

دانشگاه تهران
پردیس دانشکده‌های فنی
دانشکده علوم و فنون نوین
گروه شبکه



استنتاج درخت فیلوزنی تومور سرطانی با استفاده از داده‌های تکسلولی و تغییرات تعداد تکرار

پایان‌نامه برای دریافت درجه کارشناسی ارشد در رشته مهندسی فناوری اطلاعات
گرایش سامانه‌های شبکه‌ای

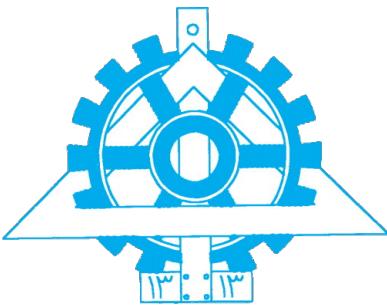
افشین بزرگپور

اساتید راهنما

دکتر سامان هراتی‌زاده و دکتر ابوالفضل مطهری

۱۴۰۰ مرداد

سُبْحَانَ رَبِّ الْجَمَلِ



دانشگاه تهران
پردیس دانشکده‌های فنی
دانشکده علوم و فنون نوین
گروه شبکه



استنتاج درخت فیلوزنی تومور سرطانی با استفاده از داده‌ای تک‌سلولی و تغییرات تعداد تکرار

پایان‌نامه برای دریافت درجهٔ کارشناسی ارشد در رشتهٔ مهندسی فناوری اطلاعات
گرایش سامانه‌های شبکه‌ای

افشین بزرگ‌پور

اساتید راهنما

دکتر سامان هراتی‌زاده و دکتر ابوالفضل مطهری

۱۴۰۰ مرداد



دانشگاه تهران
پردیس دانشکده‌های فنی
دانشکده علوم و فنون نوین

گواهی دفاع از پایان‌نامه کارشناسی ارشد

هیأت داوران پایان نامه کارشناسی ارشد آفای / خانم افشین بزرگ پور به شماره دانشجویی ۸۳۰۵۹۶۰۰۵ در رشته مهندسی فناوری اطلاعات - گرایش سامانه های شبکه ای را در تاریخ با عنوان «استنتاج درخت فیلوژنی تومور سرطانی با استفاده از داده های تک سلولی و تغییرات تعداد تکرار»

با نمره نهایی	به عدد	به حروف

و درجه ارزیابی کرد.

ردیف	مشخصات هیأت	نام و نام خانوادگی	مرتبه دانشگاهی	دانشگاه یا مؤسسه	امضا
۱	استاد راهنمای	دکتر سامان هراتی‌زاده	استادیار	دانشگاه تهران	
۲	استاد راهنمای	دکتر ابوالفضل مطهری	استادیار	دانشگاه تهران	
۳	استاد داور داخلی	دکتر داور داخلی	دانشیار	دانشگاه تهران	
۴	استاد مدعو	دکتر داور خارجی	دانشیار	دانشگاه داور خارجی	
۵	نماینده تحصیلات تکمیلی، دانشکده	دکتر نماینده	دانشیار	دانشگاه تهران	

نام و نام خانوادگی معاون آموزشی و تحصیلات

تکمیلی پردازش دانشکده‌های فنی:

تاریخ و امضا:

نام و نام خانوادگی معاون تحصیلات تکمیلی و

پژوهشی دانشکده / گروه:

تاریخ و امضا:

تعهدنامه اصالت اثر

باسم‌هه تعالی

اینجانب افشنین بزرگ‌پور تأیید می‌کنم که مطالب مندرج در این پایان‌نامه حاصل کار پژوهشی اینجانب است و به دستاوردهای پژوهشی دیگران که در این نوشه از آن‌ها استفاده شده است مطابق مقررات ارجاع گردیده است. این پایان‌نامه قبل‌اً برای احراز هیچ مدرک هم‌سطح یا بالاتری ارائه نشده است.

نام و نام خانوادگی دانشجو: افشنین بزرگ‌پور
تاریخ و امضاء دانشجو:

کلیه حقوق مادی و معنوی این اثر
متعلق به دانشگاه تهران است.

تقدیم به:

همسر و فرزندانم

و

پدر و مادرم

قدردانی

سپاس خداوندگار حکیم را که با لطف بی کران خود، آدمی را به زیور عقل آراست.
در آغاز وظیفه خود می دام از زحمات بی دریغ اساتید راهنمای خود، جناب آقای دکتر ... و ...، صمیمانه
تشکر و قدردانی کنم که در طول انجام این پایان نامه با نهایت صبوری همواره راهنمای و مشوق من بودند و قطعاً
بدون راهنمایی های ارزنده ایشان، این مجموعه به انجام نمی رسید.
از جناب آقای دکتر ... که زحمت مشاوره، بازبینی و تصحیح این پایان نامه را تقبل فرمودند کمال امتنان را
دارم.

با سپاس بی دریغ خدمت دوستان گران مایه ام، خانم ها ... و آقایان ... در آزمایشگاه ...، که با همفکری مرا
صمیمانه و مشفقانه یاری داده اند.

و در پایان، بوسه می زنم بر دوستان خداوندگاران مهر و مهربانی، پدر و مادر عزیزم و بعد از خدا، ستایش می کنم
وجود مقدس شان را و تشکر می کنم از خانواده عزیزم به پاس عاطفه سرشار و گرمای امیدبخش وجودشان، که
بهترین پشتیبان من بودند.

افشین بزرگ پور

۱۴۰۰ مرداد

چکیده

این راهنمای نمونه‌ای از قالب پژوهش، پایان‌نامه و رساله دانشگاه تهران می‌باشد که با استفاده از کلاس-*tehran* و بسته زی پرشین در *LATeX* تهیه شده است. این قالب به گونه‌ای طراحی شده است که مطابق با دستورالعمل نگارش و تدوین پایان‌نامه کارشناسی ارشد و دکتری، مورخ ۹۳/۰۶/۰۳ پردازی دانشکده‌های فنی دانشگاه تهران باشد و حروف چینی بسیاری از قسمت‌های آن، مطابق با استاندارد قالب‌های فارسی پایان‌نامه در لاتک، به طور خودکار انجام می‌شود.

چکیده بخشی از پایان‌نامه است که خواننده را به مطالعه آن علاقمند می‌کند و یا از آن می‌گریزاند. چکیده باید ترجیحاً در یک صفحه باشد. در نگارش چکیده نکات زیر باید رعایت شود. متن چکیده باید مزین به کلمه‌ها و عبارات سلیس، آشنا، بامعنی و روشن باشد. بگونه‌ای که با حدود ۳۰۰ تا ۵۰۰ کلمه بتواند خواننده را به خواندن پایان‌نامه راغب نماید. چکیده، جدای از پایان‌نامه باید به تنها یی گویا و مستقل باشد. در چکیده باید از ذکر منابع، اشاره به جداول و نمودارها اجتناب شود. تمیز بودن مطلب، نداشتن غلط‌های املایی یا دستور زبانی و رعایت دقت و تسلسل روند نگارش چکیده از نکات مهم دیگری است که باید درنظر گرفته شود. در چکیده پایان‌نامه باید از درج مشخصات مربوط به پایان‌نامه خودداری شود. چکیده باید منعکس‌کننده اصل موضوع باشد. در چکیده باید اهداف تحقیق مورد توجه قرار گیرد. تأکید روی اطلاعات تازه (یافته‌ها) و اصطلاحات جدید یا نظریه‌ها، فرضیه‌ها، نتایج و پیشنهادها متمرکز شود. اگر در پایان‌نامه روش نوینی برای اولین بار ارائه می‌شود و تا به حال معمول نبوده است، با جزئیات بیشتری ذکر شود. شایان ذکر است چکیده فارسی و انگلیسی باید حتماً به تأیید استاد راهنما رسیده باشد.

کلمات کلیدی در انتهای چکیده فارسی و انگلیسی آورده می‌شود. محتوای چکیده‌ها بر اساس موضوع و گرایش تحقیق طبقه‌بندی می‌شود و به همین جهت وجود کلمات شاخص و کلیدی، مراکز اطلاعاتی را در طبقه‌بندی دقیق و سریع پایان‌نامه یاری می‌دهد. کلمات کلیدی، راهنمای نکات مهم موجود در پایان‌نامه هستند. بنابراین باید در حد امکان کلمه‌ها یا عباراتی انتخاب شود که ماهیت، محتوا و گرایش کار را به وضوح روشن نماید.

واژگان کلیدی حداقل ۵ کلمه یا عبارت، متناسب با عنوان، قالب پایان‌نامه، لاتک

فهرست مطالب

ث

فهرست تصاویر

خ

فهرست جداول

د

فهرست الگوریتم‌ها

ر

فهرست برنامه‌ها

۱

فصل ۱: مقدمه

۵

فصل ۲: مبانی تحقیق

۵

۱.۲ تنوع ژنتیکی

۸

۲.۲ تکامل تومور^۱

۹

۳.۲ تکنولوژی‌های توالی‌بایی و فراوانی تغییرات آلل^۲

۱۰

۴.۲ ناهمگنی ژنومی تومور

۱۳

۵.۲ بازسازی زیر کلونال

۱۵

۶.۲ تغییرات تعداد کپی

۱۷

۷.۲ جهش‌های ساده بدندی

۱۸

۸.۲ ترک آللی^۳

۱۹

۹.۲ مقدمه‌ای بر مدل‌سازی احتمالی

¹Tumor Evolution

²Variant allele frequency

³Allelic dropout

۲۰	۱.۹.۲ زنجیره مارکوف مونت کارلو ^۴
۲۲	۱۰.۲ یادگیری ماشین ^۵ و یادگیری تقویتی ^۶
۲۱	۱۱.۲ شبکه‌های عصبی بازگشتی
۳۲	۱.۱۱.۲ شبکه عصبی بازگشتی چیست؟
۳۳	۲.۱۱.۲ مزایای شبکه عصبی بازگشتی ^۷
۳۳	۳.۱۱.۲ معایب شبکه عصبی بازگشتی
۳۴	۴.۱۱.۲ کاربردهای شبکه عصبی بازگشتی
۳۴	۵.۱۱.۲ انواع شبکه عصبی بازگشتی
۳۵	۶.۱۱.۲ حافظه‌ی کوتاه‌مدت بلند (LSTM)
۴۳	۱۲.۲ یادگیری تقویتی
۴۳	۱.۱۲.۲ مقدمه و بیشینه تاریخی
۴۵	فصل ۳: روش‌های پیشین
۴۵	۱.۳ مقدمه
۴۵	۲.۳ روش ساخت درخت تکاملی ^۸ با استفاده از داده‌های توالی‌یابی تک سلولی ^۹
۴۷	۱.۲.۳ مدل کیم و سایمون[۴۷]
۴۹	۲.۲.۳ پایگاه داده:
۵۰	۳.۲.۳ معیار ارزیابی:
۵۰	۴.۲.۳ الگوریتم [۸۴]: Bitphylogeny
۵۱	۵.۲.۳ پایگاه داده:
۵۲	۶.۲.۳ معیار ارزیابی:
۵۳	۳.۳ الگوریتم Scite:[۴۶]
۵۶	۱.۳.۳ پایگاه داده:

⁴Markov Chain Monte Carlo (MCMC)⁵Machine learning⁶Reinforcement learning⁷Recurrent Neural Network⁸Tumor evolutionary tree inference⁹Single cell sequencing

۵۷	الگوریتم [۶۰]:Onconem	۴.۳
۵۸	پایگاه داده: ۱.۴.۳	
۵۹	الگوریتم [۶۰]:Sasc	۵.۲
۶۰	پایگاه داده: ۱.۵.۳	
۶۰	الگوریتم [۶۶]:Scarlet	۶.۲
۷۱	الگوریتم [۹]:Deepphylo	۷.۲
۷۷	جمع‌بندی ۱.۰.۷.۳	
۸۱	فصل ۴: روش پیشنهادی	
۸۱	مقدمه ۱.۴	
۸۱	معرفی دادگان ورودی ۲.۴	
۸۲	روش پیشنهادی برای مدیریت داده‌های از دست رفته ۳.۴	
۸۲	روش محاسبه استاتیک ۱.۳.۴	
۸۵	تصادفی ۱.۱.۳.۴	
۸۶	روش پیشنهادی اول (درخت‌بازی) ۴.۴	
۸۶	پیش‌پردازش ۱.۴.۴	
۸۷	فصل ۵: نتایج تجربی	
۸۷	پایگاه داده‌های ورودی ۱.۵	
۸۷	پایگاه داده مصنوعی ۱.۱.۵	
۸۸	ساخت درخت تصادفی ۱.۱.۱.۵	
۹۱	تبديل درخت به ماتریس ژن-سلول ۲.۱.۱.۵	
۹۳	اضافه کردن نویز به ماتریس ژن-جهش ۳.۱.۱.۵	
۹۵	پایگاه داده حقیقی ۲.۱.۵	
۹۶	روش پیشنهادی بدست آوردن درخت فیلوزنی ۲.۵	
۹۷	استفاده از شبکه ژن‌ها برای یافتن درخت فیلوزنی ۱.۲.۵	

۱.۱.۲.۵	استفاده از یک درخت نمونه نابهینه و تغییر اتصالات تا رسیدن به درخت فیلوزنی بنهینه	۹۷
۳.۰.۵	نتایج تجربی	۱۰۰
۱.۳.۵	نتایج بر روی پایگاه داده مصنوعی	۱۰۰
۲.۳.۵	نتایج بر روی داده‌های حقیقی	۱۰۴
۱.۲.۳.۵	نتایج بهینه‌سازی درخت ژنی	۱۰۴
۴.۰.۵	گام‌های آتی	۱۰۵
۱.۴.۵	بهبود در ساخت درخت اولیه	۱۰۵
۲.۴.۵	افزایش سرعت همگرایی MCMC	۱۰۶
۱.۲.۴.۵	تنوع در گام‌ها با استراتژی معقول	۱۰۶
۲.۲.۴.۵	قرار دادن احتمال وزن‌دار به ازای هر انتخاب	۱۰۷
	فصل ۶: بحث و نتیجه‌گیری	۱۰۹
	مراجع	۱۱۱

فهرست تصاویر

۳	دو مدل برای ناهمگونی تومور	۱.۱
۶	مارپیچ دوگانه دی ان ای	۱.۲
۷	همانندسازی دی ان ای	۲.۲
۸	جهش تک نوکلئوتیدی	۳.۲
۸	تغییرات ساختاری	۴.۲
۹	درخت فیلوژنیک تومور	۵.۲
۱۰	تشخیص تغییر بدنه تک نوکلئوتیدی از طریق خوانش هم ترازی	۶.۲
۱۵	درخت کلون تومور	۷.۲
۱۹	نمایی از تطابق ژنتیکی	۸.۲
۲۴	معماری یک شبکه عصبی کانولوشنی	۹.۲
۲۶	عملیات کانولوشن ^{۱۰} در یک شبکه عصبی کانولوشنی ^{۱۱} با کرنل ^{۱۲} 5×5	۱۰.۲
۲۷	(a) تابع فعالیت ReLU ^{۱۳} و (b) تابع فعالیت سیگموید ^{۱۴}	۱۱.۲
۲۸	تابع max-pooling بر روی آرایه دو بعدی کوچک 2×2 و $m = s = 2$	۱۲.۲
۲۹	لایه حذف تصادفی ^{۱۵} با $\sigma = 0.5$	۱۳.۲
۳۲	یک نمونه باز شده شبکه عصبی بازگشتی	۱۴.۲
۳۵	ساختار شبکه عصبی بازگشتی یک به یک	۱۵.۲

¹⁰Convolution

¹¹Convolutional neural network

¹²Kernel

¹³Activation function

¹⁴Sigmoid

¹⁵Dropout

۱۶.۲ ساختار شبکه عصبی بازگشتی یک به چند	۳۵
۱۷.۲ ساختار شبکه عصبی بازگشتی چند به یک	۳۶
۱۸.۲ ساختار شبکه عصبی بازگشتی چند به چند	۳۶
۱۹.۲ ساختار LSTM	۳۷
۲۰.۲ مازول‌های تکرار شونده در شبکه‌های عصبی بازگشتی استاندارد فقط دارای یک لایه هستند.	۳۸
۲۱.۲ مازول‌های تکرار شونده در LSTM‌ها دارای ۴ لایه هستند که با هم در تعامل می‌باشند.	۳۹
۲۲.۲ اشکال از راست به چپ به ترتیب برابر هستند با: کپی کردن، وصل کردن، بردار انتقال، عملیات نقطه به نقطه، یک لایه‌ی شبکه عصبی.	۴۰
۲۳.۲ سلول حالت در مازول LSTM	۴۰
۲۴.۲ نمایی از نحوه تاثیر و ورود اطلاعات به سلول حالت	۴۰
۲۵.۲ قدم اول در پاک کردن اطلاعات از سلول حالت در وضعیت ورودی	۴۱
۲۶.۲ قدم دوم در اضافه کردن اطلاعات جدید به سلول حالت	۴۲
۲۷.۲ بهروزرسانی اطلاعات در سلول حالت	۴۲
۲۸.۲ قدم نهایی برای تولید خروجی مازول LSTM	۴۳
۱.۳ عنوان	۴۸
۲.۳ عنوان	۴۹
۳.۳ عنوان	۵۲
۴.۳ عنوان	۵۳
۵.۳ عنوان	۵۵
۶.۳ عنوان	۵۸
۷.۳ عنوان	۶۲
۸.۳ عنوان	۶۲
۹.۳ عنوان	۶۳
۱۰.۳ عنوان	۶۴
۱۱.۳ عنوان	۶۵

16 Likelihood

۱۰۵	۱۳.۵ نمودار تغییر انرژی در طی گام‌های مختلف
۱۰۶	۱۴.۰ بهترین درخت یافته شده و خروجی الگوریتم برای مقاله SCITE
۱۰۸	۱۵.۰ درخت بدست آمده در مقاله SCITE

فهرست جداول

۷۶	ffffffffff	۱.۳
۷۸	Comparison	۲.۳
۸۵	اندیس‌های به کار رفته در روابط روش پیشنهادی اول	۱.۴
۸۶	پارامترهای مدل ریاضی	۲.۴
۸۶	متغیرهای مدل ریاضی	۳.۴

فهرست الگوریتم‌ها

فهرست برنامه‌ها

فصل ۱

مقدمه

تومور^۱ از رشد غیر طبیعی سلول با احتمال حمله یا گسترش به سایر قسمت‌های بدن تشکیل می‌شود. تومورهای بدخیم^۲ معمولاً سرطان^۳ نامیده می‌شوند. سرطان علل مختلفی از جمله تغییرات ژنتیکی، آلودگی محیط زیست یا انتخاب‌های نادرست در سبک زندگی دارد. یک تومور ممکن است از زیرجمعیت‌های سلولی با تغییرات ژنومی مشخص تشکیل شده باشد، این پدیده ناهمگنی تومور^۴ نامیده می‌شود. ناهمگنی تومور احتمالاً برای درمان سرطان و کشف نشانگر زیستی، به ویژه در روش‌های درمانی هدفمند، تأثیراتی خواهد داشت [۳۲]. درمان‌های فعلی، سرطان را به عنوان یک بیماری همگن درمان می‌کنند [۷۶].

داروهای هدفمند در برابر زیرجمعیت‌های تک یا چند سلولی با انکوژن^۵ جهش‌یافته که آن‌ها را هدف قرار می‌دهند، تولید شده اند، در حالی که آن دسته از زیرجمعیت‌های سلولی که هیچ گونه تاثیری از داروهای به واسطه جهش خود، نمی‌گیرند بدون درمان باقی مانده و ممکن است منجر به عود مجدد تومور یا عدم درمان تومور می‌شوند [۳۲]. این زیرجمعیت‌های سلولی بدون درمان ممکن است منجر به پیشرفت تومور پس از درمان دارویی شوند [۳۲]. به عنوان مثال، رشد مجدد سلول‌های تومورزا در سرطان روده بزرگ^۶ سرطان پستان و گلیوبالستوم^۷ پس از تابش یا درمان سیکلوفسفامید مشاهده شده است [۷۶]. بنابراین، مطالعه روند رشد تومور و ناهمگنی آن تأثیرات زیادی بر تشخیص و درمان سرطان دارد.

تومورها می‌توانند خوش‌خیم، بدخیم و دارای رفتاری نامشخص یا ناشناخته باشند [۲]. تومورهای خوش‌خیم

¹Tumor

²Malignant tumor

³Cancer

⁴Tumor heterogeneity

⁵Oncogene

⁶Colorectal carcinoma

⁷Glioblastomas

شامل فیبروییدهای رحمی^۸ و خالهای ملانوسیتیک^۹ است. آنها محدود و محلی^{۱۰} هستند و به سرطان تبدیل نمی‌شوند [۴]. تومورهای بالقوه بدخیم^{۱۱} شامل سرطان در محل^{۱۲} هستند. آنها به سایر بافت‌ها حمله نکرده و از بین نمی‌روند اما ممکن است به سرطان تبدیل شوند [۳]. تومورهای بدخیم را معمولاً سرطان می‌نامند. آنها به بافت اطراف حمله کرده و از بین می‌روند، ممکن است متاستاز^{۱۳} ایجاد کنند و اگر درمان نشوند یا به درمان پاسخ ندهند، کشنده خواهد بود [۳].

ناهمگنی تومور توضیح می‌دهد که تومور بیش از یک نوع سلول شامل می‌شود. انواع مختلف سلول‌های داخل تومور دارای ویژگی‌های مورفولوژیکی و فیزیولوژیکی متمایزی مانند گیرنده‌های سطح سلول، تکثیر^{۱۴} و رگ‌زایی^{۱۵} هستند. ناهمگنی تومور می‌تواند بین تومورها (ناهمگنی بین توموری) و یا درون تومورها (ناهمگنی درون توموری) رخ دهد. به طور گسترده‌ای پذیرفته شده است که توسعه تومور یک روند تکاملی است [۱۲]، و پیشرونده^{۱۶} معمولاً از یک سلول منشأ می‌گیرند و گروهی از سلول‌ها را تشکیل می‌شوند که در نهایت یک توده را شکل می‌دهند.

دو مدل برای ناهمگنی تومور وجود دارد (شکل ۱.۱). یک مدل تشکیل سرطان از طریق سلول‌های بنیادی بوده که قابلیت ارث‌بری ندارند و مدل دیگر تشکیل سرطان از طریق تکامل کلونی^{۱۷} بوده که قابلیت ارث‌بری دارد. [۱۲]. مفهوم سلول‌های بنیادی سرطانی بیان می‌کند که رشد و پیشرفت بسیاری از تومورها توسط کسری کمی از سلول‌ها کنترل می‌شود و اکثر سلول‌های موجود در تومور محصولات تمایز غیر طبیعی سلول‌های بنیادی سرطانی هستند [۱۲]. بنابراین، برای توصیف و از بین بردن سلول‌های بدخیم در تومورها، لازم است که بر بخش کوچکی از سلول‌های تومورزا تمرکز کنیم [۴۱]. مفهوم تکامل کلونی بیان می‌کند که تومور از یک سلول طبیعی ژنتیکی بوجود می‌آید که به تعداد زیادی سلول تبدیل می‌شود. در این تکامل، جهش‌های تصادفی به طور مداوم تولید می‌شوند و در نهایت تومور حاصل میلیارد‌ها سلول بدخیم است که حاصل از تجمع تعداد زیادی جهش است [۳۸]. تکامل تومور به عنوان توالی پیدرپی گسترش کلونی توصیف می‌شود، که در آن در هر حالت جدید یک رویداد جهش اضافی ایجاد می‌شود [۱۲].

یکی از توالی‌های پی در پی گسترش کلونی، یک مدل خطی از جانشینی کلونی است، جایی که جهش‌های متوالی پیدرپی باعث ایجاد توالی خطی از مجموعه‌های گسترش کلون می‌شوند و منجر به رشد کلون می‌شوند

⁸Uterine fibroid

⁹Melanocytic nevi

¹⁰Local

¹¹Potentially malignant tumor

¹²Carcinoma In Situ

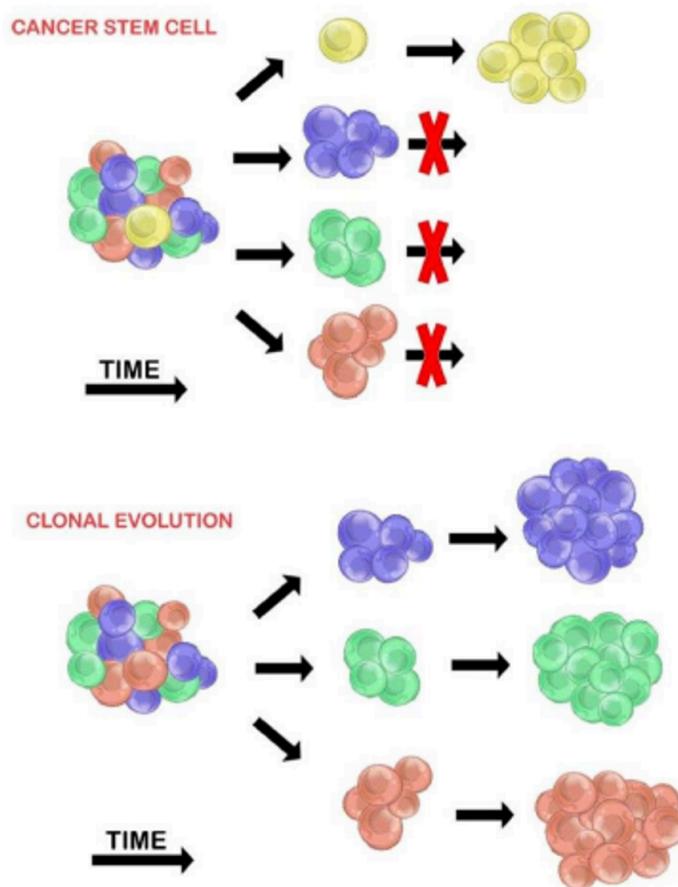
¹³Metastases

¹⁴Proliferative

¹⁵Angiogenic

¹⁶Spontaneous

¹⁷Clonal



شکل ۱.۱: دو مدل برای ناهمگونی تومور

[۱۲]. مورد دیگر یک مدل چند کلونی از پیشرفت تومور است، که در آن یک سلول منفرد از طریق مکانیزم تقسیم به چندین زیرکلون گسترش می‌یابد [۴۹]. این مدل بیش از مدل خطی با ناهمگونی تومور مرتبط است. جهش‌های اکتسابی منجر به افزایش بی ثباتی ژنومی با هر نسل متوالی می‌شود [۱۷].

تومورهای ناهمگن^{۱۸} که متشکل از چندین کلون هستند، می‌توانند حساسیت‌های مختلفی را نسبت به داروهای سمیت سلولی^{۱۹} در نشان دهند. علاوه بر این، می‌زان ناهمگنی تومور می‌تواند خود به عنوان نشانگر زیستی^{۲۰} مورد استفاده قرار گیرد زیرا هر چقدر می‌زان ناهمگنی تومور بیشتر باشد، احتمال حضور کلون‌های مقاوم در برابر درمان بیشتر است [۵۹]. دلایل حساسیت‌های مختلف می‌تواند تعاملات بین کلون‌ها باشد که ممکن است اثر درمانی را مهار یا تغییر دهد [۱۲]. تومورهایی با ناهمگنی زیاد، با احتمال بیشتری از کلون‌های گوناگون تشکیل

¹⁸Heterogenetic¹⁹Cytotoxic²⁰Biomarker

شده است که به درمان مقاوم هستند و ممکن است منجر به عدم موفقیت در درمان شوند. روش‌های نوین درمان تومورها با هدف شخصی‌سازی برنامه‌های درمانی از طریق هدف قرار دادن جمعیت‌های سلولی توموری موجود در یک بیمار، توسعه می‌یابند [۳۱]. ناهمگنی‌های توموری یکی از عوامل اصلی مقاومت در برابر دارو است و بنابراین، یک عامل بالقوه در شکست درمان محسوب می‌شود. [۳۱]. تومورها می‌توانند از راه‌های مختلف به طور همزمان به مقاومت دارویی دست یابند، بنابراین هدف قرار دادن فقط یک مکانیسم مقاومت برای غلبه بر نارسایی درمانی، می‌تواند مزیت درمان‌های هدفمند را محدود کند [۱۴]. بنابراین، ناهمگنی تومور می‌تواند برای درک توسعه تومور، پیچیدگی ایجاد کند و توسعه روش‌های موفقیت آمیز را با چالش روپرور کند [۳۱]. مطالعه ناهمگنی تومور می‌تواند منجر به پیشرفت و توسعه روش‌های درمانی شخصی‌سازی شده شوند و درک ما را از روابط عملکردی بین کلون‌ها در طول درمان افزایش دهنده [۱۴]. برای مطالعه ناهمگنی تومور، بسیاری از ابزارهای محاسباتی موثر برای تجزیه و تحلیل اطلاعات کلونی تومور و تاریخچه تکامل آن تولید شده است. این ابزارها با استفاده از داده‌های تغییرپذیری ژنتیکی، تولید شده توسط فناوری‌های توالی یابی نسبتاً دقیق، قادر هستند تا ترکیب‌های کلونی تومور و رابطه اجداد بین کلون‌ها نتیجه دهند. این اطلاعات برای درک پیشرفت تومور و کمک به پیشرفت‌های درمانی کارآمد مهم است.

