

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



ژنومیک محاسباتی با R

Population Genomics with R

آلتونا آکالین، وردا فرانک، بورا اویار و جانان رونن

مسعود علی پناه

(دانشیار، گروه تولیدات گیاهی، دانشگاه تربت حیدریه)

ایمان یوسفی جوان

(استادیار، گروه تولیدات گیاهی، دانشگاه تربت حیدریه)

(۱۴۰۴)

عنوان و نام پدیدآور	: ژنومیک محاسباتی با R/نویسندگان آلتونا آکالین... [و دیگران]؛ مترجمان مسعود علی پناه، ایمان یوسفی جوان.
مشخصات نشر	: تربت حیدریه: دانشگاه تربت حیدریه، انتشارات، سال ۱۴۰۳
مشخصات ظاهری	: ق، ۷۱۲ص:، مصور، نمودار
شابک	: ۹۷۸-۶۰۰-۸۳۳۵-۴۰-۵
وضعیت فهرست نویسی	: فیبا
یادداشت	: عنوان اصلی: Computational Genomics with R, 2020.
یادداشت	: نویسندگان آلتونا آکالین، وردا فرانک، بورا اوپار، جاناتان رونن.
یادداشت	: کتابنامه: ص: ۶۹۴-۷۱۱.
موضوع	: ژنتیک پزشکی Medical genetics آر (زبان برنامه نویسی کامپیوتر) R (Computer program language) ژنتیک -- روش های آماری Genetics -- Statistical methods
شناسه افزوده	: آکالین، آلتونا
شناسه افزوده	: Akalin, Altuna
شناسه افزوده	: علی پناه، مسعود، ۱۳۴۸ - مترجم
شناسه افزوده	: یوسفی جوان، ایمان، مترجم
شناسه افزوده	: دانشگاه تربت حیدریه
رده بندی کنگره	: RB۱۵۵
رده بندی دیویی	: ۰۴۲/۶۱۶
شماره کتابشناسی ملی	: ۹۸۲۰۳۱۳

این اثر مشمول قانون حمایت از مؤلفان و مصنفان و هنرمندان است. هر کس تمام یا قسمتی از این اثر را بدون اجازه ناشر، نشر، پخش یا عرضه کند مورد پیگرد قانونی قرار خواهد گرفت.



ژنومیک محاسباتی با R

نویسنده: آلتونا آکالین، وردا فرانک، بورا اوپار، جاناتان رونن.

مترجم: مسعود علی پناه، ایمان یوسفی جوان.

چاپ اول

بها: تومان

نشانی ناشر: تربت حیدریه، کیلومتر هفت جاده مشهد، دانشگاه دولتی تربت حیدریه

مسئولیت کلیه مطالب این کتاب به عهده نگارنده می باشد. دانشگاه تربت حیدریه هیچگونه مسئولیتی در قبال صحت و سقم مطالب ندارد.

پیشگفتار الف

فصل ۱

- ۱- مقدمه‌ای بر ژنومیک ۱
- ۱-۱ ژن‌ها، DNA و اصل مرکزی ۱
- ۲-۱ عناصر تنظیم ژن ۸
- ۳-۱ شکل دادن به ژنوم: جهش DNA ۱۹
- ۴-۱ روش‌های آزمایشی با کارایی بالا در ژنومیک ۲۱
- ۵-۱ تصویرسازی و مخازن داده برای ژنومیک ۲۷

فصل ۲

- ۲- مقدمه‌ای بر R برای تجزیه و تحلیل داده‌های ژنومی ۳۱
- ۱-۲-۱ مراحل تجزیه و تحلیل (ژنومیک) داده‌ها ۳۱
- ۲-۲-۲ شروع به کار R ۳۷
- ۳-۲ محاسبات در R ۴۰
- ۴-۲ ساختار داده ۴۰
- ۵-۲ انواع داده‌ها ۴۷
- ۶-۲ خواندن و نوشتن داده‌ها ۴۸
- ۷-۲ رسم نمودار با گرافیک پایه در R ۵۰

- ۲-۸ رسم در R با ggplot2 ۵۷
- ۲-۹ توابع و ساختارهای کنترل (برای، اگر/دیگر و غیره) ۶۴
- ۲-۱۰ تمرین ۷۲

فصل ۳

- ۳-آمار برای ژنومیک ۸۷
- ۳-۱ نحوه خلاصه کردن مجموعه نقاط داده: ایده پشت توزیع‌های آماری ۸۷
- ۳-۲ نحوه آزمایش تفاوت بین نمونه‌ها ۱۰۳
- ۳-۳ رابطه بین متغیرها: مدل‌های خطی و همبستگی ۱۱۶
- ۳-۴ تمرینات ۱۴۰

فصل ۴

- ۴-تجزیه و تحلیل داده‌های اکتشافی با یادگیری ماشینی بدون ناظر ۱۴۷
- ۴-۱ خوشه‌بندی: گروه‌بندی نمونه‌ها بر اساس شباهت آن‌ها ۱۴۷
- ۴-۲ تکنیک‌های کاهش ابعاد: تجسم مجموعه داده‌های پیچیده ۱۷۲
- ۴-۳ تمرینات ۱۹۸

فصل ۵

- ۵-مدل‌سازی پیش‌بینی با یادگیری ماشینی نظارت شده ۲۰۱
- ۵-۱ مدل‌های یادگیری ماشینی چگونه تنظیم می‌شوند؟ ۲۰۲
- ۵-۲ مراحل در یادگیری ماشینی با ناظر ۲۰۵
- ۵-۳ مورد استفاده: نوع فرعی بیماری از داده‌های ژنومیک ۲۰۶
- ۵-۴ پیش‌پردازش داده‌ها ۲۰۸

