونویسی را تنظیم می کنند. آن ها نقش کلیدی در تنظیم بیان ژن دارند و پاسخ سلول به محرک های محیطی و سیگنال های درون سلولی را هماهنگ می کنند.

ت) مهارکننده های رونویسی با مسدود کردن فعالیت پروتئین های کلیدی درگیر در رونویسی ژن های مرتبط با رشد و تکثیر سلول های سرطانی، از پیشرفت سرطان جلوگیری می کنند. یک مثال از این داروها اکتینومایسین D است که با اتصال به DNA و جلوگیری از حرکت RNA پلی مراز، مانع از رونویسی mRNA می شود. این مکانیسم منجر به مهار رشد سلول های سرطانی می گردد.

ث) فرآیندی در RNA پردازش است که طی آن اگرژن ها به روش های مختلف به هم متصل می شوند یا اینترون ها حذف می گردند. این فرایند امکان تولید چندین نوع mRNA از یک ژن را فراهم می کند، که منجر به ساخت پروتئین های مختلف از یک توالی ژنی واحد می شود. به این ترتیب، تنوع پروتئینی بدون افزایش تعداد ژن ها حاصل می شود.

۲. (۱۰ نمره) به سؤالات زیر در مورد تحلیل داده های ریزآرایه (Microarray) پاسخ دهید:

الف) تحلیل داده های ریزآرایه چیست و چه کاربردی در بیوانفورماتیک دارد؟ توضیح دهید چگونه این فناوری می تواند بیان ژن ها را اندازه گیری کند.

ب) اصول اساسی طراحی یک آزمایش ریزآرایه را توضیح دهید و بیان کنید چرا انتخاب نمونه ها اهمیت دارد.

پ) اصول عملکرد پروب ها در ریزآرایه چیست و چرا انتخاب آن ها در طراحی آزمایش مهم است؟

ت) مراحل کلی تحلیل داده های ریزآرایه چیست؟ برای هر مرحله توضیح مختصری ارائه دهید.

ث) تفاوت اصلی بین آرایه های یک رنگ (Single-Color Array) و دو رنگ (Dual-Color Array) چیست و چه مزایا یا معایبی دارند؟

ج) مزایای استفاده از ریزآرایه در مقایسه با روش های سنتی تحلیل بیان ژن چیست؟

چ) چه نوع اطلاعاتی از داده های ریزآرایه قابل استخراج است و چگونه این اطلاعات به درک عملکرد ژن ها کمک می کند؟

ح) یکی از مشکلات معمول در تحلیل داده های ریزآرایه، نویزهای پس زمینه است. چه روش هایی برای کاهش این نویزها وجود دارد؟

خ) در تحلیل داده های ریزآرایه، چرا نرمال سازی ضروری است؟ دو روش نرمال سازی را توضیح دهید.

د) چگونه از داده های ریزآرایه برای شناسایی ژن های دیفرانسیلی بیان شده (Differentially Expressed Genes) استفاده می شود؟ مراحل این فرایند را شرح دهید.

ذ) محدودیت های ریزآرایه در مقایسه با تکنیک های جدیدتر مانند RNA-Seq چیست؟

حل.

الف) ریزآرایه یک فناوری است که به کمک آن می توان بیان همزمان هزاران ژن را اندازه گیری کرد. ریزآرایه ها مجموعه ای از پروب های DNA هستند که به یک سطح جامد متصل شده اند. نمونه های RNA یا DNA تک رشته ای به پروب ها هیبرید می شوند و سطح بیان ژن ها با اندازه گیری سیگنال های فلورسانس تعیین می شود. این تکنیک برای شناسایی سطح بیان ژن ها در شرایط مختلف استفاده می شود و این امکان را می دهد که بتوان ژن های دیفرانسیلی بیان شده (DEGs) را در شرایط بیمار و سالم شناسایی کرد.

ب) مراحل طراحی شامل انتخاب نمونه ها، استخراج RNA با کیفیت بالا، برچسب گذاری فلورسانس، هیبریداسیون و تحلیل سیگنال های فلورسانس است. انتخاب دقیق نمونه ها برای اطمینان از نتایج قابل اعتماد و نماینده از شرایط زیستی بسیار مهم است.

