

به نام خدا

تمرین سری سوم
درس مقدمه ای بیوانفورماتیک
دکتر علی شریفی زارچی

فرزان رحمانی
۴۰۳۲۱۰۷۲۵

سوال اول

الف) تفاوت بین رونویسی در پروکاریوت ها و یوکاریوت ها

۱. مکانیسم های اساسی:

- پروکاریوت ها: رونویسی و ترجمه به طور همزمان در سیتوپلاسم اتفاق می افتد. یک RNA پلیمرز منفرد رونویسی را انجام می دهد.
- یوکاریوت ها: رونویسی در هسته، و ترجمه در سیتوپلاسم اتفاق می افتد. سه RNA پلیمرز مجزا (I، II، III) برای انواع مختلف RNA وجود دارد.

۲. ساختار ژن:

- پروکاریوت ها: ژن ها اغلب به اپرون (operons) سازمان دهی می شوند و به ژن های متعدد اجازه می دهند تا به عنوان یک mRNA واحد (پلی سیسترونیک) رونویسی شوند.
- یوکاریوت ها: ژن ها به طور معمول مونوسیسترونیک (monocistronic) هستند که به صورت جداگانه با اینترون ها و اگزون ها رونویسی می شوند.

۳. پیچیدگی Regulation:

- پروکاریوت ها: Regulation عمدتاً در سطح رونویسی با استفاده از سرکوب کننده ها و فعال کننده های متصل به اپراتورها است.
- یوکاریوت ها: Regulation پیچیده تر است و شامل بازسازی کروماتین، تقویت کننده ها، خاموش کننده ها و فاکتورهای رونویسی می شود.

۴. تأثیر بر تنظیم ژن (Gene Regulation):

- پروکاریوت ها: تنظیم ساده تر اجازه می دهد تا به تغییرات محیطی واکنش سریع نشان دهد.
- یوکاریوت ها: تنظیم پیچیده، کنترل دقیق، بافت خاص و رشد بیان ژن را امکان پذیر می کند.

ب) اپی ژنتیک و تأثیر آن بر رونویسی ژن

۱. متیلاسیون DNA:

- متیلاسیون باقیمانده های سیتوزین (معمولاً در جزایر CpG) با مسدود کردن اتصال فاکتور رونویسی یا جذب پروتئین های سرکوب کننده، رونویسی را سرکوب می کند.
- مثال: Hypermethylation ژن های سرکوبگر تومور می تواند منجر به سرطان شود.

۲. اصلاحات هیستون (Histone Modifications):

- استیلایسیون (Acetylation): افزودن گروه های استیل به هیستون ها ساختار کروماتین را شل می کند و دسترسی به ژن و رونویسی را افزایش می دهد.
- متیلایسیون (Methylation): بسته به مکان و زمینه می تواند رونویسی را فعال یا سرکوب کند.
- فسفوریلاسیون (Phosphorylation) و یوبی کوئیتیناسیون (ubiquitination): سایر تغییرات موثر بر دینامیک کروماتین.

۳. نقش در Regulation:

- تغییرات اپی ژنتیک هویت سلول را ایجاد می کند و به سیگنال های محیطی و رشدی پاسخ های پویا می دهد.
- پ) نقش فاکتور های رونویسی اختصاصی

۱. عملکرد در آغاز:

- فاکتورهای رونویسی خاص به توالی های تنظیم کننده DNA (پیش گیرنده، تقویت کننده یا همان promoters, enhancers) متصل می شوند و ماشین های رونویسی عمومی را به کار می گیرند.
- آنها به RNA پلیمراز در تشکیل کمپلکس شروع رونویسی در پروموتور کمک می کنند.

۲. تنظیم بیان ژن:

- عمل به عنوان فعال کننده یا سرکوب کننده:
 - فعال کننده ها (Activators): رونویسی را با به کارگیری فعال کننده های فعال یا تثبیت اتصال RNA پلیمراز تقویت می کنند.
 - سرکوبگرها (Repressors): رونویسی را با مسدود کردن اتصال فعال کننده یا به کارگیری corepressors مهار می کنند.

۳. مثال: فاکتور رونویسی p53 به عنوان یک سرکوب کننده تومور عمل می کند و ژن های دخیل در توقف چرخه سلولی و آپوپتوز (apoptosis) را تنظیم می کند.

ت) مهارکننده های رونویسی به عنوان داروهای ضد سرطان

۱. مکانیسم عمل:

- مهارکننده ها فرآیندهای رونویسی را برای سرکوب بیان oncogenes یا القای آپوپتوز (apoptosis) در سلول های سرطانی هدف قرار می دهند.

۲. مثال: Actinomycin D (Dactinomycin)

- در کمپلکس شروع رونویسی به DNA متصل می شود و از طول شدن توسط RNA پلیمراز جلوگیری می کند.
- در درمان تومور ویلمز (Wilms' tumor)، rhabdomyosarcoma و سایر سرطان ها استفاده می شود.

۳. منطق و نقش آنها:

- سلول های سرطانی برای تکثیر سریع به شدت به برنامه های رونویسی خاصی متکی هستند و رونویسی را به یک هدف قابل دوام تبدیل می کنند.

ث) مفهوم Alternative Splicing

۱. تعریف:

- Alternative splicing یک فرآیند پس از رونویسی است که در آن یک pre-mRNA منفرد به روش های مختلف برای تولید ایزوفرم های مختلف mRNA بالغ متصل می شود.

۲. مکانیسم ها:

- شامل پرش اگزون، مکان‌های پیوند جایگزین 5' یا 3' (alternative 5' or 3' splice sites)، حفظ اینترون و اگزون‌های متقابل منحصر به فرد (mutually exclusive exons).

۳. تأثیر بر تنوع پروتئین:

- چندین پروتئین را از یک ژن تولید می‌کند و به تنوع عملکردی بدون افزایش اندازه ژنوم کمک می‌کند.
- مثال: ژن انسانی tropomyosin می‌تواند ایزوفرم‌های خاص بافت با عملکردهای متمایز تولید کند.

۴. اهمیت بیولوژیکی:

- برای توسعه، تمایز، و سازگاری ضروری است.
- خطا در alternative splicing می‌تواند منجر به بیماری‌هایی مانند سرطان یا اختلالات عصبی (neurodegenerative) شود.

سوال دوم

الف) تجزیه و تحلیل داده‌های ریزآرایه و کاربرد آن

۱. تعریف:

تجزیه و تحلیل داده‌های ریزآرایه شامل پردازش و تفسیر داده‌های آزمایش‌های ریزآرایه است که سطح بیان هزاران ژن را به طور همزمان اندازه‌گیری می‌کند.

۲. کاربرد در بیوانفورماتیک:

برای مطالعه بیان ژن، شناسایی نشانگرهای زیستی (biomarkers)، درک مکانیسم‌های بیماری و طبقه‌بندی بیماری‌ها استفاده می‌شود.

۳. اندازه‌گیری بیان ژن:

ریزآرایه‌ها از نمونه‌های نشاندار cDNA یا RNA استفاده می‌کنند که با پروب‌های مکمل روی یک تراشه هیبرید می‌شوند. سیگنال‌های فلورسنت ساطع شده توسط پروب‌های متصل، میزان بیان ژن‌های مربوطه را کمیت‌بندی می‌کنند.

ب) طراحی یک آزمایش ریزآرایه

۱. اصول اساسی:

- تعریف اهداف (Define objectives): ژن‌ها، شرایط یا مسیرهای مورد علاقه را شناسایی کنیم.
- کنترل‌ها و تکرارهای مناسب را برای اطمینان از قابلیت اطمینان آماری انتخاب کنیم.
- تنوع فنی و بیولوژیکی را در نظر بگیریم.

۲. اهمیت انتخاب نمونه (Sample Selection):

- relevance را تضمین می‌کند و عوامل مخدوش‌کننده را به حداقل می‌رساند.
- نمایش دقیق جمعیت یا شرایط مورد مطالعه، کیفیت و قابلیت تفسیر داده‌ها را افزایش می‌دهد.

پ) اصول و اهمیت پروب‌ها

۱. عملکرد پروب‌ها:

توالی‌های کوتاه DNA یا RNA روی تراشه ریزآرایه با توالی‌های هدف مکمل در نمونه hybridize می‌شوند.

۲. اهمیت انتخاب پروب:

- Specificity: از cross-hybridization برای تشخیص دقیق ژن اجتناب کنیم.
- Sensitivity: پروب ها باید transcript های کم فراوانی را تشخیص دهند.
- Coverage: همه ژن ها یا ایزوفرم های مورد نظر را نشان می دهد.

ت) مراحل تجزیه و تحلیل داده های ریزآرایه

۱. پیش پردازش (Preprocessing):

تصحیح پس زمینه، کاهش نویز، و افزایش سیگنال.

۲. نرمال سازی یا Normalization:

تغییرات فنی را برای اطمینان از مقایسه بین نمونه ها تنظیم کنیم.

۳. تجزیه و تحلیل آماری:

با استفاده از آزمون های آماری، ژن های با بیان متفاوت (DEGs) را شناسایی کنیم.

۴. تحلیل عملکردی:

DEG ها را به مسیرهای بیولوژیکی یا هستی شناسی ژن نگاشت کنیم.