- ۵-۵ تقسیم داده‌ها ۲۱۴
- ۵-۶ پیش‌بینی نوع فرعی با k - نزدیک‌ترین همسایه‌ها ۲۱۷
- ۵-۷ ارزیابی عملکرد مدل ما ۲۱۸
- ۵-۸ تنظیم مدل و اجتناب از پردازش بیش از حد ۲۲۵
- ۵-۹ اهمیت متغیر ۲۳۶
- ۵-۱۰ نحوه برخورد با عدم تعادل طبقاتی ۲۴۰
- ۵-۱۱ تعامل با پیش‌بینی‌کننده‌های وابسته ۲۴۱
- ۵-۱۲ درختان و جنگل‌ها: جنگل‌های تصادفی در عمل ۲۴۲
- ۵-۱۳ رگرسیون لجستیک و منظم‌سازی ۲۴۸
- ۵-۱۴ سایر الگوریتم‌های ناظر ۲۵۸
- ۵-۱۵ پیش‌بینی متغیرهای پیوسته: رگرسیون با یادگیری ماشین ۲۷۰
- ۵-۱۶ تمرینات ۲۷۶

فصل ۶

- ۶-۶ عملیات در فواصل ژنومی و محاسبات ژنوم ۲۷۹
- ۶-۱ عملیات در فواصل ژنومی با بسته GenomicRanges ۲۷۹
- ۶-۲ پرداختن به توالی‌های نقشه‌برداری شده با توان عملیاتی بالا ۲۹۲
- ۶-۳ پرداختن به امتیازات پیوسته روی ژنوم ۲۹۴
- ۶-۴ فواصل ژنومی با اطلاعات بیشتر: کلاس SummarizedExperiment ۳۰۰
- ۶-۵ تصویرسازی و خلاصه‌کردن فواصل ژنومی ۳۰۶
- ۶-۶ تمرینات ۳۱۹

فصل ۷

- ۷-۱ بررسی کیفیت، پردازش و هم‌ترازی خوانش توالی با توان بالا..... ۳۲۵
- ۷-۱ فرمت‌های FASTA و FASTQ..... ۳۲۶
- ۷-۲ بررسی کیفیت در خواندن توالی..... ۳۲۸
- ۷-۳ فیلتر کردن و اصلاح خوانش‌ها..... ۳۳۴
- ۷-۴ نقشه‌برداری/هم‌تراز کردن خوانش با ژنوم..... ۳۳۷
- ۷-۵ پردازش بیشتر خوانش‌های هم‌تراز شده..... ۳۴۰
- ۷-۶ تمرینات..... ۳۴۰

فصل ۸

- ۸-۱ تجزیه و تحلیل RNA-seq..... ۳۴۳
- ۸-۱ بیان ژن چیست؟..... ۳۴۳
- ۸-۲ روش‌های تشخیص بیان ژن..... ۳۴۵
- ۸-۳ تجزیه و تحلیل بیان ژن با استفاده از فناوری‌های توالی‌یابی با توان بالا..... ۳۴۵
- ۸-۴ سایر کاربردهای RNA-seq..... ۴۰۲
- ۸-۵ تمرین..... ۴۰۳

فصل ۹

- ۹-۱ تجزیه و تحلیل ChiP-seq..... ۴۰۹
- ۹-۱ برهم‌کنش‌های تنظیمی DNA - پروتئین..... ۴۰۹
- ۹-۲ اندازه‌گیری برهم‌کنش‌های DNA - پروتئین با ChiP-seq..... ۴۱۰
- ۹-۳ عواملی که بر کیفیت آزمایش و تجزیه و تحلیل ChiP-seq تأثیر می‌گذارد..... ۴۱۳

- ۴-۹ پیش پردازش داده‌های تراشه ۴۱۸
- ۵-۹ کنترل کیفیت تراشه ۴۲۱
- ۶-۹ فراخوانی پیک ۴۵۲
- ۷-۹ کشف موتیف ۴۹۱
- ۸-۹ بعد چه باید کرد؟ ۴۹۶
- ۹-۹ تمرینات ۴۹۸

فصل ۱۰

- ۱۰ تجزیه و تحلیل متیلاسیون DNA با استفاده از داده‌های توالی‌یابی بی سولفیت ۵۰۱
- ۱-۱۰ متیلاسیون DNA چیست؟ ۵۰۱
- ۲-۱۰ تجزیه و تحلیل داده‌های متیلاسیون DNA ۵۰۲
- ۳-۱۰ پردازش داده‌های خام و وارد کردن داده‌ها به R ۵۰۴
- ۴-۱۰ فیلتر کردن داده‌ها و تجزیه و تحلیل اکتشافی ۵۰۶
- ۵-۱۰ استخراج مناطق جالب: متیلاسیون و تقسیم‌بندی افتراقی ۵۱۸
- ۶-۱۰ حاشیه‌نویسی DMRs/DMCها و بخش‌ها ۵۳۲
- ۷-۱۰ سایر بسته‌های R که می‌توانند برای آنالیز متیلاسیون استفاده شوند ۵۳۴
- ۸-۱۰ تمرینات ۵۳۵

فصل ۱۱

- ۱۱ تجزیه و تحلیل امیکس چند گانه ۵۳۹

۱-۱۱	مورد استفاده: داده‌های امیکس چندگانه از سرطان کولورکتال	۵۴۰
۲-۱۱	مدل‌های متغیر پنهان برای ادغام چند امیکس	۵۴۸
۳-۱۱	روش‌های فاکتورسازی ماتریسی برای یکپارچه‌سازی داده‌های امیکس چندگانه بدون ناظر	۵۴۹
۴-۱۱	خوشه‌بندی با استفاده از عوامل پنهان	۵۶۹
۵-۱۱	تفسیر بیولوژیکی عوامل نهفته	۵۷۳
۶-۱۱	تمرینات	۵۸۰
۵۸۵	منابع	