پ) پروب‌ها قطعات کوتاهی از DNA هستند که به توالی مکمل RNA یا DNA هدف متصل می‌شوند. انتخاب پروب‌های اختصاصی برای جلوگیری از اتصال غیراختصاصی و بهبود دقت آزمایش حیاتی است. معمولاً الیگونوکلوئوتیدها به دلیل قابلیت طراحی دقیق برای هر ژن خاص و ثبات بالا به‌عنوان پروب استفاده می‌شوند. در واقع، هنگام هیبریداسیون، فقط RNA یا cDNA مکمل به پروب‌ها متصل می‌شود و سیگنال فلورسانس تولید می‌کند.

ت) ۱. جمع‌آوری داده‌های خام: ثبت سیگنال‌های فلورسانس. ۲. نرمال‌سازی داده‌ها: حذف نویزها و قابل‌مقایسه کردن نمونه‌ها. ۳. شناسایی ژن‌های دیفرانسیلی بیان‌شده: تحلیل آماری برای یافتن ژن‌های با تغییر معنادار. ۴. تحلیل مسیرها: بررسی مسیرهای زیستی فعال. ۵. تفسیر نتایج: مقایسه با داده‌های زیستی موجود.

ث) عملکرد کلی آرایه‌های یک رنگ و دو رنگ مشابه است. در ریزآرایه‌های دو رنگ، دو نمونه بیولوژیکی (نمونه آزمایشی و نمونه کنترل) با رنگ‌های مختلف فلورسنت، معمولاً Cy3 (سیانین ۳) و Cy5 (سیانین ۵)، برچسب‌گذاری می‌شوند.

پس از هیبریداسیون رقابتی، اندازه‌گیری فلورسانس به طور جداگانه برای هر رنگ انجام شده و نشان‌دهنده فراوانی هر ژن در یک نمونه (نمونه آزمایشی، Cy5) نسبت به نمونه دیگر (نمونه شاهد، Cy3) است. داده‌های هیبریداسیون به صورت نسبت سیگنال‌های فلورسنت Cy5/Cy3 در هر پروب گزارش می‌شود. در این روش، مزیت اصلی امکان مقایسه مستقیم بین دو نمونه روی یک اسلاید است؛ اما عیب آن احتمال تداخل رنگ و نیاز به دقت بالا در انجام آزمایش است.

در مقابل، در ریزآرایه‌های یک رنگ، هر نمونه به طور جداگانه برچسب‌گذاری شده و در یک ریزآرایه مستقل هیبرید می‌شود. این روش مقدار مطلق فلورسانس را برای هر پروب به دست می‌دهد. به طور دقیق‌تر، فقط یک نمونه روی هر اسلاید قرار می‌گیرد که باعث وضوح بیشتر داده‌ها و حذف تداخل رنگ می‌شود.

ج) امکان تحلیل همزمان هزاران ژن، صرفه‌جویی در زمان و هزینه، و دقت بالا در شناسایی بیان ژن‌های با سطح پایین از مزایای ریزآرایه است.

چ) داده‌های ریزآرایه سطح بیان ژن، ژن‌های دیفرانسیلی بیان‌شده، و شبکه‌های تنظیمی را نشان می‌دهند. این اطلاعات برای درک عملکرد ژن‌ها و مسیرهای زیستی مفید است.

ح) روش‌هایی مانند نرمال‌سازی پس‌زمینه (RMA)، حذف پروب‌های ضعیف و فیلترگذاری داده‌ها به کاهش نویز کمک می‌کند.

خ) نرمال‌سازی داده‌ها برای کاهش نویز، تنظیم مقیاس داده‌ها، و مقایسه‌پذیری بین نمونه‌ها ضروری است. دو روش نرمال‌سازی:

۱. درون‌آرایه‌ای: تنظیم تفاوت رنگ‌های فلورسانس (مانند Lowess).

۲. بین‌آرایه‌ای: یکسان‌سازی توزیع سیگنال‌ها در آرایه‌ها (مانند Quantile Normalization).

د) فرآیند شامل شناسایی ژن‌های دیفرانسیلی بیان‌شده از طریق تحلیل آماری (مانند t-test یا ANOVA)، تعیین ژن‌هایی با adjusted p-value پایین‌تر از یک آستانه مشخص، و تهیه لیستی از ژن‌های هدف برای تحلیل بیشتر است. ژن‌هایی که سطح بیان آن‌ها در شرایط مختلف تفاوت معنی‌دار دارند شناسایی می‌شوند. این ژن‌ها می‌توانند به‌عنوان ژن‌های دیفرانسیلی بیان‌شده معرفی شوند.

