

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



تجزیه و تحلیل داده‌های RNA-seq

همراه با لوح فشرده آموزش گام به گام

دکتر علیرضا سیفی

عضو هیئت علمی دانشگاه فردوسی مشهد

مهندس محمدرضا رضائی

سرشناسه:
عنوان و نام پدیدآور:

سیفی، علیرضا، ۱۳۵۶-
تجزیه و تحلیل داده‌های RNA-seq همراه با لوح فشرده آموزش گام‌به‌گام/ علیرضا سیفی،
محمدرضا رضایی؛ ویراستار علمی امین میرشمسی کاخکی؛ ویراستار ادبی هانیه اسدپور فعال
مشهد.

مشخصات نشر:

مشهد: دانشگاه فردوسی مشهد، ۱۴۰۰.

مشخصات ظاهری:

ص. ۱۲۸

فروست:

انتشارات دانشگاه فردوسی مشهد؛ ۷۹۰.

شابک:

ISBN: 978-964-386-483-5

وضعیت فهرست‌نویسی:

فیبا.

موضوع:

آر. ان. ا.

موضوع:

نوکلئوتیدها -- توالی

شناسه افزوده:

رضایی، محمدرضا، ۱۳۷۳-

شناسه افزوده:

میرشمسی کاخکی، امین، ۱۳۶۵- ویراستار

شناسه افزوده:

دانشگاه فردوسی مشهد، انتشارات.

رده‌بندی کنگره:

QP۶۲۳

رده‌بندی دیویی:

۵۷۲/۸۸

شماره کتابشناسی ملی:

۷۶۱۴۲۶۹

RNA

Nucleotide sequence

تجزیه و تحلیل داده‌های RNA-seq

همراه با لوح فشرده آموزش گام‌به‌گام

پدیدآورندگان: دکتر علیرضا سیفی؛ مهندس محمدرضا رضایی

ویراستار علمی: دکتر امین میرشمسی کاخکی

ویراستار ادبی: هانیه اسدپور فعال مشهد

مشخصات: وزیری، ۲۰۰ نسخه، چاپ اول، پاییز ۱۴۰۰

چاپ و صحافی: چاپخانه دقت

بها: ۳۲۰/۰۰۰ ریال

حق چاپ برای انتشارات دانشگاه فردوسی مشهد محفوظ است.



انتشارات
۷۹۰

مراکز پخش:

فروشگاه و نمایشگاه کتاب پردیس: مشهد، میدان آزادی، دانشگاه فردوسی مشهد، جنب سلف یاس

تلفن: ۳۸۸۰۲۶۶۶ - ۳۸۸۳۳۷۲۷ (۰۵۱)

مؤسسه کتابیران: تهران، خیابان کارگر جنوبی، خیابان لبافی‌نژاد، بین خیابان فروردین و اردیبهشت،

شماره ۲۳۸، تلفن: ۶۶۴۹۴۴۰۹ - ۶۶۴۸۴۷۱۵ (۰۲۱)

مؤسسه دانشیران: تهران، خیابان انقلاب، خیابان منیری جاوید (اردیبهشت) نبش خیابان نظری، شماره ۱۴۲

تلفکس: ۶۶۴۰۰۲۲۰ - ۶۶۴۰۰۱۴۴ (۰۲۱)

<http://press.um.ac.ir>

Email: press@um.ac.ir

فهرست مطالب

پیشگفتار	۷
فصل ۱. تاریخچه پیشرفت‌ها در ژنتیک، بیولوژی مولکولی و بیوانفورماتیک	۹
نسل دوم روش‌های توالی‌یابی	۱۲
نسل سوم روش‌های توالی‌یابی	۱۵
ابزارهای آنالیز داده‌های NGS	۱۶
منابع	۱۹
فصل ۲. آشنایی با اصول کامپیوتر برای پژوهشگران علوم زیستی	۲۱
انتخاب سیستم کامپیوتری مناسب برای آنالیز داده‌های NGS	۲۲
حافظه‌های ذخیره‌سازی ثانویه	۲۴
ریزپردازنده مرکزی یا CPU	۲۵
حافظه اصلی یا RAM	۲۶
منابع	۲۶
فصل ۳. سیستم عامل لینوکس	۲۷
تاریخچه نرم‌افزارهای با دسترسی آزاد و پیدایش لینوکس	۲۷
نصب سیستم عامل Ubuntu	۲۸
دستورات در Ubuntu	۲۹