۵. تفسیر داده ها (Data Interpretation):

ادغام یافته ها با دانش بیولوژیکی برای به دست آوردن نتایج معنادار.

ث) آرایه های یک رنگ در مقابل آرایه های دو رنگ

۱. Single-Color Array:

- از یک برچسب فلورسنت در هر نمونه استفاده می کند.
- مزیت: طراحی آزمایشی ساده تر.
- نقطه ضعف: به آرایه های بیشتری برای مقایسه نیاز دارد.

۲. Dual-Color Array:

- از دو برچسب فلورسنت برای مقایسه دو نمونه در یک آرایه استفاده می کند.
- مزیت: مقایسه مستقیم، تنوع فنی (technical variability) را کاهش می دهد.
- عیب: سوگیری بالقوه رنگ (Potential dye bias).

ج) مزایای ریزآرایه نسبت به روشهای سنتی

۱. توان عملیاتی بالا (High Throughput):

هزاران ژن را به طور همزمان تجزیه و تحلیل کنیم.

۲. کمی (Quantitative):

سطوح بیان نسبی (relative expression levels) را فراهم می کند.

۳. مقرون به صرفه (Cost-Effective):

ارزان تر از sequencing برای مطالعات در مقیاس بزرگ.

۴. همه کاره:

قابل استفاده برای موجودات و شرایط مختلف.

چ) اطلاعات قابل استخراج و ارتباط آنها با عملکرد ژن ها

۱. اطلاعات قابل استخراج:

- پروفایل های بیان ژن (Gene expression profiles)
- ژن های بیان شده متفاوت (Differentially expressed genes)
- شبکه های هم بیان ژن. (Gene co-expression networks)

۲. Relevance:

- نشانگرهای زیستی را شناسایی کنیم، عملکرد ژن را درک کنیم و مسیرهای بیولوژیکی (biological pathways) را مشخص کنیم.

ح) کاهش نویز پس زمینه

روش ها:

- پردازش تصویر: برای فلورسانس غیر اختصاصی (non-specific fluorescence) تنظیم کنیم.
- فیلتر کردن پروب: پروب های بی کیفیت یا نویزی را برداریم.
- تکنیک های عادی سازی: technical variability را کاهش دهیم.
- از positive control و negative control استفاده کنیم.

خ) نرمال سازی در تجزیه و تحلیل داده های ریزآرایه

۱. ضرورت:

تغییرات در آماده سازی نمونه (variations in sample preparation)، hybridization و scanning را تصحیح می کند.

۲. روش ها:

- Global Normalization: سیگنال ها به یک میانگین یا میانه یکنواخت مقیاس میکنند.
- Quantile Normalization: توزیع مقادیر expression ها را در میان آرایه ها تراز (align) می کند.

د) شناسایی ژن های بیان شده متفاوت (Differentially Expressed Genes)

۱. مراحل:

- پیش پردازش و نرمال سازی داده ها.
- انجام تست های آماری (به عنوان مثال، آزمون t، ANOVA).
- اصلاحات آزمایشی متعدد (multiple testing corrections) اعمال کنیم (به عنوان مثال، FDR).
- رتبه بندی ژن ها بر اساس fold change و اهمیت آنها و تحلیل p-value.

۲. خروجی:

- فهرست ژن هایی که در شرایط خاص، تغییر بیان یافته اند.

ذ) محدودیت های ریزآرایه در مقابل RNA-Seq

۱. محدودیت ها:

- طراحی پروب استاتیک (Static Probe Design): فقط توالی های از پیش تعریف شده را تشخیص می دهد.
- Limited Sensitivity: با transcripts های دارای فراوانی کم دست و پنجه نرم می کند.
- محدوده دینامیکی پایین (Lower Dynamic Range): دقت کمتر در تعیین کمیت سطوح بیان (Less precise quantification of expression levels).

۲. مزایای RNA-Seq:

- تشخیص Unbiased رونوشت ها (transcripts).
- شناسایی رونوشت های novel و رویدادهای alternative splicing.
- sensitivity و محدوده دینامیکی بیشتر

سوال سوم

الف) آزمون آماری مورد استفاده

برای بررسی تفاوت بین دو گروه از آزمون t مستقل (independent t-test) استفاده می شود. independent t-test مناسب است زیرا:

- داده ها از دو گروه مستقل (افراد سالم و بیماران) تشکیل شده است.
- هدف مقایسه میانگین سطوح بیان GLUT4 بین این گروه ها است.

ب) فرضیه ها

- فرض صفر (H_0 : Null Hypothesis): تفاوت معنی داری در سطح بیان ژن GLUT4 بین افراد سالم و بیماران مبتلا به دیابت نوع ۲ وجود ندارد.
- فرض مقابل (H_1 : Alternative Hypothesis): تفاوت معنی داری در سطح بیان ژن GLUT4 بین افراد سالم و بیماران مبتلا به دیابت نوع ۲ وجود دارد.

ج) نتایج آزمون تی (t-test)

- آماره t (T-statistic): $t = 12.61$
- p-value: $p = 9.998 \times 10^{-10}$

در ادامه راه حل و توضیح گام به گام نحوه محاسبه آماره t و p-value برای مقایسه سطح بیان ژن GLUT4 بین افراد سالم و بیماران آورده شده است:

مرحله ۱: داده ها را شناسایی کنیم.

سطوح بیان برای دو گروه عبارتند از:

- نمونه های سالم: [12.3, 12.2, 11.7, 12.0, 12.4, 11.9, 12.1, 12.5, 11.8, 12.3]
- نمونه های بیمار: [11.0, 10.7, 10.8, 11.1, 10.9, 10.6, 10.7, 11.0, 10.8, 10.5]

مرحله ۲: فرمول T-Test را تعریف میکنیم.

فرمول آماره t (t-statistic) در آزمون t مستقل (independent t-test) به صورت زیر است:

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

که:

- \bar{X}_1, \bar{X}_2 : میانگین دو گروه اند.

- s_1^2, s_2^2 : واریانس نمونه دو گروه اند.

- n_1, n_2 : اندازه نمونه از دو گروه اند.

مرحله ۳: محاسبه میانگین گروه

$$\bar{X}_1 = \frac{\text{Sum of healthy samples}}{n_1} = \frac{12.3 + 12.2 + \dots + 12.3}{10} = 12.12$$

$$\bar{X}_2 = \frac{\text{Sum of patient samples}}{n_2} = \frac{11.0 + 10.7 + \dots + 10.5}{10} = 10.81$$

مرحله ۴: واریانس ها را محاسبه می کنیم.

واریانس به صورت زیر محاسبه می شود:

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

برای نمونه های سالم (s_1^2):

$$s_1^2 = \frac{(12.3 - 12.12)^2 + (12.2 - 12.12)^2 + \dots + (12.3 - 12.12)^2}{10 - 1} = \frac{0.636}{9} = 0.071$$

برای نمونه های بیمار (s_2^2):

$$s_2^2 = \frac{(11.0 - 10.81)^2 + (10.7 - 10.81)^2 + \dots + (10.5 - 10.81)^2}{10 - 1} = \frac{0.329}{9} = 0.037$$

مرحله ۵: مقادیر را در فرمول قرار میدهیم.

$$t = \frac{12.12 - 10.81}{\sqrt{\frac{0.071}{10} + \frac{0.037}{10}}}$$

$$t = \frac{1.31}{\sqrt{0.0071 + 0.0037}} = \frac{1.31}{\sqrt{0.0108}} = \frac{1.31}{0.1039} = 12.61$$

مرحله ۶: محاسبه درجه آزادی

با استفاده از فرمول آزمون Welch t برای درجات آزادی (Welch's t-test formula for degrees of freedom):

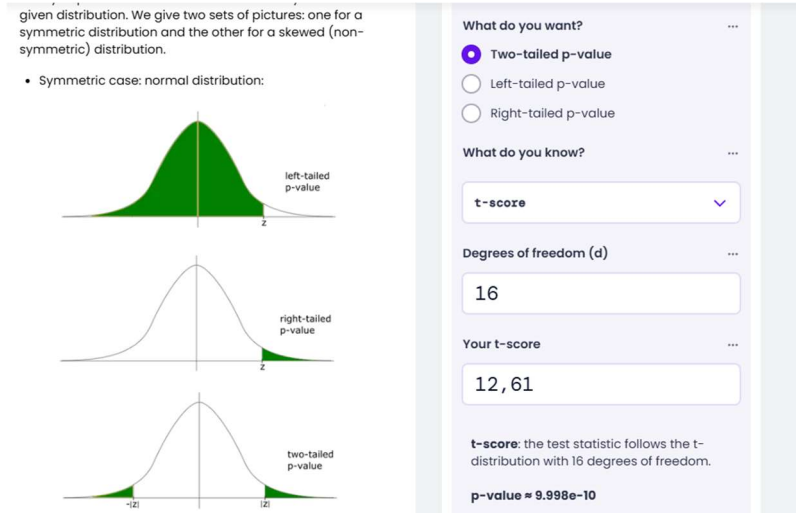
$$df = \frac{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}\right)^2}{\frac{\left(\frac{s_1^2}{n_1}\right)^2}{n_1 - 1} + \frac{\left(\frac{s_2^2}{n_2}\right)^2}{n_2 - 1}}$$
$$df = \frac{\left(\frac{0.071}{10} + \frac{0.037}{10}\right)^2}{\frac{\left(\frac{0.071}{10}\right)^2}{9} + \frac{\left(\frac{0.037}{10}\right)^2}{9}} = 16.377 \text{ (approximately 16)}$$

مرحله ۷: P-Value را پیدا کنیم.