پیشگفتار

هدف این کتاب ارائه مبانی تجزیه و تحلیل داده‌ها برای ژنومیک است. ما این کتاب را بر اساس دوره‌های ژنومی محاسباتی که هر سال ارائه می‌دهیم، تهیه کرده‌ایم. ما همواره مخاطبان بین‌رشته‌ای با پیشینه‌ای از فیزیک، زیست‌شناسی، پزشکی، ریاضی، علوم کامپیوتر یا سایر زمینه‌های کمی داشته‌ایم. ما می‌خواهیم این کتاب نقطه شروعی برای دانشجویان ژنومی محاسباتی و راهنمایی برای تجزیه و تحلیل بیشتر داده‌ها در موضوعات خاص‌تر در ژنومیک باشد. به همین دلیل است که ما سعی کردیم موضوعات مختلفی از برنامه‌نویسی گرفته تا زیست‌شناسی ژنوم پایه را پوشش دهیم. از آنجایی که این رشته بین‌رشته‌ای است، برای افراد با پیشینه‌های مختلف به نقاط شروع متفاوتی نیاز دارد. یک زیست‌شناس ممکن است بخش‌های مربوط به زیست‌شناسی ژنوم اولیه را نادیده بگیرد و با برنامه‌نویسی R شروع کند، در حالی که یک دانشمند کامپیوتر ممکن است بخواهد با زیست‌شناسی ژنوم شروع کند. به همین ترتیب، یک فرد باتجربه‌تر ممکن است بخواهد در صورت نیاز به انجام نوع خاصی از تحلیل، اما بدون تجربه قبلی، به این کتاب مراجعه کند.

نسخه آنلاین این کتاب تحت مجوز Creative Commons Attribution-NonCommercial-ShareAlike 4.0 مجوز بین‌المللی^۱ دارد.

این کتاب برای چه کسانی است؟

این کتاب شامل جنبه‌های عملی و نظری ژنومیک محاسباتی است. زیست‌شناسی و پزشکی بیش از هر زمان دیگری داده تولید می‌کنند؛ بنابراین، ما باید افراد بیشتری را با مهارت‌های تجزیه و تحلیل داده‌ها و درک ژنومیک محاسباتی آموزش دهیم. از آنجایی که ژنومیک محاسباتی بین‌رشته‌ای است، هدف این کتاب برای زیست‌شناسان، دانشمندان علوم پزشکی،

¹ <https://creativecommons.org/licenses/by-nc-sa/4.0/>

دانشمندان کامپیوتر و افرادی با سایر زمینه‌های کمی است. ما این کتاب را برای مخاطبان زیر نوشتیم:

- زیست‌شناسان و دانشمندان پزشکی که داده‌ها را تولید می‌کنند و خود مشتاق تجزیه و تحلیل آن هستند.
- دانش‌آموزان و محققانی که به طور رسمی شروع به تحقیق در مورد ژنومیک محاسباتی می‌کنند یا از آن استفاده می‌کنند، دانش گسترده‌ای در حوزه خاص ندارند، اما حداقل درک سطح مبتدی در زمینه‌های کمی به‌عنوان مثال، ریاضی، آمار، دارند.
- محققان باتجربه به دنبال دستورالعمل‌ها یا روش‌های سریع برای شروع انجام تجزیه و تحلیل داده‌های خاص مربوط به ژنومیک محاسباتی هستند.

خروجی چه خواهد بود؟

این منبع مهارت‌ها را توصیف می‌کند و روش‌هایی را ارائه می‌دهد که به خوانندگان کمک می‌کند تا داده‌های ژنومیک خود را تجزیه و تحلیل کنند.
بعد از خواندن:

- اگر با R آشنا نیستید، اصول اولیه R را دریافت کرده و مستقیماً به سراغ استفاده‌های تخصصی از R برای ژنومیک محاسباتی خواهید رفت.
- فواصل ژنومی و عملیات روی آن‌ها مانند همپوشانی را درک خواهید کرد.
- شما می‌توانید از R و کتابخانه گسترده بسته‌های نرم‌افزاری آن برای انجام تجزیه و تحلیل توالی، مانند محاسبه محتوای GC برای بخش‌های معینی از ژنوم یا یافتن مکان‌های اتصال فاکتور رونویسی استفاده کنید.
- شما با تکنیک‌های تصویرسازی مورد استفاده در ژنومیک، مانند نقشه‌های حرارتی، نمودارهای متاژن و تصویرسازی مسیر ژنومی آشنا خواهید شد.

- شما با تکنیک‌های یادگیری نظارت شده و بدون نظارت که در مدل‌سازی داده‌ها و تجزیه و تحلیل اکتشافی داده‌های با ابعاد بالا مهم هستند آشنا خواهید شد.
- شما با تجزیه و تحلیل مجموعه داده‌های مختلف توالی‌یابی با کارایی بالا (RNA-seq، ChIP-seq، BS-seq و ادغام چندگانه امیکس) که عمدتاً از ابزارهای مبتنی بر R استفاده می‌کنند، آشنا خواهید شد.