ذ) محدودیت‌های ریزآرایه شامل عدم حساسیت به توالی‌های ناشناخته، دقت پایین در اندازه‌گیری ژن‌های با بیان کم، و نیاز به طراحی پروب‌ها است. مقایسه مزایای این روش شامل هزینه کمتر، توان عملیاتی بالا و مناسب بودن برای ژنوم‌های با اطلاعات کامل است.

۳. (۲۰ نمره) در این تمرین به بررسی اثر دیابت نوع ۲ بر سطح بیان یک ژن خاص می‌پردازیم. ژن GLUT۴ نقش کلیدی در انتقال گلوکز به داخل سلول‌ها و ارتباط مستقیم با حساسیت به انسولین دارد. کاهش بیان این ژن ممکن است یکی از عوامل مهم در ایجاد یا پیشرفت دیابت نوع ۲ باشد. می‌خواهیم بدانیم آیا سطح بیان ژن

GLUT4 در بیماران مبتلا به دیابت نوع ۲ نسبت به افراد سالم تغییر کرده است یا خیر. برای این منظور، سطح بیان ژن GLUT4 در دو گروه از افراد اندازه گیری شده است. گروه اول شامل نمونه های سالم و گروه دوم شامل نمونه های بیمار مبتلا به دیابت نوع ۲ است. سطح بیان این ژن در این نمونه ها پس از نرمال سازی به شرح زیر گزارش شده است:

۱۲/۳	۱۲/۲	۱۱/۷	۱۲/۰	۱۲/۴	۱۱/۹	۱۲/۱	۱۲/۵	۱۱/۸	۱۲/۳
------	------	------	------	------	------	------	------	------	------

جدول ۱: سطح بیان ژن GLUT4 در نمونه های سالم

۱۱/۰	۱۰/۷	۱۰/۸	۱۱/۱	۱۰/۹	۱۰/۶	۱۰/۷	۱۱/۰	۱۰/۸	۱۰/۵
------	------	------	------	------	------	------	------	------	------

جدول ۲: سطح بیان ژن GLUT4 در نمونه های بیمار

- الف) مشخص کنید که کدام آماره برای بررسی تفاوت بین دو گروه استفاده می شود و چرا مناسب است؟
 ب) فرض صفر (H_0) و فرض مقابل (H_1) در این آزمایش را تعریف کنید.
 ج) آزمون t را انجام دهید. آماره t را محاسبه کنید و مقدار p -value را به دست آورید.
 د) آیا مقدار p -value کمتر از سطح معنی داری (۰/۰۵) است؟ توضیح دهید که آیا فرض صفر رد می شود یا خیر.
 ه) براساس نتایج به دست آمده، توضیح دهید که آیا دیابت نوع ۲ تأثیری بر سطح بیان ژن GLUT4 داشته است؟

حل.

الف) آماره t مناسب است. چون دو گروه مستقل هستند، متغیر کمی است، هدف مقایسه میانگین هاست و با فرض نرمال بودن داده ها.

- ب) فرض صفر: میانگین سطح بیان ژن GLUT4 در دو گروه (سالم و بیمار) یکسان است.
- فرض مقابل: میانگین سطح بیان ژن GLUT4 در دو گروه (سالم و بیمار) متفاوت است.

ج)

$$\bar{X} = \frac{\sum x_i}{n},$$

$$s^2 = \frac{\sum (\bar{X} - x_i)^2}{n - 1},$$

$$t = \frac{\bar{X}_2 - \bar{X}_1}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}},$$

$$df = \frac{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}\right)^2}{\frac{\left(\frac{s_1^2}{n_1}\right)^2}{n_1 - 1} + \frac{\left(\frac{s_2^2}{n_2}\right)^2}{n_2 - 1}}.$$

$$\bar{X}_1 = \frac{11 + 10/7 + 10/8 + 11/1 + 10/9 + 10/6 + 10/7 + 11 + 10/8 + 10/5}{10} = 10/81,$$

$$\bar{X}_2 = \frac{12/3 + 12/2 + 11/7 + 12 + 12/4 + 11/9 + 12/1 + 12/5 + 11/8 + 12/3}{10} = 12/12.$$

$$s_1^2 = \frac{(10/81 - 11)^2 + (10/81 - 10/7)^2 + \dots + (10/81 - 10/5)^2}{9}$$