با استفاده از یک جدول توزیع t یا نرم افزار آماری یا کد پایتون یا یک وب سایت مانند <https://www.omnicalculator.com/statistics/p-value>، مقدار p برای موارد زیر محاسبه می شود (توجه: مقدار p-value به صورت two-tailed محاسبه شده است که دو برابر one-tailed یا همان right-tailed است):

$$t = 12.61$$
$$\text{Degrees of freedom} = 16$$

همان طور که میبینیم مقدار p-value بسیار کوچک است ($p = 9.998 \times 10^{-10}$).



```
1 from scipy.stats import t
2
3 # Given values
4 t_statistic = 12.61
5 df = 16 # Degrees of freedom
6
7 # Calculate the two-tailed p-value
8 p_value = 2 * t.sf(abs(t_statistic), df)
9 # p_value_right_tailed = t.sf(abs(t_statistic), df)
10
11 print(f"P-value: {p_value:.10e}")
```

P-value: 9.9984873396e-10

نتایج نهایی

- آماره t (T-statistic): $t = 12.61$
- p-value: $p = 9.998 \times 10^{-1}$

(د) تصمیم گیری

مقدار p-value (9.998×10^{-1}) بسیار کوچکتر از سطح معنی داری مانند $\alpha = 0.05$ است. بنابراین:

- فرض صفر (H_0 : Null Hypothesis) رد می شود.
- فرض مقابل (H_1 : Alternative Hypothesis) درست است و در نتیجه تفاوت معنی داری در سطح بیان ژن GLUT4 بین افراد سالم و بیماران مبتلا به دیابت نوع ۲ وجود دارد.

$$9.998 \times 10^{-10} = p < \alpha = 0.05$$

(ه) نتیجه گیری

بر اساس نتایج، تفاوت آماری معنی داری در سطح بیان ژن GLUT4 بین افراد سالم و بیماران مبتلا به دیابت نوع ۲ وجود دارد چرا که مقدار p-value حاصل از t-test بسیار کوچکتر از سطح معنی داری مانند $\alpha = 0.05$ است. این امر نشان می دهد که دیابت نوع ۲ تأثیر قابل توجهی بر کاهش سطح بیان ژن GLUT4 دارد. به بیان دیگر چون مقدار p-value بسیار کوچکتر از سطح معنی داری مانند $\alpha = 0.05$ است، پس این اتفاق نمی تواند شانسی رویداده باشد و فرض مقابل (H_1 : Alternative Hypothesis) درست است یعنی تفاوت معنی داری در سطح بیان ژن GLUT4 بین افراد سالم و بیماران مبتلا به دیابت نوع ۲ وجود دارد.

سوال چهارم

الف) محاسبه نسبت شانس (Odds Ratio)

تعریف نسبت شانس (OR)

نسبت شانس، شانس وقوع یک رویداد در یک گروه را نسبت به گروه دیگر مقایسه می کند. برای این مطالعه:

$$\text{Odds Ratio} = \frac{\text{Odds in Aspirin group}}{\text{Odds in Ibuprofen group}}$$

محاسبه گام به گام

- شانس در گروه آسپرین:

$$\text{Odds}_{\text{Aspirin}} = \frac{\text{Number of Heart Attacks}}{\text{Number of No Heart Attacks}} = \frac{80}{220} = 0.3636$$

- شانس در گروه ایبوپروفن:

$$\text{Odds}_{\text{Ibuprofen}} = \frac{\text{Number of Heart Attacks}}{\text{Number of No Heart Attacks}} = \frac{120}{180} = 0.6667$$

- نسبت شانس (OR):

$$\text{OR} = \frac{\text{Odds}_{\text{Aspirin}}}{\text{Odds}_{\text{Ibuprofen}}} = \frac{0.3636}{0.6667} = 0.5455$$

شانس متقابل (reciprocal) یا معکوس آن هم برابر مقدار زیر است:

$$\text{OR}_{\text{Ibuprofen vs Aspirin}} = \frac{1}{0.5455} = 1.83$$

بنابراین، افرادی که ایبوپروفن مصرف می کنند، $1/83$ برابر بیشتر از افرادی که آسپرین مصرف می کنند، دچار حمله قلبی می شوند. به بیان دیگر، افرادی که آسپرین مصرف می کنند، $0/5455$ برابر کمتر از افرادی که ایبوپروفن مصرف می کنند، دچار حمله قلبی می شوند.

تفسیر مفصل تر نسبت شانس ها:

شانس در هر گروه

- شانس در گروه آسپرین: $0/3636$
 - به این معنی که به ازای هر ۱ عدم حمله قلبی، تقریباً $0/3636$ حمله قلبی در گروه آسپرین وجود دارد.
 - متناوباً، می توانید آن را تقریباً ۱ حمله قلبی در هر $2/75$ عدم حمله قلبی تفسیر کنید.
- شانس در گروه ایبوپروفن: $0/6667$
 - به این معنی که به ازای هر ۱ عدم حمله قلبی، تقریباً $0/6667$ حمله قلبی در گروه ایبوپروفن وجود دارد.
 - متناوباً، نشان دهنده ۱ حمله قلبی در هر $1/5$ عدم حمله قلبی است.

نسبت شانس (OR)

- نسبت شانس برای آسپرین در مقابل ایبوپروفن: $0/5455$
 - احتمال حمله قلبی در گروه آسپرین $0/5455$ برابر کمتر از گروه ایبوپروفن است.
 - این نشان می دهد که آسپرین با کاهش خطر حمله قلبی در مقایسه با ایبوپروفن همراه است.

نسبت شانس متقابل (Reciprocal Odds Ratio)

- نسبت شانس برای ایبوپروفن در مقابل آسپرین: $1/83$
 - احتمال حمله قلبی در گروه ایبوپروفن $1/83$ برابر بیشتر از گروه آسپرین است.

این نسبت ها به چه معناست؟

۱. اثربخشی داروها:

- به نظر می رسد آسپرین در کاهش خطر حملات قلبی در مقایسه با ایبوپروفن موثرتر باشد.

۲. پیامدهای بالینی:

- نسبت شانس کمتر ($OR < 1$) برای آسپرین نشان دهنده اثر محافظتی بالقوه آن است.
- نسبت معکوس ($OR > 1$) برای ایبوپروفن نشان می دهد که ممکن است خطر بیشتری نسبت به آسپرین داشته باشد.

۳. تصمیم گیری:

- این یافته ها می توانند بر تصمیم گیری های درمانی تأثیر بگذارند و از آسپرین به جای ایبوپروفن برای پیشگیری از حمله قلبی، با اعتبارسنجی آماری بیشتر (مثلاً p-value) تأثیر بگذارند.

ب) فرضیه ها و نتیجه محاسبات

- فرضیه صفر (H_0 : Null Hypothesis): هیچ تفاوتی در خطر حمله قلبی بین دو دارو وجود ندارد.

- فرضیه جایگزین (H_1 : Alternative Hypothesis): آسپرین در مقایسه با ایبوپروفن خطر حمله قلبی را کاهش می دهد.

از آنجایی که نسبت شانس مخالف صفر و کوچکتر از ۱ است ($OR = 0.5455$, $OR_{Ibuprofen\ vs\ Aspirin} = 1.83$)، خطر کمتر حمله قلبی را برای کسانی که آسپرین مصرف می کنند نشان می دهد که از فرضیه جایگزین (H_1) پشتیبانی می کند. به بیان دیگر فرضیه صفر رد می شود و فرضیه جایگزین تایید میشود البته باید مقدار p-value نیز برای تایید چک شود.

توضیحات مفصل تر در ادامه آمده است:

وقتی فرضیه صفر (H_0) درست باشد

- فرضیه صفر (H_0) بیان می کند که:

احتمال حمله قلبی برای مصرف کنندگان آسپرین و ایبوپروفن یکسان است

- از نظر ریاضی این یعنی:

$$\text{Odds Ratio (OR)} = 1$$

- اگر $OR = 1$ باشد، هیچ تفاوتی بین دو گروه از نظر خطر حمله قلبی وجود ندارد.

وقتی $OR \neq 1$

- $OR > 1$:

- احتمال حمله قلبی در گروه آسپرین بیشتر از گروه ایبوپروفن است.
- این نشان می دهد که آسپرین ممکن است خطر را نسبت به ایبوپروفن افزایش دهد.

- $OR < 1$:

- احتمال حمله قلبی در گروه آسپرین کمتر از گروه ایبوپروفن است.
- این نشان می دهد که آسپرین خطر حمله قلبی را در مقایسه با ایبوپروفن کاهش می دهد، همانطور که در این مطالعه مشاهده شد ($OR = 0.5455$).