ساختار کتاب

این کتاب با این ایده طراحی شده است که درک عملی و مفهومی روش‌های تجزیه و تحلیل داده‌ها به همان اندازه، اگر نه مهم‌تر از درک نظری، مانند استخراج دقیق معادلات در آمار یا یادگیری ماشین مهم است. به همین دلیل است که ابتدا سعی می‌کنیم توضیحی مفهومی از مفاهیم، ارائه دهیم سپس سعی می‌کنیم بخش‌های اساسی فرمول‌های ریاضی را برای درک دقیق‌تر ارائه دهیم. در این مفهوم، ما همیشه کد یک کار تجزیه و تحلیل داده خاص را نشان داده و توضیح می‌دهیم. ما همچنین به خوانندگان که مایل‌اند درک نظری عمیق‌تری از روش‌ها یا مفاهیم مرتبط با تجزیه و تحلیل داده‌ها به دست آورند، منابع اضافی مانند کتاب‌ها، وبسایت‌ها، سخنرانی‌های ویدئویی و مقالات علمی ارائه می‌کنیم.

فصل ۱: "مقدمه‌ای بر ژنومیکس" به معرفی مفاهیم اساسی در زیست‌شناسی ژنوم و ژنومیک می‌پردازد. درک این مفاهیم برای ژنومیک محاسباتی مهم است.

فصل ۲: "مقدمه‌ای بر R برای تجزیه و تحلیل داده‌های ژنومی" علاوه بر پارادایم‌های رایج تجزیه و تحلیل داده‌ها که در تجزیه و تحلیل داده‌های ژنومی مشاهده می‌کنیم، مهارت‌های پایه‌ای لازم R برای دنبال کردن کتاب را فراهم می‌کند. فصل ۳: "آمار برای ژنومیکس"، فصل ۴: "تجزیه و تحلیل داده‌های اکتشافی با یادگیری ماشین بدون ناظر" و فصل ۵: "مدل‌سازی پیش‌بینی با یادگیری ماشینی با ناظر" مهارت‌های کمی لازم را که فرد هنگام تجزیه و تحلیل داده‌های ژنومیک با ابعاد بالا به آن نیاز دارد، معرفی می‌کند.

فصل ۶: "عملیات بر روی فواصل ژنومی و محاسبات ژنومی" ابزارهای اساسی در مواجهه با فواصل ژنومی و ارتباط آنها با یکدیگر بر روی ژنوم را معرفی می‌کند. علاوه بر این، این فصل انواع روش‌های تصویرسازی داده‌های ژنومی را معرفی می‌کند. مهارت‌های معرفی‌شده در این فصل، مهارت‌های کلیدی هستند که برای کار با داده‌های ژنومی پردازش‌شده که از طریق پایگاه‌های عمومی داده قابل‌دسترس هستند مانند Ensembl و مرورگر UCSC، موردنیاز هستند.

فصل‌های بعدی به تجزیه و تحلیل خاص داده‌های توالی‌یابی با توان بالا و ادغام انواع مختلف مجموعه داده‌ها می‌پردازد. فصل ۷: "بررسی کیفیت، پردازش و ترازسازی خوانش‌های توالی با توان بالا" بررسی‌های کیفیتی را که باید در مورد خواندن توالی انجام شود و روش‌های مختلف برای پردازش بیشتر آن‌ها معرفی می‌کند. در فصل‌های ۸، ۹ و ۱۰ به تجزیه و تحلیل RNA-seq، تجزیه و تحلیل ChIP-seq و تجزیه و تحلیل BS-seq پرداخته می‌شود. فصل آخر، فصل ۱۱: "تحلیل چندگانه امیکس" به روش‌هایی برای ادغام مجموعه داده‌های امیکس متعدد می‌پردازد.

اکثر فصل‌ها دارای تمرین‌هایی هستند که برخی از نکات مهم معرفی شده در فصل‌ها را تقویت می‌کند. تمرینات در دسته‌های مبتدی، متوسط و پیشرفته طبقه‌بندی می‌شوند. اگر در یک موضوع خاص به خوبی آشنا هستید، ممکن است بخواهید از تمرینات سطح مبتدی صرف نظر کنید.

به طور خلاصه، این کتاب یک راهنمای جامع برای ژنومیک محاسباتی است. برخی از بخش‌ها به خاطر مخاطبان بین‌رشته‌ای گسترده و تکمیل اطلاعات آن‌ها وجود دارد و همه بخش‌ها به یک اندازه برای همه خوانندگان این مخاطب گسترده مفید نیستند.

اطلاعات نرم‌افزار و قراردادها

نام بسته‌ها و کد درون‌خطی و نام فایل‌ها با فونت ماشین‌تحریر (به‌عنوان مثال methylKit) قالب‌بندی می‌شوند. نام توابع با پرانتز دنبال می‌شود (به‌عنوان مثال genomation: (ScoreMatrix)). عملگر دو نقطه‌ای: به معنای دسترسی به یک شیء از یک بسته است.

کنوانسیون عملگر انتساب

به طور سنتی، \rightarrow عملگر انتساب ترجیحی است. با این حال، در سراسر کتاب ما از $=$ و \rightarrow به‌عنوان عملگر انتساب به جای یکدیگر استفاده می‌کنیم. بسته‌های موردنیاز برای اجرای کد کتاب