$$= \frac{0/0361 + 0/0121 + 0/0001 + 0/0841 + 0/0081 + 0/0441}{9}$$

$$+ \frac{0/0121 + 0/0361 + 0/0001 + 0/0961}{9} = 0/0366.$$

$$s_2^2 = \frac{(12/12 - 12/3)^2 + (12/12 - 12/2)^2 + \dots + (12/12 - 12/3)^2}{9}$$

$$= \frac{0/0324 + 0/0064 + 0/1764 + 0/0144 + 0/0784 + 0/0484}{9}$$

$$+ \frac{0/0004 + 0/1444 + 0/1024 + 0/0324}{9} = 0/0707.$$

$$t = \frac{12/12 - 10/81}{\sqrt{\frac{0/0366}{10} + \frac{0/0707}{10}}}$$

$$= \frac{1/31}{0/1036} = 12/65.$$

$$df = \frac{\left(\frac{0/0366}{10} + \frac{0/0707}{10}\right)^2}{\frac{\left(\frac{0/0366}{10}\right)^2}{9} + \frac{\left(\frac{0/0707}{10}\right)^2}{9}}$$

$$= \frac{0/0001151329}{\frac{0/0000133956}{9} + \frac{0/0000499849}{9}} = 16/35.$$

برای $t = 12/65$ و $df = 16/35$:

$$p = 2 \cdot (1 - CDF(|t|, df)).$$

مقدار p محاسبه شده برابر است با:

$$p = 7/29 \times 10^{-10}.$$

د) مقدار p -value محاسبه شده از سطح معناداری $0/05$ خیلی کمتر است. در نتیجه فرض صفر رد می‌شود.

ه) بله، با توجه به اینکه فرض صفر رد شد، پس میانگین سطح بیان ژن 4GLUT در دو گروه (سالم و بیمار) متفاوت است.

۴. (۲۰ نمره) در این سوال به بررسی ارتباط بین دو دارو و خطر حمله قلبی می‌پردازیم. یک مطالعه بالینی طراحی شده است تا ارتباط بین مصرف دو داروی Aspirin و Ibuprofen و کاهش خطر ابتلا به حمله قلبی (Myocardial Infarction) بررسی شود. این مطالعه شامل گروهی از افراد است که به طور تصادفی به دو دسته تقسیم شده‌اند: گروه اول از داروی Aspirin و گروه دوم از داروی Ibuprofen استفاده کرده‌اند. در پایان دوره‌ی مطالعه، محققان نتایج زیر را ثبت کردند:

گروه	داروی Aspirin	داروی Ibuprofen	مجموع
دچار حمله قلبی نشده‌اند	۲۲۰	۱۸۰	۴۰۰
دچار حمله قلبی شده‌اند	۸۰	۱۲۰	۲۰۰
مجموع	۳۰۰	۳۰۰	۶۰۰

فرضیات مطالعه:

H_0 : هیچ ارتباطی بین مصرف دارو و خطر ابتلا به حمله قلبی وجود ندارد. به عبارت دیگر، شانس وقوع حمله قلبی برای دو دارو یکسان است.

H_1 : شانس وقوع حمله قلبی برای افرادی که داروی Aspirin مصرف کرده‌اند کمتر از افرادی است که داروی Ibuprofen مصرف کرده‌اند.

(الف) نسبت شانس (Odds Ratio) را محاسبه کنید. با استفاده از داده‌های جدول، نسبت شانس برای وقوع حمله قلبی در گروه‌های مصرف‌کننده دو دارو را محاسبه کنید. توضیح دهید که این نسبت چه مفهومی دارد.

(ب) نتیجه محاسبات خود را در چارچوب فرضیه‌های صفر و جایگزین تفسیر کنید. آیا نسبت شانس به نفع داروی خاصی است؟

(ج) از آزمون دقیق فیشر استفاده کنید تا p-value را برای این داده‌ها محاسبه کنید. توضیح دهید که این آزمون چگونه کار می‌کند و چرا در این مطالعه به کار می‌رود.

(د) نتیجه آزمون را در چارچوب فرضیه‌ها تحلیل کنید. آیا مصرف داروی Aspirin تأثیر معناداری در کاهش خطر حمله قلبی دارد؟

حل.