نکات کلیدی در مورد نسبت شانس (OR)

۱. چگونه OR خطر را منعکس می کند:

- $OR = 1$: تفاوتی وجود ندارد (H_0 درست است).
- $OR < 1$: گروه اول (آسپرین) شانس کمتری برای نتیجه (حمله قلبی) دارند.
- $OR > 1$: گروه اول (آسپرین) شانس بیشتری برای نتیجه دارد.

۲. اهمیت آماری:

- حتی اگر $OR < 1$ یا $OR > 1$ باشد، باید از نظر آماری معنی دار باشد ($p\text{-value} < 0.05$) تا H_0 با اطمینان رد شود.

توجه به عدم تفسیر نادرست از نسبت شانس

- OR همیشه مثبت است زیرا شانس به صورت مقادیر مثبت بیان می شود.
- $OR = 0$ نامعتبر است زیرا به این معنی است که شانس یک گروه برای نتیجه صفر است که واقع بینانه نیست.
- یک OR کوچک (به عنوان مثال، نزدیک به ۰/۱) یک اثر محافظتی بسیار قوی از حمله قلبی را نشان می دهد.

در این مطالعه $OR = 0.5455$:

- نشان دهنده اثر محافظتی آسپرین در مقایسه با ایبوپروفن است.
- خطر حمله قلبی برای کسانی که آسپرین مصرف می کنند در مقایسه با ایبوپروفن تقریباً ۴۵ درصد کمتر است.

ارتباط با فرضیه ها:

- اگر $p\text{-value}$ از آزمایش دقیق فیشر کمتر از ۰/۰۵ باشد، H_0 را رد می کنیم و می پذیریم که آسپرین به طور قابل توجهی خطر حمله قلبی را کاهش می دهد.

ج) آزمون دقیق فیشر: آزمون دقیق فیشر یک روش آماری است که برای تعیین اینکه آیا ارتباط معنی داری بین دو متغیر طبقه بندی شده (categorical variable) در جدول احتمالی (contingency table) وجود دارد یا خیر.

- P-value محاسبه شده: $p = 0.00071$

ابتدا راجع به خود آزمون، نحوه کارکرد آن و اینکه چرا در این مطالعه استفاده میشود توضیح میدهم و سپس وارد محاسبات می شویم:

هدف این آزمون: این آزمون تعیین می کند که آیا ارتباط معنی داری بین دو متغیر طبقه بندی شده در جدول احتمالی وجود دارد یا خیر. احتمال دقیق مشاهده داده های داده شده (یا شدیدتر: or more extreme) را تحت فرضیه صفر محاسبه می کند.

- آزمون دقیق فیشر برای بررسی اهمیت ارتباط بین دو متغیر طبقه بندی شده (categorical variables) در جدول احتمالی 2×2 استفاده می شود.
- برای اندازه های نمونه کوچک یا زمانی که فرکانس های مورد انتظار در جدول پایین است ایده آل است چرا که باعث می شود تست های دیگری مانند تست Chi-square کمتر قابل اعتماد باشند.
- احتمال دقیق مشاهده داده ها را تحت فرضیه صفر محاسبه می کند.

چگونه کار می کند؟

۱. فرضیه های صفر و جایگزین:

- H_0 : هیچ ارتباطی بین این دو متغیر وجود ندارد (به عنوان مثال، داروی مصرف شده و احتمال حمله قلبی).
- H_1 : بین این دو متغیر ارتباط وجود دارد (به عنوان مثال، داروی مصرف شده بر احتمال حمله قلبی تأثیر می گذارد).

۲. محاسبه دقیق احتمال:

- آزمون فیشر احتمال دقیق (p-value) مشاهده داده ها را تحت فرضیه صفر محاسبه می کند.
- از توزیع فوق هندسی (hypergeometric distribution) برای محاسبه احتمال تمام نتایج ممکن به صورت افراطی یا شدیدتر از آنچه مشاهده شده است (extreme or more extreme than the observed one)، با توجه به مجموع سطرها و ستون ها استفاده می کند.

۳. داده های ورودی:

این آزمون از یک جدول احتمالی 2×2 استفاده می کند.

۴. فرمول کلیدی:

$$P = \frac{\binom{a+b}{a} \binom{c+d}{c}}{\binom{n}{a+c}}$$

که:

- a, b, c, d : تعداد سلول ها در جدول اند.
- n : تعداد کل مشاهدات است.

به هر جدول ممکن یک احتمال اختصاص داده می شود و آزمون فیشر احتمالات همه جداول را به صورت افراطی یا شدیدتر از جدول مشاهده شده جمع می کند.

چرا از آزمون دقیق فیشر در این مطالعه استفاده می شود؟

۱. حجم نمونه کوچک: حجم کل نمونه (۶۰۰ شرکت کننده) متوسط است، اما برای برخی از پیامدها (به عنوان مثال، ۸۰ حمله قلبی در گروه آسپرین)، فرکانس سلولی نسبتاً کوچک است.

۲. نتایج دقیق: بر خلاف روش‌های تقریبی مانند آزمون Chi-square، آزمون فیشر مقادیر p دقیقی را ارائه می‌کند که وقتی اندازه نمونه یا فرکانس‌های مورد انتظار پایین باشد قابل اعتمادتر هستند.

۳. نتیجه باینری: داده‌ها شامل متغیرهای باینری هستند (حمله قلبی: بله/خیر) و ب = ه خوبی در جدول ۲×۲ قرار می‌گیرند و آزمایش فیشر را مناسب می‌کند.

مزایای تست دقیق فیشر

- دقیق برای نمونه‌های کوچک: نتایج دقیق را بدون توجه به حجم نمونه ارائه می‌دهد.
- غیر پارامتریک: هیچ فرضی در مورد توزیع داده‌ها ندارد.
- مقاوم: ایده آل برای مواردی که فرکانس‌های مورد انتظار پایین است.

محدودیت‌ها

- از نظر محاسباتی برای جداول یا مجموعه داده‌های بزرگ فشرده است (اما قابل مدیریت برای جداول ۲×۲ است).
- در این مطالعه، آزمایش فیشر بسیار مفید است زیرا به طور دقیق رابطه بین مصرف آسپرین و ایبوپروفن و خطر حمله قلبی را ارزیابی می‌کند و حتی با توزیع داده‌های مشاهده‌شده نتایج قابل اعتمادی را تضمین می‌کند.

محاسبات به شرح زیر است:

۱. فرضیه صفر (H_0): هیچ ارتباطی بین نوع دارو (آسپرین یا ایبوپروفن) و خطر حمله قلبی وجود ندارد.

۲. درک معنای "Extreme or More Extreme"

در یک جدول ۲×۲، extreme به جدول‌هایی اشاره دارد که در آنها:

- انحراف از توزیع مورد انتظار تحت فرضیه صفر به اندازه یا بیشتر از آنچه مشاهده شد است.
- به عنوان مثال:

- داده‌های مشاهده‌شده: آسپرین در مقایسه با ایبوپروفن حملات قلبی را کاهش می‌دهد.
- داده‌های شدیدتر بیشتر به نفع آسپرین است یا تفاوت آشکارتری بین گروه‌ها خواهد داشت.

۳. فرمول آزمون دقیق فیشر

احتمال هر پیکربندی جدول ۲×۲ خاص با استفاده از توزیع ابر هندسی محاسبه می‌شود:

$$P = \frac{\binom{a+b}{a} \binom{c+d}{c}}{\binom{n}{a+c}}$$

- a، b، c، d: مقادیر در جدول اند.
 - n: تعداد کل مشاهدات است ($n = a+b+c+d$).
- در صفحه بعد حل گام به گام آمده است.

۴. حل گام به گام

داده های مشاهده شده عبارتند از:

	Aspirin	Ibuprofen	Total
No Heart Attack	a=220	c=180	a+c=400
Heart Attack	b=80	d=120	b+d =200
Total	a+b=300	c+d=300	n=a+b+c+d=600

احتمال مشاهده این جدول خاص:

$$P = \frac{\binom{a+b}{a} \binom{c+d}{c}}{\binom{n}{a+c}}$$

$$P = \frac{\binom{220+80}{220} \binom{180+120}{180}}{\binom{600}{400}}$$

محاسبات فاکتوریل ها:

$$\binom{a+b}{a} = \frac{(300)!}{(220)!(80)!}, \quad \binom{c+d}{c} = \frac{(300)!}{(180)!(120)!}, \quad \binom{n}{a+c} = \frac{(600)!}{(400)!(200)!}$$

جداول Extreme تر:

- آزمون دقیق فیشر احتمالات (P) را برای همه جداول جمع می کند که اختلاف نسبت ها به اندازه جدول مشاهده شده بزرگتر و بزرگتر است.
- جداول افراطی (Extreme tables) مواردی را شامل میشود که $a > 220$ باشد.

فرآیند محاسبه:

- تمام جداول extreme ممکن را با مجموع ردیف ها و ستون های یکسان (marginal sums) برمی شماریم.
- احتمال هر جدول را با استفاده از فرمول های پرهندسی محاسبه میکنیم.
- احتمالات همه جداول را به صورت extreme یا شدیدتر از جدول مشاهده شده جمع میکنیم.