این کتاب در درجه اول در مورد استفاده از بسته‌های R برای تجزیه و تحلیل داده‌های ژنومیک است، بنابراین اگر می‌خواهید تجزیه و تحلیل را در این کتاب تکرار کنید، باید بسته‌های مربوطه را در هر فصل با استفاده از `install.packages` یا `BiocManager::install` نصب کنید. در هر فصل، هنگامی که از توابع مورد نیاز از بسته‌های مربوطه استفاده می‌کنیم، بسته‌های لازم را با تابع `library()` یا `require()` بارگذاری می‌کنیم. با نگاه کردن به فراخوان‌ها، می‌توانید ببینید چه بسته‌هایی برای کد آن قطعه یا فصل مورد نیاز است. اگر نیاز به نصب تمام بسته‌های الحاقی برای کتاب دارید، می‌توانید دستور زیر را اجرا کنید و در حین انتظار یک فنجان چای بنوشید.

```
if (requireNamespace("BiocManager", quietly = TRUE))
  install.packages("BiocManager")
BiocManager::install(c('qvalue','plot3D','ggplot2','pheatmap','cowplot',
  'cluster', 'NbClust', 'fastICA', 'NMF','matrixStats',
  'Rtsne', 'mosaic', 'knitr', 'genomation',
  'ggbio', 'Gviz', 'DESeq2', 'RUVSeq',
  'gProfileR', 'ggfortify', 'corrplot',
  'gage', 'EDASeq', 'citr', 'formatR',
  'svglite', 'Rqc', 'ShortRead', 'QuasR',
  'methylKit', 'FactoMineR', 'iClusterPlus',
  'enrichR', 'caret', 'xgboost', 'glmnet',
  'DALEX', 'kernlab', 'pROC', 'nnet', 'RANN',
  'ranger', 'GenomeInfoDb', 'GenomicRanges',
  'GenomicAlignments', 'ComplexHeatmap', 'circlize',
  'rtracklayer', 'BSgenome.Hsapiens.UCSC.hg38',
  'BSgenome.Hsapiens.UCSC.hg19', 'tidyr',
  'AnnotationHub', 'GenomicFeatures', 'normr',
  'MotifDb', 'TFBSTools', 'rGADEM', 'JASPAR2018'
```

داده‌های مورد استفاده در این کتاب

ما در سراسر کتاب به داده‌های بسته‌های مختلف R و Bioconductor تکیه می‌کنیم. برای مجموعه داده‌هایی که با آن بسته‌ها ارسال نمی‌شوند، ما بسته `compGenomRData` خودمان را ایجاد کردیم. می‌توانید این بسته را از طریق `devtools::install_github` (" `compgenomr/compGenomRData`") نصب کنید. ما از تابع `system.file()` برای دریافت مسیر فایل‌ها استفاده می‌کنیم. ما متوجه شدیم که بسیاری از کاربران بی‌تجربه در مورد این عملکرد گیج شده‌اند. این تابع فقط خروجی‌های مسیر کامل فایل نصب شده با بسته داده را می‌دهد.

تمرین‌های موجود در کتاب

در پایان هر فصل مجموعه‌ای از تمرین‌ها وجود دارد. تمرین‌ها در بخش‌های موضوعی که از بخش‌های اصلی فصل پیروی می‌کنند، جدا شده‌اند. علاوه بر این، هر تمرین بر اساس سختی آن به‌عنوان "مبتدی"، "متوسط" و "پیشرفته" طبقه‌بندی می‌شود. تمرین‌های سطح مبتدی را معمولاً می‌توان با تغییردادن کد موجود در فصل انجام داد. تمرین‌های سطح پیشرفته معمولاً به ترکیبی از کدها از بخش‌ها یا فصل‌های مختلف نیاز دارند. سطح متوسط جایی در این بین است. راه‌حل‌های تمرینات در آدرس زیر موجود است.

<https://github.com/compgenomr/exercises>

حرف‌های تکراری

این کتاب با R 4.0.0 و بسته‌های زیر گردآوری شده است. ما فقط بسته‌های اصلی و نسخه‌های آن‌ها را و نه ملحقات آن را فهرست کرده‌ایم.

```
## qvalue_2.20.0 | plot3D_1.3 | ggplot2_3.3.1 | pheatmap_1.0.12
## cowplot_1.0.0 | cluster_2.1.0 | NbClust_3.0 | fastICA_1.2.2
## NMF_0.23.0 | matrixStats_0.56.0 | Rtsne_0.15 | mosaic_1.7.0
## knitr_1.28 | genomation_1.20.0 | ggbio_1.36.0 | Gviz_1.32.0
```

² <https://github.com/compgenomr/compGenomRData>

```
## DESeq2_1.28.1 | RUVSeq_1.22.0 | gProfileR_0.7.0 | ggfortify_0.4.10
## corrplot_0.84 | gage_2.37.0 | EDASeq_2.22.0 | citr_0.3.2
## formatR_1.7 | svglite_1.2.3 | Rqc_1.22.0 | ShortRead_1.46.0
## QuasR_1.28.0 | methylKit_1.14.2 | FactoMineR_2.3 |
iClusterPlus_1.24.0
## enrichR_2.1 | caret_6.0.86 | xgboost_1.0.0.2 | glmnet_4.0
## DALEX_1.2.1 | kernlab_0.9.29 | pROC_1.16.2 | nnet_7.3.14
## RANN_2.6.1 | ranger_0.12.1 | GenomeInfoDb_1.24.0 |
GenomicRanges_1.40.0
## GenomicAlignments_1.24.0 | ComplexHeatmap_2.4.2 | circlize_0.4.9 |
rtracklayer_1.48.0
## tidyr_1.1.0 | AnnotationHub_2.20.0 | GenomicFeatures_1.40.0 |
normr_1.14.0
## MotifDb_1.30.0 | TFBSTools_1.26.0 | rGADEM_2.36.0 |
JASPAR2018_1.1.1
## BSgenome.Hsapiens.UCSC.hg38_1.4.3 |
BSgenome.Hsapiens.UCSC.hg19_1.4.3
```

سپاسگزاری‌ها

مایلم از انجمن‌های R و Bioconductor برای توسعه و نگهداری کتابخانه‌ها برای تجزیه و تحلیل داده‌های ژنومی تشکر کنم. بدون تلاش و فداکاری مستمر آن‌ها، نوشتن چنین کتابی ممکن نخواهد بود.