در ادامه مفاهیم حوزه تحقیق مثل مدل‌های ناهمگنی توموری، روش‌های مختلف توالی‌یابی، روش‌های مختلف ساخت درخت فیلوزنی تومور، مباحث مرتبط به یادگیری عمیق و یادگیری تقویتی به اختصار توضیح داده شد. در فصل سوم تحقیق پیشرو، به بررسی الگوریتم‌هایی که با استفاده از داده‌های توالی‌یابی تکسولی، درخت فیلوزنی تومور را استباط کرده‌اند پرداخته شد. هر یک از این روش‌ها برای ساخت درخت فیلوزنی به همراه دادگان مورد استفاده، مورد ارزیابی قرار گرفت و در انتهای فصل سوم مقایسه‌های بین روش‌های مختلف صورت گرفت. در فصل چهارم روش پیشنهادی استباط درخت فیلوزنی بر مبنای یادگیری تقویتی و داده‌های توالی‌یابی تکسولی به تفصیل بیان شده و در فصل پایانی نتایج بدست آمده و مقایسه آن با نتایج پیشین، گزارش شده است. در پایان موضوعات پیشنهادی که در کارهای آتی در راستای ادامه این پژوهش می‌تواند مورد بررسی قرار گیرند، توضیح داده شد.

فصل ۲

مبانی تحقیق

در این فصل ابتدا مفاهیم مورد نیاز جهت تعریف مسئله مانند مدل‌های ناهمگنی تومور، روش‌های یافتن درخت تکاملی تومور، روش‌های توالی‌یابی داده مورد بررسی قرار می‌گیرند. در ادامه مدل‌های مورد استفاده برای استتباط درخت تکاملی تومور معرفی می‌شوند. در پایان مفاهیم مرتبط با یادگیری ماشینی، یادگیری عمیق و یادگیری تقویتی به منظور استتباط درخت تکاملی تومور با رویکرد مبتنی بر داده^۱ توضیح داده می‌شوند.

۱.۲ تنوع ژنتیکی

دی‌ان‌ای^۲ یک مولکول بیولوژیکی است که توسط نوکلئوتیدها^۳ پلیمری شده است. در دی‌ان‌ای چهار نوع نوکلئوتید وجود دارد: آدنین^۴، (A) تیمین^۵، (T) سیتوزین^۶ (C) و گوانین^۷ (G). دی‌ان‌ای اساس توالی اسیدهای آمینه است که پروتئین را تشکیل می‌دهد. یک مولکول دی‌ان‌ای از دورشته تشکیل شده است. که در موازات^۸ هم و درجهت‌های مخالف قرار دارد و ساختاری از مارپیچ دوتایی ایجاد می‌کند. هر نوع نوکلئوتید روی یک رشته با نوع دیگری از نوکلئوتید در رشته دیگر مرتبط است: A با T؛ C با G (شکل ۱.۲) [۹]. این به عنوان قانون پایه

¹Data driven

²DNA

³Nucleotid

⁴Adenine

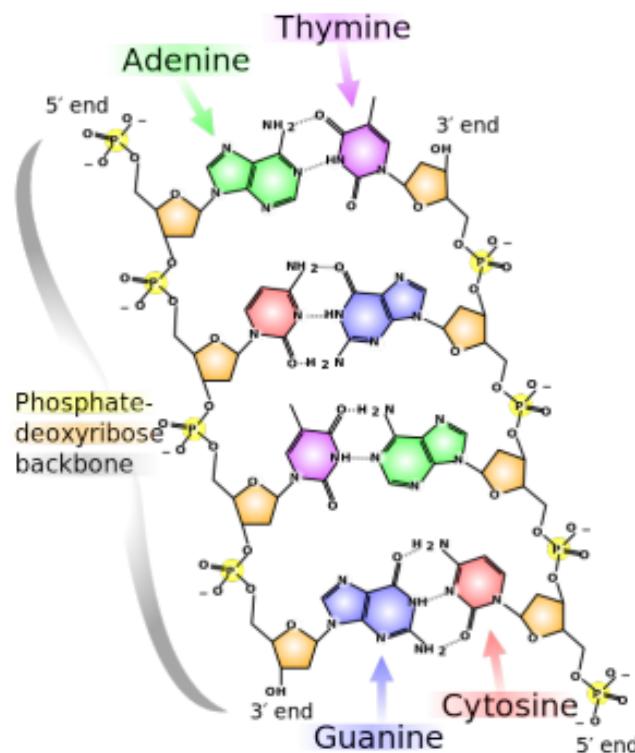
⁵Thymine

⁶Cytosine

⁷Guanine

⁸Antiparallel

جفت شدن نوکلئوتیدها در هر رشته از دی ان ای شناخته می شود.



شکل ۱.۲: مارپیچ دوگانه دی ان ای

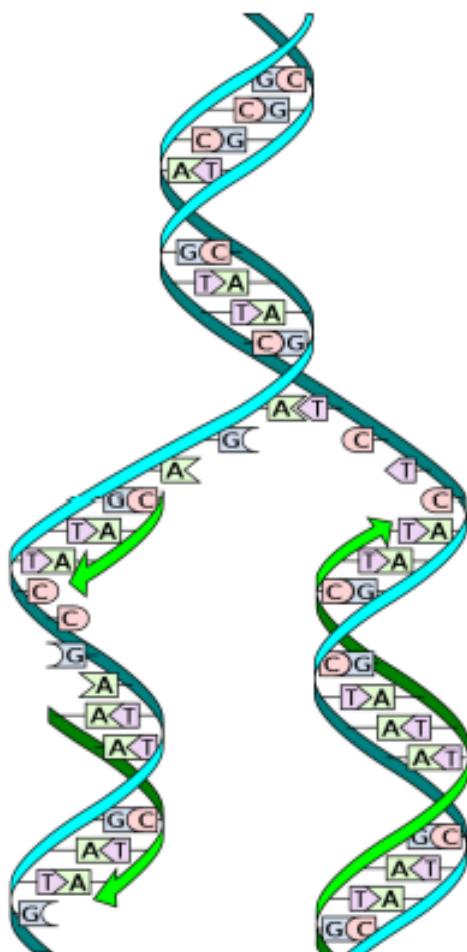
همانند سازی دی ان ای فرآیند تولید دو مولکول دی ان ای یکسان از مولکول دی ان ای اصلی است. وقتی تکثیر شروع می شود، دو رشته یک مولکول دی ان ای از یکدیگر جدا می شوند و هر رشته به عنوان الگویی برای ساخت نمونه مشابه خود عمل می کند. نوکلئوتیدها در هر موقعیت از یک رشته با نوع دیگری از نوکلئوتید مبتنی بر قانون پایه جفت شدن، به منظور سنتز همتای این رشته، متصل می شود. پس از همانند سازی، مولکول دی ان ای اصلی به دو مولکول یکسان تبدیل می شود (شکل ۲.۲) [۶].

ژن ناحیه‌ای از دی ان ای است و به عنوان مولکول واحد وراثت شناخته می شود. ژن‌های متعددی در ساختار دی ان ای با عملکردهای متفاوت وجود دارد. جهش به تغییر دائمی توالی هسته‌ای ژنوم اتلاق می شود. جهش‌ها می توانند در حین فرآیند تکثیر دی ان ای و با جفت‌گیری اشتباه در قسمت‌های مختلف دی ان ای ایجاد می شود. انواع مختلفی از جهش‌ها مانند جهش تک نوکلئوتیدی^۹ (جهش نقطه‌ای)^{۱۰} (شکل ۳.۲) و تغییرات ساختاری^{۱۱}

⁹Single nucleotide mutation

¹⁰Point mutation

¹¹Single variant



شکل ۲.۲: همانندسازی دی‌ان‌ای

شامل درج^{۱۲}، حذف^{۱۳} و برگشت^{۱۴} (شکل ۴.۲) وجود دارد. جهش‌های سلولی می‌توانند به بنا بر دلایلی چون مواد شیمیایی، سمیت یا ویروس ایجاد شوند. جهش در یک ژن می‌تواند محصولات آن را تغییر دهد (مانند ایجاد پروتئین متفاوت) یا از عملکرد صحیح ژن جلوگیری کند [۶].

¹²Insertion¹³Deletion¹⁴reversion

original sequence:

ACTTGGTCAGAATTCCCAGGTGTCA

point mutation:

ACTTGGTCATAATTCCCAGGTGTCA

شکل ۳.۲: جهش تکنوکلئوتیدی

insertion:

ACTTGGTCAGAATTCCCAGGTGTCA
↓
ACTTGGTCAGATAGGCATTCCCAGGTGTCA

deletion:

ACTTGGTCAG~~GAATT~~CCCAGGTGTCA
ACTTGGTCACCCAGGTGTCA

reversion:

ACTTGGTCAGAATTCCCAGGTGTCA
X X
ACTTGGTC~~TTAAGA~~CCCAGGTGTCA

شکل ۴.۲: تغییرات ساختاری

۲۰.۲ تکامل تومور

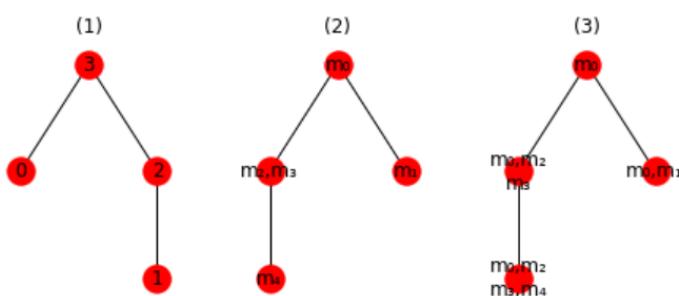
جهشی که در هر سلول از بدن اتفاق می‌افتد، به استثنای سلول‌های جنسی (اسپرم و تخمک)، جهش جسمی^{۱۵} نامیده می‌شود [۱]. تجمع جهش بدنی در طول زندگی یک فرد می‌تواند منجر به رشد کنترل نشده مجموعه‌ای از سلول (تومور) شود [۵۸] و می‌تواند باعث شکل‌گیری سرطان یا بیماری‌های دیگر شود [۱]. بدلیل تجمع سلول‌های گوناگون، بیش از یک نوع سلول در تومور وجود خواهد داشت. به گروه‌های سلول با مجموعه‌ای از جهش مشخص، کلون یا جمعیت سلولی تومور گفته می‌شود. کلون‌های موجود در تومور از نظر فیلوزنیک با هم مرتبط هستند و رابطه آنها را می‌توان با یک درخت فیلوزنیک نشان داد [۱۲]. درخت فیلوزنیک رابطه تکاملی بین

¹⁵somatic

کلون و ترتیب وقوع هر جهش را نشان می‌دهد. به عنوان مثال، شکل ۵.۲ :

- یک درخت فیلوزنیک از یک تومور با چهار کلون با برچسب ۰ تا ۳ را نشان می‌دهد.
- جهش جدیدی را نشان می‌دهد که در هر کلون در طول تکامل این تومور رخ داده است.

همچنین هر کلون جهشی را در مسیر از کلون بالایی به سمت خود به ارث می‌برد. به عنوان مثال، کلون ۰ جهش‌های m_1, m_0 دارد. کلون ۱ دارای جهش m_0, m_2 ، کلون ۲ دارای جهش m_0, m_3 و کلون ۳ دارای جهش m_0, m_1, m_2, m_3 است.



شکل ۵.۲: درخت فیلوزنیک تومور

۳.۲ تکنولوژی‌های توالی‌یابی و فراوانی تغییرات آلل

تعیین توالی دی‌ان‌ای روشنی برای تشخیص ترتیب دقیق نوکلئوتیدها در یک رشته دی‌ان‌ای است. روش توالی‌یابی نسل بعدی^{۱۶} از تعدادی فناوری مدرن توالی تشکیل شده است که امکان تعیین هزینه و زمان توالی‌یابی را به طور موثر فراهم می‌کند. با استفاده از نمونه بیولوژیکی به عنوان ورودی این تکنولوژی‌ها، توالی‌های کوتاه نوکلئوتیدی تولید می‌شود (که به آن خوانش^{۱۷} گفته می‌شود). سپس خوانش با استفاده از الگوریتم هم‌ترازی^{۱۸} متنوعی مانند الگوریتم تبدیل Burrows-Wheeler با ژنوم مرجع تراز می‌شوند. پس از ترازبندی، می‌توان با جمع‌آوری خوانش‌های همپوشانی^{۱۹}، توالی اجماعی^{۲۰} ایجاد کرد (شکل ۶.۲). در موقعیتی از توالی اجماع به دلیل همپوشانی خوانش‌ها، ممکن است بیش از یک نوع خوانش از نوکلئوتید تراز شده وجود داشته باشد (تعداد

¹⁶Next generation sequencing

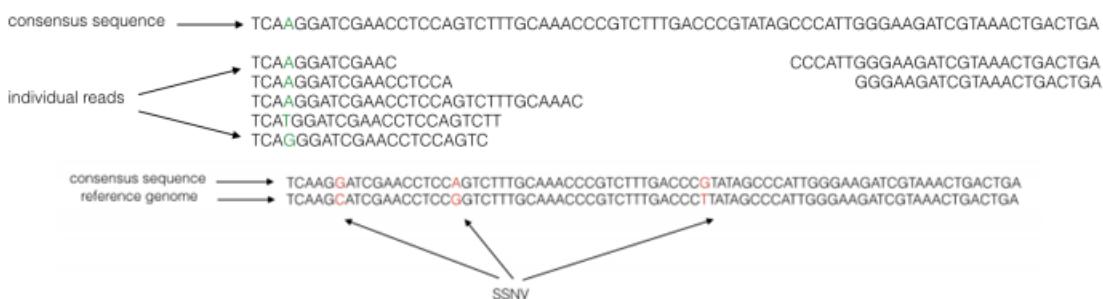
¹⁷Read

¹⁸Alignment

¹⁹Overlapping read

²⁰Consensus

کل قرائت مرتبط با یک نوع جهش، را پوشش خوانش^{۲۱} نامیده می‌شود). نوکلئوتید موجود در این موقعیت به عنوان رایج‌ترین نوکلئوتید تراز شده، مشخص می‌شود. به عنوان مثال، در شکل ۶.۲، سه آدنین، (A) یک گوانین (G) و یک تیمین (T) در موقعیت سوم توالی اجمع تراز می‌شوند، سپس نوکلئوتید در آن موقعیت به عنوان آدنین (A) تعیین می‌شود. پس از ایجاد توالی اجمع، نوکلئوتیدهای موجود در آن توالی، که متفاوت از ژنوم مرتع هستند، شناسایی شده و به عنوان تغییرات بدنی تک نوکلئوتیدی^{۲۲} شناخته می‌شود. با استفاده از نمونه‌های متعدد استخراج شده از یک نمونه تومور، ما می‌توانیم تغییرات بدنی تک نوکلئوتیدی را در هر نمونه با فناوری تعیین توالی‌بایی تشخیص دهیم. نسبت تعداد سلول‌های موجود در یک نمونه حاوی تغییرات بدنی تک نوکلئوتیدی به کل سلول‌ها، فراوانی تغییرات آلل یک تغییر بدنی تک نوکلئوتیدی در این نمونه نامیده می‌شود. مقادیر فراوانی تغییرات آلل برای هر تغییر بدنی تک نوکلئوتیدی در هر نمونه تومور قابل محاسبه است. ابزارهای زیادی برای بازسازی درخت فیلوجنتیک تومور از مقادیر فراوانی تغییرات آلل تومور به عنوان ورودی الگوریتم استفاده می‌کنند.



شکل ۶.۲: تشخیص تغییر بدنی تک نوکلئوتیدی از طریق خوانش هم‌ترازی

۴.۲ ناهمگنی ژنومی تومور

سرطان بیماری‌ای است که بدلیل ایجاد ناهنجاری‌های اساسی در فرآیندهای بنیادی سلول مانند تکثیر^{۲۳}، تمایز^{۲۴} و مرگ^{۲۵} سلول ایجاد می‌شود [۴۰]. این ناهنجاری منجر به رشد کنترل نشده تومور و به کارگیری بافت غیرسرطانی برای حمایت از این رشد می‌شود. علت اصلی این تغییرات جهش است. جهش یک اصطلاح گسترده است که چندین دسته از تغییرات ژنتیکی را پوشش می‌دهد. هنگام حاملگی، یک جنین دارای یک ژنوم خاص

²¹Read coverage

²²Somatic single nucleotide variation

²³Replication

²⁴Differentiation

²⁵Death

و منحصر به فرد است. این ژنوم که به ژنوم جوانهزنی^{۲۶} معروف است، می‌تواند با ژنوم انسانی مرجع مقایسه شود. ژنوم انسانی مرجع یک نمونه از ژنوم انسان است و از دی‌ان‌ای چند نفر تشکیل شده است. تفاوت بین ژنوم جوانهزنی و ژنوم مرجع به عنوان جهش ژنوم جوانهزنی شناخته می‌شود. جهش‌های جوانهزنی می‌توانند مسئول افزایش خطر ابتلا به سرطان باشند [۷۳]^{۲۷}، اما بnderت خود مسئول مستقیم توسعه تومور هستند.

عموماً^{۲۸} تومورها در اثر جهش‌های اکتساب شده پس از لقاح، که معروف به جهش‌های بدنی هستند، ایجاد می‌شوند. جهش‌های بدنی نتیجه اشتباهات در تکثیر دی‌ان‌ای [۱۱]^{۲۹}، قرار گرفتن در معرض جهش‌های با منشأ داخلی یا خارجی یا واردشدن توالی‌های دی‌ان‌ای با منشأ بیرونی بدلیل قرار گرفتن در معرض ویروس است [۷۸]^{۳۰}. غالباً در سرطان، جهش‌های بدنی باعث ایجاد اختلال در روند تکثیر دی‌ان‌ای یا ترمیم آن می‌شوند و حتی جهش‌های بدنی بیشتری ایجاد می‌کنند [۷۴]^{۳۱}. نظریه کلونی بودن سرطان [۵۸]^{۳۲} سرطان را به عنوان یک تک سلولی با منشأ غیرجنسي در نظر می‌گیرد که در اثر تولید مثل فراوان، یک توده متشكل از کلون‌های سلولی گوناگون را ایجاد می‌کند. در این مدل سلولهای توموری با یکدیگر در رقابت هستند و جهش‌های بدنی که مزیت رشد را ایجاد می‌کنند در جمعیت سلول‌های توموری از نسبت بیشتری برخوردار خواهند بود. جهش‌های بدنی که باعث رشد تومور شده و از سلولی به سلولی دیگر منتقل می‌شوند به عنوان جهش‌های راننده^{۳۳} شناخته می‌شوند. اولین سلولی که دارای جهش راننده بوده و آن را به جهش‌های بعدی منتقل می‌کند به عنوان سلول بنیانگذار شناخته می‌شود. همه فرزندان این سلول بنیانگذار، جهش راننده و هر جهش دیگری را که سلول بنیانگذار قبل از به دست آوردن جهش راننده بدست آورده است، دارند. این جهش‌های دیگر، که مزیتی برای رشد و گسترش تنوع توموری ندارند، به عنوان جهش‌های مسافر^{۳۴} شناخته می‌شوند. شایان ذکر است که تعریف جهش راننده و مسافر به زمینه ژنتیکی و محیطی بستگی دارد. به عنوان مثال، شیمی درمانی داروهای سمیت سلولی (سیتوتوکسیک) می‌تواند باعث تغییر جهش از مسافر به جهش راننده شود و عامل اصلی مقاومت در برابر درمان باشد. همچنین جهش‌ها را می‌توان بر اساس نوع تغییری که در دی‌ان‌ای ایجاد می‌شود، به طبقات متمایز تقسیم کرد. حذف و تغییر تک نوکلئوتیدها^{۳۵} جهش‌هایی هستند که یک پایه در ژنوم را به پایه دیگری تغییر می‌دهند. ایندل^{۳۶} درج یا حذف یک بخش دی‌ان‌ای است که می‌تواند کوتاه یا طولانی باشد. از ایندل کوتاه و تغییرات تک نوکلئوتیدی در مجموع به عنوان جهش‌های ساده بدنی^{۳۷} یاد می‌شود. در همه قسمت‌های یک ژنوم، از جمله کل کروموزوم‌ها، قابلیت حذف یا کپی شدن قسمتی از ژنوم وجود دارد. تغییرات شماره کپی به جهشی اتلاق می‌شود که منجر به حذف یا کپی شدن قسمتی از ژنوم می‌شود. تغییرات شماره کپی^{۳۸} نوعی تغییر ساختاری هستند که شامل وارونگی (وقتی

²⁶Germline genome

²⁷Driver mutation

²⁸Passenger mutation

²⁹Single nucleotide variants (SNV)

³⁰Indel

³¹Single Somatic Mutation

³²Copy number alteration

قسمت بزرگی از ژنوم معکوس شده باشد) و انتقال متعادل (جایی که دو بخش ژنومی مکان‌های خود را با یکدیگر تعویض می‌کنند) می‌باشند^[۷۴]. این گونه‌های مختلف جهش مستقل از یکدیگر نیستند و می‌توانند در رابطه با یکدیگر اتفاق یافته‌اند (به عنوان مثال یک جهش می‌تواند منجر به تقویت یک وارونگی شود).

تکنیک توالی‌بایی نسل بعدی این امکان را فراهم کرده است تا با صرف هزینه بسیار کم و با استفاده از یک نمونه توموری، توالی‌بایی از دی‌ان‌ای صورت پذیرد و همین امر منجر به تحول گسترده‌ای در زمینه مطالعه تکامل تومور شده زیر امکان نمونه-برداری در تعداد بسیار بالا را از تومور فراهم می‌کند. نمونه‌گیری در حجم بالا این امکان را فراهم آورده است تا ناهمگنی تومور از نقطه منظر ژنتیکی مورد بررسی قرار گیرد و پاسخ به درمان بیماران سرطانی با جزئیات بیشتری مورد ارزیابی قرار گیرد.

تقریباً همه نمونه‌های استخراج شده از تومور ترکیبی از سلول‌ها با ژنتیپ‌های مختلف را شامل می‌شود. یک نمونه توموری به ندرت فقط شامل بافت سرطانی است زیرا شامل سلول‌های غیر سرطانی از استرومای اطراف^[۳۳] یا سلول‌های ایمنی نفوذی^[۳۴] است. مطالعات ژنومیک نشان داده است که حتی در میان سلول‌های سرطانی، غالباً زیرجمعیت‌های متعدد سرطانی نیز وجود دارد. به عنوان مثال، در یک مطالعه مهم در سال ۲۰۱۲، گرلینگر و همکارانش^[۳۵] توالی‌بایی ژنوم و تغییرات شماره کپی را از طریق نمونه‌های مکانی مجزا استخراج شده از سرطان کلیه اولیه و نقاط متاستاز ثانویه بدست آورده‌اند. با بررسی این نمونه‌های متعدد، مشخص شد که یک ناهمگنی ژنتیکی قابل توجهی در تومور وجود دارد. تعداد بسیار زیادی از جهش‌های شناسایی شده در همه سلول‌های توموری مشاهده نشده‌اند و این بدان معناست که این جهش‌ها بیش از آن‌که یک ناحیه کلونی باشند، به صورت یک ناحیه زیر کلونی بوده‌اند. با استفاده از روش‌های پردازش غیراتوماتیک، تغییرات تک نوکلئوتیدی‌ها و تغییرات شماره کپی بر اساس نمونه‌هایی که از آن استخراج شده‌اند، به خوش‌های مجزا دسته‌بندی شده و یک درخت فیلورژنی به آن‌ها نسبت داده شد. بازسازی درخت فیلورژنیک تومور این امکان را فراهم آورد تا سیر تکاملی تومور با استفاده از شاخه‌های مختلف درخت فیلورژنی شامل جهش‌هایی با عملکرد یکسان از سه ژن متفاوت مورد بررسی قرار گیرد.

در همان سال، یک مطالعه مهم دیگر، "تاریخچه زندگی ۲۱ سرطان پستان"^[۵۷]، حضور ITH را نیز نشان داد. در این مطالعه آنها توالی‌بایی کامل ژنوم را در عمق متوسط ۱۸۸X بر روی تومور پستان PD4120a انجام دادند. این عمق اجازه می‌دهد تا جمعیت‌های شیوع تا ۵٪ کم باشد. آنها مشاهده کردند که تغییرات تک نوکلئوتیدی‌ها در تعداد کمی از خوش‌های مجزا مشاهده می‌شوند که با توجه به کسر نوع آلل (VAF) آنها مشاهده می‌شود، نسبت خواندن‌ها در یک مکان متفاوت شامل آلل نوع. علاوه بر این، آنها توانستند نشان دهنده که برخی از این خوش‌های مجزارانمی‌توان با جهش‌های موجود در تمام جمعیت‌های سرطانی توضیح داد، که این نشان دهنده

³³Surrounding stroma

³⁴Infiltrating immune cell

حضور تغییرات تک نوکلئوتیدی‌های تحت کلونال است. در همان زمان، آنها دریافتند که بسیاری از جهش‌ها در تمام سلول‌های سرطانی موجود در نمونه وجود دارد، که نشان می‌دهد جد مشترک اخیر نسبتاً دیر در زمان تکامل رشد کرده است. مشاهده اینکه جهش‌های زیر کلونال به جای توزیع یکنواخت یا مطابق قانون قدرت در خوش‌های متمایز پیدا شده است، شواهدی را نشان می‌دهد که این جهش‌های زیرکلونالی بیش از آنکه ناشی از تکامل خشی یا مصنوعات فنی باشد، در زیرمجموعه‌های متمایز ناشی از فشارهای انتخابی یافت می‌شود. نویسندهای همچنین با تأیید اینکه جهش‌های زیر کلونال محدود به تغییرات تک نوکلئوتیدی نیستند، توانستند حضور تغییرات شماره کپی‌های کلونال و زیرکلونال را تأیید کنند. نویسندهای یک الگوریتم خوشه‌بندی غیرپارامتریک (یک مدل مخلوط فرآیند دیریشله (DPMM)) را با استدلال قابل توجه دستی برای استنباط فیلوزنی شاخه‌ای از چهار زیر جمعیت سرطانی در آن نمونه منفرد تومور ترکیب کردند. درک معماری ژنتیکی این زیر جمعیت‌ها می‌تواند به مطالعه زیست‌شناسی سرطان کمک کند و نشان داده شده است که در پیش‌بینی بقا در بسیاری از انواع سرطان مفید است [۳]. به عنوان مثال، زیر جمعیت‌های مختلف، که توسط مجموعه جهش‌های جسمی حمل شده تعریف می‌شوند، توانایی‌های مختلفی در مقاومت در برابر درمان و متاستاز دارند. برای انجام این کار، باید از یک یا تعداد کمی از نمونه‌های تومور فله، ژنوتیپ‌های موجود در نمونه را شناسایی کرد. این مسئله، تحت عنوان بازسازی ساب کلونال، موضوع اصلی این پایان‌نامه است. مطالعات پیشگام که نشان داد ITH برای انجام این بازسازی به استدلال دستی قابل توجهی نیاز دارد. استدلال دستی کند، مستعد خطأ است و به تخصص قابل توجهی نیاز دارد. مزایای بازسازی کاملاً خودکار بدیهی است. این بخش پیش زمینه مشکل بازسازی زیر کلونال، چگونگی پرداختن به آن برای انواع مختلف جهش، خصوصیات اصلی الگوریتم‌های بازسازی زیر کلونال و خلاصه‌ای از کارهای موجود در این زمینه را توصیف می‌کند.

۵.۲ بازسازی زیر کلونال

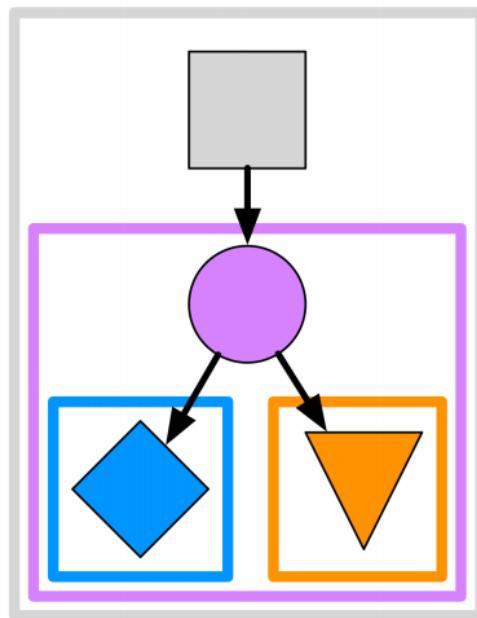
بازسازی ساب کلونال سعی دارد ژنوتیپ‌های موجود در تومور را از تعداد کمی از نمونه‌های توالی دی‌ان‌ای از آن تومور استنباط کند. تعداد ژنوتیپ‌های موجود در تومور از قبل مشخص نیست. این ژنوتیپ‌های زیر کلونال به طور معمول با جهش‌هایی که در مقایسه با ژنوم خط جوانه‌ای دارند، توصیف می‌شوند. ژنوم جوانه‌زنی علاوه بر نمونه‌های تومور، با تعیین توالی یک نمونه غیرسرطانی تعیین می‌شود. در حال حاضر در هنگام تعریف این جمعیت از دو نوع جهش به طور معمول استفاده می‌شود: جهش‌های ساده بدنی‌های متشکل از تعویض‌ها و درج / حذف کوچک (ایندل) و CNA حاصل از تغییرات ساختاری بزرگتر. مشاهده انواع جهش‌های دیگر، مانند مجموعه گسترده‌ای از SV‌ها که شامل بازآرایی هستند، مشاهده آنها دشوارتر است و روش‌های شناسایی آنها در

مراحل اولیه رشد است.