۴. محاسبه P-Value

مقدار P-Value در این مطالعه ($P = 0.00071$) با جمع کردن احتمالات جدول مشاهده شده و همه جداول افراطی تر (more extreme tables) به دست می آید. این محاسبات به دلیل پیچیدگی شمارش تمام جداول ممکن، معمولاً به صورت برنامه نویسی شده ای انجام می شوند. برای محاسبه میتوانیم از کد پایتون یا نرم افزار های آماری یا وب سایت های مختلف استفاده کنیم. همان طور که در زیر دیده میشود:

Easy Fisher Exact Test Calculator

Success! The Fisher exact test statistic and statement of significance appear beneath the table. Blue means you're dealing with dependent variables; red, independent.

Results			
	No Heart Attack	Heart Attack	Marginal Row Totals
Aspirin	220	80	300
Ibuprofen	180	120	300
Marginal Column Totals	400	200	600 (Grand Total)

The Fisher exact test statistic value is 0.0007. The result is significant at $p < .05$.

Start Again

```

1 from scipy.stats import fisher_exact
2 import numpy as np
3
4 contingency_table = np.array([[80, 220], [120, 180]])
5
6 odds_ratio, p_value_fisher = fisher_exact(contingency_table)
7
8 odds_ratio, p_value_fisher

```

(0.5454545454545454, 0.0007120141417653778)

۵. چرا از تست دقیق فیشر استفاده می شود؟

- احتمالات دقیق: بر خلاف تست های تقریبی (مثلاً Chi-square)، تست فیشر مقادیر p دقیق را محاسبه می کند و آن را برای نمونه های کوچک یا تعداد نمونه های کم قابل اعتماد می کند.
- تناسب سناریو: این مطالعه شامل یک جدول 2×2 با داده های طبقه بندی شده و تعداد نسبتاً کوچک در برخی سلول ها است (به عنوان مثال، $d = 120$ ، $b = 80$)، که آزمایش دقیق فیشر را مناسب می کند.

خلاصه

مقدار p به صورت زیر محاسبه می شود:

۱. محاسبه احتمال جدول مشاهده شده.
 ۲. جمع کردن احتمالات برای همه جداول به صورت افراطی یا افراطی تر (extreme or more extreme) تحت H_0 .
 ۳. استفاده از احتمال مجموع برای ارزیابی اینکه آیا فرضیه صفر باید رد شود یا خیر.
- در این مطالعه، مقدار p-value بسیار کوچک (0.00071) و کوچکتر از سطح معنی دار $\alpha = 0.05$ قویاً نشان می دهد که آسپرین در مقایسه با ایبوپروفن خطر حمله قلبی را کاهش می دهد.

(د) تجزیه و تحلیل نتایج و فرضیه

- سطح معنی داری مانند $\alpha = 0.05$ در نظر می گیریم.
- مقدار p-value (0.00071) بسیار کوچکتر از α است که منجر به رد فرضیه صفر (H_0) می شود. این نشان می دهد که آسپرین به طور قابل توجهی خطر حمله قلبی را در مقایسه با ایبوپروفن کاهش می دهد.
- در واقع فرضیه جایگزین (Alternative Hypothesis: H_1) تایید میشود. یعنی این اتفاق شانس نبوده و افرادی که ایبوپروفن مصرف می کنند، بیشتر از افرادی که آسپرین مصرف می کنند، دچار حمله قلبی می شوند. به بیان دیگر، داروی آسپرین تاثیرگذار بوده و افرادی که آسپرین مصرف می کنند، کمتر از افرادی که ایبوپروفن مصرف می کنند، دچار حمله قلبی می شوند.
- همچنین در بخش های قبل هم دیدیم که با توجه به $OR = 0.5455$ مصرف کنندگان آسپرین در مقایسه با مصرف کنندگان ایبوپروفن شانس کمتری برای حمله قلبی دارند.

نتیجه گیری:

آسپرین در مقایسه با ایبوپروفن از نظر آماری تاثیر معناداری در کاهش خطر حمله قلبی دارد.

سوال پنجم الف) تکنیک RNA-Seq و مزایای آن

۱. RNA-Seq چیست؟

- RNA-Seq (توالی یابی RNA یا RNA sequencing) یک تکنیک توالی یابی با توان عملیاتی بالا است که از توالی یابی نسل بعدی (NGS: next-generation sequencing) برای تجزیه و تحلیل transcriptome یا رونوشت (همه مولکول های RNA، از جمله mRNA، ncRNA و غیره) از نمونه استفاده می کند.

۲. چگونه به مطالعه بیان ژن کمک می کند:

- بیان ژن را با اندازه گیری فراوانی رونوشت های RNA تعیین می کند.
- رونوشت های جدید، ایزوفرم ها و انواع مختلف اسپلایس را شناسایی می کند.
- تغییرات پویا در بیان ژن در شرایط یا نقاط زمانی را ثبت می کند.

۳. مزایای نسبت به Microarray:

- تشخیص بی طرفانه (Unbiased Detection): نیازی به دانش قبلی از توالی ندارد.
- Higher Sensitivity and Specificity: رونوشت ها و ایزوفرم های کم فراوانی را تشخیص می دهد.
- Wide Dynamic Range: ژن های با بیان بالا و کم را با دقت بیشتری تعیین می کند.
- کشف رونوشت های جدید (Discovery of Novel Transcripts): ژن های unannotated و رویدادهای splicing را شناسایی می کند.

ب) داده های خام در مقابل داده های پردازش شده در RNA-Seq

۱. داده های خام (Raw Data):

- حاوی توالی خام read ها در قالب FASTQ.
- شامل base call quality scores است.
- مفید برای تجزیه و تحلیل مجدد یا خطوط لوله پردازش جایگزین (alternative processing pipelines).

۲. داده های پردازش شده (Processed Data):

- از داده های خام پس از کنترل کیفیت، هم ترازی و کمی سازی (quality control, alignment, and quantification) به دست می آید.
- شامل تعداد عبارات، داده های نرمال شده، یا رونوشت های حاشیه نویسی شده (expression counts, normalized data, or annotated transcripts) است.
- آماده برای تفسیر آماری یا بیولوژیکی پایین دستی (downstream statistical or biological interpretation).

۳. اهمیت:

- داده های خام: تکرارپذیری (reproducibility) و انعطاف پذیری (flexibility) در تجزیه و تحلیل را تضمین می کند.
- داده های پردازش شده: بینش و مقایسه سریع را تسهیل می کند.

ج) پایگاه داده SRA

۱. تعریف:

- The Sequence Read Archive (SRA) یک مخزن عمومی است که توسط NCBI مدیریت می شود و داده های توالی یابی خام از مطالعات مختلف از جمله RNA-Seq را ذخیره می کند.

۲. داده های ذخیره شده:

- توالی خام read ها (Raw sequencing reads) در قالب FASTQ.
- Metadata در مورد نمونه ها، روش های توالی یابی و شرایط.

۳. دسترسی به داده های RNA-Seq:

- پایگاه داده SRA را از طریق وب سایت NCBI با کلمات کلیدی، ارگانیسم، نوع مطالعه یا شماره های الحاق (keywords, organism, study type, or accession numbers) قابل جستجو است.
- داده ها را با استفاده از NCBI SRA Toolkit میتوان دانلود کرد.

(د یافتن یک مطالعه مرتبط با سرطان در پایگاه داده SRA

۱. مراحل جستجو:

- از پایگاه داده NCBI SRA (<https://www.ncbi.nlm.nih.gov/sra>) دیدن میکنیم.
- جستجوی عباراتی مانند «RNA-Seq breast cancer.» انجام میدهیم.
- نتایج را بر اساس نوع مطالعه، ارگانیسم یا ارتباط (study type, organism, or relevance) فیلتر میکنیم.
- یک مطالعه مناسب را به دلخواه شناسایی میکنیم و شماره الحاق (accession number) آن را یادداشت میکنیم.