۳۳ کلیدهای ترکیبی در خط فرمان لینوکس

۳۴ منابع

فصل ۴. محیط نرم‌افزاری و زبان برنامه‌نویسی R ۳۵

۳۷ نصب R

۳۹ توابع R

۴۰ ساده‌سازی با ایجاد متغیر

۴۰ پوشه کاری و تغییر آن

۴۱ داده‌ها در R

۴۲ ساختار داده‌ها در R

۴۵ روش‌های وارد کردن داده‌ها به R

۴۷ عملیات ریاضی روی وکتورها

۴۷ مقایسات منطقی

۴۸ عملگرهای منطقی

۴۹ ترسیم نمودارهای آماری در R

۵۰ نمودار هیستوگرام

۵۱ بسته نرم‌افزاری ggplot2

۵۵ برخی دیگر از توابع مهم R

۵۶ نوشتن برنامه در R

۵۸ منابع

فصل ۵. RNA-seq: انقلابی در مطالعه ترنسکرپتوم ۵۹

۵۹ روش‌های ارزیابی بیان ژن

۶۱ طراحی آزمایشات RNA-seq

۶۵ روش‌های نوین RNA-seq

۶۶ استفاده از داده‌های RNA-seq موجود در NCBI

منابع ۶۸

فصل ۶. ارزیابی کیفیت و پردازش داده‌های ایلومینا ۷۰

فایل fastq ۷۰

کنترل کیفیت خوانش‌ها با استفاده از ابزار FASTQC ۷۴

نصب نرم‌افزار FASTQC ۷۴

پردازش توالی‌ها با استفاده از نرم‌افزار Trimmomatic ۷۹

منابع ۸۳

فصل ۷. هم‌ردیف کردن خوانش‌ها با توالی مرجع ۸۶

مسیرهای مختلف آنالیز داده‌های RNA-seq ۸۶

روش‌های هم‌ردیف کردن خوانش‌ها با توالی مرجع ۸۸

هم‌ردیفی روی ژنوم مرجع با استفاده از HISAT2 ۸۹

هم‌ردیف کردن خوانش‌ها با استفاده از BWA ۹۰

جدول تعداد خوانش‌های هم‌ردیف شده ۹۲

فیلتر کردن جدول خوانش‌های هم‌ردیف شده ۹۲

منابع ۹۳

فصل ۸. آنالیز بیان ژن‌ها ۹۴

بررسی اولیه خوانش‌های هم‌ردیف شده ۹۴

شناسایی ژن‌های با بیان متفاوت با استفاده از DESeq2 ۹۶

ترسیم نقشه حرارتی ۹۸

بررسی بیان ژن‌ها با استفاده از StringTies-Ballgown ۱۰۰

تجزیه و تحلیل‌های پس از شناسایی ژن‌های با بیان متفاوت ۱۰۲

تهیه ترنسکرپتوم مرجع ۱۰۲

منابع ۱۰۴

۱۰۵	فصل ۹. ملاحظات آماری در تجزیه و تحلیل آزمایشات RNA-seq
۱۰۵	طرح آزمایشی در مطالعات RNA-seq
۱۰۶	آزمایش بدون تکرار زیستی
۱۰۷	نرمال سازی تعداد خوانش‌های هم‌ردیف شده
۱۰۹	مدل سازی خوانش‌های RNA-seq
۱۱۰	یادآوری
۱۱۱	یادآوری
۱۱۲	منابع
۱۱۳	پیوست: نصب لینوکس مجازی روی سیستم عامل ویندوز
۱۳۴	نمایه

پیشگفتار

فناوری‌های نوین توالی‌یابی DNA که با عنوان متداول NGS شناخته می‌شوند، تحول شگرفی در رشته‌های مختلف علوم زیستی ایجاد کرده و امکان مطالعه سریع‌تر و جامع‌تر اساس ژنتیکی و مولکولی پدیده‌های زیستی را فراهم آورده‌اند. هم‌راستا با موفقیت‌های چشمگیر سخت‌افزاری در فناوری‌های توالی‌یابی DNA، پیشرفت‌های اساسی نیز در تولید ابزارهای تجزیه و تحلیل داده‌های حاصل از این فناوری‌ها حاصل شده است.