به طور متوسط، حتی در شرایط ایده آل، هر سلول در هر بخش یک جهش پیدا می‌کند [۱۱]، به همین ترتیب، بیشتر سلول‌های تومور ژنتیک منحصر به فردی خواهند داشت. بنابراین، به طور دقیق، اکثر سلول‌های تومور می‌توانند به طور بالقوه نمایانگر زیرجمعیت منحصر به فرد خود باشند. با این حال، به طور عملی، جهش‌هایی که مختص سلول‌های منفرد است یا فقط تعداد کمی از سلول‌ها آنها را به اشتراک می‌گذارد، در حین فراخوانی نوع شناسایی نمی‌شوند. تماس متغیر در بخش ۳.۵.۲ بیشتر مورد بحث قرار گرفته است. بعلاوه، سلول‌هایی که بخش عمده‌ای از جهش‌های خود را به اشتراک می‌گذارند، خصوصاً جهش‌های راننده، صفات مشابهی دارند. به همین ترتیب، من قرارداد گسترهای را اتخاذ کرده و یک زیر جمعیت را به عنوان تمام سلول‌هایی که دارای زیر مجموعه یکسان جهش‌های بدنه در هنگام فراخوانی نوع هستند، تعریف می‌کنم.

یک گام مهم در بازسازی ساب کلونال محاسبه شیوع سلولی تبارهای زیرکلونال و سپس، در نهایت، زیرجمعیت‌های سرطانی است. شیوع سلولی یک زیرجمعیت، نسبت سلول‌های نمونه توالی شده متعلق به آن است. غالباً، شیوع سلولی با تقسیم بر خلوص نمونه، یعنی نسبت سلول‌های سرطانی در نمونه، به بخش سلول‌های سرطانی، نسبت سلول‌های سرطانی، تبدیل می‌شود. هر سلول دقیقاً به یک زیرمجموعه تعلق دارد، بنابراین این شیوع باید در یک جمع باشد. به طور کلی، سلول‌های غیر سرطانی در یک زیرمجموعه واحد قرار می‌گیرند. با این حال، از آنجا که جهش‌ها اغلب در زیرجمعیت‌های متعدد وجود دارند، شیوع سلولی بسیاری از زیرجمعیت‌ها را نمی‌توان مستقیماً از جهش‌های آن استنباط کرد. برای پرداختن به این موضوع، ما یک نسب زیرکلونال برای یک جهش به عنوان مجموعه زیرجمعیت‌هایی که در آن وجود دارد، تعریف می‌کنیم. به طور رسمی، دودمان‌های زیرکلونال از زیر جمعیت بنیانگذار تشکیل می‌شود (جایی که جهش برای اولین بار ظاهر می‌شود) و همه زیرجمعیت‌های بعدی آن (که وراثت جهش) علاوه بر جهش‌های خاص خود، این زیرمجموعه‌های فرزندی حاوی تمام جهش‌های موجود در نژاد تعریف کننده زیر جمعیت هستند (به جز در صورت حذف محل منبع جهش، برای جزئیات بیشتر به فصل ۳ مراجعه کنید). نسب مربوط به یک زیر درخت (یا کلاد) از درخت کلون تومور است. شیوع سلولی یک تبار مجموع شیوع سلولی زیرجمعیت‌هایی است که متعلق به آن تبار هستند. از آنجا که سلول‌ها می‌توانند در چندین نژاد زیرکلونال وجود داشته باشند، شیوع نسب در یک جمع نیست.

شکل ۷.۲ تصویری از یک درخت کلون نمونه را ارائه می‌دهد. گره‌های موجود در درخت، همانطور که در بالا تعریف شد، نشان دهنده زیر جمعیت است. فلاش‌ها از جمعیت والدین به سمت فرزندانشان هدایت می‌شوند. دودمانهای زیرکلونال به صورت مستطیل نشان داده می‌شوند و با توجه به زیرمجموعه بنیادی آنها که در ریشه تیغه یافت می‌شوند، رنگی هستند.



شکل ۷.۲: درخت کلون تومور

۶.۲ تغییرات تعداد کپی

بیشتر ژنوم انسان دیپلوبتید است، به این معنی که دو نسخه از توالی دی‌ان‌ای ما در سلول‌های ما وجود دارد، یکی از پدر و دیگری از مادر. تغییرات شماره کپی این تغییر را می‌دهند، یا با تغییر در تعداد نسخه‌ها (مثلاً از طریق تکثیر کل ژنوم)، نسبت کپی‌های مادر به پدر (مثلاً از دست دادن خشی هتروزیگوتیه در تعداد کپی‌ها، جایی که برای همان منطقه یک ژنوم والدین تکثیر می‌شود و دیگری حذف شده است) یا هر دو (به عنوان مثال کپی کروموزوم مادر). بیشتر این تغییرات (به استثنای تکثیر کل ژنوم) دامنه محدودی از ژنوم را تحت تأثیر قرار می‌دهد، اما می‌تواند از تأثیر یک ژن تا یک کروموزوم کامل باشد. این بخش از ژنوم تغییر یافته به عنوان یک بخش شناخته می‌شود.

تغییرات شماره کپی می‌توانند تعداد کپی کل یک بخش و / یا تعداد نسبی دو کروموزوم والدین را تغییر دهند. هر یک از این تغییرات توسط توالی‌یابی ژنومی هسته قابل تشخیص است. تغییر در تعداد کپی کل یک بخش را می‌توان تشخیص داد زیرا نسبت خواندن آن نقشه به آن بخش بین خط جوانه زنی و نمونه تومور متفاوت خواهد بود. بخش از یک قطعه نسبت ورود خوانده شده است که به یک قطعه در یک نمونه غیر سرطانی ترسیم شده است به نسبت خوانده شده که به یک نمونه سرطانی ترسیم شده است. از نسبت نسبت‌ها برای محاسبه این واقعیت استفاده می‌شود که تعداد کل قرائت‌ها اغلب بین توالی‌یابی سرطانی و غیرسرطانی متفاوت

است، در مناطق مختلف ژنوم عمق خواندن بیشتر یا پایین‌تر ناشی از محتوای GC یا نقشه برداری وجود دارد و تردستی یک تومور با بافت طبیعی متفاوت است. تکرر یک ژنوم، میانگین تعداد کپی از هر کروموزوم است که برای طول کروموزوم نرمال می‌شود.

با تغییر در کسر آلل می‌توان عدم تعادل در تعداد نسخه‌های مادری و پدری این بخش را تشخیص داد. در مناطق دیپلولئید ژنوم‌ها، اگر یک بازه بین کپی‌های مادر و پدر متفاوت باشد، موقعیت هتروزیگوت نامیده می‌شود. جهش‌های تک پایه، خط جوانه زنی همچنین به عنوان چند شکلی تک هسته‌ای نامیده می‌شوند. وقتی یک ژنوم توالی‌یابی شود، حدود نیمی از قرائت آن مکان هتروزیگوت حاوی هر یک از بازها خواهد بود، در نتیجه کسر آلل ۵۰ است. این امر تازمانی که نسبتی برابر با نسخه‌های مادرانه و پدری وجود داشته باشد، صادق خواهد بود. اگر این نسبت تغییر کند، کسر آلل تمام پولیمورفیسم تک هسته‌ای در بخش آسیب دیده تغییر می‌کند. پولیمورفیسم تک هسته‌ای هتروزیگوت به طور متوسط هر ۱۵۰۰ باز [۱۵] رخ می‌دهد و بنابراین برای بخش‌های طولانی بسیاری از پولیمورفیسم تک هسته ای هتروزیگوت تحت تأثیر قرار می‌گیرند. توزیع کسر آلل S تمام پولیمورفیسم تک هسته‌ای در بخش، حالت دوگانه‌ای پیدا می‌کند که هر حالت نشان دهنده نسبت نسخه‌های آن بخش از هر والد است.

فراخوانی CNA چالش برانگیز است زیرا با مشاهده مستقل هر بخش، مسئله هنوز مشخص نشده است. حتی با فرض اینکه هر بخش فقط توسط یک CNA تحت تأثیر قرار گیرد، CNA موسوم به سه پارامتر (نسبت سلولهای حاوی CNA، تعداد کپی‌های مادر و تعداد کپی‌های پدری) وجود دارد و فقط دو مشاهده برای توضیح وجود دارد (و کسر آلل)

همه روش‌ها با فرض اینکه تعداد کمی از نژادهای زیرکلونال مسئول بیشتر یا تمام تغییرات شماره کپی هستند، این ابهام را برطرف می‌کنند. روشی که توسط الگوریتم باتبرگ [۵۷] به کار رفته است، به بیشتر تغییرات شماره کپی وابسته به یک نژاد زیر کلونال منفرد و شایع به نام تبار کلونال متکی است. تحت این روش، شیوع این تبار، همراه با تعداد کپی اصلی و جزئی در تمام تغییرات تعداد کپیکلونال، می‌تواند با یک فرآیند دو مرحله‌ای تخمین زده شود. در گام اول، این روش با فرض شیوع نژاد کلون f_c آغاز می‌شود. شیوع تبار کلونال در بیشتر موارد با خلوص نمونه تومور برابر است. با توجه به شیوع کلونال، هر بخش پس از آن فقط دو متغیر برای توضیح دارد (تعداد کپی بزرگ و جزئی). از آنجا که هر بخش دارای دو مشاهدات است، اکنون مسئله هنوز به درستی تعیین نشده است و بهترین کپی اصلی و مینور متناسب است. سپس، ترکیب کلی مقدار Φ فرض شده با ترکیب مناسب در تمام بخشها تعیین می‌شود. الگوریتم با بهینه سازی این تناسب بهترین مقدار Φ را انتخاب می‌کند. سپس برای هر بخش، شماره کپی اصلی و جزئی با بهینه سازی متناسب بودن قطعه با بهترین مقدار Φ انتخاب می‌شود. این روش فرض می‌کند که تمام تغییرات شماره کپی به نژاد کلونال تعلق دارند، که همیشه درست نیست. در مرحله بعدی، بخش‌هایی که حاوی تغییرات تعداد کپیتخت کلونال هستند با جستجوی بخش‌هایی با اطلاعات مناسب

ضعیف با استفاده از Φ_c استنباط شده مشخص می‌شوند. در این بخش‌ها، روش به طور همزمان و مستقل از هر بخش دیگر، عدد Φ_i و عدد کپی بزرگ و جزئی را استنباط می‌کند.

از آنجا که سه متغیر وجود دارد و تنها دو مشاهده وجود دارد، راه حل‌های بسیاری با تناسب داده برابر وجود دارد که از نظر زیست شناختی برای این تغییرات تعداد کپی زیر کلونال قابل قبول است. این ابهام با انتخاب راه حلی که نزدیکترین شماره به شماره نسخه طبیعی است برطرف می‌شود، اما تعدادی از موارد متداول وجود دارد که این ابتکار عمل ناموفق است. سپس این روش‌ها انتساب تغییرات تعداد کپی زیرکلونال به دودمان و تمام استنباط‌های فیلوزنیک را برای روش‌های پایین دست رها می‌کنند.

رویکرد عمدۀ دیگر این است که فرض کنیم همه تغییرات شماره کپی از تعداد کمی تبار ساب کلونال به وجود می‌آیند. الگوریتم‌هایی که از این روش استفاده می‌کنند به طور مشترک شیوع این نژادها و تعداد کپی بزرگ و جزئی را برای هر بخش استنتاج می‌کنند (به عنوان مثال THetA [۸۰]، TITAN [۸۱] و BIC) یا اثبات دودمانهای زیر کلونال معمولاً با استفاده از احتمال جرمیه شده‌ای مانند معیار اطلاعات بیزی (BIC) یا انواع BIC تعیین می‌شود (به عنوان مثال THetA از BIC اصلاح شده با پارامتر مقیاس گذاری استفاده می‌کند [۸۱]). بنابراین این روش‌ها هم تغییرات شماره کپیرا فراخوانی می‌کنند و هم آنها را به دودمان‌های زیرکلونال اختصاص می‌دهند. هیچ روش موجود این دودمان‌ها را در یک درخت فیلوزنیک قرار نمی‌دهد

۷.۲ جهش‌های ساده بدنی

جهش‌های ساده بدنی جهش‌های کوچکی هستند که می‌توانند مستقیماً از طریق توالی‌یابی و نسبت کروموزوم‌های موجود در نمونه حاوی آنها از تعداد قرائت‌های حاوی جهش و تعداد کل خوانده‌ها در آن مکان، مشاهده شوند. نسبت قرائت حاوی جهش به کل قرائت به عنوان VAF جهش شناخته می‌شود. جهش‌های ساده بدنی‌ها معمولاً با بررسی مشترک ترازها و یک نمونه غیرسرطانی خوانده می‌شوند. این استنباط مشترک برای جداسازی انواع بدنی و ژرمنیال مورد نیاز است.

این فرایند به دلیل انواع مختلف خطاهای و تعصبات که در داده‌های NGS وجود دارد، دشوار می‌شود [۳۳]. یک مشکل اساسی در تشخیص جهش‌های ساده بدنی این است که به نظر می‌رسد خطاهای توالی جهش‌های ساده بدنی شیوع کمی دارند. به طور خاص، در Illumina Hiseq2000 که به طور گسترده استفاده می‌شود، از هر ۱۰۰۰ پایه یکی از آنها دارای یک خط است (به طور معمول یک تعویض) [۶۱]. به همین ترتیب، در طول سه میلیارد پایه ژنوم انسانی، یک احتمال غیر قابل اغماض وجود دارد که در بعضی موقعیت‌ها، چندین بار خواندن دقیقاً شامل خطاهای توالی دقیقاً در همان موقعیت‌ها است. به نظر می‌رسد این خطاهای شیوع کم جهش‌های ساده

بدنی دارند. تمایز بین این خطاهای وشیوع کم واقعی جهش‌های ساده بدنه شامل یک معامله بین حساسیت و ویژگی و در حالت ایده‌آل، یک مدل نویز بسیار دقیق است. حل این مشکل امتداد طبیعی کار گسترهای است که در زمینه فراخوانی جهش‌های جوانهزنی انجام شده است و الگوریتم‌های زیادی برای انجام این کار وجود دارد (به عنوان مثال [۲۱، ۲۲])

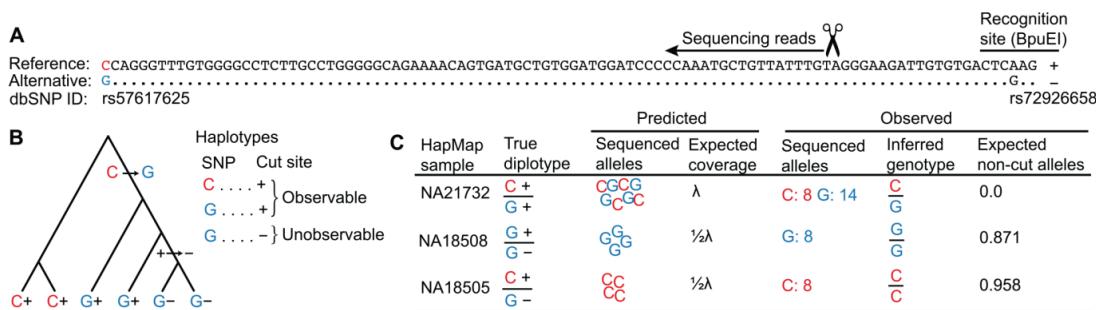
۸.۲ ترک آللی

اگرچه روش‌های تعیین توالی با بازدهی بالا [۴۴] ارزان هستند، اما تحت تاثیر مقدار بایاس هستند و مارکرهای ژنتیکی ای تولید می‌کنند که تقریباً به طور تصادفی در کل ژنوم تقسیم می‌شوند. این روشها با موفقیت در نگاشت ۳۵ صفات [۵۸، ۳۷]، ساخت مپ پیوندی [۶۳، ۳۰]، اسکن انتخاب [۲۲، ۸۱]، و برآورد تنوع ژنتیکی [۱۹] استفاده شده است. یکی از این روش‌ها، تعیین ژنوتیپ براساس توالی [۵] (GBS) است. در GBS، هدف توالی‌یابی فقط با اتصال آداتورهای توالی به محل‌های برش آنژیم محدود کننده، به کمتر از ۵٪ از ژنوم کاهش می‌یابد (شکل زیر). قرائت GBS همچنین می‌تواند به صورت کانکت‌های کوتاه مونتاژ شود، که بدون نیاز به توالی ژنوم فراخوانی یک نوع تغییر تک هسته‌ای (تغییرات تک نوکلئوتیدی) را امکان پذیر می‌کند [۳۹]. از این رو، GBS یک روش محبوب در سیستم‌های غیر مدلی است که به طور معمول فاقد منابعی مانند مجموعه ژنوم و ریزآرایه‌ها است.

بر خلاف توالی‌یابی کل ژنوم (WGS)، GBS مستعد ابتلا به خطاهای مختلف تماس به دلیل محدودیت چندشکلی‌های سایت است (کاهش آللیک). کاهش آللیک در GBS می‌تواند برنامه‌هایی را که به فراخوانی دقیق تغییرات نادر، از جمله تخمین طیف فرکانس سایت در ژنتیک جمعیت متکی هستند، را دچار اختلال کند. یک رویکرد آماری سیستماتیک برای تشخیص کاهش آللیک در داده‌های توالی GBS، اجرا شده و در بسته نرم افزاری منبع باز GBStools وجود دارد. این روش مبتنی بر این واقعیت است که کاهش آللیک متناسب با تعداد آلل‌های سایت محدود کننده بدون برش که در آنجا حمل می‌کند، میزان خوانش نمونه را در یک سایت خاص کاهش می‌دهد. بنابراین GBStools پوشش هر نمونه را در یک سایت خاص به عنوان یک متغیر تصادفی پواسون مورد استفاده قرار می‌دهد که از توزیع با میانگین λ (آللیک‌های بدون برش صفر)، توزیع با میانگین $\lambda/2$ (یک آللیک بدون برش)، یا با میانگین صفر (دو آللیک بدون برش). GBStools حداقل احتمال پارامتر λ را با استفاده از تعداد واقعی آللیک‌های بدون برش در هر نمونه که به عنوان متغیرهای نهفته (مشاهده نشده) در نظر رفته می‌شود و از طریق حداقل‌رساندن مقدار چشم انتظاری (EM)، محاسبه می‌کند. از مقادیر مورد انتظار این متغیرهای نهفته می‌توان برای تخمین اینکه کدام نمونه‌ها یک آللیک بدون برش دارند استفاده کرد. به طور همزمان، GBStools

³⁵Mapping

فرکانس سایت آلل‌های SNP مرجع قابل مشاهده و جایگزین، φ_1 و φ_2 ، و آلیک بدون برش، φ_3 ، که در آن $\varphi_1 + \varphi_2 + \varphi_3 = 1$ برآورد می‌کند و در نهایت، آزمون نسبت احتمال با مقایسه فرضیه صفر $= \varphi_3$ با فرضیه $> \varphi_3$ جایگزین می‌کند. GBStools در اجرای فعلی خود نمی‌تواند ژنتیپ‌های واقعی پنهان شده توسط کاوش آلیک را استنباط کند، اما می‌توان با فیلتر کردن سایت‌هایی که نسبت احتمال آنها زیاد است خطای را حذف کند.



شکل ۸.۲: نمایی از تطابق ژنتیکی

در شکل بالا، آلل BpuEI بدون برش ناشی از SNP rs72926658 با برچسب “-” و آلل برش با “+” برچسب گذاری شده است. آلل “-” در هاپلوتیپ با آلل G مشتق شده بوجود آمده و باعث شده تا برخی از آلل‌های G توسط GBS قابل مشاهده نباشند. نمونه‌های نشان داده شده دارای سه دیپلوتیپ هتروزیگوت است. نتایج توالی با پیش‌بینی‌ها مطابقت داشت و نمونه NA18505 به اشتباه هموزیگوت نامیده می‌شد، اما انتظار می‌رود تعداد آلل‌های کاوشی محاسبه شده توسط GBStools (0.958) با تعداد واقعی (۱) مطابقت داشته باشد، و آن را به عنوان یک تماس اشتباه احتمالی مشخص کند.

۹.۲ مقدمه‌ای بر مدل‌سازی احتمالی

وظیفه اصلی یادگیری ماشین، یادگیری از داده‌ها است، کاری که به عنوان استنباط شناخته می‌شود. برای یادگیری از داده‌ها، باید فرضیاتی را مطرح کرد. توصیف رسمی فرضیات صورت گرفته به عنوان یک مدل ذکر می‌شود. یک مدل احتمالی مفروضات ارائه شده را تعریف می‌کند که اطلاعات آموخته شده را با استفاده از متغیرهای تصادفی و توزیع‌های احتمال به داده‌های مشاهده شده پیوند می‌دهد. توزیع‌های احتمال توابع ریاضی هستند که یک رویداد را ورودی می‌کنند و احتمال آن واقعه را بیرون می‌آورند. توزیع احتمال می‌تواند تابعی بیش از واقعه باشد و این متغیرهای اضافی به عنوان پارامترهای توزیع شناخته می‌شوند [۳۹]. رویکرد بیزی در یادگیری ماشین شامل استنباط احتمالی مقادیر پارامترهای منوط به مشاهدات است [۴۰]. چهار مولفه دارد:

- احتمال: احتمال مشاهده داده‌ها است، مشروط به تنظیم پارامتر ($P(\text{data} | \text{parameters})$)
 - پارامترهای احتمال
 - پارامترهای قبلی
 - داده‌های مشاهده شده

پارامترها خود مجموعه‌ای از متغیرهای تصادفی هستند که از توزیع قبلی ($P(\text{parameters})$) گرفته شده‌اند، که باورهای ما را در مورد احتمال حالت‌های مختلف پارامتر در غیاب مشاهده مشاهده می‌کند. این اصطلاحات با استفاده از قانون بیز با هم ترکیب می‌شوند:

$$P(\text{parameters} | \text{data}) = P(\text{data} | \text{parameters}) * P(\text{parameters}) / P(\text{data}) \quad •$$

$$\text{Posterior} \propto \text{likelihood} * \text{prior} \quad •$$

پس زمینه توزیع پارامترهای مشروط به مشاهده داده‌ها است و خروجی اصلی استنتاج بیزی است. از توزیع پسین می‌توان برای انجام کارهایی مانند پیش‌بینی مشاهدات آینده استفاده کرد.

۱.۹.۲ زنجیره مارکوف مونت کارلو

برای انجام استنتاج بیزی^{۳۶}، ما اغلب می‌خواهیم در توزیع پسین ادغام شده، پیش‌بینی کنیم یا خلاصه‌هایی پیدا کنیم، به عنوان مثال میانگین پارامتر پسین. به طور کلی، انجام چنین ادغامی (جمع بندی در مورد متغیرهای گسته) از نظر تحلیلی غیرقابل حل است. با این حال، می‌توان چنین ادغام‌هایی را با استفاده از نمونه‌هایی که از قسمت پسین ترسیم شده‌اند تقریبی داد:

$$E[f] = \int f(x)p(x)dx \approx 1/N \sum_{1..N} f(x_i) \quad (1.2)$$

که در آن x_i نمونه i از $p(x)$ و $f(x)$ به ترتیب توزیع و عملکرد مورد نظر ما است. به ندرت می‌توان مستقیماً از توزیع پسین نمونه برداری کرد. برای تولید موثر نمونه‌ها از توزیع، حتی در ابعاد بالا، می‌توان از تکنیک زنجیره مارکوف مونت کارلو استفاده کرد. زنجیره مارکوف مونت کارلو یک زنجیره مارکوف می‌سازد که در آن توزیع

³⁶Bayesian

تعادل توزیع پسین است. سپس مقادیر زنجیره می‌تواند به عنوان نمونه از پسین با توجه به همگرایی کافی به توزیع تعادل مورد استفاده قرار گیرد. برای انجام زنجیره ماکوف مونت کارلو، تاز زمانی که بتوان $p(x) \propto p$ را محاسبه کرد، نیازی به محاسبه (x) نیست. این زنجیره ماکوف مونت کارلو را قادر می‌سازد تا از محاسبه ثابت‌های نرمال سازی، که اغلب غیرقابل حل هستند، خودداری کند. یک زنجیره مارکوف به عنوان یک سری متغیرهای تصادفی تعریف می‌شود که دارای ویژگی استقلال شرطی زیر هستند:

$$p(z^{N+1} | z^1..z^N) = p(z^{N+1} | z^N) \quad (2.2)$$

نمونه‌ای از الگوریتم زنجیره ماکوف مونت کارلو الگوریتم Metropolis-Hastings (MH) است [۴۲]. الگوریتم MH از حالت دلخواه Z^t شروع می‌شود. سپس یک حالت پیشنهادی z از توزیع پروپوزال $q(z|z^t)$ ترسیم می‌شود. این حالت پیشنهادی z با احتمال زیر پذیرفته می‌شود:

$$\min(1, \hat{p}(z^*) q(z^t | z^*) / \hat{p}(z^t) q(z^* | z^t)) \quad (3.2)$$

می‌توان نشان داد که الگوریتم MH تعادل دقیق را برآورده می‌کند و از این رو، $p(x)$ توزیع تعادل است [۱۳]. در حالی که توازن دقیق برای اثبات اینکه در محدوده نمونه‌های بی‌نهایت زنجیره به توزیع مورد نظر همگراست کافی است، اما در عمل فقط تعداد محدودی از نمونه‌ها را می‌توان ترسیم کرد. واضح است که نمونه‌های ابتدای زنجیره، که از یک مکان دلخواه در فضای حالت شروع می‌شوند، بعید است از توزیع تعادل باشد. این نمونه‌ها به عنوان نمونه‌های سوختنی کنار گذاشته می‌شوند. هرچه همگرایی زنجیره مارکوف سریعتر باشد، نمونه‌های کمتری باید کنار گذاشته شوند و می‌توان از تعداد بیشتری برای محاسبه انتظارات استفاده کرد. با بررسی اثری از مقادیر مهم پارامتر یا احتمال همگرایی می‌توان نظارت کرد، اما این امر ممکن است چند حالت را از دست بدهد. متاسفانه دانستن اینکه آیا همگرایی حاصل شده است غیرممکن است، فقط گاهی اوقات می‌توان همگرایی را رد کرد [۳۵]. گذشته از همگرایی، یکی دیگر از خصوصیات اصلی یک زنجیره مارکوف میزان اختلاط زنجیره است. با توجه به n نمونه مستقل از توزیع، واریانس میانگین پارامتر برآورده σ_n است که σ انحراف استاندارد توزیع خلفی پارامتر است. نمونه‌های گرفته شده از زنجیره مارکوف مستقل نیستند، زیرا به وضعیت فعلی زنجیره بستگی دارند (یعنی فقط از نظر شرطی مستقل هستند). برای تخمین اندازه نمونه موثر یک زنجیره مارکوف، یعنی تعداد نمونه‌های مستقل با همان خطای استاندارد همان زنجیره، می‌توان از معادله زیر استفاده کرد:

$$ESS = \frac{n}{1 + 2 \sum_{j=1}^{\infty} \rho_j} \quad (4.2)$$

حاصل جمع بی نهایت محاسبه ESS را می‌توان با استفاده از برآوردگر پریودوگرام کوتاه تطبیقی [۷۱] Sokal تخمین زد.