۲. مطالعه انتخابی:

- لینک: [https://www.ncbi.nlm.nih.gov/sra/SRX27146400\[accn\]](https://www.ncbi.nlm.nih.gov/sra/SRX27146400[accn])

- Accession Number :SRX27146400

- SRA :SRP552834

- Run Accession :SRR31785012

- توضیحات مطالعه:

The Study of Ru Doped Carbon Quantum Dots (Ru CDs) Therapy on Mouse breast cancer Cell Line 4T1
Abstract: Research on the anti-tumor activities of Ru doped carbon quantum dots (Ru CDs) nanomaterials

NIH National Library of Medicine
National Center for Biotechnology Information

SRA SRA Advanced Search Help

Full Send to: Related information BioProject BioSample Taxonomy Recent activity Turn Off Clear

SRX27146400: RNA-Seq of mus musculus cancer cell line 4T1with nanoparticles treatments
1 ILLUMINA (Illumina HiSeq 4000) run: 21.6M spots, 6.5G bases, 1.9Gb downloads

Design: RNA-seq of Ru-CDs-GOx@HA+H2O2 group 1
Submitted by: Yunnan University
Study: The Study of Ru Doped Carbon Quantum Dots (Ru CDs) Therapy on Mouse breast cancer Cell Line 4T1
PRJNA1200808 • SRP552834 • All experiments • All runs
hide Abstract
Research on the anti-tumor activities of Ru doped carbon quantum dots (Ru CDs) nanomaterials

Sample: SAMN45915661 • SRS23601538 • All experiments • All runs
Organism: Mus musculus

Library:
Name: D1
Instrument: Illumina HiSeq 4000
Strategy: RNA-Seq
Source: TRANSCRIPTOMIC
Selection: cDNA
Layout: PAIRED

Runs: 1 run, 21.6M spots, 6.5G bases, 1.9Gb

Run	# of Spots	# of Bases	Size	Published
SRR31785012	21,551,544	6.5G	1.9Gb	2024-12-19

ID: 36630219

ه) دانلود داده های خام با استفاده از NCBI SRA Toolkit

۱. دستورات:

- NCBI SRA Toolkit را نصب کنیم:

```
sudo apt-get install sra-toolkit
```

- دانلود داده های خاص:

```
prefetch Run_Accession_ID # Replace with specific Run Accession ID
```

- تبدیل SRA به FASTQ:

```
fastq-dump --split-files Run_Accession_ID
```

۲. مثال برای Accession انتخابی:

```
prefetch SRR31785012
```

```
fastq-dump --split-files SRR31785012
```

```
1 |prefetch SRR31785012 # Replace with specific Run Accession ID

2024-12-26T16:36:56 prefetch.2.11.3: Current preference is set to retrieve SRA Normalized Format files with full base quality scores.
2024-12-26T16:36:56 prefetch.2.11.3: 1) Downloading 'SRR31785012'...
2024-12-26T16:36:56 prefetch.2.11.3: SRA Normalized Format file is being retrieved, if this is different from your preference, it may be due to current fi
2024-12-26T16:36:56 prefetch.2.11.3: Downloading via HTTPS...
2024-12-26T16:37:51 prefetch.2.11.3: HTTPS download succeed
2024-12-26T16:38:04 prefetch.2.11.3: 'SRR31785012' is valid
2024-12-26T16:38:04 prefetch.2.11.3: 1) 'SRR31785012' was downloaded successfully
2024-12-26T16:38:04 prefetch.2.11.3: 'SRR31785012' has 0 unresolved dependencies

[7] 1 |fastq-dump --split-files SRR31785012

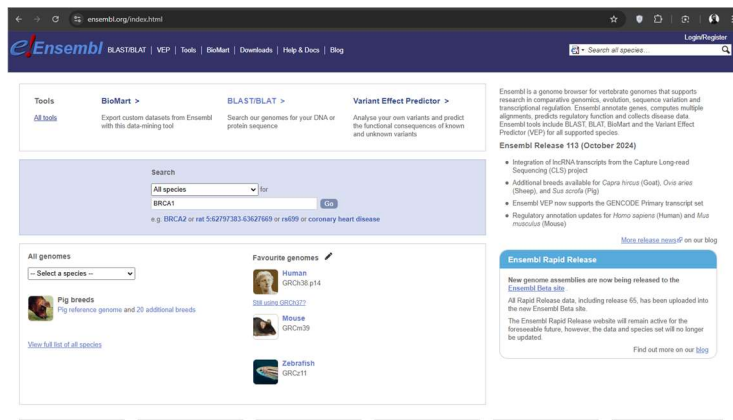
Read 21551544 spots for SRR31785012
Written 21551544 spots for SRR31785012
```

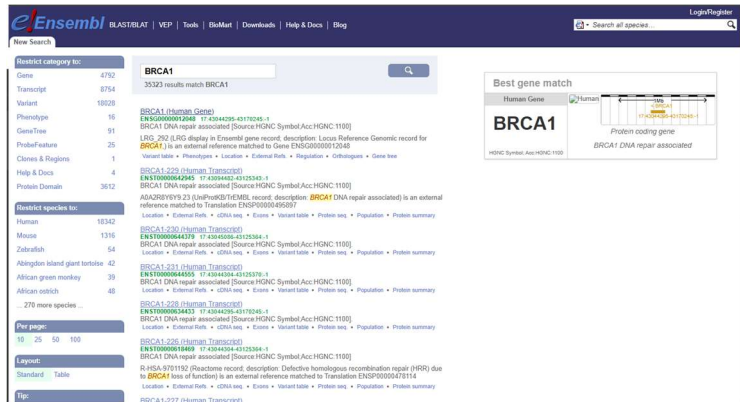
و) استفاده از پایگاه داده Ensembl

۱. یافتن یک ژن:

- از (<https://www.ensembl.org>) دیدن میکنیم.

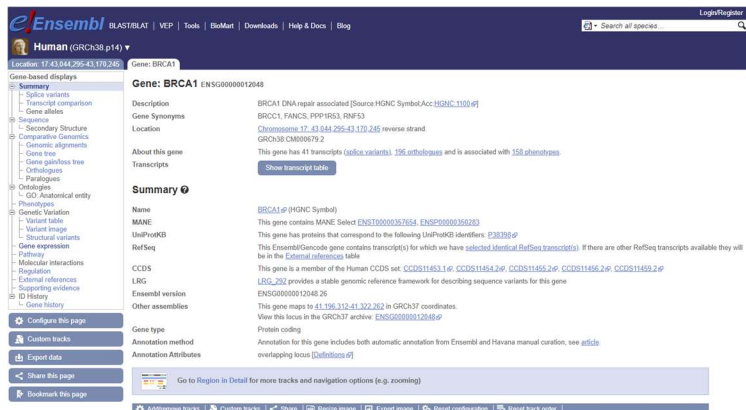
- "BRCA1" را با استفاده از نوار جستجو، جستجو میکنیم.



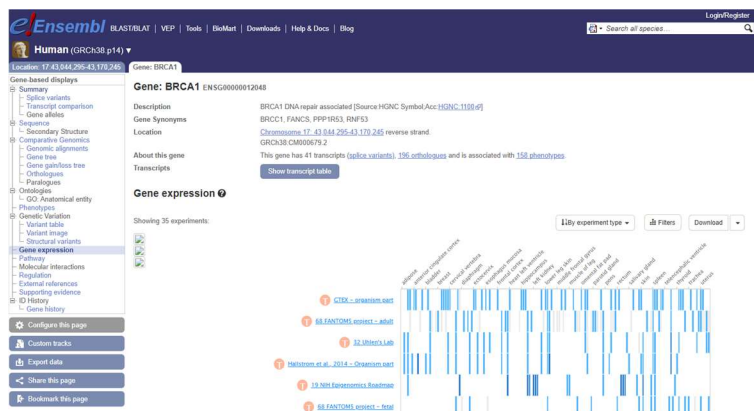


۲. داده های بیان ژن:

- ابتدا صفحه مربوط به ژن را باز میکنیم.

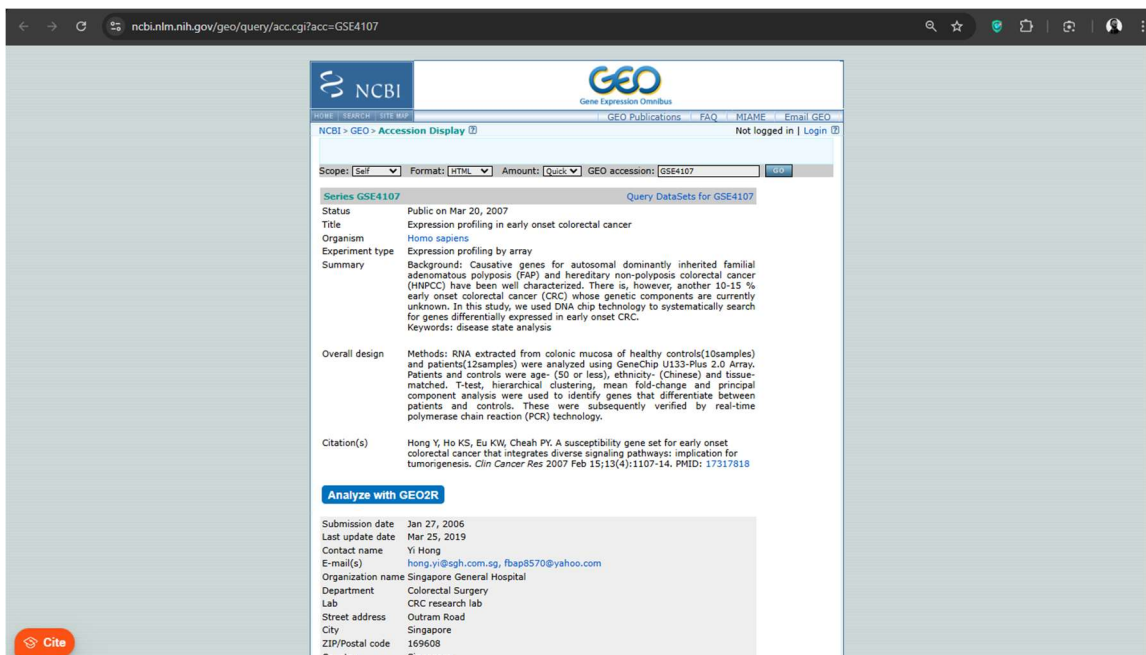


- حال به بخش "Gene expression" صفحه ژن می رویم.