(الف)

$$\text{Ratio Odds} = \frac{1 \text{ group in event of Odds}}{2 \text{ group in event of Odds}}$$

$$(\text{Aspirin}) \text{ Odds} = \frac{(\text{Aspirin}) \text{ occurred attack Heart}}{(\text{Aspirin}) \text{ occur not did attack Heart}} = \frac{80}{220} = 0.364$$

$$(\text{Ibuprofen}) \text{ Odds} = \frac{(\text{Ibuprofen}) \text{ occurred attack Heart}}{(\text{Ibuprofen}) \text{ occur not did attack Heart}} = \frac{120}{180} = 0.667$$

$$\text{OR} = \frac{(\text{Aspirin}) \text{ Odds}}{(\text{Ibuprofen}) \text{ Odds}} = 0.546$$

نسبت شانس (Ratio Odds) یا (OR) یک معیار آماری است که برای مقایسه شانس وقوع یک رخداد در دو گروه مختلف استفاده می‌شود. این مفهوم به ویژه در مطالعات پزشکی و اپیدمیولوژی برای بررسی ارتباط میان یک عامل (مثل مصرف دارو) و یک پیامد (مثل بروز بیماری) به کار می‌رود.

تفسیر نسبت شانس:

- $\text{OR} = 1$: شانس وقوع رخداد در دو گروه یکسان است (هیچ ارتباطی وجود ندارد).
- $\text{OR} < 1$: شانس وقوع رخداد در گروه اول بیشتر از گروه دوم است.
- $\text{OR} > 1$: شانس وقوع رخداد در گروه اول کمتر از گروه دوم است.

(ب)

مصرف داروی خاص (Aspirin) تأثیری بر کاهش یا افزایش خطر وقوع حمله قلبی ندارد. $OR = 1$: H_0

مصرف داروی Aspirin تأثیر دارد $OR \neq 1$: H_1

نتیجه محاسبات: مقدار OR محاسبه شده برابر با ۰/۵۴۶ است. این مقدار کمتر از ۱ است، به این معنا که احتمال وقوع حمله قلبی با مصرف داروی Aspirin نسبت به Ibuprofen کاهش می یابد. پس مصرف Aspirin به نسبت Ibuprofen به نفع کاهش خطر حمله قلبی است.

(ج) آزمون دقیق فیشر احتمال مشاهده جدول توافقی فعلی (یا جدولی با نسبت های شدیدتر) را در صورت درست بودن فرضیه صفر محاسبه می کند. این آزمون می تواند به طور دقیق بررسی کند که آیا رابطه معناداری بین دو متغیر (مثلاً استفاده از دارو و بروز نتیجه خاص) وجود دارد یا خیر.

(د) مقدار p-value به دست آمده ۰/۰۰۰۷۱۲ است که از ۰/۰۵ کمتر است. یعنی فرض صفر رد می شود. پس می توان گفت که Ibuprofen نسبت به Aspirin در کاهش شانس حمله قلبی مؤثرتر است.

۵. (۲۰ نمره) به سؤالات زیر در مورد RNA-Seq و پایگاه های مرتبط پاسخ دهید:

(الف) تکنیک RNA-Seq چیست و چگونه به مطالعه بیان ژن کمک می کند؟ مزایای این روش نسبت به تکنیک های سنتی مانند Microarray چیست؟

(ب) تفاوت میان داده های raw و processed در RNA-Seq چیست؟ چرا هر دو نوع داده برای تحلیل ها اهمیت دارند؟

(پ) پایگاه داده SRA چیست و چه نوع داده هایی در آن ذخیره می شود؟ چگونه می توان به داده های RNA-Seq در این پایگاه دسترسی پیدا کرد؟

(ت) یک مطالعه مرتبط با سرطان در پایگاه SRA پیدا کنید که داده های RNA-Seq را ارائه داده باشد. Accession Number آن را مشخص کنید و مراحل جستجوی خود را توضیح دهید.

(ث) با استفاده از ابزار NCBI SRA Toolkit داده های خام مطالعه ای که در بخش قبل یافتید را دانلود کنید. دستورات مورد استفاده برای این کار را بیان کنید.

(ج) پایگاه داده Ensembl را بازدید کنید و یک ژن خاص مرتبط با بیماری (مثلاً BRCA1 برای سرطان پستان) را پیدا کنید. داده های بیان ژن مرتبط با این ژن را جستجو و توضیح دهید چه اطلاعاتی ارائه می شود.

حل.

(الف) تکنیک RNA-Seq روشی پیشرفته برای تحلیل بیان ژن است که با استفاده از توالی یابی نسل جدید (Next-Generation Sequencing) انجام می شود. در این روش، RNA استخراج شده از نمونه ها به cDNA تبدیل و سپس توالی یابی می شود تا میزان بیان ژن ها با دقت بسیار بالا مشخص شود. مزیت RNA-Seq نسبت به تکنیک های سنتی مانند Microarray در این است که RNA-Seq نیازی به پروب های از پیش طراحی شده ندارد، امکان شناسایی ترانسکرپت های جدید، انواع مختلف RNA و بیان ژن ها در مقیاس وسیع را فراهم می کند و حساسیت بیشتری در اندازه گیری بیان ژن دارد.