۱۰.۲ یادگیری ماشین و یادگیری تقویتی

آنالیز داده‌های بالینی یک حوزه مهم تحقیقاتی در انفورماتیک، علوم کامپیوتر و پزشکی است که توسط محققان شاغل در دانشگاه‌ها، صنعت و مراکز بالینی انجام می‌شود. یکی از بزرگترین چالش‌ها در تعزیز و تحلیل داده‌های پزشکی، استخراج و تعزیز و تحلیل داده‌ها از تصاویر است. در چند سال اخیر روش‌های یادگیری ماشین انقلابی بزرگ در بینایی کامپیوتر^{۳۷} به وجود آورده است که راه حل‌های جدید و کارآمدی را در مورد خیلی از مسائل و مشکلات موجود در آنالیز تصاویر که مدت زمان طولانی است حل نشده باقی مانده‌اند معرفی می‌کنند. برای اینکه این انقلاب وارد حوزه آنالیز تصاویر پزشکی شود شیوه و روش‌های اختصاصی ای باید طراحی شوند تا خاص بودن تصاویر پزشکی را در نظر گیرند. سیستم‌های کامپیوتری هوشمند چندین دهه است که در دنیا جایگاه برجسته‌ای پیدا کرده‌اند. در حال حاضر، به خاطر تکنیک‌های جدید هوش مصنوعی^{۳۸}، قابلیت پردازش کامپیوتری بالا و رشد گسترده تصویربرداری و ذخیره‌سازی دیجیتالی داده، کاربرد هوش مصنوعی در حال انتقال به حوزه‌های گوناگون می‌باشد. در حوزه پزشکی، سیستم‌های هوش مصنوعی به منظور آشکارسازی بیماری، پیش‌بینی و به عنوان استراتژی پشتیبان در تصمیم‌گیری بالینی در حال توسعه، کاوش و ارزیابی هستند. در زمینه سرطان سینه^{۳۹} از هوش مصنوعی به منظور آشکارسازی زودهنگام و تفسیر ماموگرام‌ها^{۴۰} به منظور بهبود غربالگری سرطان پستان و کاهش تشخیص مثبت کاذب^{۴۱} استفاده می‌شود و این امکان فراهم شده است تا متخصصانی مانند رادیولوژیست‌ها^{۴۲} بتوانند بر اساس میلیون‌ها تصویر از بیماران قبلی که مشخصات مشابهی دارند، تصمیمات آگاهانه‌ای بگیرند. استفاده از هوش مصنوعی در شیوه‌های تشخیص سرطان سینه به مدالیته تصویربرداری^{۴۳} و همچنین تفسیر آسیب‌شناسی^{۴۴} نیز گسترش یافته است. یادگیری عمیق^{۴۵} که زیر شاخه‌ای از یادگیری ماشین می‌باشد یکی از تکنیک‌های هوش مصنوعی است که در انواع مختلفی از مسائل کلینیکی و پردازش تصاویر

³⁷Computer Vision

³⁸Artificial Intelligence (AI)

³⁹Breast cancer

⁴⁰Mammogram

⁴¹False positive

⁴²Radiologist

⁴³Imaging modality

⁴⁴Pathology

⁴⁵Deep learning

پزشکی شامل آشکارسازی^{۴۶}/شناسایی^{۴۷}، قطعه‌بندی^{۴۸} و تشخیص به کمک کامپیوتر^{۴۹} به کار گرفته می‌شود. یادگیری عمیق مجموعه‌ای از الگوریتم‌های ماشین است که قادر به مدل‌سازی الگوهای بطور مستقیم از داده‌های خام می‌باشد. الگوریتم‌های یادگیری عمیق از مجموعه‌ای از لایه‌های چندگانه با واحدهای پردازنده غیرخطی برای استخراج و تبدیل ویژگی استفاده می‌کنند. هر لایه از خروجی لایه قبل به عنوان ورودی استفاده می‌کند. این مفهوم با بسیاری از روش‌های دیگر یادگیری ماشین که نیاز به استخراج ویژگی دارند متفاوت است. به همین ترتیب این الگوریتم‌ها حتی در مسائلی که دانش بسیار کمی در موردشان وجود دارد، می‌توانند مورد استفاده قرار گیرند. اگرچه در دهه ۱۹۹۰ این الگوریتم‌ها در برخی از مطالعات مورد استفاده قرار گرفته‌اند، اما در چند سال اخیر شاهد نتایج بسیار چشمگیر این الگوریتم‌ها هستیم. با توجه به وجود داده‌های بیشتر و همچنین قدرت محاسباتی بالا، این روش‌ها در بسیاری از زمینه‌ها توانسته‌اند به عملکرد انسان یا بهتر از انسان دست یابند^{۵۰}. شبکه‌های عصبی مصنوعی نوع خاصی از مدل‌های یادگیری عمیق هستند که برای کار با داده‌های از نوع تصویر مناسب هستند.

شبکه‌های عصبی مصنوعی مدل‌هایی هستند که در بسیاری از زمینه‌های تحقیقاتی از جمله یادگیری ماشین کاربرد دارند. یک شبکه عصبی مصنوعی از واحدهای ساده‌ای به نام نورون^{۵۱} تشکیل شده است که در یک سیستم پیچیده سازمان یافته‌اند. هر نورون بر اساس ورودی‌های خود، یک خروجی (فعال‌سازی^{۵۲}) را محاسبه می‌کند که می‌تواند فعالیت‌ها یا داده‌های سایر نورون باشد. متدائل‌ترین نوع شبکه عصبی، شبکه عصبی کاملاً متصل شبکه عصبی کاملاً متصل پیش‌خور^{۵۳} است. این شبکه‌ها دارای ورودی (جایی که داده‌ها وارد می‌شوند) و خروجی هستند. به طور معمول، هدف از استفاده از این مدل‌ها حل رگرسیون^{۵۴} یا طبقه‌بندی^{۵۵}، توسط تقریب فعال‌سازی خروجی با مقدار هدف، برای هر داده ورودی است. این شبکه‌ها به صورت لایه^{۵۶} متواالی سازماندهی شده‌اند که یک نورون (واحد) از لایه k تمام نورون لایه $1 - k$ را به عنوان ورودی دریافت می‌کند، ترکیبی خطی از این مقادیر را محاسبه کرده و آن را از طریق تابع غیرخطی عبور می‌دهد

محاسبه خروجی نورون i ام لایه k

$$O_{k,i} = \text{actv}(\mathbf{W}_{k,i} \cdot \mathbf{l}_{k-1} + b_{k,i}) \quad (5.2)$$

⁴⁶Detection

⁴⁷Recognition

⁴⁸Segmentation

⁴⁹Computer-aided diagnosis

⁵⁰Neuron

⁵¹Activation

⁵²Fully-connected feed forward neural network

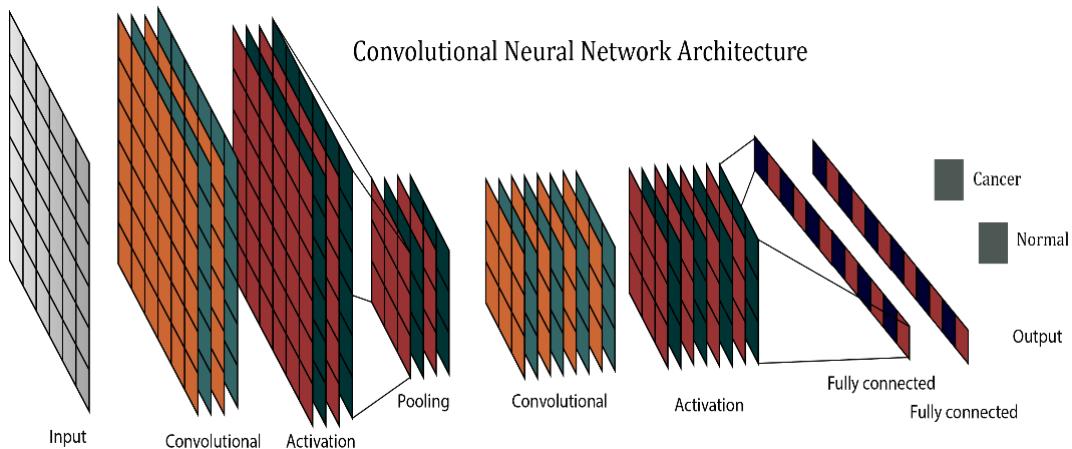
⁵³Regression

⁵⁴Classification

⁵⁵Layer

که واحد $O_{k,i}$ واحد i ام لایه k و $b_{k,i}$ پارامترهای $W_{k,i}$ عدد l_{k-1} بردار تمام فعال سازهای لایه $1 - k$ است. بردار $O_{k,i}$ و عدد $b_{k,i}$ گفته می شود که برای یک وظیفه خاص آموخته می شوند. تابع فعال سازی غیرخطی $actv$ می تواند اشکال مختلفی به خود بگیرد. هر مدل با یک لایه پنهان و تعداد مشخصی نورون اگر پارامترهای کافی داشته باشد می تواند هر تابع پیوسته ای را با خطأ دلخواه تقریب بزند [۲۵].

شبکه های عصبی کانولوشنی یک نوع شبکه عصبی مصنوعی هستند که از نورون ها، لایه ها و وزن ها تشکیل شده اند. مطالعه ای که در سال ۱۹۶۸ میلادی صورت گرفت نشان داد که قشر بینایی مغز برای پردازش اطلاعات از تصاویر از الگوی پیچیده ای استفاده می نماید [۷۷]. نواحی ادراکی که قشر بینایی در آن قرار دارد، همانند فیلترهای محلی بر روی اطلاعات تصویر اعمال می شود. سلول های ساده تر برای تشخیص ویژگی های ادراکی سطح پایین تر در نواحی ادراکی مانند لبه ها کاربرد دارند، همچنین سلول های پیچیده قادر به تشخیص ویژگی های مهم تر و اختصاصی تر و در سطوح بالاتر می باشند. تشخیص ویژگی های اختصاصی تر نتیجه و ترکیبی از ویژگی های سطح پایین می باشد. این عملکرد مغز الهام بخش شبکه های عصبی عمیق امروزی می باشد. مفهوم شبکه کانولوشن نخستین بار در سال ۱۹۸۰ توسط فکوشیما مطرح گردید [۳۴]. اما به دلیل نیاز به سخت افزار ها و پردازشگرهای گرافیکی قوی استفاده از این شبکه ها برای تشخیص تا سال ۲۰۱۲ که به شکل اختصاصی برای تشخیص تصاویر ارایه و معرفی گردیدی به تعویق افتاد [۴۸].



شکل ۹.۲: معماری یک شبکه عصبی کانولوشنی

همانطور که قبل^{۵۶} بیان شد، شبکه های عصبی کانولوشنی مدل های شبکه عصبی کاملا متصل پیش خور هستند که از لایه های زیادی تشکیل شده اند. بسیاری از این مدل ها محدودیت های پارامتر و مکانی دارند که در ادامه توضیح داده خواهد شد. با این حال، آنها در تغییراتی که بر ورودی شان اعمال می کنند تفاوت دارند. در اینجا ما

^{۵۶}Network weight

تمام لایه‌های یک شبکه کانولوشنی و توابع مورد استفاده در آموزش آن‌ها را شرح می‌دهیم. یک معماری می‌تواند یاد بگیرد که مسائل بسیار متفاوتی را حل کند تا زمانی که پارامترها برای هر یک از مسائل به خوبی بهینه شوند. لایه ورودی فقط نمایشی از داده خام است که به مدل داده می‌شود که نیاز به شکل ورودی ثابت دارد. در رایج‌ترین حالت، یک تصویر به یک آرایه 3×3 ^{۵۷} بعدی تبدیل می‌شود با ابعاد $[w, h, 3]$ که w و h عرض و ارتفاع هستند. بعد آخر به دلیل استفاده از تصاویر رنگی ^{۵۸} RGB اغلب 3×3 است. وقتی از تصاویر اشعه ایکس استفاده می‌کنیم چون دارای یک کanal ^{۶۰} هستند بعد سوم برابر با ۱ است.

این لایه اصلی ترین لایه شبکه‌های عصبی کانولوشنی است و این شبکه‌ها نام خود را از این لایه‌ها دریافت می‌کنند. وظیفه این لایه استخراج ویژگی‌ها است. این لایه عملیات کانولوشن را بر روی داده ورودی اعمال می‌کند و خروجی‌هایی به نام نقشه ویژگی ^{۶۱} از این لایه به دست می‌آید. در نتیجه تمامی نوروون‌ها در یک نقشه ویژگی، وزن‌ها و بایاس‌ها ^{۶۳} مشابه و مشترکی دارند که باعث می‌شود، ویژگی‌های تصویر در موقعیت‌های مختلف قابل شناسایی باشند. از طرف دیگر این اشتراک وزن‌ها باعث کاهش تعداد پارامترهای مورد نیاز برای آموزش می‌شود. در شبکه‌های کانولوشن اتصالات به صورت نواحی کوچک و محلی صورت می‌گیرد. به بیان دیگر هر نوروون در نخستین لایه مخفی به ناحیه کوچکی از نوروون‌های ورودی متصل می‌شود. برای مثال اگر این ناحیه 5×5 باشد این ناحیه کوچک 25×25 پیکسلی ناحیه ادراک محلی ^{۶۴} یا کرنل کانولوشن نامیده می‌شود. با توجه به شکل ^{۱۰.۲} یک تصویر ورودی 28×28 داریم که یک کرنل 5×5 بر روی پیکسل‌های ورودی از چپ به راست حرکت می‌کند هر پنجره به نورونی در لایه مخفی متصل می‌شود. بنابراین همان طور که در شکل ^{۱۰.۲} مشخص است لایه مخفی شامل یک شبکه 24×24 نورونی خواهد بود.

در شکل ^{۱۰.۲} هر نوروون لایه مخفی دارای یک بایاس و تعداد 5×5 وزن می‌باشد که به ناحیه ادراکی خود متصل شده است. تمامی نوروون‌های لایه مخفی مذکور که دارای ابعاد 24×24 هستند، دارای وزن‌ها و بایاس‌های مشترکی می‌باشند. به عبارت دیگر خروجی نوروون لایه کانولوشن $y_{w,h,m}$ در طول و عرض w, h و عمق m به صورت رابطه ^{۶۰.۲} است.

$$y_{w,h,m} = f \left(\sum_{i=(w-1)S+1}^{(w-1)S+K} \sum_{j=(h-1)S+1}^{(h-1)S+K} \sum_{k=1}^N W_{k,m}(x_{i,j,k}) + b_m \right) \quad (6.2)$$

⁵⁷ Dimension

⁵⁸ Red Green Blue

⁵⁹ X-ray

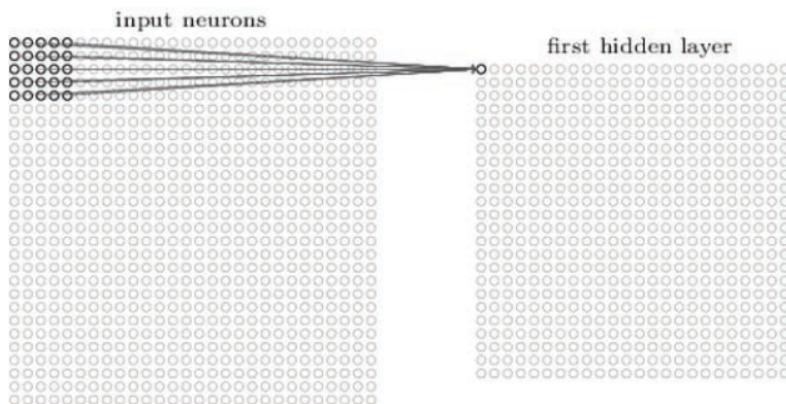
⁶⁰ Channel

⁶¹ Intensity

⁶² Feature map

⁶³ Bias

⁶⁴ Local receptive field



شکل ۱۰.۲: عملیات کانولوشن در یک شبکه عصبی کانولوشنی با کرنل 5×5

که در این رابطه f تابع فعالیت، b_m بایاس مشترک نورون‌ها، $W_{k,m}$ وزن‌های 5×5 مشترک نورون‌ها و $x_{i,j,k}$ ورودی در موقعیت k ، i, j می‌باشد. بنابراین تمامی نورون‌های واقع در لایه مخفی اول به طور دقیق ویژگی‌های مشابهی را در نواحی مختلف تصویر شناسایی می‌کنند. در نهایت خروجی لایه ورودی یا نورون‌های لایه مخفی به عنوان نقشه ویژگی شناخته می‌شوند. ابعاد مربوط به ماتریس خروجی لایه کانولوشن $D_2 \times H_2 \times W_2$ که از ماتریس ورودی با ابعاد $D_1 \times H_1 \times W_1$ است، به صورت رابطه ۷.۲ به دست می‌آید.

$$W_2 = \frac{W_1 - F + 2P}{S+1}, \quad H_2 = \frac{H_1 - F + 2P}{S+1}, \quad D_2 = K \quad (7.2)$$

در روابط ۷.۲ که بیانگر نحوه محاسبه ابعاد ماتریس خروجی کانولوشن است، F, P, S و k به ترتیب نشان‌دهنده اندازه کرنل، مدار لایه‌گذاری صفر^{۶۵}، اندازه اندازه گام^{۶۶} و تعداد فیلترها می‌باشد. طبق این روابط به ازای هر فیلتر تعداد $D_1 \times F \times F$ وزن داریم و با توجه به تعداد k فیلتر موجود، در مجموع تعداد $(D_1 \times F \times F) \times k$ وزن و k بایاس ایجاد می‌شود. بنابراین تعداد پارامترهایی که شبکه در یک لایه کانولوشن خود می‌بایست آموخته بینند زیاد است.

بکارگیری تابع فعالیت در لایه کانولوشن باعث ایجاد خصوصیات غیر خطی در خروجی می‌شود و باعث می‌شود عملکرد مدل متمایز کننده‌تر شود. این توابع با حفظ اندازه لایه، بدون نیاز به پارامترهای آموخته شده، یک عملکرد ساده عنصرگونه در مدل انجام می‌دهند. تابع تابع واحد اصلاح شده خطی^{۶۷} متداول ترین تابع مورد

⁶⁵Zero padding

⁶⁶Stride

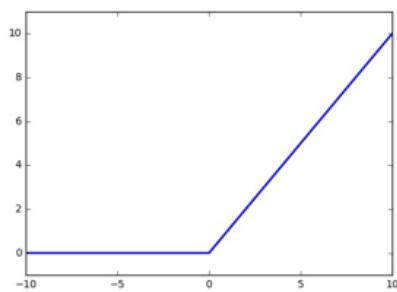
⁶⁷Rectified linear unit (ReLU)

استفاده به خاطر آسان کردن مرحله آموزش است. مثال‌های دیگر شامل تابع سیگموید و هایپربولیک^{۶۸} است.

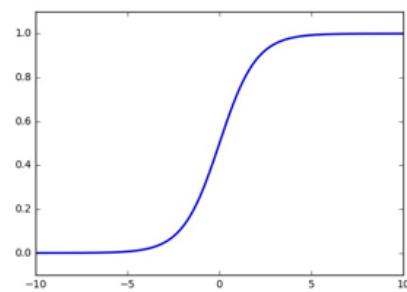
$$\text{ReLU: } r_{m,n,c} = \max\{\circ, l_{x,y,z}\} \quad (8.2)$$

$$\text{Sigmoid: } s_{m,n,c} = \frac{1}{1 + \exp(-l_{x,y,z})}$$

در یک شبکه عصبی کانولوشن معمولاً پس از هر لایه کانولوشن یک لایه pooling قرار می‌گیرد. این لایه از آن



(a) ReLU



(b) Sigmoid

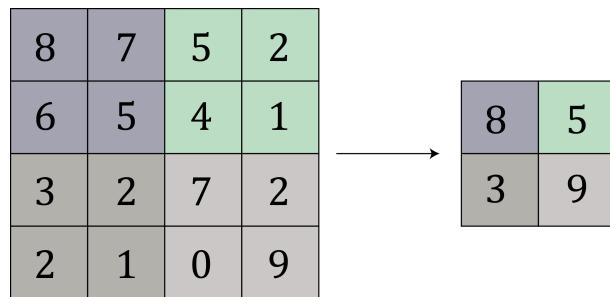
شکل ۱۱.۲: (a) تابع فعالیت ReLU و (b) تابع فعالیت سیگموید

جهت اهمیت دارد که باعث کاهش تعداد پارامترهایی می‌شود که باید آموزش بینند. بنابراین با بکارگیری این لایه ضمن کاهش محاسبات مورد نیاز در بخش آموزش، باعث کنترل بیش‌پردازش^{۶۹} احتمالی در شبکه می‌شود. این لایه بر روی هر عمق از ورودی اعمال می‌شود و اندازه آن را تغییر می‌دهد. دو تابع عملکردی معروف این لایه mean-pooling و max-pooling نام دارند که تابع اول دارای کاربرد بیشتری در شبکه‌های عصبی کانولوشنی است. طریقه عملکرد max-pooling به این صورت است که در هر پنجره بزرگترین پیکسل^{۷۰} را به خروجی می‌فرستد. این پنجره بر روی تصویر مانند تابع کانولوشن از چپ به راست و از بالا به پایین با انداه گام‌های مشخص حرکت می‌کند و نتیجه را به خروجی می‌فرستد. به دلیل اینکه این عملیات بر روی تمامی عمق‌ها اعمال می‌گردد، عمق خروجی همان عمق ورودی به لایه pooling است. یک مثال از عمل max-pooling در شکل ۱۲.۲ به نمایش گذاشته شده است.

$$\text{with } l \in [s \times x, s \times x + m], j \in [s \times y, s \times y + m], \quad R_{x,y,x} = \max\{l_{i,j,z}\} \quad (9.2)$$

⁶⁸Hyperbolic tangent⁶⁹Over-fitting⁷⁰Pixel

لایه کاملاً متصل لایه آخر یک شبکه عصبی کانولوشنی محسوب می‌شود و اتصالات کاملی با خروجی لایه قبلی



شکل ۱۲.۲: تابع max-pooling بر روی آرایه دو بعدی کوچک ۲ و $m = 2$

ایجاد می-کند. این لایه ورودی را دریافت و سپس خروجی را به صورت برداری با N مولفه تولید می‌کند که N تعداد کلاس‌هایی که شبکه باید طبقه بندی کند است. در واقع یک شبکه عصبی کانولوشنی جهت تولید یک بردار خروجی با N مولفه عددی طراحی می‌شود که هر عدد در این بردار خروجی درصد احتمال تعلق به کلاس مورد نظر را نشان می‌دهد. برای یک مسئله با تعداد k کلاس، k نورون خروجی داریم که هر احتمال را با تابع SoftMax محاسبه می‌کنند

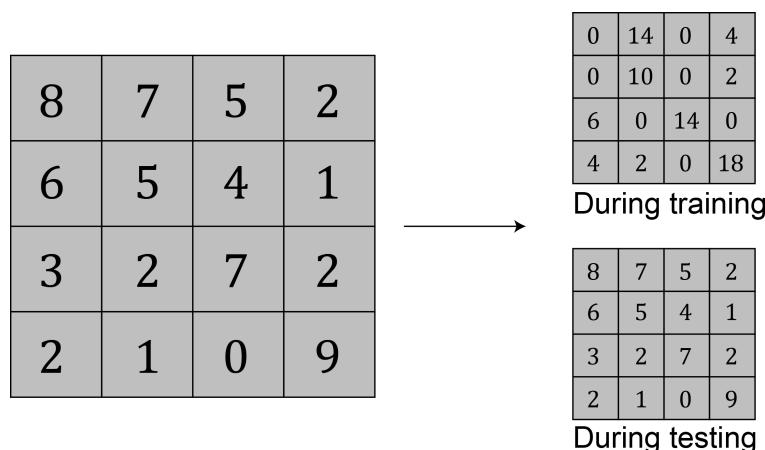
$$P(C)_j = \frac{e^{c_j}}{\sum_{k=1}^K e^{c_k}} \quad (10.2)$$

اگر دو کلاس داشته باشیم می‌توانیم از تابع SoftMax با دو خروجی استفاده کنیم یا از یک نورون استفاده کنیم و تابع سیگموید را محاسبه کنیم. برای دو کلاس احتمال توسط معادله؟؟ محاسبه می‌شود

$$P(1) = \frac{1}{1 + e^i} \quad P(0) = 1 - P(1) \quad (11.2)$$

حذف تصادفی یک روش بسیار رایج برای جلوگیری از بیش‌پردازش شبکه عصبی مصنوعی از جمله مدل‌های یادگیری عمیق است [۷۲]. ایده این تکنیک این است که با جلوگیری از هماهنگی نورون‌ها، ویژگی‌های قوی تری ایجاد شود. اجرای آن ساده است تنها نیاز به بهم چسباندن لایه‌های اضافی در شبکه معمولاً پس از توابع فعال سازی است. این مازول بطور تصادفی برخی از نقاط نقشه ویژگی ورودی را صفر می‌کند. هریک از مازول‌ها دارای یک احتمال مستقل σ برای نگهداری نقاط هستند و در صورت بروز چنین اتفاقی، توسط $\frac{1}{\sigma}$ مقیاس بندی می‌شوند. نقاطی که نگهداری نمی‌شوند بر روی صفر تنظیم می‌شوند. این لایه فقط یک پارامتر σ دارد، که برای

آموزش در فاصله [۱، ۰] قرار دارد و برای آزمایش روی ۱ قرار می‌گیرد. به طور شهودی، می‌توان این فرآیند را به عنوان حذف برخی از نورون‌های شبکه عصبی، به طور موقت، همراه با اتصالات ورودی و خروجی آن تصور کرد. مکانیزم حذف، نورون‌هایی را که به اتصالات ورودی کمتری متکی هستند را در نظر می‌گیرد. زیرا افت بک زیر مجموعه از ورودی‌ها در مقایسه با یک نورون که به بسیاری از ورودی‌ها متکی است، قابل توجه‌تر خواهد بود و به این ترتیب ویژگی‌های کلی تر مهم‌تر می‌شوند. شکل ۱۳.۲ یک مثال از لایه حذف تصادفی را نمایش می‌دهد. نرمال‌سازی دسته^۱ یک تکنیک جدید ولی خیلی کارآمد است. در طی آموزش مدل‌های عمیق، وزن‌ها در هر



شکل ۱۳.۲: لایه حذف تصادفی با $\sigma = 0.5$

تکرار^{۱۲} به روز می‌شوند. یک اثر جانبی این امر این است که در هر لایه توزیع‌های ورودی تغییر می‌کند، پدیده‌ای که به آن تغییر همبستگی داخلی^{۱۳} می‌گویند. این پدیده فرایند آموزش را کند می‌کند، به مقدار دهی دقیق‌تر وزن احتیاج دارد و مانع بهینه‌سازی^{۱۴} مدل‌های غیرخطی اشباع، مانند مماس‌های سیگموید یا هایپربولیک می‌شود. برای حل این مشکل نرمال‌سازی دسته را پیشنهاد می‌شود که مشابه با حذف تصادفی، به عنوان لایه‌ای در شبکه با رفتارهای متفاوت در حین آموزش و آزمون پیاده سازی می‌شود. برای رفع مشکل تغییر کواریانس^{۱۵} داخلی، این لایه برای هر دسته آموزش با کم کردن میانگین و تقسیم بر انحراف استاندارد^{۱۶} همه نورون‌های عمق مشابه، ورودی خود را نرمال می‌کند. به میانگین و انحراف استاندارد آمار mini-batch گفته می‌شود. برای اطمینان از اینکه مدل می‌تواند دقیقاً همان تابع را با یا بدون نرمال‌سازی دسته عادی نشان دهد، دو وزن جدید قابل تمرین^۷

^{۱۱}Batch normalization

^{۱۲}Iteration

^{۱۳}Internal covariate shift

^{۱۴}Optimization

^{۱۵}Covariance

^{۱۶}Standard deviation

و β اضافه می‌شوند که خروجی را اندازه‌گیری و جبران می‌کنند. بنابراین خروجی به صورت معادله ۱۲.۲ است.

$$\begin{aligned} I_c &= \gamma \left(\frac{I_c - \text{mean}(I_c)}{\text{std}(I_c)} \right) + \beta & : \text{در طی آموزش} \\ I_c &= \gamma \left(\frac{I_c - u_c}{v_c} \right) + \beta & : \text{در طی آزمایش} \end{aligned} \quad (12.2)$$

که u_c و v_c متوسط‌های در حال اجرا (I_c) و $\text{mean}(I_c)$ و $\text{std}(I_c)$ هستند. نشان داده شده است که نرمال‌سازی دسته باعث آهنگ یادگیری بالاتر می‌شود و مدل در تکرارهای کمتری همگرا خواهد شد. این روش دارای اثر رگولاژیشن^{۷۷} است. مدل با استفاده از تابع هزینه^{۷۸} یاد می‌گیرد. این روشی است برای ارزیابی اینکه تا چه میزان خوب یک الگوریتم داده‌های مشاهده شده را می‌تواند مدل سازی کند. اگر پیش‌بینی‌ها بیش از حد از نتایج واقعی منحرف شوند، تابع هزینه مقدار بالایی خواهد داشت. به تدریج، با کمک برخی توابع بهینه سازی، تابع هزینه می‌آموزد تا خطأ در پیش‌بینی را کاهش دهد.

بهینه سازی مهمترین بخش در الگوریتم‌های یادگیری عمیق است. این کار با تعریف تابع هزینه شروع می‌شود و با به حداقل رساندن آن با استفاده از یک روش بهینه سازی به پایان می‌رسد. فرض کنید یک مجموعه داده D با تعداد I تصویر داریم. این تصاویر می‌توانند ضایعه باشند یا نباشند، بنابراین دارای برچسب $\{0, 1\}$ هستند. باید مدلی بسازیم که با توجه به یک تصویر ورودی I_i ، یک احتمال (I_i) p تولید کند که تا حد ممکن به برچسب مربوط به آن تصویر (y_i) نزدیک باشد. برای این منظور الگوریتم‌های بهینه سازی متفاوتی وجود دارد مانند^{۷۹} SGD و Adadelta.

به حداقل رساندن تابع هزینه با کاهش گرادیان تقریباً رایج ترین الگوریتم برای بهینه سازی شبکه‌های عصبی است. اگر تابع هزینه آنتروپی متقاطع دودویی^{۸۰} باشد و بخواهیم محاسبه کنیم که (I_i) p تا چه حد خوب می‌تواند برچسب y_i را تقریب بزند از معادله ۱۳.۲ استفاده می‌شود.

$$L = \frac{1}{|\mathcal{D}|} \sum_i^{|D|} \left(y_i \log(P(I_i)) + (1 - y_i) \log(1 - P(I_i)) \right) \quad (13.2)$$

احتمال برای یک ورودی به وزن‌های آن (θ) بستگی دارد و با (I, θ) p نمایش داده می‌شود. با توجه به θ می‌توان $L(\theta)$ را با اجرای مدل بر روی مجموعه داده به دست آورد.