چالش اصلی در تجزیه و تحلیل داده‌های NGS حجم بالا و پیچیدگی این داده‌هاست که بهره‌برداری مناسب از آن‌ها نیازمند استفاده از ابزارهای بیوانفورماتیکی کارا و دقیق است. این ابزارهای بیوانفورماتیکی معمولاً توسط متخصصان علوم ریاضی، کامپیوتر و بیوانفورماتیک ساخته می‌شوند و در غالب موارد بهره‌برداری از آن‌ها برای متخصصان علوم زیستی که اطلاعات کافی از علوم کامپیوتری ندارند، چالش برانگیز است. معمولاً در مراکز تحقیقاتی پویا در دنیا یک یا چند متخصص بیوانفورماتیک عهده‌دار تجزیه و تحلیل داده‌های NGS تولیدشده توسط محققان زیست‌شناسی هستند و بنابراین تمام محققان نیازمند فراگیری این روش‌های تجزیه و تحلیل نیستند. لیکن در ایران به دلیل اینکه هنوز توسعه یافتگی مطلوبی در مراکز تحقیقاتی رخ نداده است و کماکان تشکیل تیم‌های تحقیقاتی با تخصص‌های گوناگون میسر نیست، محققان علوم زیستی نیاز دارند که اشراف نسبی بر روش‌های تجزیه و تحلیل NGS داشته باشند.

مطالعه ترنسکرپتوم (کل محتوای RNA سلول) با استفاده از روش‌های توالی‌یابی RNA (که با عنوان RNA-seq شناخته می‌شود) یکی از قدرتمندترین روش‌های مطالعه پدیده‌های زیستی برای شناسایی مکانیسم‌های مولکولی و ژنتیکی کنترل‌کننده این پدیده‌هاست. ابزارهای بیوانفورماتیکی بسیار متنوعی برای تجزیه و تحلیل داده‌های RNA-seq ارائه شده است و ابزارهای جدید نیز با سرعت در حال توسعه و ارائه است. هدف از گردآوری کتاب حاضر فراهم آوردن مجموعه‌ای است برای محققان علوم زیستی که الزاماً اطلاعات گسترده‌ای از بیوانفورماتیک ندارند، ولی نیازمند تجزیه و تحلیل داده‌های RNA-seq هستند.

در این کتاب مباحث نظری مورد نیاز در مورد RNA-seq و اصول پایه و مقدماتی کامپیوتر توضیح داده می‌شود. سپس روش‌های استاندارد تجزیه و تحلیل داده‌های RNA-seq که در حال حاضر بیشترین کاربرد را در بین پژوهشگران دارند، معرفی می‌شوند. با استفاده از داده‌های واقعی کل مسیرهای تجزیه و تحلیل به صورت گام‌به‌گام توضیح داده می‌شوند، به نحوی که مخاطب با در اختیار داشتن این کتاب بتواند کل مسیر تجزیه و تحلیل RNA-seq را روی کامپیوتر شخصی خود تمرین کند. کلیه نرم‌افزارها، داده‌های مورد استفاده

و خروجی‌های موردانتظار از اجرای هر کدام از مراحل تجزیه و تحلیل در قالب یک لوح فشرده همراه کتاب در اختیار خوانندگان قرار می‌گیرد. تلاش شده است که مطالب کتاب و لوح فشرده همراه کتاب به نحوی تنظیم شود که علاقه‌مندان به فراگیری روش‌های تجزیه و تحلیل RNA-seq با صرف حداقل زمان، اصول این روش‌ها را دریابند و بتوانند از آن‌ها در پروژه‌های پژوهشی خود استفاده کنند.

مزیت عمده این کتاب این است که ابزارهایی که در آن معرفی می‌شوند ابزارهای استاندارد، رایگان و با متن باز هستند. استفاده از این ابزارها به پرداخت هزینه و یا عدول از اصول اخلاقی مربوط به حقوق مالکیت معنوی نیاز ندارد. چنانچه تجزیه و تحلیل با این ابزارها فراگرفته شود، استفاده از ابزارهای جدیدتری که به سرعت در حال توسعه و توزیع در جامعه علمی هستند، بسیار ساده‌تر خواهد بود و پژوهشگر همواره به جدیدترین ابزارهای تجزیه و تحلیل RNA-seq دسترسی خواهد داشت.

هرچند تلاش زیادی شده است که اشکالات نگارشی و علمی در متن وجود نداشته باشد، با این حال سپاسگزار خواهیم بود چنانچه خوانندگان محترم اشکالات احتمالی و یا پیشنهادهای سازنده خود را برای ارتقای کیفی کتاب به نحو مقتضی به این جانب منعکس کنند.

علیرضا سیفی

زمستان ۱۳۹۹