همچنین از تمامی مربیان، همکاران و کارفرمایان گذشته و حال خود تشکر می‌کنم. تعامل با آن‌ها انگیزه نوشتن چنین کتابی و سازماندهی و تدریس دوره‌های عملی ژنومیک محاسباتی را فراهم کرد.

می‌خواهم از جان کیمل، ویراستار Chapman & Hall/CRC که در انتشار این کتاب به من کمک کرد، تشکر کنم. کار با او لذت‌بخش بود. او سخاوتمندانه موافقت کرد که به من اجازه دهد نسخه آنلاین این کتاب را نگه دارم، تا بتوانم پس از چاپ آن را به‌روز کنم. این یک سفر طولانی برای من بوده است. نوشتن بخش‌هایی از این کتاب را از اوایل سال ۲۰۱۳ شروع کردم. اگر ودران فرانکه، بورا اویار و جاناتان رونن نبودند، حتی بیشتر طول می‌کشید. آن‌ها با مهربانی موافقت کردند که در نگارش فصل‌های گمشده مشارکت کنند و کار بزرگی انجام دادند. من از کمک‌های آن‌ها سپاسگزارم.

افراد زیر با مهربانی برای رفع اشتباهات تایپی و کد و پیشنهادهای مختلف کمک کردند: توماس شالچ، الکس گوسدشان، رودریگو اوگاوا، فی ژائو، جاناتان کیت، جانانی راوی، کریستین شودوما، ساموئل اسلدزیسکی، دانیاهامو و سروش نیکومب.

آلتونا آکالین

برلین، آلمان

درباره نویسندگان

دکتر آلتونا آکالین^۳ ساختار کتاب را سازماندهی کرد، بیشتر کتاب را نوشت و بقیه را ویرایش کرد. او یک دانشمند بیوانفورماتیک و رئیس بیوانفورماتیک و پلتفرم علوم داده‌های امیکس در مؤسسه زیست‌شناسی سیستم‌های پزشکی برلین، مرکز Max Delbrück در برلین است. او از سال ۲۰۰۲ در حال توسعه روش‌های محاسباتی برای تجزیه و تحلیل و یکپارچه‌سازی مجموعه‌های داده‌های ژنومیک در مقیاس بزرگ است. او علاقه‌مند به استفاده از یادگیری ماشینی و آمار برای کشف الگوهای مرتبط با متغیرهای مهم بیولوژیکی مانند وضعیت و نوع بیماری است. او در ایالات متحده آمریکا، نروژ، ترکیه، ژاپن و سوئیس زندگی کرد تا کارهای تحقیقاتی و آموزش مرتبط با ژنومیک محاسباتی را دنبال کند. هدف اساسی کار فعلی او استفاده از امضاهای مولکولی پیچیده برای ارائه سیستم‌های پشتیبانی تصمیم برای تشخیص بیماری و کشف نشانگرهای زیستی است. وی علاوه بر تلاش‌های پژوهشی و مدیریت یک آزمایشگاه علمی، از سال ۲۰۱۵، در دوره‌های ژنومی محاسباتی در برلین با شرکت کنندگانی از سراسر جهان، سازماندهی و تدریس می‌کند. این کتاب بیشتر نتیجه مطالبی است که برای آن‌ها و تلاش‌های قبلی تدریس در کالج پزشکی ویل کورنل در نیویورک و مؤسسه فردریش میشر در بازل، سوئیس تهیه شده است.

دکتر آکالین و نویسندگان مشارکت‌کننده زیر چندین دهه تجربه ترکیبی در تجزیه و تحلیل داده‌ها برای ژنومیک دارند. آن‌ها توسعه‌دهندگان بسته‌های Bioconductor مانند methylKit^۴، genomation^۵، RCAS^۶ و netSmooth^۷ هستند. علاوه بر این، آن‌ها نقش کلیدی در توسعه تمام و کمال مراحل تجزیه و تحلیل داده‌های ژنومیک برای RNA-seq، Bisulfite-seq، ChiP-seq و RNA-seq تک سلولی به نام PiGx^۸ ایفا کرده‌اند.

^۳ <https://github.com/al2na>

^۴ <https://bioconductor.org/packages/release/bioc/html/methylKit.html>

^۵ <https://bioconductor.org/packages/release/bioc/html/genomation.html>

^۶ <https://bioconductor.org/packages/release/bioc/html/RCAS.html>

^۷ <https://bioconductor.org/packages/release/bioc/html/netSmooth.html>

^۸ <http://bioinformatics.mdc-berlin.de/pigx/>

نویسندگان مشارکت کننده

دکتر بورا اویار^۹ در فصل ۸، "تحلیل RNA-seq" مشارکت داشت. او آموزش بیوانفورماتیک خود را در دانشگاه سابانچی (استانبول/ترکیه) آغاز کرد و از آنجا مدرک کارشناسی خود را گرفت. بعداً از دانشگاه سایمون فریزر (ونکوور/کانادا) مدرک کارشناسی ارشد گرفت و سپس از آزمایشگاه بیولوژی مولکولی اروپا در هایدلبرگ/آلمان مدرک دکترا گرفت. از سال ۲۰۱۵، او به عنوان دانشمند بیوانفورماتیک در پلتفرم بیوانفورماتیک و پلتفرم علوم داده Omics در مؤسسه زیست‌شناسی سیستم‌های پزشکی برلین مشغول به کار است. او از طریق تحقیق، همکاری، خدمات و توسعه روش تجزیه و تحلیل داده‌ها به پلت فرم بیوانفورماتیک کمک کرده است. علاقه اصلی تحقیقاتی فعلی او ادغام انواع مختلف مجموعه داده‌های Omics برای کشف نشانگرهای زیستی پیش‌آگهی/تشخیصی سرطان‌ها است.