⁷⁷Regularization

⁷⁸Cost function

⁷⁹Stochastic gradient descent

⁸⁰Binary cross-entropy

بکپروپگیشن^{۱۱} اساس آموزش شبکه عصبی است. این عمل تنظیم-دقیق وزن‌های یک شبکه عصبی بر اساس میزان خطا^{۱۲} در هر دوره^{۱۳} قبلی است که این امر با محاسبه مشتق‌های تابع خطا بر اساس وزن‌ها $\nabla_{\theta} L(\theta)$ در زمان آموزش امکان پذیر است. تنظیم مناسب وزن‌ها باعث کاهش میزان خطا می‌شود. در فرایند بکپروپگیشن ابتدا ورودی در سراسر شبکه انتشار داده می‌شود سپس $L(\theta)$ محاسبه شده و در نهایت این خطا از طریق تمام وزن‌ها در شبکه رو به عقب منتشر می‌شود. مشتق تابع هزینه از خروجی توسط معادله^{۱۴.۲} محاسبه می‌شود.

$$\frac{\partial L}{\partial P} = \frac{\partial \left(- (y_i \log(p) + (1-y) \log(1-P)) \right)}{\partial P} = \frac{P-y}{P(1-P)} \quad (14.2)$$

همچنین محاسبه مشتق تابع هزینه L از ورودی i به صورت معادله^{۱۵.۲} محاسبه می‌شود.

$$\frac{\partial L}{\partial i} = \frac{\partial L}{\partial P} \frac{\partial P}{\partial i} = P - y \quad (15.2)$$

همچنین محاسبه مشتق تابع هزینه بر اساس وزن‌های لایه آخر w به صورت،

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial P} \frac{\partial P}{\partial i} \frac{\partial i}{\partial w} = (P - y)a \quad (16.2)$$

می‌باشد که a در آن برابر با ترکیب خطی از ورودی‌های لایه آخر است. این کار را می‌توان به راحتی به لایه‌های قبلی تعمیم داد، بنابراین می‌توان $\nabla_{\theta} L(\theta)$ را محاسبه کرد.

۱۱.۲ شبکه‌های عصبی بازگشتی

قبل از آشنا شدن با شبکه‌های عصبی بازگشتی بهتر است مروی بر مفهوم شبکه عصبی داشته باشیم. شبکه‌های عصبی مجموعه‌ای از الگوریتم‌ها هستند که شباهت نزدیکی به مغز انسان داشته و به منظور تشخیص الگوهای طراحی شده‌اند. شبکه‌ی عصبی داده‌های حسی را از طریق ادرارک ماشینی، برچسب زدن یا خوشه بندی ورودی‌های خام تفسیر می‌کند. شبکه می‌تواند الگوهای عددی را شناسایی کند؛ این الگوها بردارهایی هستند که همه‌ی داده‌های دنیای واقعی (تصویر، صدا، متن یا سری‌های زمانی) برای تفسیر باید به شکل آن‌ها درآیند. شبکه‌های

^{۱۱}Back-propagation

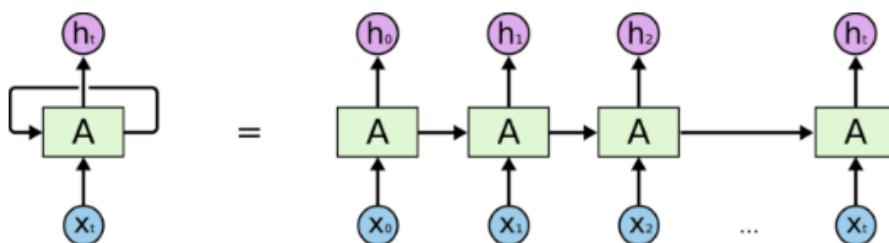
^{۱۲}Loss

^{۱۳}epoch

عصبي مصنوعی از تعداد زیادي مؤلفه‌ی پردازشی (نورون) تشکیل شده‌اند که اتصالات زیادي بینشان وجود دارد و برای حل یک مسئله با یکدیگر همکاری دارند. شبکه‌ی عصبی مصنوعی معمولاً تعداد زیادي پردازشگر دارد که به صورت موازی کار می‌کنند و در ردیف‌هایی کنار هم قرار می‌گیرند. ردیف اول، همچون عصب‌های بینایی انسان در پردازش بصری، اطلاعات ورودی‌های خام را دریافت می‌کند. سپس هر کدام از ردیف‌های بعدی، به جای ورودی خام، خروجی ردیف قبلی را دریافت می‌کنند؛ در پردازش بصری نیز نورون‌هایی که از عصب بینایی فاصله دارند، سیگنال را از نورون‌های نزدیک‌تر می‌گیرند. ردیف آخر خروجی کل سیستم را تولید می‌کند.

۱۰.۱۱.۲ شبکه عصبی بازگشتی چیست؟

شبکه‌ی عصبی بازگشتی شکلی از شبکه‌ی عصبی پیشخور است که یک حافظه‌ی داخلی دارد. شبکه عصبی بازگشتی ذاتاً بازگشتی است، زیرا یک تابع یکسان را برای همه‌ی داده‌های ورودی اجرا می‌کند، اما خروجی داده‌ی (ورودی) فعلی به محاسبات ورودی قبلی بستگی دارد. خروجی بعد از تولید، کپی شده و مجدداً به شبکه‌ی بازگشتی فرستاده می‌شود. این شبکه برای تصمیم‌گیری، هم ورودی فعلی و هم خروجی که از ورودی قبلی آموخته شده را در نظر می‌گیرد. شبکه عصبی بازگشتی برخلاف شبکه‌های عصبی پیشخور می‌توانند از حالت (حافظه‌ی) درونی خود برای پردازش دنباله‌هایی از ورودی‌ها استفاده کنند. این خاصیت باعث می‌شود در مسائلی همچون تشخیص دست خط زنجیره‌ای یا تشخیص گفتار کاربرد داشته باشند. در سایر شبکه‌های عصبی، ورودی‌ها از یکدیگر مستقل هستند، اما در شبکه عصبی بازگشتی ورودی‌ها به هم مرتبط می‌باشند. به شکل ۱۴.۲ توجه کنید، این شبکه ابتدا X_1 را از دنباله‌ی ورودی‌ها گرفته و خروجی h_1 را تولید می‌کند که همراه با X_1 ورودی گام بعدی



An unrolled recurrent neural network.

شکل ۱۴.۲: یک نمونه بازشده شبکه عصبی بازگشتی

محسوب خواهند شد. یعنی X_1 ورودی گام بعدی هستند. به همین صورت h_1 بعدی همراه با X_1 ورودی گام بعدی خواهند بود. شبکه عصبی بازگشتی بدین طریق می‌تواند هنگام آموزش زمینه را به خاطر داشته باشد.

فرمول حالت^{۱۴} کنونی به صورت رابطه ۱۷.۲ خواهد بود که در آن،

$$h_t = f(h_{t-1}, x_t) \quad (17.2)$$

خواهد بود که در آن h_t برابر است با،

$$h_t = \tanh(W_{hh}h_{t-1} + W_{hx}x_t) \quad (18.2)$$

در این فرمول W وزن، h تکبردار نهان، W_{hh} وزن حالت نهان قبلی، W_{hx} وزن حالت ورودی کنونی و \tanh تابع فعالیت است که با استفاده از تابعی غیرخطی، خروجی را فشرده می‌کند تا در بازهی $[-1, 1]$ جای گیرند. در نهایت حالت خروجی Y_t از طریق رابطه ۱۹.۲ بدست می‌آید،

$$y_t = W_{hy}h_t \quad (19.2)$$

که در آن W_{hy} برابر وزن در حالت تولید شده را نشان می‌دهد.

۲.۱۱.۲ مزایای شبکه عصبی بازگشتی

شبکه عصبی بازگشتی می‌تواند دنباله‌ای از داده‌ها را به شکلی مدل‌سازی کند که هر نمونه وابسته به نمونه‌های قبلی به نظر برسد. شبکه عصبی بازگشتی را می‌توان با لایه‌های پیچشی نیز به کار برد تا گسترهی همسایگی پیکسلی را افزایش داد.

۳.۱۱.۲ معایب شبکه عصبی بازگشتی

- گرادیان کاهشی و مشکلات ناشی از آن
- آموزش بسیار دشوار
- ناتوانی در پردازش دنباله‌های طولانی از ورودی در صورت استفاده از تابع فعالیت ReLU یا \tanh

⁸⁴State

۴.۱۱.۲ کاربردهای شبکه عصبی بازگشته

- شرح نویسی عکس^{۸۵}: شبکه عصبی بازگشته با تحلیل حالت کنونی عکس، برای شرح نویسی عکس به کار می‌رود
- پیش‌بینی سری‌های زمانی^{۸۶}: هر مسئله سری زمانی مانند پیش‌بینی قیمت یک سهام در یک ماه خاص، با شبکه عصبی بازگشته قابل انجام است
- پردازش زبان طبیعی^{۸۷}: کاوش متن و تحلیل احساسات می‌تواند با استفاده از شبکه عصبی بازگشته انجام شود
- ترجمه ماشینی^{۸۸}: شبکه شبکه عصبی بازگشته می‌تواند ورودی خود را از یک زبان دریافت و آن را به عنوان خروجی به زبان دیگری ترجمه کند

۵.۱۱.۲ انواع شبکه عصبی بازگشته

به طور کلی ۴ نوع شبکه عصبی بازگشته داریم:
 یک به یک (one to one): این نوع شبکه عصبی به عنوان شبکه عصبی وانیلی نیز شناخته می‌شود و برای مسائل یادگیری ماشین که یک ورودی و یک خروجی دارند به کار می‌رود.
 یک به چند (one to many): این شبکه عصبی بازگشته دارای یک ورودی و چند خروجی است. یک نمونه آن، شرح نویسی عکس است.

چند به یک (many to one): این نوع از شبکه عصبی بازگشته، دنباله ایی از ورودی‌ها را می‌گیرد و یک خروجی تولید می‌کند. تحلیل احساسات مثال خوبی از این نوع شبکه است که یک جمله را به عنوان ورودی می‌گیرد و آن را با احساس مثبت یا منفی طبقه بندی می‌کند.

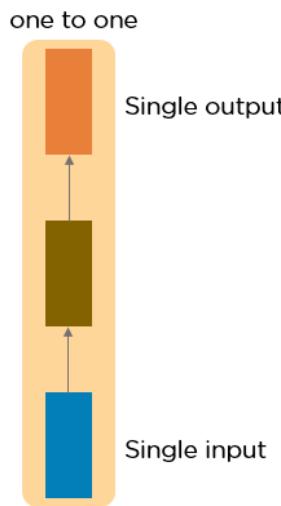
چند به چند (many to many): دنباله ایی از ورودی‌ها را می‌گیرد و دنباله ایی از خروجی‌ها را تولید می‌کند. ترجمه ماشینی نمونه ایی از این نوع شبکه است.

⁸⁵Image Captioning

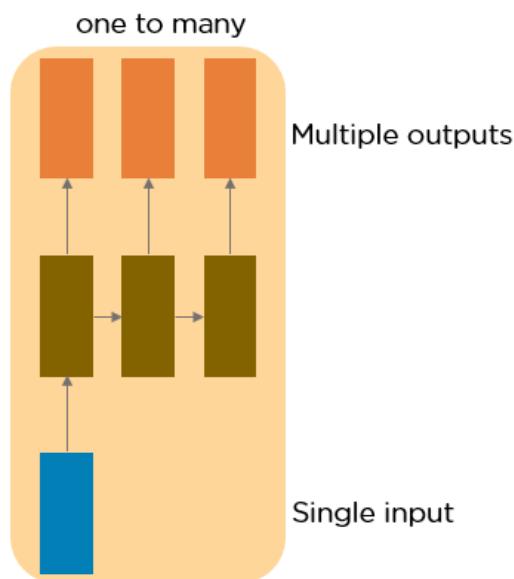
⁸⁶Time Series Prediction

⁸⁷Natural Language Processing

⁸⁸Machine Translation



شکل ۱۵.۲: ساختار شبکه عصبی بازگشتی یک به یک

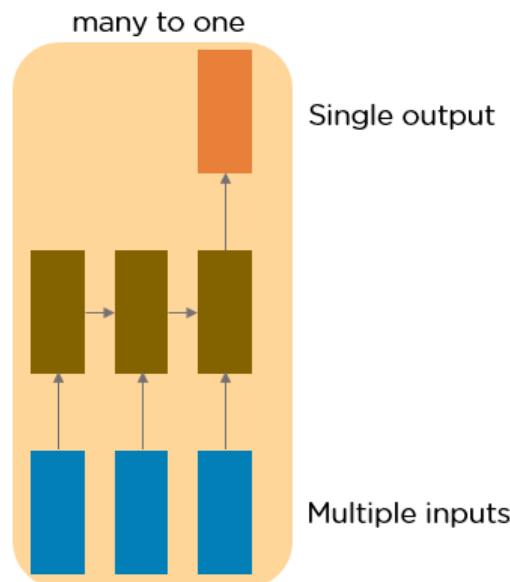


شکل ۱۶.۲: ساختار شبکه عصبی بازگشتی یک به چند

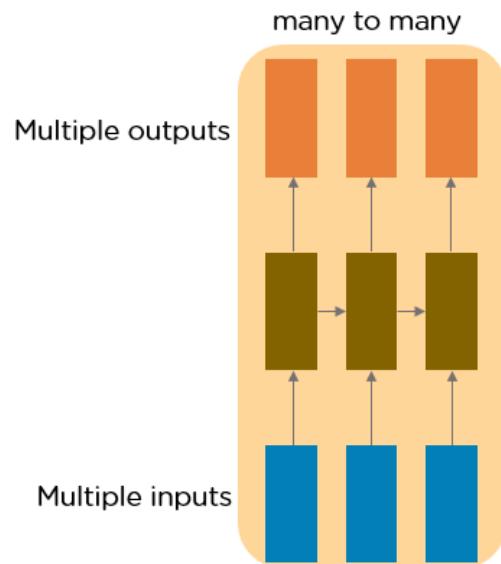
۶.۱۱.۲ حافظه‌ی کوتاه‌مدت بلند (LSTM)

شبکه‌های حافظه‌ی کوتاه‌مدت بلند^{۸۹} یا LSTM نسخه‌ی تغییریافته‌ای از شبکه‌های عصبی بازگشتی هستند که یادآوری داده‌های گذشته در آن‌ها تسهیل شده است. مشکل گرادیان کاهشی که در شبکه عصبی بازگشتی وجود داشت نیز در این شبکه‌ها حل شده است. شبکه‌های LSTM برای مسائل رده‌بندی، پردازش و پیش‌بینی سری‌های

⁸⁹Long Short Term Memory (LSTM)



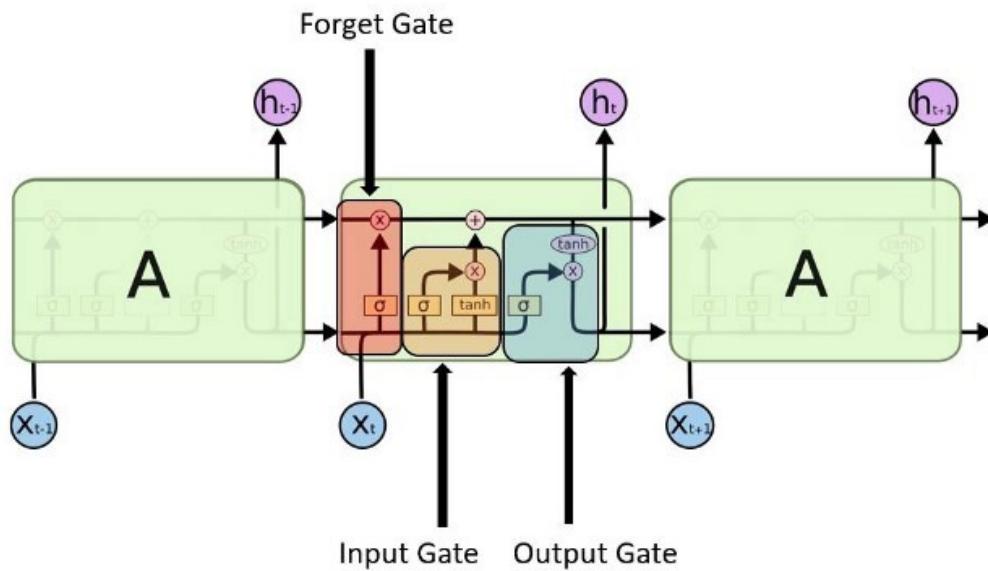
شکل ۱۷.۲: ساختار شبکه عصبی بازگشتی چند به یک



شکل ۱۸.۲: ساختار شبکه عصبی بازگشتی چند به چند

زمانی با استفاده از برقسپهای زمانی مدت‌های نامعلوم مناسب هستند. این شبکه‌ها مدل را با استفاده از انتشار رو به عقب آموزش می‌دهند.

همان‌طور که در شکل ۱۹.۲ نمایش داده شده است، در یک شبکه‌ی LSTM سه دریچه وجود دارد:



شکل ۱۹.۲: ساختار LSTM

دریچه‌های LSTM

۱) دریچه‌ی ورودی: با استفاده از این دریچه می‌توان دریافت کدام مقدار از ورودی را باید برای تغییر حافظه به کار برد.تابع سیگموید تصمیم می‌گیرد مقادیر بین ۰ و ۱ اجازه‌ی ورود دارند و تابع \tanh با ضریب دهی (بین ۱ - تا ۱+) به مقادیر، در مورد اهمیت آن‌ها تصمیم می‌گیرد.

$$\begin{aligned} i_t &= \sigma(W_i \cdot [h_{t-1}, x_t] + b_i) \\ \tilde{C}_t &= \tanh(W_C \cdot [h_{t-1}, x_t] + b_C) \end{aligned} \quad (۲۰.۲)$$

۲) دریچه‌ی فراموشی: از طریق این دریچه می‌توان جزئیاتی را که باید از بلوک حذف شوند، تشخیص داد. تصمیم‌گیری در این مورد بر عهده‌ی تابع سیگموید است. این تابع با توجه به حالت قبلی h_{t-1} و ورودی محتوا X_t ، عددی بین ۰ تا ۱ به هر کدام از اعداد موجود در حالت سلولی C_{t-1} اختصاص می‌دهد؛ نشان‌دهنده‌ی حذف

آن عدد و ۱ به معنی نگه داشتن آن است.

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) \quad (21.2)$$

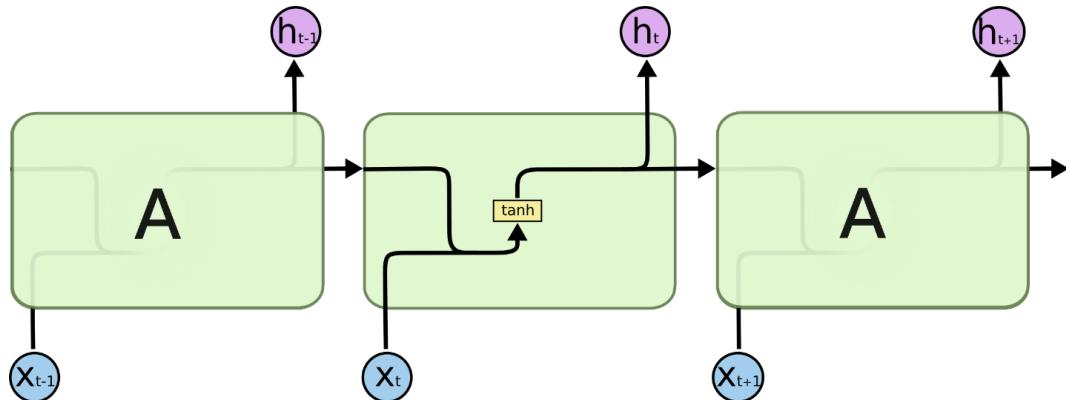
۳) دریچه‌ی خروجی: ورودی و حافظه‌ی بلوک برای تصمیم‌گیری در مورد خروجی مورد استفاده قرار می‌گیرند. تابع سیگموئید تصمیم می‌گیرد مقادیر بین ۰ و ۱ اجازه‌ی ورود دارند و تابع \tanh با ضریب‌دهی (بین -۱ تا +۱) به مقادیر و ضرب آن‌ها در خروجی تابع سیگموئید در مورد اهمیت آن‌ها تصمیم‌گیری می‌کند.

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o) \quad (22.2)$$

$$h_t = o_t * \tanh(C_t)$$

در حقیقت هدف از طراحی شبکه‌های LSTM، حل کردن مشکل وابستگی بلندمدت بود. به این نکته مهم توجه کنید که به یاد سپاری اطلاعات برای بازه‌های زمانی بلند مدت، رفتار پیش‌فرض و عادی شبکه‌های LSTM است و ساختار آن‌ها به صورتی است که اطلاعات خیلی دور را به خوبی یاد می‌گیرند که این ویژگی در ساختار آن‌ها نهفته است.

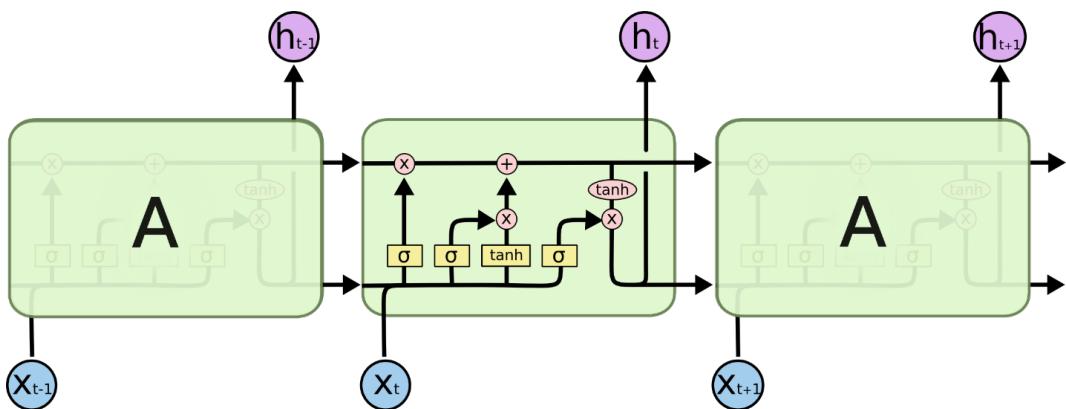
همه شبکه‌های عصبی بازگشتی به شکل دنباله‌ای (زنگیره‌ای) تکرار شونده از مازول‌های (واحدهای) شبکه‌های عصبی هستند. در شبکه‌های عصبی بازگشتی استاندارد، این مازول‌های تکرار شونده ساختار ساده‌ای دارند، برای مثال تنها شامل یک لایه تائزانتِ هایپربولیک (\tanh) هستند.



شکل ۲۰.۲: مازول‌های تکرار شونده در شبکه‌های عصبی بازگشتی استاندارد فقط دارای یک لایه هستند.

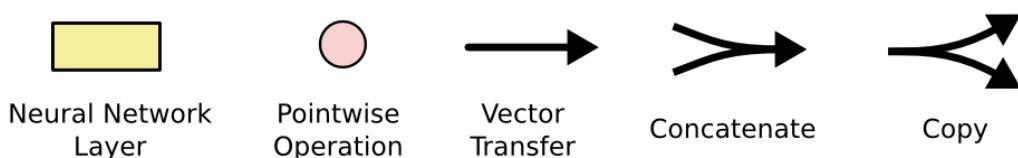
شبکه‌های LSTM نیز چنین ساختار دنباله یا زنجیره‌مانندی دارند ولی مازول تکرار شونده ساختار متفاوتی دارد. به جای داشتن تنها یک لایه شبکه عصبی، ۴ لایه دارند که طبق ساختار ویژه‌ای با یکدیگر در تعامل و ارتباط

هستند. در ادامه قدم به قدم ساختار شبکه‌های حافظه‌ی کوتاه‌مدت بلند را توضیح خواهیم داد. اما در ابتدا معنی



شکل ۲۱.۲: مراحلهای تکرار شونده در LSTM‌ها دارای ۴ لایه هستند که با هم در تعامل می‌باشند.

هر کدام از شکل و علامت‌هایی را که از آن‌ها استفاده خواهیم کرد توضیح می‌دهیم. در شکل ۲۲.۲، هر خط



شکل ۲۲.۲: اشکال از راست به چپ به ترتیب برابر هستند با: کپی کردن، وصل کردن، بردار انتقال، عملیات نقطه به نقطه، یک لایه‌ی شبکه عصبی.

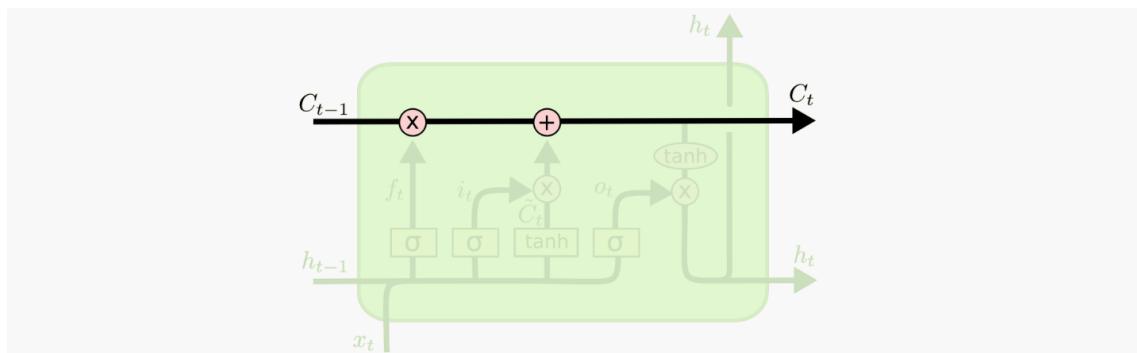
یک بردار را به صورت کامل از خروجی یک گره به ورودی گره دیگر انتقال می‌دهد. دایره‌های صورتی نمایش دهنده عملیات‌های نقطه به نقطه مانند «جمع کردن دو بردار» هستند. مستطیل‌های زرد، لایه‌های شبکه‌های عصبی هستند که شبکه پارامترهای آن‌ها را یاد می‌گیرد. خط‌هایی که با هم ادغام می‌شوند نشان‌دهنده الحاق^{۹۰} و خط‌هایی که چند شاخه می‌شوند نشان‌دهنده‌ای این موضوع است که محتوای آن‌ها کپی و به بخش‌های مختلف ارسال می‌شود.

عنصر اصلی LSTM‌ها سلول حالت^{۹۱} است که در حقیقت یک خط افقی است که در بالای شکل ۲۳.۲ قرار دارد. سلول حالت را می‌توان به صورت یک تسمه نقاله تصور کرد که از اول تا آخر دنباله یا همان زنجیره با تعاملات خطی جزئی در حرکت است (یعنی ساختار آن بسیار ساده است و تغییرات کمی در آن اتفاق می‌افتد).

LSTM این توانائی را دارد که اطلاعات جدیدی را به سلول حالت اضافه یا اطلاعات آن را حذف کنید. این کار

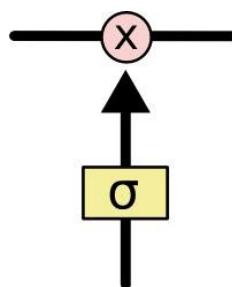
⁹⁰Concatenation

⁹¹Cell state



شکل ۲۳.۲: سلول حالت در مازول LSTM

توسط ساختارهای دقیقی به نام دروازه‌ها^{۹۲} انجام می‌شود. دروازه‌ها راهی هستند برای ورود اختیاری اطلاعات. آن‌ها از یک لایه شبکه عصبی سیگموید به همراه یک عملگر ضرب نقطه به نقطه تشکیل شده‌اند.



شکل ۲۴.۲: نمایی از نحوه تاثیر و ورود اطلاعات به سلول حالت

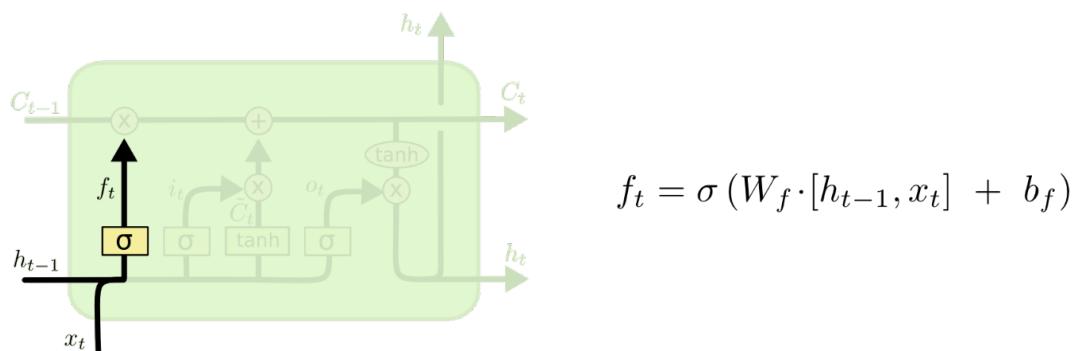
خروجی لایه سیگموید عددی بین صفر و یک است، که نشان می‌دهد چه مقدار از ورودی باید به خروجی ارسال شود. مقدار صفر یعنی هیچ اطلاعاتی نباید به خروجی ارسال شود در حالی که مقدار یک یعنی تمام ورودی به خروجی ارسال شود!