- اطلاعاتی مختلفی در این بخش ارائه می شود مانند:

- Transcript table
- سطوح بیان (Expression Levels): در بافت ها و شرایط مختلف.
- Disease Associations: پیوند بین بیان ژن و بیماری های خاص.
- داده های مقایسه ای (Comparative Data): پروفایل های بیان بین گونه ها.



خلاصه مطالعه به شرح زیر است:

Background: Causative genes for autosomal dominantly inherited familial adenomatous polyposis (FAP) and hereditary non-polyposis colorectal cancer (HNPCC) have been well characterized. There is, however, another 10-15 % early onset colorectal cancer (CRC) whose genetic components are currently unknown. In this study, we used DNA chip technology to systematically search for genes differentially expressed in early onset CRC.

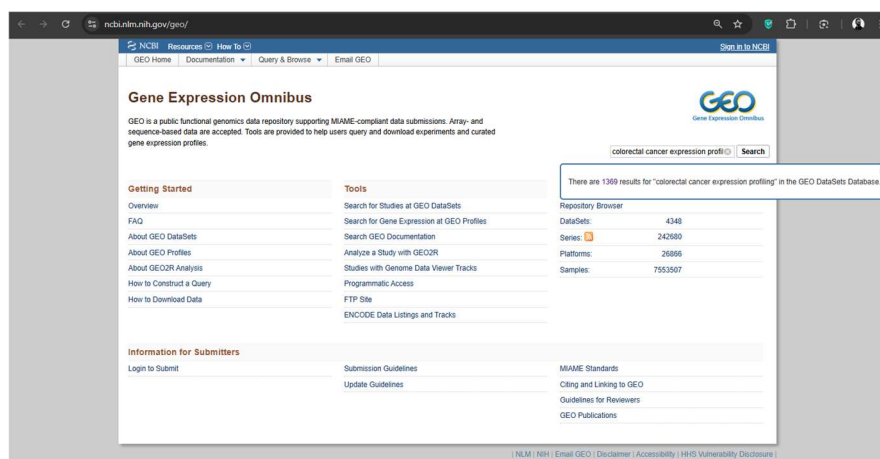
Keywords: disease state analysis

(ب) توضیح مراحل جستجو و انتخاب داده ها.

۱. دسترسی به پایگاه داده GEO: به صفحه اصلی GEO (<https://www.ncbi.nlm.nih.gov/geo/>) میرویم.

۲. جستجو برای مجموعه داده ها:

- از نوار جستجو برای وارد کردن کلمات کلیدی مرتبط با علاقه تحقیقاتی خود مانند "colorectal cancer expression profiling" استفاده میکنیم.
- فهرست مجموعه داده هایی که با معیارهای جستجوی شما مطابقت دارند را مرور میکنیم و با فیلترهای گوناگون به مجموعه داده دلخواه می رسیم.



۳. یک مجموعه داده را انتخاب میکنیم به عنوان مثال GSE4107:

- روی مجموعه داده مورد علاقه کلیک میکنیم تا اطلاعات دقیق آن شامل Summary، Overall design و Sample های موجود را مشاهده کنیم.
- مطمئن می شویم که مجموعه داده شامل نمونه‌های کنترل و بیمار (control and patient samples) مرتبط با مطالعه ما باشد.

۴. به Accession Number توجه داریم: هر مجموعه داده دارای یک شماره دسترسی (Accession Number) منحصر به فرد (به عنوان مثال، GSE4107) است که برای ارجاع و تجزیه و تحلیل استفاده می شود.

(ج) ژن هایی را با تفاوت بیان معنی دار بین گروه کنترل و بیمار با استفاده از GEO2R شناسایی میکنیم.

۱. دسترسی به GEO2R: به GEO2R tool موجود در این لینک (<https://www.ncbi.nlm.nih.gov/geo/geo2r/>) میریم.

۲. شماره دسترسی (Accession Number) را وارد میکنیم: "GSE4107" را در کادر GEO accession وارد کرده و روی "Set" و "Analyze" کلیک میکنیم. یا اینکه در صفحه مطالعه مربوطه روی "Analyze with GEO2R" کلیک میکنیم.

Series GSE4107

Query DataSets for GSE4107

Status	Public on Mar 20, 2007
Title	Expression profiling in early onset colorectal cancer
Organism	Homo sapiens
Experiment type	Expression profiling by array
Summary	Background: Causative genes for autosomal dominantly inherited familial adenomatous polyposis (FAP) and hereditary non-polyposis colorectal cancer (HNPCC) have been well characterized. There is, however, another 10-15 % early onset colorectal cancer (CRC) whose genetic components are currently unknown. In this study, we used DNA chip technology to systematically search for genes differentially expressed in early onset CRC. Keywords: disease state analysis
Overall design	Methods: RNA extracted from colonic mucosa of healthy controls(10samples) and patients(12samples) were analyzed using GeneChip U133-Plus 2.0 Array. Patients and controls were age- (50 or less), ethnicity- (Chinese) and tissue-matched. T-test, hierarchical clustering, mean fold-change and principal component analysis were used to identify genes that differentiate between patients and controls. These were subsequently verified by real-time polymerase chain reaction (PCR) technology.
Citation(s)	Hong Y, Ho KS, Eu KW, Cheah PY. A susceptibility gene set for early onset colorectal cancer that integrates diverse signaling pathways: implication for tumorigenesis. <i>Clin Cancer Res</i> 2007 Feb 15;13(4):1107-14. PMID: 17317818

Analyze with GEO2R

۳. گروه های نمونه را تعریف میکنیم:

- روی «Define group» کلیک میکنیم و دو گروه ایجاد میکنیم، به عنوان مثال، «Control» و «Patient».

- بر اساس اطلاعات نمونه مجموعه داده، نمونه های مناسب را به هر گروه اختصاص می دهیم.

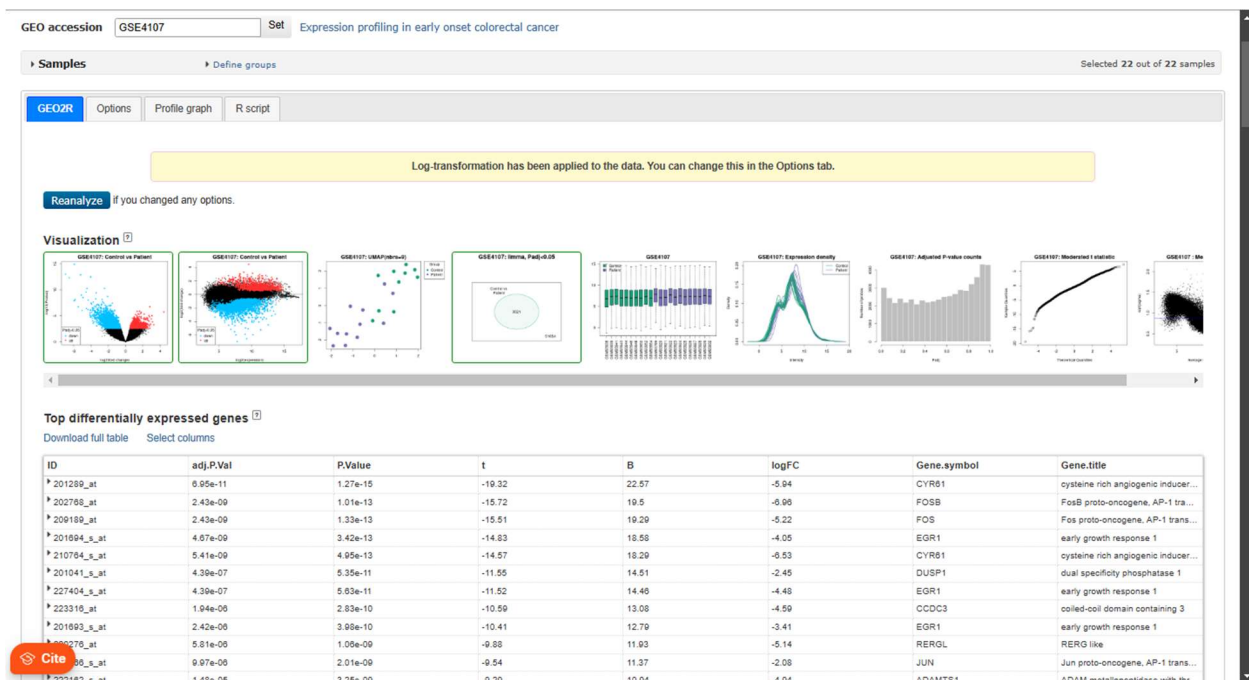
NCBI GEO2R interface for GSE4107. The 'Define groups' section shows a list of samples with checkboxes for 'Control (10 samples)' and 'Patient (12 samples)'. The 'Patient' group is selected. The table below shows the details of the selected samples, including GSM accession, sample name, description, sex, age, and tissue type.