دکتر ودران فرانکه^{۱۰} در فصل ۹، "تحلیل CHIP-seq" مشارکت کرد. وی دکترای خود را از دانشگاه زاگرب دریافت کرد. کار او بر روی بیو ژنز و عملکرد مولکول‌های RNA کوچک در طول جنین‌زایی اولیه و ایجاد پرتوانی متمرکز بود. قبل از دکترا، او به عنوان محقق علمی زیر نظر بوریس لنهاارد در دانشگاه برگن، نروژ، با تمرکز بر اصول عملکردهای تقویت کننده ژن کار می‌کرد. او تحقیقات خود را در بستر علم داده‌های بیوانفورماتیک و Omics در مؤسسه زیست‌شناسی سیستم پزشکی برلین ادامه می‌دهد. او ابزارهایی را برای یکپارچه‌سازی داده‌های چندگانه Omics، با تمرکز بر توالی‌یابی RNA تک سلولی، و اپی ژنومیک توسعه می‌دهد. دانش یکپارچه او از فیزیولوژی سلولی همراه با مهارت او در تجزیه و تحلیل داده‌ها او را قادر می‌سازد تا راه‌حل‌های خلاقانه‌ای برای مشکلات زیستی دشوار بیابد.

⁹ <https://bioconductor.org/packages/release/bioc/html/methylKit.html>

¹⁰ <https://github.com/frenkiboy>

دکتر جاناتان رونن¹¹ در فصل ۱۱، "تحلیل چندگانه امیکس" مشارکت داشت. دکتر رونن مدرک کارشناسی ارشد خود را در رشته مهندسی کنترل از دانشگاه علم و صنعت نروژ در سال ۲۰۱۰ دریافت کرد. سپس به عنوان توسعه‌دهنده نرم‌افزار در اسلو، بروکسل و مونیخ مشغول به کار شد. در آن زمان، او همچنین در تیم مؤسس www.holderdeord.no بود، وب‌سایتی که رای‌های پارلمان نروژ را به وعده‌های داده شده در اعلامیه‌های حزبی پیوند می‌دهد. در سال‌های ۲۰۱۴-۲۰۱۵، او به عنوان دانشمند داده در آزمایشگاه رسانه‌های اجتماعی و مشارکت سیاسی دانشگاه نیویورک کار کرد. در آن زمان، او همچنین www.lahadam.co.il را راه‌اندازی کرد، وب‌سایتی که پست‌های فیس‌بوک سیاستمداران اسرائیلی را ردیابی می‌کرد. او در سال ۲۰۲۰ مدرک دکترای خود را در زیست‌شناسی محاسباتی گرفت، جایی که او ابزارهایی را برای انتساب RNA-seq تک‌سلولی با استفاده از پیشین‌ها و تجزیه و تحلیل یکپارچه داده‌های چندگانه Omics با استفاده از یادگیری عمیق منتشر کرد.

¹¹ <https://github.com/jonathanronen>

فصل ۱

مقدمه‌ای بر ژنومیک

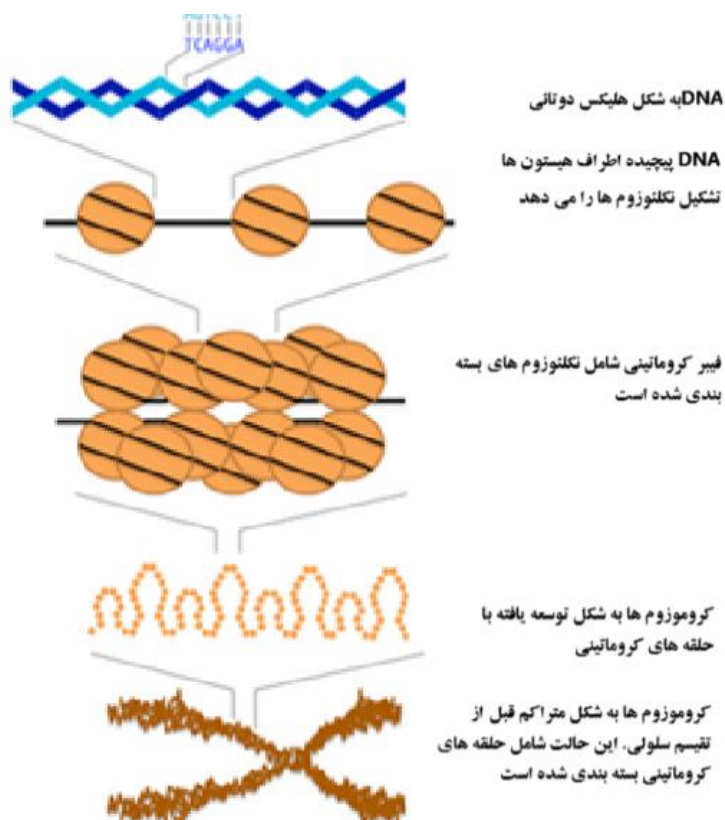
هدف این فصل ارائه برخی از اصول مورد نیاز برای درک زیست‌شناسی ژنوم در اختیار خواننده است. به هیچ وجه، این یک مرور کلی از موضوع نیست، بلکه فقط خلاصه‌ای است که به خواننده غیر زیست‌شناس کمک می‌کند تا مفاهیم بیولوژیکی تکرار شونده در ژنومیک محاسباتی را درک کند. خوانندگانی که در زیست‌شناسی ژنوم و سنجش‌های کمی ژنومی مدرن به خوبی مسلط هستند، باید از این فصل بگذرند یا آن را مرور کنند.