LSTM دارای ۳ دروازه مشابه برای کنترل مقدار سلول حالت است که در ادامه به بررسی قدم به قدم آن‌ها از لحظه ورود تا خروج اطلاعات خواهیم پرداخت.

قدم اول در LSTM تصمیم در مورد اطلاعاتی است که می‌خواهیم آن‌ها را از سلول حالت پاک کنیم. این تصمیم توسط یک لایه سیگموید به نام «دوازه فراموشی»^{۹۳} انجام می‌شود. این دروازه با توجه به مقادیر h_{t-1} و x_t ، برای هر عدد، مقدار صفر یا یک را در سلول حالت C_{t-1} به خروجی می‌برد. مقدار یک یعنی به صورت کامل مقدار حال حاضر سلول حالت C_{t-1} را به C_t انتقال داده شود و مقدار صفر یعنی به صورت کامل اطلاعات سلول حالت

⁹²Gate⁹³Forget gate

کنونی C_{t-1} را پاک شود و هیچ مقداری از آن به C_t برد نشود. باید به مثال قبلی مان که یک مدل زبانی‌ای بود که در آن تلاش داشتیم کلمه بعدی را بر اساس همه کلمه‌های قبلی حدس بزنیم، برگردیم. در چنین مسأله‌ای، سلول حالت ممکن است در بردارنده جنسیت فاعل کنونی باشد، که با توجه به آن می‌توانیم تشخیص دهیم از چه ضمیری باید استفاده کنیم. زمانی که یک فاعل جدید در جمله ظاهر می‌شود، می‌بایست جنسیت فاعل قبلی حذف شود.



شکل ۲۵.۲: قدم اول در پاک کردن اطلاعات از سلول حالت در وضعیت ورودی

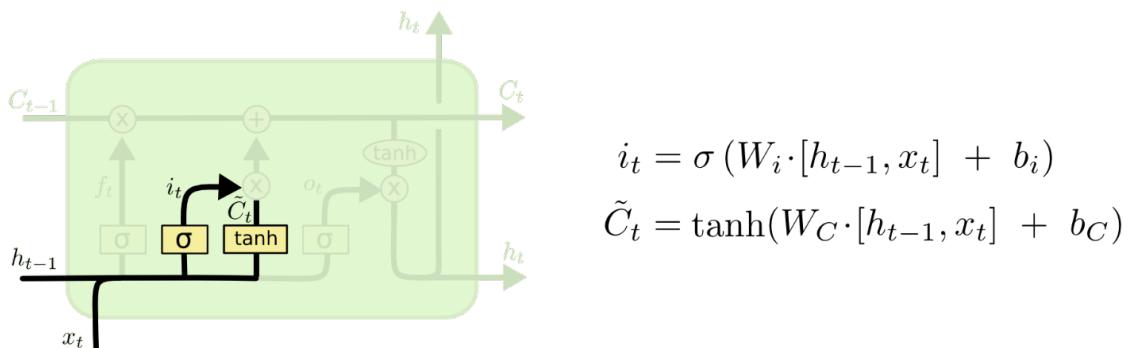
قدم بعدی است که تصمیم بگیریم چه اطلاعات جدیدی را می‌خواهیم در سلول حالت ذخیره کنیم. این تصمیم دو بخشی است. ابتدا یک لایه سیگموید به نام دروازه ورودی^{۹۴} داریم که تصمیم می‌گیرد چه مقادیری به روز خواهند شد. مرحله بعدی یک لایه تانژانت هایپربولیک است که برداری از مقادیر به نام \tilde{C}_t می‌سازد که می‌توان آن را به سلول حالت اضافه کرد. در مرحله بعد، ما این دو مرحله را با هم ترکیب می‌کنیم تا مقدار سلول حالت را به روز کنیم.

در مثال مدل زبانی‌ای که پیش‌تر داشتیم، قصد داریم جنسیت فاعل جدید را به سلول حالت اضافه کنیم تا جایگزین جنسیت فاعل قبلی شود که در مرحله قبلی تصمیم گرفتیم آن را فراموش کنیم.

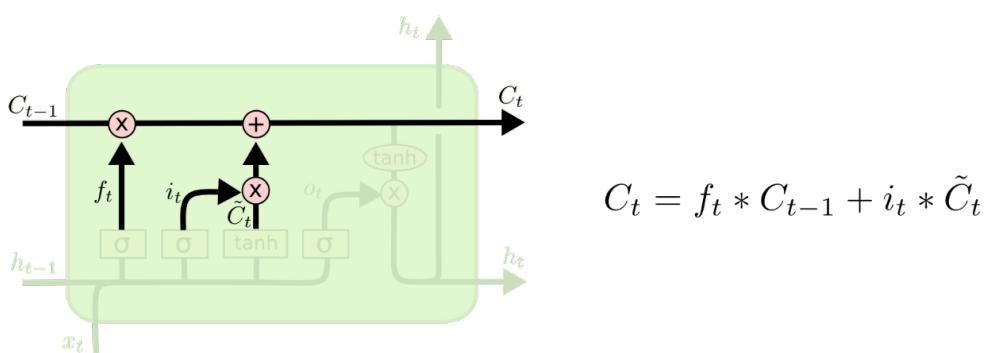
حال زمان آن فرا رسیده است که سلول حالت قدیمی یعنی C_{t-1} را سلول حالت جدید یعنی C_t به روز کنیم. در مراحل قبلی تصمیم گرفته شد که چه کنیم و در حال حاضر تنها لازم است تصمیماتی را که گرفته شد عملی کنیم. ما مقدار قبلی سلول حالت را در f_t ضرب می‌کنیم که یعنی فراموش کردن اطلاعاتی که پیش‌تر تصمیم گرفتیم آن‌ها را فراموش کنیم. سپس $i_t * \tilde{C}_t$ را به آن اضافه می‌کنیم. در حال حاضر مقادیر جدید سلول حالت با توجه به تصمیماتی که پیش‌تر گرفته شده بود بدست آمدند. در مثال مدل زبانی، اینجا دقیقاً جائی است که اطلاعاتی که در مورد جنسیت قبلی داشتیم را دور می‌ریزیم و اطلاعات جدید را اضافه می‌کنیم.

در نهایت باید تصمیم بگیریم قرار است چه اطلاعاتی را به خروجی ببریم. این خروجی با در نظر گرفتن مقدار سلول حالت خواهد بود، ولی از فیلتر مشخصی عبور خواهد کرد. در ابتدا، یک لایه سیگموید داریم که تصمیم

^{۹۴}Input gate



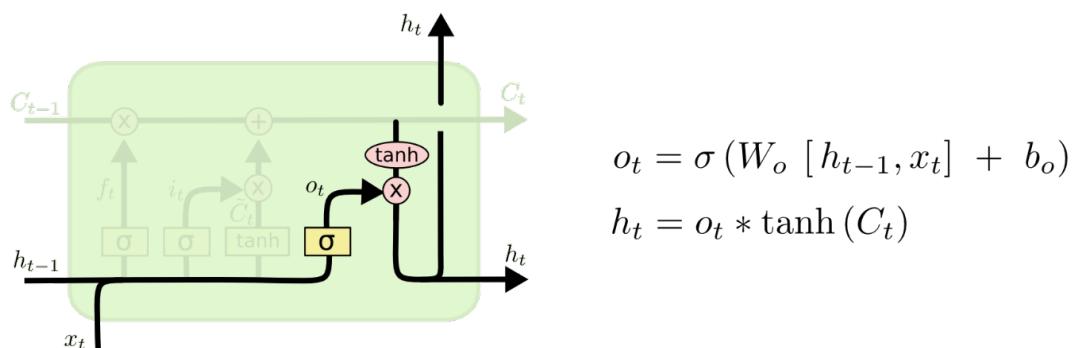
شکل ۲۶.۲: قدم دوم در اضافه کردن اطلاعات جدید به سلول حالت



شکل ۲۷.۲: بهروز رسانی اطلاعات در سلول حالت

می‌گیرد چه بخشی از سلول حالت قرار است به خروجی برده شود. سپس مقدار سلول حالت (پس از بهروز شدن در مراحل قبلی) را به یک لایه تانژانت هایپربولیک (تا مقادیر بین -1 و $+1$ باشند) می‌دهیم و مقدار آن را در خروجی لایه سیگموید قبلی ضرب می‌کنیم تا تنها بخش‌هایی که مدنظرمان است به خروجی برود.

در مثال مدل زبانی، با توجه به اینکه تنها فاعل را دیده است، در صورتی که بخواهیم کلمه بعدی را حدس بزنیم، ممکن است بخواهد اطلاعاتی در ارتباط با فعل را به خروجی ببرد. برای مثال ممکن است اینکه فاعل مفرد یا جمع است را به خروجی ببرد، که ما با توجه به آن بدانیم فعل به چه فرمی خواهد بود.



شکل ۲۸.۲: قدم نهایی برای تولید خروجی مازول LSTM

۱۲.۲ یادگیری تقویتی

۱.۱۲.۲ مقدمه و بیشینه تاریخی

ادوارد ثورندایک^{۹۵} پدر روانشناسی مدرن در سال ۱۸۷۴ میلادی در ایالت ماساچوست آمریکا متولد شد. وی در اوایل قرن ۲۰ میلادی آزمایشی انجام داد که باعث ارائه قانون اثر شد. او برای این آزمایش، گربه ای را در جعبه ای موسوم به جعبه معملاً قرار داد. هر کوشش درستی، از این گربه برای نجات از جعبه صورت می‌گرفت، باعث میشد ثورندایک به عنوان پاداش به او غذا بدهد. به تدریج گربه به کارهای درست خود پی برد و آنها را تکرار کرد، تا جایی که دیگر هیچ کار اشتباہی نمی‌کرد و بالاخره موفق به خروج از جعبه شد. ثورندایک در سال ۱۹۱۲ به ریاست انجمن روانشناسان، در سال ۱۹۱۷ به عضویت انجمن علوم، در سال ۱۹۳۴ به ریاست انجمن علوم پیشرفته نایل آمد و در سال ۱۹۴۷ در سن ۷۴ سالگی، بدرود حیات گفت. در سال ۲۰۰۲ رتبه ای از برترین روانشناسان تاریخ ارائه شد که ثورندایک جزء ۱۵ روانشناس برتر تاریخ قرار گرفت. می‌توان مهمن ترین کشف وی را، اثبات وجود یادگیری تقویتی در روانشناسی دانست.

شاید ریچارد بلمن^{۹۶} (مختصر الگوریتم بلمن-فورد) را بتوان اولین کسی دانست که یادگیری تقویتی را وارد هوش مصنوعی ساخت. در اوایل دهه ۱۹۵۰ بلمن مسئله ای با عنوان «کنترل بهینه» را مطرح ساخت که با استفاده از روش‌های پویا در برنامه ریزی پویا کنترل کننده‌ها را به سمت نتیجه بهینه رهنمون می‌شد. در اواخر دهه ۵۰ میلادی مینسکی در پایان نامه دکتری خود روش‌های محاسبات آزمون و خطاب توسط مفهوم یادگیری تقویتی را مطرح نمود و الگوریتم‌های یادگیری تقویتی را پایه ریزی کرد. در کل دهه ۵۰ میلادی را میتوان دهه تشکیل الگوریتم‌های محاسباتی اولیه یادگیری تقویتی دانست. در دهه ۶۰ میلادی اولین کابرد‌های یادگیری تقویتی

⁹⁵Edward Thorndike

⁹⁶Richard E. Bellman

به وقوع پیوستند. در اولین تلاش‌ها فارلی و کلارک، از یادگیری تقویتی برای تشخیص الگو استفاده کردند بدین صورت که هر بار برنامه نتیجه بهتری به دست می‌آمد او را تشویق می‌کردند. در اواخر دهه ۶۰ میلادی، یادگیری نظارتی از یادگیری تقویتی، مشتق شد. در یادگیری نظارتی طراح نتیجه نهایی را در دست دارد و از هوش مصنوعی می‌خواهد هر بار مسیر بین ورودی و نتیجه را طراحی کرده و هر بار که برنامه، مسیر بهتری به دست می‌آورد، تشویق می‌شود. همچنین طراح نظارت مستقیم بر عملکرد عامل دارد.

فصل ۳

روش‌های پیشین

۱.۳ مقدمه

در فصل گذشته به معرفی مفاهیم و موضوعات مرتبط با این حوزه پرداخته شد. در ادامه در این فصل با توجه به اطلاعاتی که کسب کرده‌اید به معرفی و بررسی روش‌هایی که مرتبط با موضوع این پایان‌نامه است پرداخته خواهد شد و نتایج آن‌ها را برای فرض‌های و داده‌های ورودی خود مشاهده خواهیم نمود. در این بین تا جایی که ممکن باشد به بررسی نقاط قوت و ضعف آن‌ها نیز خواهیم پرداخت و در انتهای این فصل یک جدول مقایسه بین روش‌هایی که تا به حال معرفی شده‌اند را ارائه خواهیم داد.

۲.۳ روش ساخت درخت تکاملی با استفاده از داده‌های توالی‌یابی تک

سلولی

سرطان نامی است که به مجموعه‌ای از بیماری‌ها اطلاق می‌شود که از تکثیر مهار نشده سلول‌ها پدید می‌آیند. تحقیقات انجام شده نشان می‌دهد که سرطان در واقع یک فرآیند تکاملی از جهش‌های ژنتیکی، شامل حذف و تغییر تعداد کپی، حذف و تغییر تک نوکلئوتید‌ها، بازسازی و جایگزینی ژن‌ها در سلول‌های توموری است. در واقع تومور زمانی ایجاد می‌شود که یک سلول جهش یافته بتواند با عبور از سیستم دفاعی بدن زندگی کرده و تکثیر شود به گونه‌ای که نسبت مرگ به تولید آن گونه ایجاد شده بسیار کوچک تر از $1 \ll \alpha$

باشد. با پیشرفت تومور، ناهنجاری‌های ژنتیکی مختلف منجر به افزایش گروه‌های جمعیتی ناهمگنی به نام کلون می‌شود. فرآیند تکاملی همه این کلون‌ها را می‌توان با یک درخت فیلوزنی و آنالیز فیلوزنیتیکی از چندین کلون سلولی سرطانی مدل‌سازی کرد که می‌تواند مطالعه انواع تومور را تسهیل کند. ساختار و الگو‌های درون این درخت میزان وابستگی بین گونه‌های خاص را با توجه به تعداد و فواصل بین اجداد مشترک‌شان تعیین می‌کند. درخت‌های فیلوزنی عملکرد بارزی در توصیف فرآیند توسعه تومور دارند که بهتر از دیگر الگوریتم‌های مشابه عمل می‌کنند. تحلیل توپولوژی درخت درخت پیشرفت تومور نشان می‌دهد که مسیر توسعه تومور در طول مراحل مختلف تشکیل تومور، تا حد زیادی تغییر می‌کند و البته نتایج مبتنی بر درخت بهتر از نتایج داده‌های بدست آمده از طریق روش‌های دیگر در تشخیص تومور می‌باشد. در حال حاضر ظهور تکنولوژی‌های بر پایه دی‌ان‌ای یک سلول منفرد، با هدف افزایش دانش از جنبه‌های مختلف بیولوژی سرطان، شامل بررسی زیرساخت کلونال، ردیابی تکامل تومور، شناسایی زیرکلون‌های نادر و درک ریزمحیط‌های سرطانی در پیشرفت تومور، به یاری محققان این حوزه آمده و بالاترین وضوح را از تاریخچه سرطان (درخت فیلوزنی) فراهم کرده است. در واقع از آنجایی که در روش‌توالی‌یابی تک سلولی گونه‌های مختلف از ابتدا از هم جدا می‌شوند، از نقطه منظر از دست دادن تنوع در زیرجمعیت بافت مورد آزمایش نداریم و به همین دلیل دقت این روش نیز نسبت به روش‌انبوه بالاتر می‌باشد. در کنار مزایا این روش، معایبی چون، هزینه بالا، از دست دادن سلول‌ها، جهش ثانویه در هنگام کشت، از دست دادن میزان فراوانی درون تومور حقیقی و زمان‌گیر بودن فرآیند نمونه گیری اشاره کرد.

در ابتدا استفاده از روش‌های توالی‌یابی انبوه بدلیل اینکه حجم بالایی از اطلاعات در اثر این توالی‌یابی ایجاد می‌شود، از محبوبیت بیشتری برخوردار بود اما با پیشرفت تکنولوژی و ظهور روش‌های نوینی چون توالی‌یابی تک‌سلولی این مهم دچار تغییر شد. در روش توالی‌یابی انبوه، نمونه‌برداری بر روی تعداد بسیار زیادی سلول (از محدوده‌ی هزار تا میلیون سلول) صورت می‌گرفت و حجم بالای داده‌ها و امکان تفکیک پایین نواحی ناهمگن، اطلاعات کافی از ساختار درون تومور و ناهمگنی‌های درون توموری بدست نمی‌داد. در مقابل، در روش توالی‌یابی تک‌سلولی، اگر‌چه میزان هزینه نمونه‌برداری افزایش قابل‌توجهی داشت و یا میزان اطلاعات از دست رفته و نویز موجود در داده‌های توالی‌یافته بالا بود [۲۰، ۲۳].

اما در این روش رزولوشن یا قدرت تفکیک جهش‌های گوناگون از یکدیگر بسیار بالا بود و تشخیص نواحی ناهمگنی تومور و تفکیک زیرکلون‌ها از یکدیگر بسیار راحت‌تر از گذشته صورت می‌گرفت. در این فصل روش‌هایی را مورد بررسی قرار می‌دهیم که ساخت درخت فیلوزنی تومور و ناهمگنی‌های درون توموری را از طریق داده‌های توالی‌یابی تک‌سلولی مورد بررسی قرار می‌دهد.

۱۰.۲.۳ مدل کیم و سایمون [۴۷]

این مدل در سال ۲۰۱۴ با تمرکز بر ساخت درخت فیلوزنی^۱ از طریق رابطه ترکیبی میان جهش‌های ایجاد شده در داده‌های توالی‌یابی تک‌سلولی دی‌ان‌ای ارائه گردید. بررسی رابطه‌ی ترتیبی هر یک از جهش‌های رخ داده با یکدیگر، این امکان را فراهم می‌آورد تا اطلاعاتی در مورد نحوه تشکیل کلون‌ها و ترتیب زمانی رخدادن جهش‌های گوناگون بدست آید. همچنین امکان محاسبه نسبت زمانی سپری شده میان جهش‌های اولیه موجود در داده‌های توالی‌یابی تک‌سلولی تا نزدیک ترین جد مشترک وجود دارد. استنباط درخت فیلوزنی از طریق لگوریتم کیم و سایمون، بر مبنای منطق بیزی است، یعنی از این منطق به منظور تعیین رابطه ترتیبی بین هر دو جهش گوناگون استفاده شده است. در ادامه مقدار بیشینه درستنمایی درخت استنباط شده بر مبنای احتمال ترتیبی دو به‌دوی بین هر دو جهش مختلف در دو جایگاه از یک دنباله، که از طریق ژنولوژی متفاوت با جهش‌های گوناگون در گره‌های درخت به هم مرتبط می‌شوند، محاسبه می‌شود. سرانجام مقادیر بیشینه‌ی احتمالات با شرط کمینه کردن میزان تفاوت با داده‌های مشاهده شده محاسبه می‌گردد.

از نکات قوت این الگوریتم در نظر گرفتن خطای توالی‌یابی و ترک آلل^۲ است. این عدم قطعیت در داده‌ها از طریق محاسبه ماکزیمم درست‌نمایی ترتیبی هر یک از جهش‌ها بدست خواهد آمد. به عنوان مثال در نظر بگیرید که هفت زوج مرتب از جهش‌های یک دی‌ان‌ای موجود است. برای سادگی بیشتر مولفه اول را با x و مولفه دوم را با y نشان داده‌می‌شود. داده‌های نمونه‌گیری شده از این دی‌ان‌ای در جدول زیر نشان داده شده‌است:

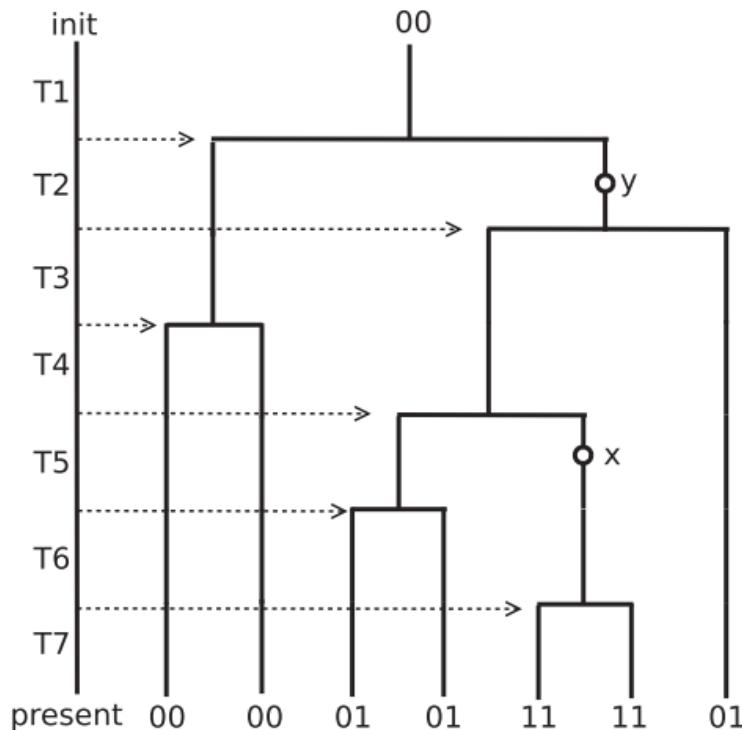
	۷	۶	۵	۴	۳	۲	۱	Sample
X mutation	۰	۱	۱	۰	۰	۰	۰	
Y mutation	۱	۱	۱	۱	۱	۰	۰	

در این جدول صفر بیانگر عدم وجود جهش و یک بیانگر وجود جهش است. تعداد رخداد جهش‌ها با فرض عدم وجود خطأ در توالی‌یابی داده‌ها، برابر یک در نظر گرفته می‌شود، یعنی در هر موقعیت تنها یکبار جهش رخ داده است. همچنین ترتیب زمانی رخداد جهش‌ها یک ترتیب جزئی است، به این معنی که زوج (۱،۱) بیانگر این است که یا جهش x مقدم بوده است یا جهش y . زوج (۰،۰) بیانگر آن است که جهش x وجود نداشته است ولی جهش y وجود داشته و با فرض اینکه هیچ جهشی از بین نمی‌رود، در نتیجه می‌توان استنباط کرد که y نسبت به x قدیمی‌تر است و به عنوان یکی از اجداد x در درخت فیلوزنی تومور قرار می‌گیرد. در نتیجه با استفاده از جدول داده‌های نمونه‌برداری شده، استنباط یک رابطه زمانی میان جهش‌های صورت گرفته امکان

¹Phylogeny tree

²Allele dropout

پذیر است. شکل ۱.۳ یک درخت فیلوزنیک تومور را نشان می‌دهد که از داده‌های جدول بالا استنباط شده است. در این همه هفت نمونه به عنوان برگ‌های درخت مشاهده می‌شود و ریشه درخت زوج (۰،۰) می‌باشد به این معنی که در ابتدا هیچ جهشی رخ نداده است. محور عمودی بیانگر سیر زمانی تکامل تومور است که به تعداد نمونه‌ها تقسیم شده است.



شکل ۱.۳: عنواننتنتنتنتنتنتنت

برای استنباط درخت فیلوزنی تومور، الگوریتم کیم و سایمون از سه بخش اصلی تشکیل شده است. طبق قضیه بیز برای محاسبه هر یک از این سه احتمال به مقادیر درست‌نمایی نیاز داریم. مقدار احتمال رخداد طبق رابطه زیر محاسبه می‌گردد:

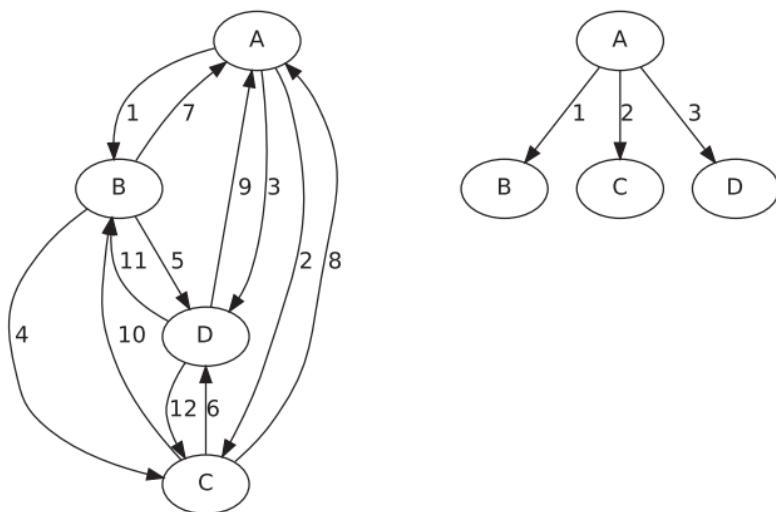
$$P(x) = yyyyyyyyyyyyyyyyyyyyyyyyyyyyy$$

طبق این رابطه و با توجه به اینکه رابطه زمانی میان جهش‌های x و y دارای ۳ حالت

$$x \not\leftrightarrow y \quad y \rightarrow x, \quad x \rightarrow y$$

است، مقدار احتمال محاسبه شده از رابطه فوق به ازای یکی از این سه حالت بیشینه است و به ازای آن حالت یک مسیر جهت‌دار در درخت فیلوزنی قرار خواهد گرفت. طبق آنچه گفته شد یک گرف جهت دار فیلوزنی

بلغه مشابه آنچه در شکل زیر نشان داده شده است استنبط خواهد شد. در نهایت از این گراف جهت دار، یک درخت به طوری که روابط میان چهش-ها از آن استنباط شود ساخته خواهد شد.



شكل ٢.٣: عنوان

در ابتدا یال‌های گرف از طریق رابطه زیر وزن دهی می‌شود:

$-loqzzzzzzzzzzzzzzzzzzzzzzzzzz$

که در آن ($x \sim y$) رابطه بین جهش-های x و y است و D نمونه یا سمپل-های موجود در داده است.
بهترین درخت \hat{T} از طریق کمینه-کردن وزن-های گراف بدست می-آید.

$$y = xxxxxxxxxxxxxxxxxxxxxxx$$

در شکل بالا محتمل-ترین درخت فیلوزنی با بیشینه درست‌نمایی بر اساس اطلاعات نمونه-برداری شده بدست می-آید. در این شکل گراف اولیه و درخت متناظر آن مشاهده می-شود. مجموع همه وزن-ها در درخت نهایی، با شرط کمینه سازی برای هفت است که این مقدار کمترین مقدار ممکن است.

یاگاه داده: ۲۰۲۳

در این مقاله از پایگاه داده تولی-یابی تک سلولی هو- و همکاران [۴۳] استفاده شده است. این مجموع داده از توالی-یابی تک-سلولی دی ان ای نمونه-های توموری یک نوع خاص از سرطان خون^۳ جمع-آوری شده است. این مجموعه داده شامل ۵۸ سلول منفرد و ۱۸ نوع حشر، بکتا است. اطلاعات کاما، در مورد این پایگاه

³Thrombocythemia

داده از جمله، نام و نوع جهش‌های موجود در دیتابیس، نوع روش نمونه‌برداری و اطلاعاتی دیگر در پایگاه داده COSMIC در دسرس عموم قرار دارد. ماتریس ژنوتاپی این پایگاه داده شامل سه مقدار صفر، یک و دو می‌باشد که در آن صفر بیانگر عدم وجود جهش، یک بیانگر جهش هتروزیگوت و دو نمایانگر جهش هموزیگوت است. یکی از معایب این پایگاه داده نرخ بالای خطای توالی‌یابی تک‌سلولی و بالا-بودن نرخ داده‌های از دست رفته (در حدود ۴۵ درصد کل داده‌ها) می‌باشد. همین امر سبب می‌شود تنوع درخت فیلوزنی نسبت داده شده به این پایگاه داده زیاد باشد. در واقع با در نظر گرفتن حالت‌های مختلف روابط دو-به-دوی جهش‌های گوناگون، می‌توان درخت‌های جهشی متنوعی از داده‌ها استنباط کرد.

۳.۲.۳ معیار ارزیابی:

ارزیابی درختهای جهشی گوناگون از طریق روش LOOCV^۴ صورت میگیرد. این روش همانند روش ارزیابی‌های متقابل^۵ با K قسمت میباشد با این تفاوت که در آن k برابر تعداد جهش‌ها (تعداد ستونهای ماتریس ژنوتاپ) می‌باشد. در هر یک از درخت‌های استنباط شده، یکبار یک جهش حذف شده و میزان دقت مدل محاسبه می‌گردد. سپس این کار برای همه جهش‌های موجود تکرار می‌شود و در نهایت میانگین دقت مدل در حالتهای مختلف محاسبه می‌شود و به عنوان دقت نهایی مدل گزارش می‌شود.