Sample Type	GSM Accession	Sample Name	Description	Sex	Age	Tissue
Patient	GSM93923	3894M patient's mucosa	mucosa collected and archived within 30 mins after surgery	F	33	Mucosa
Patient	GSM93924	3894M patient's mucosa	mucosa collected and archived within 30 mins after surgery	M	33	Mucosa
Patient	GSM93925	3894M patient's mucosa	mucosa collected and archived within 30 mins after surgery	F	48	Mucosa
Patient	GSM93926	3894M patient's mucosa	mucosa collected and archived within 30 mins after surgery	M	36	Mucosa
Patient	GSM93927	3950M patient's mucosa	mucosa collected and archived within 30 mins after surgery	M	44	Mucosa
Patient	GSM93928	4145M patient's mucosa	mucosa collected and archived within 30 mins after surgery	F	47	Mucosa
Patient	GSM93929	4180M patient's mucosa	mucosa collected and archived within 30 mins after surgery	M	41	Mucosa
Patient	GSM93932	3446M patient's mucosa	mucosa collected and archived within 30 mins after surgery	F	44	Mucosa
Control	GSM93938	NM14 healthy control	biopsies from healthy individuals undergoing colonoscopic examination, snap-frozen immediately, stored at -80°C	F	45	Mucosa
Control	GSM93939	NM16 healthy control	biopsies from healthy individuals undergoing colonoscopic examination, snap-frozen immediately, stored at -80°C	F	44	Mucosa
Control	GSM93941	NM18 healthy control	biopsies from healthy individuals undergoing colonoscopic examination, snap-frozen immediately, stored at -80°C	F	41	Mucosa
Control	GSM93943	NM18 healthy control	biopsies from healthy individuals undergoing colonoscopic examination, snap-frozen immediately, stored at -80°C	M	28	Mucosa
Control	GSM93944	NM20 healthy control	biopsies from healthy individuals undergoing colonoscopic examination, snap-frozen immediately, stored at -80°C	M	27	Mucosa
Control	GSM93946	NM22 healthy control	biopsies from healthy individuals undergoing colonoscopic examination, snap-frozen immediately, stored at -80°C	M	41	Mucosa
Control	GSM93948	NM24 healthy control	biopsies from healthy individuals undergoing colonoscopic examination, snap-frozen immediately, stored at -80°C	M	44	Mucosa
Control	GSM93950	NM3 healthy control	biopsies from healthy individuals undergoing colonoscopic examination, snap-frozen immediately, stored at -80°C	F	48	Mucosa

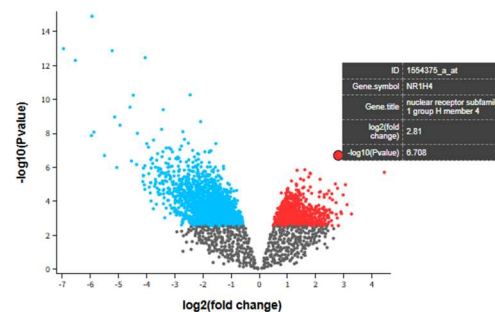
۴. تجزیه و تحلیل را انجام می‌دهیم:

- پس از تخصیص نمونه‌ها به گروه‌ها، برای شناسایی ژن‌های بیان شده متفاوت بین دو گروه، روی "Analyze" کلیک می‌کنیم.
- ۵. بررسی نتایج:

- GEO2R جدولی از ژن‌ها را با آماره‌هایی مانند adjusted p-values و log fold changes نمایش می‌دهد که نشان‌دهنده اهمیت و بزرگی تفاوت‌های بیان است. همچنین نمودارهای متنوعی برای visualization ارائه می‌دهد. در زیر نتیجه آن آمده است.
- ژن‌های با تفاوت بیان معنی‌دار بین گروه کنترل و بیمار: به عنوان مثال در volcano plot ژن‌های قرمز و آبی ژن‌های با تفاوت بیان معنی‌دار بین گروه کنترل و بیمار هستند که هر چه بالاتر باشند و از نقطه ۰ فاصله بیشتری داشته باشند تفاوت بیان معنی‌دارتری دارند. همچنین در جدول داده شده این ژن‌ها به ترتیب اهمیت مرتب شده‌اند که همان‌طور که می‌بینید ژن‌های اولیه در جدول دارای adjusted p-values کوچک و log fold changes بزرگ هستند که به معنای تفاوت بیان معنی‌دار بین گروه کنترل و بیمار است.



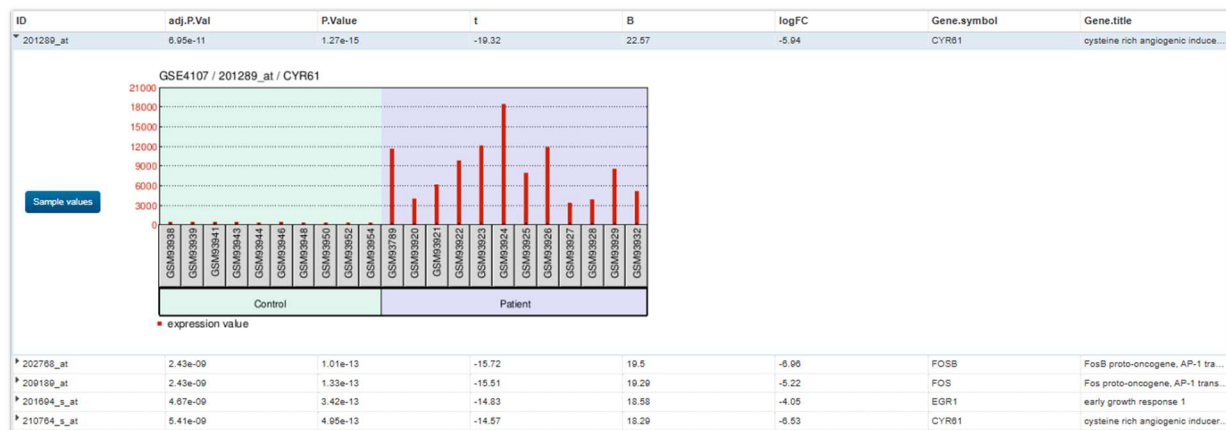
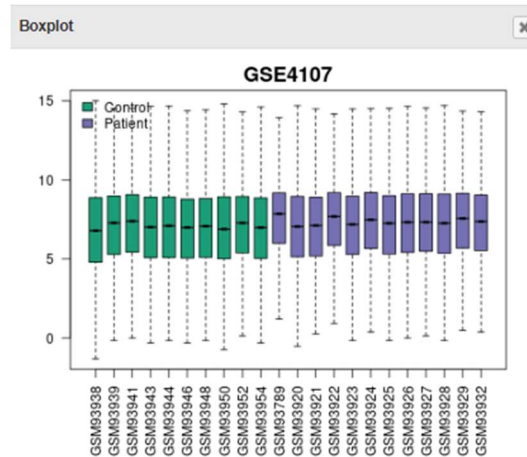
Volcano plot
GSE4107: Expression profiling in early onset
colorectal cancer
Control vs Patient, Padj<0.05



Select contrast to display

Control vs Patient

Download significant genes



د) نتایج را در جدولی شامل ژن های مهم و p-value ارائه می‌دهیم.

در زیر جدولی از ژن های انتخاب شده از مجموعه داده GSE4107 وجود دارد که تفاوت های بیانی قابل توجهی را بین نمونه های کنترل و بیمار نشان می دهد. همان طور که میبینید تمام ژن ها را مشاهده نمیکنید بلکه فقط ۱۰ ژن با کمترین adj.P.Val آمده اند (۱۰ ژن ابتدایی در جدول خروجی ابزار GEO2R):

ID	adj.P.Val	P.Value	t	logFC	Gene.symbol
201289_at	6.95E-11	1.27E-15	-19.316086	-5.94112291	CYR61
202768_at	2.43E-09	1.01E-13	-15.7194503	-6.95512411	FOSB
209189_at	2.43E-09	1.33E-13	-15.514348	-5.22397202	FOS
201694_s_at	4.67E-09	3.42E-13	-14.8278778	-4.0543598	EGR1
210764_s_at	5.41E-09	4.95E-13	-14.5653311	-6.53075681	CYR61
201041_s_at	4.39E-07	5.35E-11	-11.5465574	-2.45413476	DUSP1
227404_s_at	4.39E-07	5.63E-11	-11.516511	-4.47587421	EGR1
223316_at	1.94E-06	2.83E-10	-10.5931505	-4.59230998	CCDC3
201693_s_at	2.42E-06	3.98E-10	-10.4060636	-3.41023566	EGR1
220276_at	5.81E-06	1.06E-09	-9.8761837	-5.13632928	RERGL

همچنین جدول کامل در فایل GSE4107.top.table.tsv در پیوست موجود است. در این جدول که از پراهمیت به کم اهمیت مرتب شده است در سطرهای اولیه میتوانید مهم ترین ژن ها را مشاهده کنید. در این جدول تمام ژن های با adjusted p-value کوچک تر از ۰/۰۵ از کوچک به بزرگ آمده اند.

توجه: adjusted p-value اصلاحات آزمایشی متعدد (multiple testing corrections) را به حساب می آورد، و log fold change نشان دهنده جهت و میزان تغییر بیان است (مقادیر منفی نشان دهنده کاهش در بیماران در مقایسه با گروه کنترل است).

برای نمایش بصری استفاده از GEO2R، آموزش زیر مفید بود که از آن استفاده کردم:

<https://www.youtube.com/watch?v=9RyWjzSnaEO>

پایان