۱-۱ ژن‌ها، DNA و اصل مرکزی

یک مفهوم اصلی که بارها و بارها مطرح می‌شود "ژن" است. قبل از اینکه بتوانیم آن را توضیح دهیم، باید چند مفهوم دیگر را معرفی کنیم که برای درک مفهوم ژن مهم هستند. بدن انسان از میلیاردها سلول تشکیل شده است. این سلول‌ها در کارهای مختلفی تخصص دارند. به عنوان مثال، در کبد سلول‌هایی وجود دارد که به تولید آنزیم‌هایی برای شکستن سموم کمک می‌کنند. در قلب، سلول‌های ماهیچه‌ای تخصصی وجود دارد که باعث تپش قلب می‌شود. با این حال، همه این انواع مختلف سلول‌ها از یک جنین تک‌سلولی می‌آیند. تمام دستورالعمل‌های ساخت انواع مختلف سلول در آن سلول وجود دارد و با هر تقسیم آن سلول، این دستورالعمل‌ها به سلول‌های جدید منتقل می‌شوند. این دستورالعمل‌ها می‌توان در یک رشته - یک مولکول DNA، یک پلیمر ساخته شده از واحدهای تکرار شونده به نام نوکلئوتید رمزگذاری کرد. چهار نوکلئوتید موجود در مولکول‌های DNA، آدنین، گوانین، سیتوزین و تیمین (با چهار حرف رمزگذاری شده‌اند: A، C، G و T) در یک توالی خاص، اطلاعات را برای زندگی ذخیره می‌کنند. DNA به شکل دو مارپیچ سازماندهی شده است که در آن دو پلیمر مکمل با یکدیگر درهم می‌پیچند و به شکل مارپیچ آشنا می‌پیچند.

۱-۱-۱ ژنوم چیست؟

توالی کامل DNA یک موجود زنده که حاوی تمام اطلاعات ارثی است، ژنوم نامیده می‌شود. ژنوم شامل تمام اطلاعات برای ساخت و نگهداری ارگانیسم است. ژنوم‌ها در اندازه‌ها و ساختارهای مختلف هستند. ژنوم ما تنها یک بخش خالص DNA نیست. در سلول‌های یوکاریوتی، DNA به‌دور پروتئین‌ها (هیستون‌ها) پیچیده می‌شود و ساختارهای مرتبه بالاتری مانند نوکلئوزوم‌ها را تشکیل می‌دهند که کروماتین‌ها و کروموزوم‌ها را می‌سازند (شکل ۱-۱ را ببینید).



شکل ۱-۱: ساختار کروموزوم در حیوانات.

بسته به ارگانسیم ممکن است چندین کروموزوم وجود داشته باشد. با این حال، در برخی از گونه‌ها (مانند اکثر پروکاریوت‌ها) DNA به شکل دایره‌ای ذخیره می‌شود. اندازه ژنوم بین گونه‌ها نیز متفاوت است. ژنوم انسان دارای ۴۶ کروموزوم و بیش از ۳ میلیارد جفت پایه است، در حالی که ژنوم گندم دارای ۴۲ کروموزوم و ۱۷ میلیارد جفت پایه است. هم اندازه ژنوم و هم تعداد کروموزوم‌ها بین موجودات مختلف متغیر است. توالی ژنوم موجودات با استفاده از فناوری توالی‌یابی به دست می‌آید. با این فناوری، قطعاتی از توالی DNA از ژنوم به دست می‌آید که خوانش نامیده می‌شود. تکه‌های بزرگ‌تری از توالی ژنوم بعداً با اتصال قطعات اولیه به قطعات بزرگ‌تر با استفاده از همپوشانی خوانش‌ها به دست می‌آید. جدیدترین فناوری‌های توالی‌یابی، توالی‌یابی ژنوم را ارزان‌تر و سریع‌تر کرده است. این فناوری‌ها تولید خوانش بیشتر، خوانش طولانی‌تر و خوانش دقیق‌تر می‌کنند.

هزینه تخمینی اولین ژنوم انسان ۳۰۰ میلیون دلار در سالهای ۱۹۹۹-۲۰۰۰ بود. امروزه با ۱۵۰۰ دلار می‌توان یک ژنوم انسانی با کیفیت بالا به دست آورد. از آنجایی که هزینه‌ها در حال کاهش است، محققان و پزشکان می‌توانند داده‌های بیشتری تولید کنند. این امر هزینه‌های ذخیره‌سازی داده‌ها را افزایش می‌دهد و همچنین تقاضا برای افراد واجد شرایط برای تجزیه و تحلیل داده‌های ژنومی را افزایش می‌دهد. این یکی از انگیزه‌های نگارش این کتاب بود.

۱-۱-۲ ژن چیست؟

در ژنوم، مناطق خاصی حاوی اطلاعات دقیقی وجود دارد که محصولات فیزیکی اطلاعات ژنتیکی را رمزگذاری می‌کند. منطقه‌ای در ژنوم با این اطلاعات به طور سنتی "ژن" نامیده می‌شود. با این حال، تعریف دقیق ژن هنوز در حال توسعه است. طبق کتاب‌های درسی کلاسیک زیست‌شناسی مولکولی، یک ژن بخشی از یک توالی DNA است که مرتبط به یک پروتئین منفرد یا یک مولکول RNA ساختاری و کاتالیزوری است (Alberts et al., 2002). یک تعریف مدرن این است: "یک منطقه (مناطق) که شامل تمام عناصر دنباله‌ای لازم برای رمزگذاری رونوشت عملکردی است (Eilbeck