۴.۲.۳ الگوریتم [۸۴]: Bitphylogeny

این الگوریتم در سال ۲۰۱۵ ارائه شد و مانند الگوریتم کیم و سایمون از منطق بیزی بهره می‌برد. هدف این الگوریتم در کنار ساخت درخت جهشی تومور، پیدا کردن روابط بین کلونهای مختلف درون یک تومور است. در داده‌های توالی‌یابی تک‌سلولی، بدلیل کمبود میزان نمونه‌گیری و در نتیجه محتمل بودن عدم حضور گونه‌های ژنومی جهش‌یافته در نمونه‌ها، برای تشخیص ناهمگنی‌های درون توموری باید رویکرد متفاوتی را برگزید. شاید یکی از دلایلی که هنوز از داده‌های توالی‌یابی انبوه^۶ برای استنباط درخت فیلوزنی استفاده می‌شود همین باشد. در هر صورت در این مقاله سعی بر این است تا هر ۲ چالش زیر مورد بررسی قرار گیرد:

- تشخیص زیرنواحی یا کلونهای درون یک تومور

- کشف روابط تکاملی کلونهای درون یک تومور با یکدیگر

⁴Leave one out cross validation

⁵Cross validation

⁶Bulk sequencing

ماتریس ورودی (ماتریس ژنتوتایپی) این الگوریتم تعدادی سطر و ستون است که در آن سطرها بیانگر سلولها و ستونها نمایانگر انواع جهش‌های مختلف است. این ماتریس، یک ماتریس دودویی‌ها^۷ است که در آن بودن درایه z و زیانگ آن است که در سطر z ام جهشی از نوع زام وجود ندارد. متعاقباً، اگر مقدار درایه z و ز برابر یک باشد در سطر z ام جهشی از نوع زام وجود دارد.

در این مقاله برای جستجو درختی که بیشترین تطابق با داده‌های ورودی را داشته باشد از الگوریتم زنجیره مارکوف مونت کارلو استفاده می‌شود. این الگوریتم سلول‌ها با ژنتوتایپ مشابه را درون یک گروه قرار می‌دهد و به این گروه‌ها کلون گفته می‌شود. در طی دسته‌بندی سلول‌ها کلون‌هایی ایجاد می‌شود که با احتمال زیاد توموری بوده ولی در نمونه‌گیری از بافت توموری حضور نداشته‌اند. شناسایی این گونه از کلون‌ها با توجه به روند گسترش و تکامل تومور، که به مرور زمان صورت می‌گیرد، امکان پذیر است. این الگوریتم قادر است تا یک تخمین زمانی از انتقال جهش از سطوح بالای درخت فیلوزنی به سطوح پایین‌تر را محاسبه کند. در این الگوریتم از داده‌های تغییرات تک نوکلئوتید استفاده شده است اما این روش این قابلیت را دارد تا بدون در نظر گرفتن فرض مکان‌های بینهایت برای داده‌های متیلاسیون دی‌ان‌ای استفاده شود. از نکات قوت این الگوریتم می‌توان به محاسبه رخداد هر جهش در درخت فیلوزنی تومور اشاره کرد اما این مقدار احتمال بدلیل تعداد بالای داده‌های از دست رفته و نرخ بالای خطای مثبت کاذب و منفی کاذب^۸، بیش از مقدار واقعی است.

شایان ذکر است که این الگوریتم محدودیت‌های خاص خود را دارد. به عنوان مثال، در نظر گرفتن فرض مکان‌های بینهایت برای رخداد جهش‌ها و زمان محاسباتی بسیار بالا الگوریتم زنجیره مارکوف مونت کارلو برای استنباط درخت فیلوزنی از جمله این محدودیت‌ها می‌باشد. از دیگر محدودیت‌های این الگوریتم می‌توان به عدم تشخیص کلون‌های هموژنی و هتروژنی در یک نوع جهش از یکدیگر اشاره کرد. منظور از کلون‌های هموژنی در یک جهش معین آن است که اجزای تشکیل دهنده آن با توزیع یکنواخت در کنار یکدیگر قرار گرفته‌اند و این بدان معناست که احتمال رخداد هر جهش در این توده برابر با احتمال رخداد دیگر جهش‌هاست. در مقابل، یک توده دارای خاصیت هتروژنی است اگر اجزای تشکیل دهنده آن توزیع غیریکنواخت داشته باشد و به همین امر سبب می‌شود تا بدلیل حضور سلول‌های مختلف با توزیع گوناگون، احتمال رخداد جهش‌های مختلف متفاوت باشد.

۵.۲.۳ پایگاه داده:

به منظور ارزیابی مدل استنباط کننده درخت تکاملی تومور، از دو پایگاه داده متفاوت در این مقاله استفاده شده است:

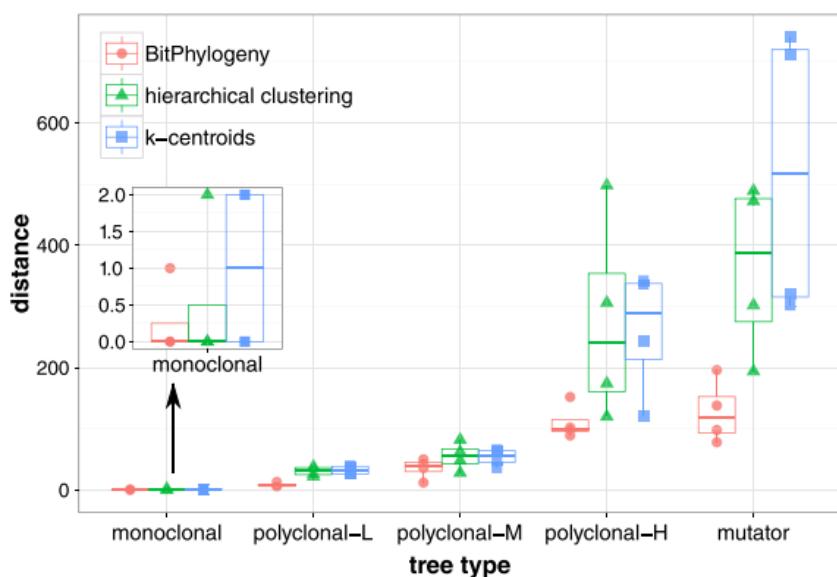
⁷Binary

⁸False negative

- دادگان مربوط به الگوهای متیلاسیون سرطان روده بزرگ
 - دادگان شبیه‌سازی شده مربوط به سرطان خون

۶۰۲۰۳ معیار ارزیابی:

به منظور ارزیابی عملکرد الگوریتم بیت‌فیلوژنی یک مقایسه بین خروجی این الگوریتم و خروجی‌های الگوریتم‌های خوشبندی k هسته‌ای^۹ و دسته‌بندی سلسله مراتبی^{۱۰} صورت گرفته است. این مقایسه از طریق محاسبه معیار بیشینه عمق درخت تکاملی استنباط شده و مقدار درست‌نمایی صورت گرفته است. نتایج گزارش شده در این مقاله گواه از پایداری^{۱۱} و دقیقت^{۱۲} بسیار بهتر الگریتم بیت‌فیلوژنی نسبت به دو الگوریتم دیگر است. در شکل ۳.۳ میزان خطای عملکرد الگوریتم بیت‌فیلوژنی نسبت به دو الگوریتم خوشبندی k هسته‌ای و دسته‌بندی سلسله مراتبی در سطوح مختلف درخت در حالت‌های تک‌کلونی و چندکلونی قابل مشاهده است.



شكل ٣.٣: عنوان

در شکل ۴.۲ به طور کلی مراحل عملکرد الگوریم بیت فیلوژنی را مشاهده می‌کنید. این شکل تومور چندکلونی A را نشان می‌دهد که به روش توالی یابی نمونه‌گیری شده است. این تومور شامل سه کلون مجزا و سلول‌های سالم

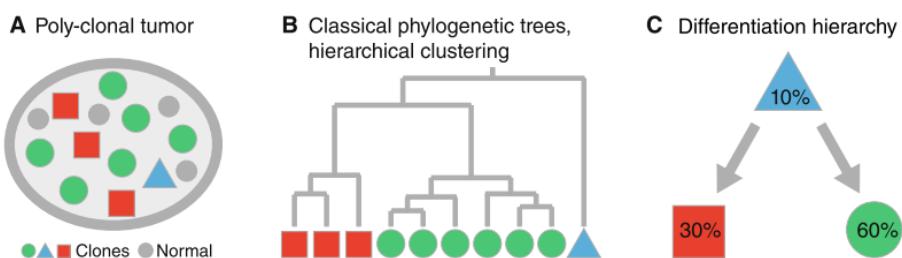
⁹K-Centroids

R-Centroids

Hierarchical ¹¹Consistency

Consistent ¹² Accuracy

(دایره‌های خاکستری رنگ) است. در تصویر میانی یک درخت بلقوه که نشان‌دهنده سیر تکاملی تومور است نشان داده شده است. در تصویر سمت راست درخت کلونی بدست آمده از درخت تکاملی تومور گفته شده با الگوریتم بیت‌فیلوژنی مشاهده می‌شود که در آن کلون‌ها و فراوانی هر یک مشهود است.



شكل ٤.٣: عنوان

٣.٣ الگوریتم Scite [٤٦]

این الگوریتم با استفاده از داده‌های توالی‌یابی تک سلولی سعی در استباط درخت فیلوزنی تومور دارد. همانطور که پیشتر نیز اشاره شد، یک تومور ناشی از تجمع تعدادی سلول با ویژگی‌های ژنی متفاوت است و این سلول‌ها سعی دارند تا این ویژگی‌های ژنی منحصر به فرد را از طریق تکثیر سلولی به سلول‌های بعدی منتقل کنند. [۱۸]

وجود سلول‌ها با جهش‌های متفاوت سبب می‌شود که تومور از زیرنواحی گوناگون، که به کلون مشهور هستند، تشکیل شود. هر چه تومور از تعداد کمتری زیرکلون تشکیل شده باشد درمان آن ساده‌تر خواهد بود. در نظر گرفتن هر کلون به صورت یک تومور جداگانه، مطالعه و بررسی هر یک از این زیرتومورها به صورت دقیق‌تر و یافتن سیر تکاملی آنها سبب می‌شود تا درمان تومور به صورت کارآمدتری انجام شود. [۱۰]

یکی از چالش‌های بزرگ در زمینه تشخیص و مطالعه کلون‌های درون‌تومور، توالی‌یابی قسمت‌های مشترک دنباله‌های دی‌ان‌ای است، زیرا شامل ترکیب‌های بسیار زیادی (در حدود میلیون‌ها ترکیب) از ژنهای سلول‌های گوناگون است. جهش‌های بدست آمده از ترکیب توالی سلول‌های مختلف، با تعداد زیرنواحی توموری (کلون) متناسب است و با استفاده از تعداد زیرنواحی می‌توان تخمین نزدیکی از جهش‌های درون یک نمونه را بدست آورد. [۵۵]

به همین دلیل به منظور شناسایی دقیق هر یک از زیرنواحی توموری (کلون‌ها) لازم است تا اطلاعات حاصل از نواحی مشترک کلون‌ها به دقت مورد تحلیل و تجزیه قرار گیرد. [۶۴]

الگوریتم Scite از طریق داده‌های توالی‌یابی تک سلولی قادر است سیر تکاملی تومور را از طریق درخت جهشی تومور که در آن ترتیب وقوع جهش‌ها مشخص است یا از طریق استنباط درخت فیلوزنیک تومور که در آن هر برگ نشان دهنده یک سلول است، نشان دهد. خروجی مدل Scite نتیجه ارزیابی بهتری در مقایسه با الگوریتم بیت‌فیلوزنی بر روی داده‌های واقعی داراست. الگوریتم Scite از طریق معیار بیشینه درست‌نمایی و احتمال رخداد هر جهش و با استفاده از ماتریس ژنتوتایپ ورودی تعیین می‌کند که کدام درخت استنباط بهتری از سیر تکاملی تومور است. در حالتی که تعداد جهش‌ها بسیار زیاد باشد یعنی تعداد ستون‌های ماتریس ژنتوتایپ ورودی زیاد باشد، ساخت درخت فیلوزنیک راحت‌تر خواهد بود، اما در حالتی که تعداد سلولها زیاد باشد (تعداد سطرهای ماتریس ژنتوتایپ بالا باشد) ساخت درخت جهشی تومور (ترتیب وقوع جهش‌ها) راحت‌تر است. به طور خلاصه اینکه کدام نوع درخت (جهشی یا فیلوزنیک) در نهایت بیان‌کننده سیر تکاملی تومور باشد به نزد جهش‌های توموری و روش توالی‌یابی داده‌ها بستگی دارد.

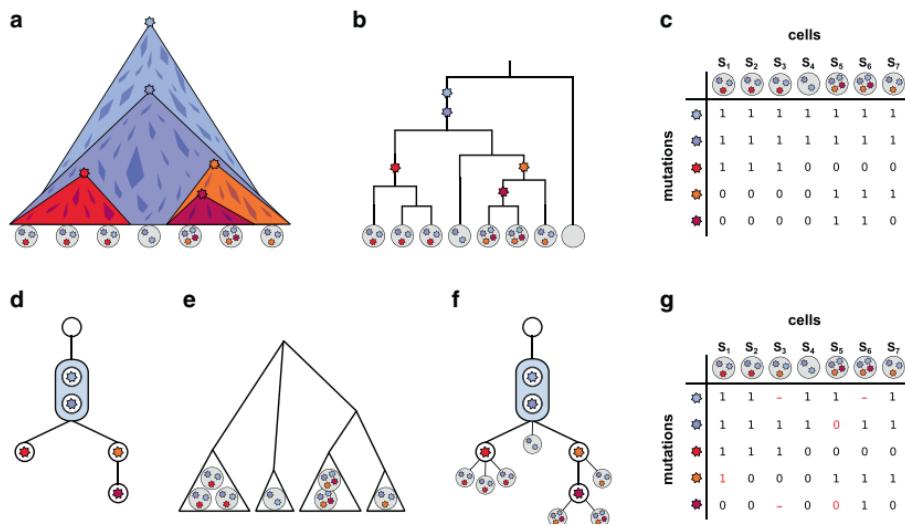
در الگوریتم Scite از دو فرض اصلی استفاده می‌شود:

- فرض مکان‌های بی‌نهایت^{۱۳} که بر طبق آن هر جهش تنها یکبار در هر موقعیت از ژنوم رخ می‌دهد.
- فرض جهش‌های نقطه‌ای یعنی مدل تکاملی تومور به جهش‌های نقطه‌ای محدود می‌شود.

مانند الگوریتم بیت‌فیلوزنی از یک ماتریس ژنتوتایپ (ماتریس دودویی‌ها که سطرها نمایانگر نمونه‌ها و ستون‌ها بیان‌گر جهش‌هایست) به عنوان ورودی الگوریتم استفاده می‌شود. موقعیت هر جهش به صورت درایه $n \times m$ مشخص می‌شود. به این صورت که مقدار صفر در سطر i ام و درایه j ام بیان‌گر آن است که جهش از نوع i در سطر j وجود ندارد. یک ماتریس ژنتوتایپ را ماتریس فیلوزنی کامل^{۱۴} گوییم هر گاه به ازای آن یک درخت فیلوزنیک متناظر باشد. در الگوریتم scite همانند الگوریتم بیت‌فیلوزنی از الگوریتم زنجیره مارکوف مونت کارلو برای استنباط درخت تکاملی تومور از داده‌های توالی‌یابی تک سلولی استفاده می‌شود، با این تفاوت که فضای جستجو برای انتخاب پارامترها بسیار محدودتر از حالت بیت‌فیلوزنی است و نرخ خطاهای داده (مثبت کاذب و منفی کاذب) برای همه جهش‌ها یکسان در نظر گرفته شده است. محدود کردن فضای جستجو برای انتخاب پارامترها از طریق نمونه‌برداری در این فضای سبب می‌شود تا بر اساس سیر زمانی جهش، بیشینه درست‌نمایی از روی توزیع احتمال پیشین نمونه‌ها بدست آید. یکی از مزایای این روش محاسبه نرخ خطاهای توالی‌یابی است. شکل ۵.۳ یک استنتاج تکاملی از داده‌های توالی‌یابی تک سلولی را نشان می‌دهد. در این شکل، a یک ماتریس فیلوزنی کامل است اما b ماتریس داده‌های واقعی است که شامل مقادیر از دست‌رفته، خطای مثبت کاذب و خطای منفی کاذب است. ماتریس داده‌های واقعی با D و ماتریس فیلوزنی کامل را با E نشان داده شده است.

¹³Infinite sites

¹⁴Perfect phylogenetic matrix



شکل ۵.۳: عنواننتنتنتنتنتنتنتنتنتنتنت

منظور از خطای مثبت کاذب این است که به عنوان مثال در یک موقعیت خاص از ماتریس E جهشی وجود ندارد (مقدار ماتریس برابر صفر است) اما در همین موقعیت مقدار یک (وجود جهش) در ماتریس D وجود دارد. نرخ خطای مثبت کاذب با α و نرخ خطای منفی کاذب با β نشان داده می‌شود. مقادیر α و β از طریق روابط زیر تعریف می‌گردند:

$$y = xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx$$

در این معادلات فرض بر استقلال نرخ خطاهای مشاهده شده است. مقدار درست‌نمایی درخت جهشی T با بردار ضمیمه θ و نرخ خطای (α, β) به صورت زیر محاسبه می‌گردد.

$$y = xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx$$

در معادله بالا E ماتریس جهشدار است که با درخت جهشی T و بردار ضمیمه θ تعریف می‌گردد. توزیع احتمال پسین به صورت زیر محاسبه می‌گردد:

$$y = xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx$$

به منظور بالا رفتن سرعت همگرایی مدل زنجیره مارکوف مونت کارلو فرض می‌شود که بردار ضمیمه θ توزیع یکنواخت دارد. در نتیجه:

$$y = xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx$$

اندازه فضای جستجو برای دو پارامتر درخت جهشی T و بردار ضمیمه θ برابر با $(n+1)^{n-1} * (n+1)^m$ می‌باشد. این فضا جستجو با فرض یکنواخت بودن توزیع بردار ضمیمه θ و طبق معادله بالا و حذف بردار ضمیمه به $(n+1)^{n-1}$ انتخاب کاهش می‌یابد. پس از همگرایی با استفاده از الگوریتم زنجیره مارکوف مونت کارلو و احتمال پسین، بهترین ترکیب درخت جهشی T با بردار ضمیمه θ با بیشینه درست‌نمایی بدست می‌آید:

$$y = xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx$$

منظور از MAP در این معادله حالتی است بیشینه درست‌نمایی رخ داده است.

۱.۳.۳ پایگاه داده:

۱

به منظور ارزیابی عملکرد الگوریتم Scite برای استنباط درخت تکاملی تومور از داده‌های توالی‌بایی تک سلولی از داده‌های واقعی و شبیه‌سازی شده، استفاده شده است. مجموعه داده‌های استفاده شده جهت ارزیابی الگوریتم عبارتند از:

- داده‌های توالی‌بایی تک سلولی از یک نمونه تومور مغز استخوان با ۵۸ سلول سرطانی و ۱۸ نوع جهش با نرخ خطای مثبت کاذب $4^{+} * 10^{-4}$ و نرخ خطای منفی کاذب $9^{+} * 10^{-4}$.
- داده‌های توالی‌بایی تک سلولی یک نوع خاص از سرطان کبد با ۱۷ سلول سرطانی و ۵۰ نوع جهش با مقادیر نرخ خطای مثبت کاذب $5^{+} * 10^{-5}$ و نرخ خطای منفی کاذب $3^{+} * 10^{-5}$ و نرخ داده‌های از دست رفته ۲۲ درصد.
- داده‌های توالی‌بایی تک سلولی نمونه‌گیری شده از سرطان سینه با ۴۷ سلول سرطانی و ۴۰ نوع جهش و با نرخ خطای ترک آلل ۷۳ درصد و نرخ خطای مثبت کاذب $6^{+} * 10^{-6}$.

شایان ذکر است که مدت زمان استنباط یک درخت فیلوزنی تا حد زیادی به پیچیدگی داده‌های ورودی بستگی دارد بطوریکه برای ساخت یک درخت با ۵۰ تا ۱۰۰ سلول، مدت زمانی در حدود چندین دقیقه طول می‌کشد. از مهمترین محدودیت‌های این الگوریتم می‌توان به فرض مکان‌های بی‌نهایت اشاره کرد، زیرا این امکان وجود دارد که در یک محل مشخص از یک دنباله دی‌ان‌ای، یک جهش مشخص چندین بار رخ دهد و یا در محل‌های مختلف از یک دنباله ژنی جهش‌های مشابه رخ دهد که این موارد در فرض مکان‌های بی‌نهایت در نظر گرفته نمی‌شود. از دیگر محدودیت‌های این روش آن است که جهش‌هایی که در همه سلول‌ها وجود دارند یا جهش‌هایی که فقط در یک سلول مشاهده شده‌اند (سطری با مقادیر تماماً یک در ماتریس ورودی) در روند استنباط درخت مورد استفاده قرار نمی‌گیرند.

٤.٣ الگوریتم [٦٠]: Onconem

توده‌های سرطانی حاصل تجمیع کلون‌های متفاوت هستند که کلون‌ها متشکل از سلول‌هایی با جهش‌های ژنتیکی مشابه می‌باشند. وجود زیرناحی با جهش‌های ژنتیکی متفاوت و جمعیت‌های سلولی گوناگون منجر به مقاوم شدن تومور در برابر درمان‌های دارویی شده و روند درمان تومور را دچار اختلال می‌کند زیرا که هر یک از داروها به طور موثر یک زیرناحیه توموری را هدف قرار می‌دهد. عدم درمان کامل زیرناحیه توموری می‌تواند منجر به عود مجدد تومور شود. به همین منظور نیاز به یک متادارمان بهینه به گونه‌ای که همه ناحیه‌های کلونی تومور را به طور موثر تحت تاثیر قرار دهد، بیش از پیش احساس می‌شود. [٦٠]

الگوریتم Onconem در سال ۲۰۱۶ با هدف یافتن تاریخچه تکاملی ناحیه‌های درون توموری با استفاده از داده‌های توالی‌بایی تک‌سلولی ارائه گردید. این الگوریتم قادر است تا ناحیه‌های درون توموری مشابه را درون یک دسته قرار دهد و برای آنها یک ژنوتایپ یکتا در نظر بگیرد. این الگوریتم بر مبنای تغییرات تک نوکلئوتیدی، درخت تکاملی تومور را استباط می‌کند و قادر به یافتن خطاهای ژنوتایپی می‌باشد. در نهایت با ارزیابی بر روی داده‌های آزمایش، مدل نهایی سنجیده شده و سلول‌ها با جهش‌های یکسان در یک گروه دسته‌بندی شده و در انتهای رابطه میان جهش‌ها و ژنوتایپ‌های مشاهده شده و مشاهده نشده (پیش‌بینی شده) مشخص می‌گردد. این الگوریتم هم می‌تواند درخت کلونال توموری و هم درخت فیلوزنیک توموری (قرارگرفتن سلول‌ها به عنوان برگهای درخت) را به عنوان خروجی بدست دهد. ورودی این الگوریتم ماتریس دودویی ژنوتایپ به همراه نرخ خطأ مثبت کاذب و نرخ خطأ منفی کاذب و نرخ خطأ داده‌های از دست رفته است. در ادامه، الگوریتم سعی می‌کند تا سلول‌ها با ژنوتایپ‌های مشابه را در یک گروه قرار دهد و در نهایت درختی که بیشترین شباهت را با دسته‌بندی صورت گرفته را دارد به عنوان درخت تکاملی تومور استباط کند. از نکات قوت این الگوریتم آن است که قادر است کلون‌هایی را که احتمال وجود آنها بالاست اما در داده‌های نمونه‌گیری شده حضور ندارند حدس بزند. این الگوریتم از دو قسمت اصلی تشکیل شده است:

- ایجاد یک مدل احتمالاتی به منظور مدل کردن جمعیت جهش‌ها بر مبنای داده‌های نویزی و روابط میان داده‌ها

- پیدا کردن درخت‌هایی با بیشترین میزان درست‌نمایی در فضای جستجو

توزیع احتمال پسین با فرض D به عنوان مجموعه داده‌های مدل به صورت زیر محاسبه می‌گردد:

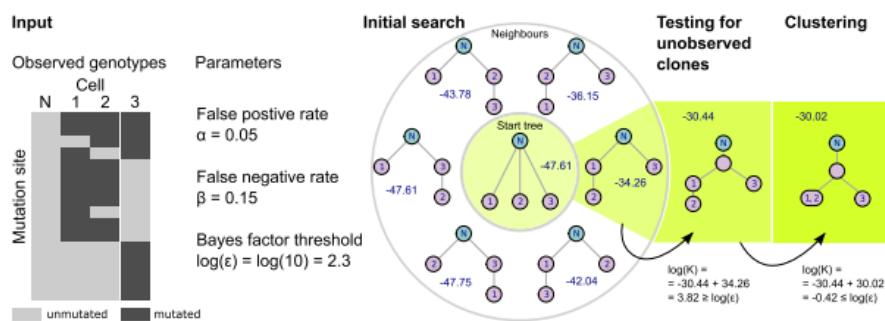
$$y = xxxxxxxxxxxxxxxxxxxxxxx$$

که در آن τ نمایانگر یک درخت جهش‌دار (که باید حتماً دودویی باشد) است که ریشه آن یک گره سالم و بدون جهش است و θ یک بردار رخداد است. در این رابطه فرض بر آن است که (τ) دارای توزیع یکنواخت

است. رابطه بالا می‌تواند به شکل زیر بازنویسی شود:

$$y = xxxxxxxxxxxxxxxxxxxxxxxxx$$

بر طبق این رابطه، برای درختی با n راس، فضای جستجو شامل n^{n-2} انتخاب است که هزینه محاسباتی بسیار بالایی برای درختانی با راس‌های بیشتر از ۹ دارد. طبق شکل ۶.۳، برای محدود کردن فضای جستجو از یک الگوریتم اکتشافی استفاده می‌شود تا اطمینان حاصل شود که خروجی الگوریتم یک نقطه بهینه محلی نباشد. نقطه قوت این الگوریتم سرعت بالای استنباط درخت برای داده‌های کم است ولی در مقابل از محدودیت‌های آن می‌توان به فرض بینهایت اشاره کرد.



شکل ۶.۳: عنواننتنتنتنتنتنتنتنتنتنتنتنتنتنتنتنتنت

رونده کلی الگوریتم Onconem در شکل بالا توضیح داده است. طبق این شکل، ماتریس دودویی ژنوتایپی به همراه نرخ خطاهای α و β به عنوان ورودی الگوریتم استفاده می‌شوند. طبق شکل بالا میزان درستنمایی اولیه برابر $61/47$ – محاسبه شده است اما از میان همه درخت‌های همسایه درخت اولیه، آن درختی که درختی که درستنمایی را دارد به عنوان درخت اولیه انتخاب می‌شود (با درستنمایی $26/36$ –). در ادامه یک گرهای که احتمال رخداد آن طبق ماتریس ورودی بالاست ولی در داده‌های ورودی وجود ندارد به درخت اضافه می‌شود. در این حالت مقدار درستنمایی به $82/82$ افزایش می‌یابد و این کلون مشاهده نشده بدلیل بزرگتر بودن مقدار درستنمایی از آستانه تعیین شده، به مدل افزوده می‌شود. در نهایت گرهای یک شاخه تا جایی که سبب کاهش میزان درستنمایی نشوند، در یک کلون تجمعی می‌شوند.

۱.۴.۳ پایگاه داده:

به منظور ارزیابی عملکرد الگوریتم Onconem از دو پایگاه داده مجزا استفاده شده است:

- داده‌های توالی‌یابی تک سلولی مربوط به سرطان مثانه که شامل ۴۴ سلول سرطانی است. در حدود ۵۵ درصد از انواع جهش‌های موجود در این پایگاه داده، اطلاعاتی در دسترس نیست یعنی بیش از نیمی از

داده‌های موجود از اطلاعات، از دست رفته^{۱۵} آند. خروجی الگوریتم Onconem برای این پایگاه داده یک درخت فیلوزنی با سه کلون اصلی می‌باشد و یک چهارم سلول‌های جهش‌یافته را شامل می‌شود.

- داده‌های مربوط به سرطان خون که در مدل کیم و سایمون و الگوریتم بیت‌فیلوزنی از آن استفاده شده بود، در این ارزیابی مورد استفاده قرار گرفت. میزان لگاریتم درست‌نمایی الگوریتم Onconem برای این مجموع داده برابر ۹۹۶۴ – گزارش شده است که بالاتر از مقداری است که الگوریتم بیت‌فیلوزنی به آن رسیده بود (۱۱۵۸۴ –).