(ب) داده های raw در RNA-Seq شامل اطلاعات خام توالی یابی شده هستند که معمولاً در قالب فایل هایی مانند FASTQ ارائه می شوند و مستقیماً از دستگاه توالی یاب به دست می آیند. داده های processed نتایج تحلیل های پردازشی روی داده های خام هستند که شامل اطلاعاتی مانند مقادیر نرمال شده بیان ژن و ترانسکرپت ها می شود. داده های خام برای بازتحلیل و بررسی کیفیت ضروری اند، در حالی که داده های پردازش شده برای تحلیل زیستی نهایی و تفسیر داده ها اهمیت دارند.

پ) پایگاه داده SRA (Sequence Read Archive) یکی از منابع اصلی ذخیره‌سازی داده‌های خام توالی‌یابی است که توسط NCBI نگهداری می‌شود. این پایگاه داده شامل توالی‌های خام از پروژه‌های RNA-Seq و سایر انواع توالی‌یابی است. برای دسترسی به داده‌های RNA-Seq در این پایگاه، می‌توان از شناسه‌های مربوط به پروژه‌ها یا نمونه‌ها استفاده کرد و داده‌ها را با ابزارهایی مانند SRA Toolkit دانلود کرد. این امکان برای محققان فراهم می‌شود تا داده‌های خام را جهت بازتحلیل یا استفاده در پروژه‌های جدید دریافت کنند.

۱. (۵۰ نمره) برای سوالات پاسخ به سوالات زیر به پایگاه داده GEO مراجعه فرمایید.

- الف) داده‌ای به دلخواه انتخاب کنید و Accession Number آن را مشخص کنید (مثال: GSE4107). سپس با استفاده از GEO2R، آن را تحلیل کنید.
- ب) مراحل جستجو و انتخاب داده را توضیح دهید.
- ج) ژن‌های با تفاوت بیان معنی‌دار بین گروه کنترل و بیمار را شناسایی کنید.
- د) نتایج را به صورت یک جدول شامل ژن‌های مهم و مقدار p-value ارائه دهید.

- حل. ۱. جستجو و انتخاب داده در GEO
۲. در نوار جستجو، کلمه کلیدی مرتبط با تحقیق (مثلاً cancer microarray) را وارد کنید یا مستقیماً Accession Number داده دلخواه خود را وارد کنید. مثال: Number Accession = GSE4107 (داده‌های مربوط به سرطان روده).
۳. داده مورد نظر را انتخاب کنید و صفحه مربوط به آن را باز کنید. در این صفحه اطلاعاتی شامل نوع مطالعه، پلتفرم استفاده شده (مانند Affymetrix)، و لینک دانلود داده‌ها وجود دارد.
۴. استفاده از ابزار GEO2R
۵. در صفحه داده، روی گزینه Analyze with GEO2R کلیک کنید.
۶. گروه‌های نمونه‌ها (مانند Case و Control) را مشخص کنید.
۷. برای این کار از بخش Define groups استفاده کنید و نمونه‌ها را به دو گروه تقسیم کنید.
۸. تنظیمات پیش فرض آماری ابزار را بررسی کنید (برای شناسایی ژن‌های دیفرانسیلی بیان شده).
۹. اجرای تحلیل و استخراج ژن‌های دیفرانسیلی بیان شده (DEGs)
۱۰. پس از اجرای GEO2R، جدولی تولید می‌شود که شامل اطلاعاتی مانند:
- Gene Symbol: نام ژن
- Log₂ Fold Change: میزان تغییر بیان ژن
- p-value: سطح معنی‌داری آماری
۱۱. فیلترسازی کنید:
- ژن‌هایی با $|\text{Log}_2 \text{ Fold Change}| > 1$ (افزایش یا کاهش بیان معنادار).
- ژن‌هایی با $p\text{-value} < 0.05$ (سطح معنی‌داری).
۱۲. جدول خروجی را دانلود کنید و گزارش کنید.

۲. (۵۰ نمره) برای حل این سوال به نوتبوک تست‌های آماری در کوئرا مراجعه کنید.