۵.۳ الگوریتم [Sasc]

سرطان ناشی از جهش‌های ژنومیک یک سلول است که این جهش‌ها به مرور زمان رشد و تکثیر می‌یابند و زیرنواحی متفاوتی را ایجاد می‌کنند. این زیرنواحی، که به آنها کلون نیز گفته می‌شود، خصویات متفاوتی دارند و در کنار هم یک توده سرطانی را تشکیل می‌دهند. بررسی تاریخچه تکاملی تومور می‌تواند کارآمدی درمان‌های موجود را بهبود بخشد و امکان عود مجدد تومور را تا حد زیادی کاهش دهد. به منظور درک بهتر تاریخچه تکاملی تومور فرض‌های گوناگونی جهت ساده‌سازی مسئله صورت می‌گیرد، مثل فرض مکان‌های بی‌نهایت که طبق آن هر جهش یکتایی تنها یکبار رخ می‌دهد. مطالعات زیادی صورت گرفته است که نشان می‌دهد در نظر گرفتن فرض مکان‌های بی‌نهایت به تنهایی برای استنباط روند تکاملی تومور کافی نیست و محدودیت‌هایی دارد، به همین منظور برای درک بهتر نواحی ناهمگن توموری باید فرض‌های دیگری را به مسئله اضافه کنیم. به همین دلیل یک فرضیه جدید تحت عنوان k -dollo^{۱۶} ارائه گردید که بر طبق آن و بر خلاف فرض مکان‌های بی‌نهایت، هر جهشی تنها یکبار رخ می‌دهد اما امکان از دست دادن این جهش به تعداد k در تاریخچه تکاملی تومور وجود دارد. الگوریتم Sasc که در سال ۲۰۱۸ ارائه گردید، از اولین الگوریتم‌هایی بود که از فرض k -dollo^{۱۶} جهت استنباط درخت تکاملی تومور بهره برد. به مانند الگوریتم Onconem، این الگوریتم به منظور محدود کردن فضای جستجو از یک الگوریتم درخت اکتشافی بهره می‌برد. الگوریتم اکتشافی استفاده شده در این روش، الگوریتم شبیه‌سازی ذوب فلزات است و هدف آن پیدا کردن بیشینه درست‌نمایی برای تابع احتمال رخداد پسین در فضا جستجو است. طبق این الگوریتم، ابتدا از طریق مجموعه‌ای از انتخاب‌های نمونه‌برداری شده از فضا جستجو یک راه حل برای مسئله ارائه می‌گردد. اگر مقدار درست‌نمایی نسبت به حالت اولیه بهبود یافته بود، با احتمال یک پذیرفته می‌شود در غیر این صورت احتمال رخداد آن حالت صفر در نظر گرفته می‌شود. این الگوریتم سعی دارد تا بیشینه درست‌نمایی ماتریس ژنوتایپ ورودی را حساب کند. ورودی این الگوریتم در کنار ماتریس ژنوتایپ، نرخ خطای مثبت کاذب،

¹⁵Missing data

نرخ خطای منفی کاذب و نرخ خطای اطلاعات از دست رفته است و بیشینه درستنمایی از رابطه زیر بدست می‌آید:

$$y = xxxxxxxxxxxxxxxxxxxxxxxxx$$

۱.۵.۳ پایگاه داده:

به منظور ارزیابی عملکرد الگوریتم Sasc از دو پایگاه داده مجزا استفاده شده است:

- داده‌های توالی‌یابی تک سلولی سرطان مثانه
- داده‌های شبیه‌سازی شده سرطان خون

خروجی الگوریتم در مقایسه با الگوریتم Scite از مقدار بیشینه درستنمایی بیشتری برای مدل کردن داده‌ها برخوردار است.

۶.۳ الگوریتم Scarlet [۶۶]

سرطان فرایند تکاملی است که سلول‌های داخل یک تومور به مرور زمان جهش‌های جسمی تجمعی شونده خواهند داشت. اگرچه تنها تعداد کمی از این جهش‌های جسمی منجر به توسعه سرطان می‌شوند، با این حال تمام جهش‌های جسمی می‌توانند به عنوان یک نشانگر زیستی برای استنباط^{۱۶} تاریخچه تکامل سرطان استفاده شوند. تکنولوژی‌های اخیر توالی‌یابی تک سلولی دی‌ان‌ای، قابلیت اندازه‌گیری جهش‌های جسمی را در سلول‌های منفرد یک تومور ممکن می‌سازند

از انجایی که ابزارهای اندازه‌گیری توالی‌یابی تک سلولی دی‌ان‌ای، اندازه‌گیری دگرگونی تک‌هسته‌ای^{۱۷} را با نرخ بالای مثبت کاذب و منفی کاذب انجام می‌دهند و همچنین دگرگونی تک‌هسته‌ای تنها نوع جهش جسمی در سرطان نیستند، جهش‌های تغییرات شماره کپی که کارشان کپی کردن یا حذف قسمت‌هایی از زنوم است، نیز قابل شناسایی هستند.

دگرگونی تک‌هسته‌ای و جهش‌های حذف و تغییر تعداد کپی^{۱۸}، نشانگرهای زیستی تکامل سرطان هستند. هر دو جهش دگرگونی تک‌هسته‌ای و تغییر تعداد کپی در طول تکامل سرطان تجمعی خواهند شد و این جهش‌ها

¹⁶Inference

¹⁷Single nucleotide variant (SNV)

¹⁸Copy number variation (CNV)

ممکن است در ژنوم همپوشانی داشته باشند. به عنوان مثال حذف یا پایان یک جهش ممکن است منجر به حذف دگرگونی تک‌هسته‌ای‌ها نیز شود. از انجایی که مدل مکان‌های بینهایت اجازه حذف دگرگونی تک‌هسته‌ای را نمی‌دهد، روش‌هایی که از این مدل استفاده می‌کنند به صورت دقیق درخت فیلوژنی را با حضور حذف و تغییر تعداد کپی‌ها بازسازی نخواهند کرد.

دگرگونی تک‌هسته‌ای و جهش‌های حذف و تغییر تعداد کپی نشانگرهای زیستی تکامل سرطان هستند. هر دو جهش دگرگونی تک‌هسته‌ای و حذف و تغییر تعداد کپی در طول تکامل سرطان تجمعی خواهند شد و این جهش‌ها ممکن است در ژنوم همپوشانی داشته باشند. به عنوان مثال حذف یا پایان یک جهش ممکن است منجر به حذف دگرگونی تک‌هسته‌ای‌ها نیز شود. از انجایی که مدل مکان‌های بینهایت اجازه حذف دگرگونی تک‌هسته‌ای را نمی‌دهد، روش‌هایی که از این مدل استفاده می‌کنند به صورت دقیق درخت فیلوژنی را با حضور حذف و تغییر تعداد کپی‌ها بازسازی نخواهند کرد.

مدل ارائه شده در این مقاله، یک مدل تکاملی است که امکان حذف هر نوع جهشی را با در نظر گرفتن حذف خطای^{۱۹} در نظر می‌گیرد. این مدل اجازه حذف دگرگونی تک‌هسته‌ای را تنها هنگامی که با شواهد داده توالی‌یابی تک‌سلولی دی‌ان‌ای از یک حذف در همان مکان هندسی^{۲۰} همراه شده باشد، می‌دهد. این مدل پایه الگوریتم اسکارلت خواهد بود که فیلوژنی تومور را از داده توالی‌یابی تک‌سلولی دی‌ان‌ای با احتساب هر دوی خطای توالی‌یابی و حذف جهش‌ها نتیجه می‌دهد.

تعداد کمی از جهش‌های سوماتیک منجر به پیشروی سرطان می‌شوند، اما تمام جهش‌های سوماتیک نشانگرهای زیستی تاریخچه تکامل تومور هستند. روش‌های غالب ساخت فیلوژنی داده توالی‌یابی تک‌سلولی دی‌ان‌ای از دگرگونی تک‌هسته‌ای‌ها به عنوان نشانگرهای زیستی استفاده می‌کنند اما در به حساب آوردن تغییر تعداد کپی، که ممکن است با دگرگونی تک‌هسته‌ای همپوشانی داشته باشد و منجر به حذف دگرگونی تک‌هسته‌ای شود، ناتوان است.

الگوریتم پیشنهادی اسکارلت، فیلوژنی تومور را از داده توالی‌یابی تک‌سلولی دی‌ان‌ای، خطای توالی‌یابی و حذف دگرگونی تک‌هسته‌ای از طریق تغییر تعداد کپی را لحاظ می‌کند. این الگوریتم عملکرد بهتری نسبت به روش‌های موجود بر روی داده‌های شبیه‌سازی شده دارد. توالی‌یابی تک‌سلولی دی‌ان‌ای از تومور بدليل افزایش بازدهی الگوریتم و کاهش هزینه ایزوله کردن، نشانه‌گذاری و توالی‌یابی سلول‌های انفرادی از محبوبیت روزافروزی برخوردار است. این روش از پیچیدگی بازسازی فیلوژنی با نمونه‌گیری‌های فراوان از سلول‌ها را جلوگیری می‌کند. شایان توجه است که نمونه‌برداری تک‌سلولی به علت خطای توالی‌یابی، نمونه‌برداری کمتر و خطای تشديد دی‌ان‌ای، خطای اطلاعات از دست رفته را بالاتر خواهد برد.

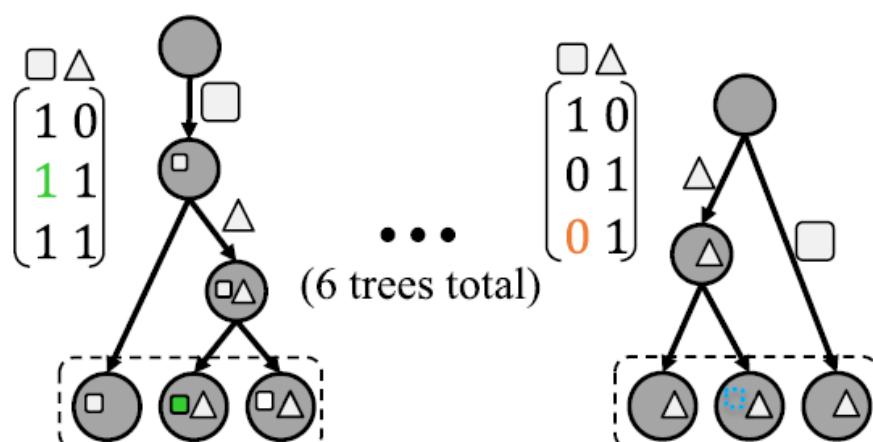
¹⁹ loss-supported²⁰ Loci

یکی از چالش‌های ساخت فیلوزنی با استفاده از داده‌های توالی‌بایی تک سلولی دی‌ان‌ای، که از دگرگونی‌های تک‌هسته‌ای به عنوان نشانگر زیستی استفاده می‌کند، نرخ بالای خطای ترک آلل (تا ۳۰ درصد) هست. به همین منظور باید بتوان مدل تکاملی فیلوزنی را همزمان با در نظر گرفتن داده‌های ورودی از بین رفته ساخت. شکل ۷.۳ نمونه ماتریس داده‌ها را نشان می‌دهد. عدد صفر به معنای عدم وجود جهش و عدد یک به معنای جهش برای سلول مورد نظر است.

mutations	
cells	□ △
○□	1 0
○△	0 1
△□	1 1

شکل ۷.۳: عنواننتنتنتنتنتنتنتنتنتنتنتنت

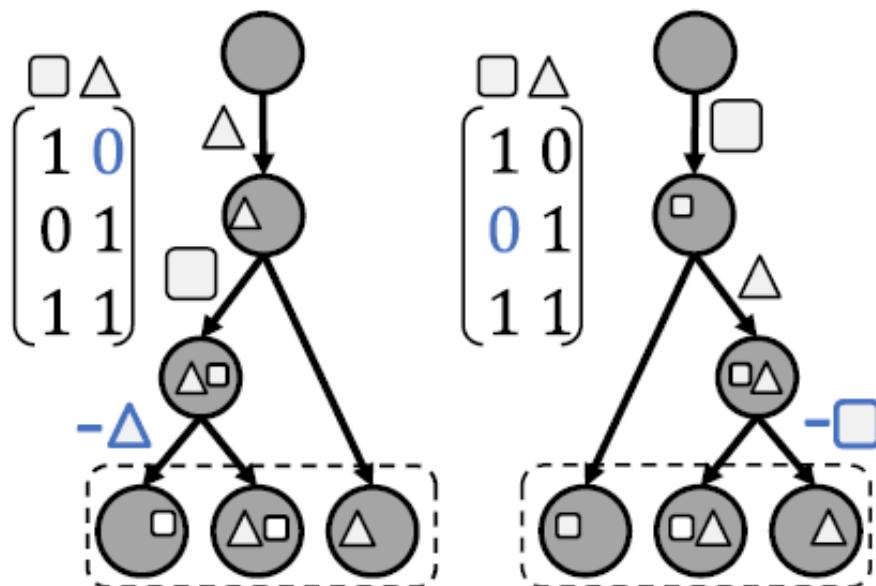
اگر داده‌های توالی‌بایی عاری از هر گونه خطا باشند، فیلوزنی کامل منحصر به فرد^{۲۱} بدست خواهد آمد. با در نظر گرفتن وجود خطا در ماتریس فیلوزنی و کاهش خطاهای موجود، می‌توان مدل فیلوزنی کامل را ساخت. از آنجایی که حالت‌های متعددی از تصحیح‌های متعدد امکان‌پذیر است، امکان استنباط چندین فیلوزنی گوناگون سازگار با داده‌ها وجود دارد. با استفاده از مدل مکان‌بی‌نهایت و اصلاح خطاهای می‌توان درخت‌های فیلوزنی‌های کامل ۸.۳ را ساخت. در اینجا ۶ درخت ممکن است.



شکل ۸.۳: عنواننتنتنتنتنتنتنتنتنتنتنتنت

²¹Perfect phylogeny

شکل ۹.۳ فیلوزنی با فرض دولو^{۲۲} را نشان می‌دهد. این مدل با شناسایی حذف جهش‌ها به منظور رفع تناقض مدل مکان‌های بی‌نهایت می‌تواند درخت‌های ۹.۳ را بسازد. هر دو مدل دولو و مکان‌های بی‌نهایت می‌توانند چندین درخت ممکن را بسازند.



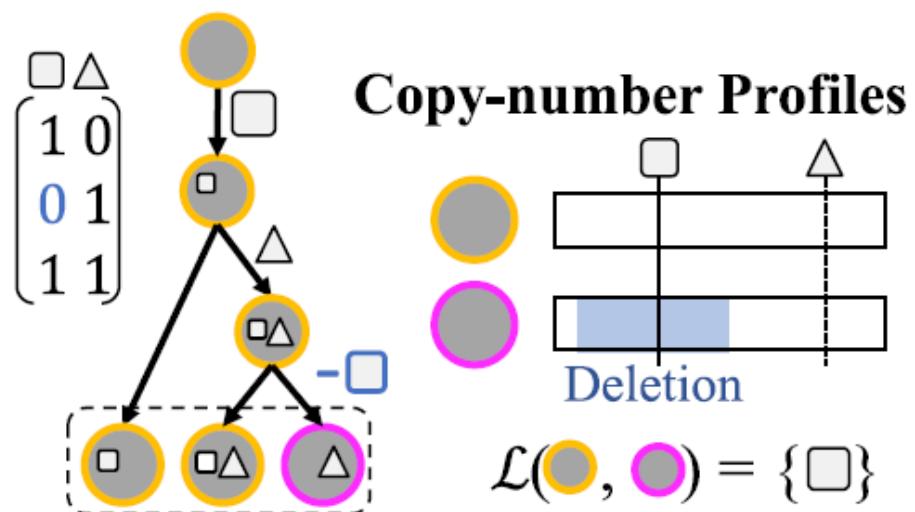
شکل ۹.۳: عنواننتننتننتننتننتننتننتن

حتی در حالت‌های ساده‌ای که خطاب وجود ندارد، استتباط چندین فیلوزنی سازگار با داده‌ها ممکن است وجود داشته باشند. در صورتی که خطاب وجود داشته باشد و عدم قطعیت در ماتریس جهش وجود داشته باشد، تعداد این درخت‌های احتمالی بسیار بیشتر خواهد شد. خطای داده توالی‌یابی تک‌سلولی دی‌ان‌ای و حذف جهش‌ها منجر به پیچیدگی مسئله و ابهام در استتباط فیلوزنی خواهد شد. به عنوان مثال با مشاهده کردن ۰ در ماتریس جهش به ۱ نمی‌توان براحتی بین خطاهای داده‌ها و حذف جهش‌ها تفاوتی قائل شد. عمدۀ محدودیت الگوریتم‌های دولو یا مکان‌های بی‌نهایت‌یابی که اجازه حذف جهش‌ها را می‌دهند این است که هیچکدام از این روش‌ها شواهد تغییر تعداد کپی در حذف جهش‌ها را در یک مکان هندسی در نظر نمی‌گیرند. مدل‌های چند حالته^{۲۳} از تکامل تومور که از داده‌های توالی‌یافته با نمونه‌های زیادی از تومور استفاده می‌کنند. این نگرش‌ها نه خطای موجود در داده توالی‌یابی تک‌سلولی دی‌ان‌ای را مدل می‌کنند و نه در ابعاد صدها یا هزاران سلول قابلیت مدل کردن را دارند. از آنجایی که حذف جهش‌ها پیچیده‌ترین قسمت در تکامل دگرگونی تک‌هسته‌ای است و مسئول اکثر تناقضات در مدل مکان‌های بی‌نهایت در داده‌های توالی‌یابی تک‌سلولی دی‌ان‌ای هستند، در نگرش ارائه شده در این

²²Dollo

²³Multi-state

الگوریتم، حذف جهش‌ها را با استفاده از داده‌های جهش‌های تغییر تعداد کپی از همان سلول‌ها محدود خواهد کرد. در نتیجه الگوریتم اسکارلت با یکپارچه کردن دگرگونی تک‌هسته‌ای و داده‌های حذف و تغییر تعداد کپی، درخت فیلوژنی را براساس داده توالی یابی تک سلولی دی‌ان‌ای می‌سازد. الگوریتم اسکارلت براساس مدل فیلوژنی با در نظر گرفتن حذف خطای است که حذف جهش‌ها را محدود به مکان‌های هندسی خواهد کرد. به عنوان مثال در این الگوریتم داده تغییر تعداد کپی گواه یک حذف است. شکل زیر مدل فیلوژنی با در نظر گرفتن حذف خطای را نشان می‌دهد که با استفاده از داده تغییر تعداد کپی سعی در محدود کردن حذف جهش‌ها دارد تا بتواند ابهام ^{۲۴} ایجاد شده را رفع کند.



مدل فیلوزنی با در نظر گرفتن حذف خطای مدلی از تکامل دگرگونی تک هسته‌ای است که جهش حداقل یکبار رخ خواهد داد (۰-۱۱) اما حذف جهش‌ها (۰-۱۱) توسط مجموعه از مقدار خطای حذف که توسط تغییر تعداد کپی‌ها تعریف می‌شوند محدود خواهد شد. برای هر جفت سلول، از مجموعه جهش‌های تغییرات تعداد کپی، مجموعه خطای به صورت ، تعریف خواهد شد. مدل فیلوزنی با در نظر گرفتن حذف خطای توسعه دهنده مدل‌های مکان‌های بین‌نهایت و دولو می‌باشد. ضمناً الگوریتم اسکارلت متکی بر مدل احتمالاتی تعداد خوانش‌ها برای هر دگرگونی تک هسته‌ای است تا خطاهای و داده‌های از بین رفته، که در توالی پایی تک سلولی دی‌ان‌ای معمول هستند، را مورد توجه قرار می‌دهد.

اسکارلت سہ ویژگی، مہم دارد:

24 conflict

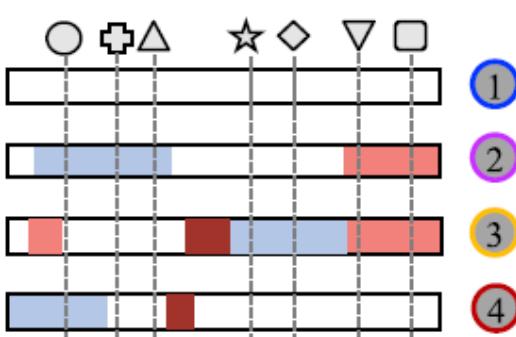
- مدل فیلوژنی با در نظر گرفتن حذف خط، حذف جهش ها را محدود به مکان هایی می کند که کاهش متناظر با آن در تعداد جهش های کپی وجود داشته باشد.
 - این الگوریتم با استفاده از مدل فیلوژنی با در نظر گرفتن حذف خط، ابتدا درخت فیلوژنی اولیه استنتاج شده را پایش و سپس از طریق داده های تغییر تعداد کپی، فیلوژنی نهایی را استنباط می کند.
 - استنتاج مبتنی بر بیشینه درست نمایی از دگرگونی های تک هسته ای با استفاده از مدل احتمالاتی تعداد خوانش های مشاهده شده در داده های توالی یابی تک سلولی دی ان ای توسط الگوریتم اسکارلت اجرا می شود.

اگر تعدادی جهش‌های حذف و تغییر تعداد کپی موجود باشند و وضعیت جهش‌های تغییر تعداد کپی را با رنگ قرمز و حذف نواحی ژنوم در طول کل ژنوم را با رنگ آبی نشان دهیم:

CNAs

Copy-number profiles

Mutations



شكل ١١.٣: عنوان

آنگاه مجموعه خطای مدل فیلوزنی با در نظر گرفتن حذف خطای توسط مجموعه های ۱۳.۳ نمایش داده خواهد شد:

الگوریتم اسکارلت به صورت مستقیم وضعیت جهش‌های حذف و تغییر تعداد کپی سلول‌های پدری را نشان نخواهد داد. به منظور غلبه بر این موضوع یک درخت پرای جهش‌های حذف و تغییر تعداد کمی زیر را در نظر

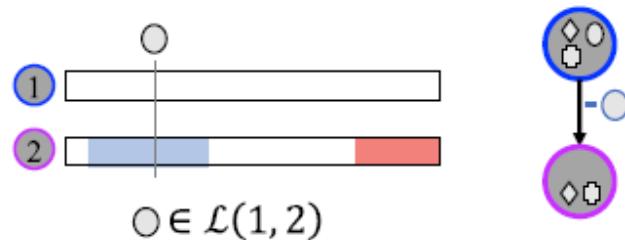
Supported losses \mathcal{L}

$$\mathcal{L}(1, 2) = \{\circ, \square, \Delta\}$$

$$\mathcal{L}(1, 4) = \{\circ\}$$

$$\mathcal{L}(2, 3) = \{\star, \diamond\}$$

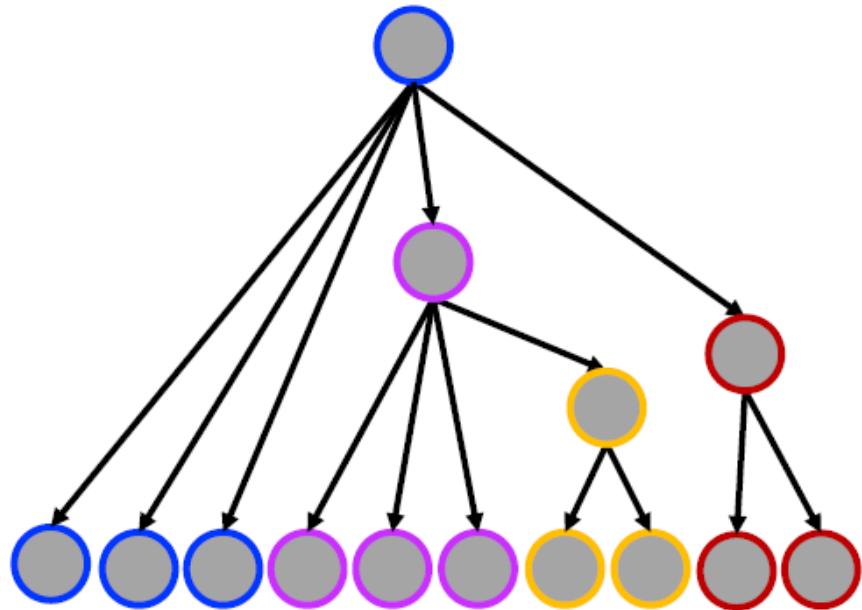
شکل ۱۲.۳: عنواننتنتنتنتنتنتنتنتنت



شکل ۱۳.۳: عنواننتنتنتنتنتنتنتنتنت

خواهد گرفت که از روی وضعیت جهش‌های حذف و تغییر تعداد کپی سلول‌های مشاهده شده، ساخته شده است.

Copy-number tree T

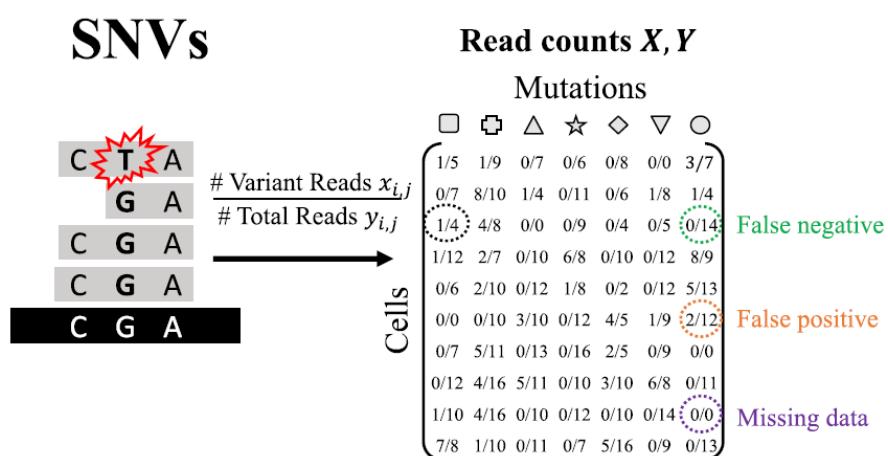


شکل ۱۴.۳: عنواننتنتنتنتنتنتنتنتنت

د ورودی برای الگوریتم اسکارلت در نظر گرفته می‌شود:

- مجموعه خطاهای ناشی از حذف جهش‌ها است، که مجموعه‌های تهی در آن نمایش داده نمی‌شوند. این مجموعه جهش‌هایی که تحت تاثیر حذف قرار می‌گیرند را نشان می‌دهند.
 - یک درخت فیلوژنی برای جهش‌های تغییر تعداد کپی، که با استفاده از آن می‌توان روابط بین سلول‌های مشاهده شده (برگها) را آنگونه که توسط وضعیت جهش‌های تغییر تعداد کپی تعیین شده، نشان داد.

برای دگرگونی‌های تک‌هسته‌ای تنوع X و مجموع Y از تعداد خوانش‌ها^{۲۶}^{۲۷} برای هر سلول و هر جهش مطابق ماتریس $15.0.3$ تهیه شده است:



١٥.٣: عنوانشون

در ادامه الگوریتم اسکارلت، روابط بین اتصال سلولها (T) را از سلول های مشاهده شده (برگها) و ماتریس B را با محدود کردن حذف چهش ها به مجموعه چهش پیشنه درست نمایی^{*}

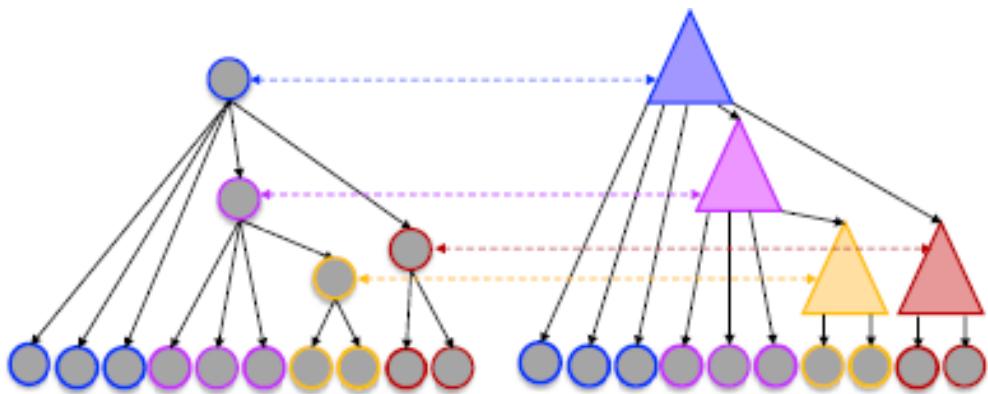
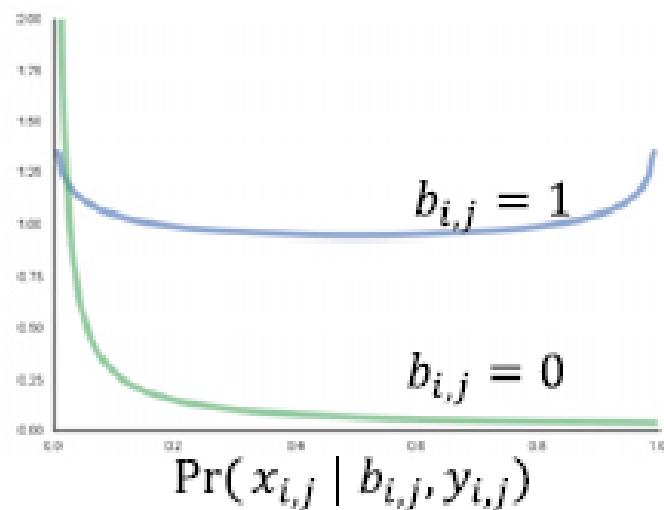
از خطاهای احتمالی حساب میکند. سپس با مقایسه T از T' و انتخاب بیشینه درستنمایی B^* را با استفاده از مدل احتمالاتی برای حضور ($b_{i,j} = 1$) و یا عدم حضور ($b_{i,j} = 0$) هر دگرگونی تک هسته‌ای در هر سلول را انجام می‌دهد.

مدل احتمالات برای توالی بایه داده:

25 Variant

26 Total

Total ²⁷Read counts

شکل ۱۶.۳: مقایسه T' از T 

شکل ۱۷.۳: مدل احتمالاتی برای توالی‌بایی داده

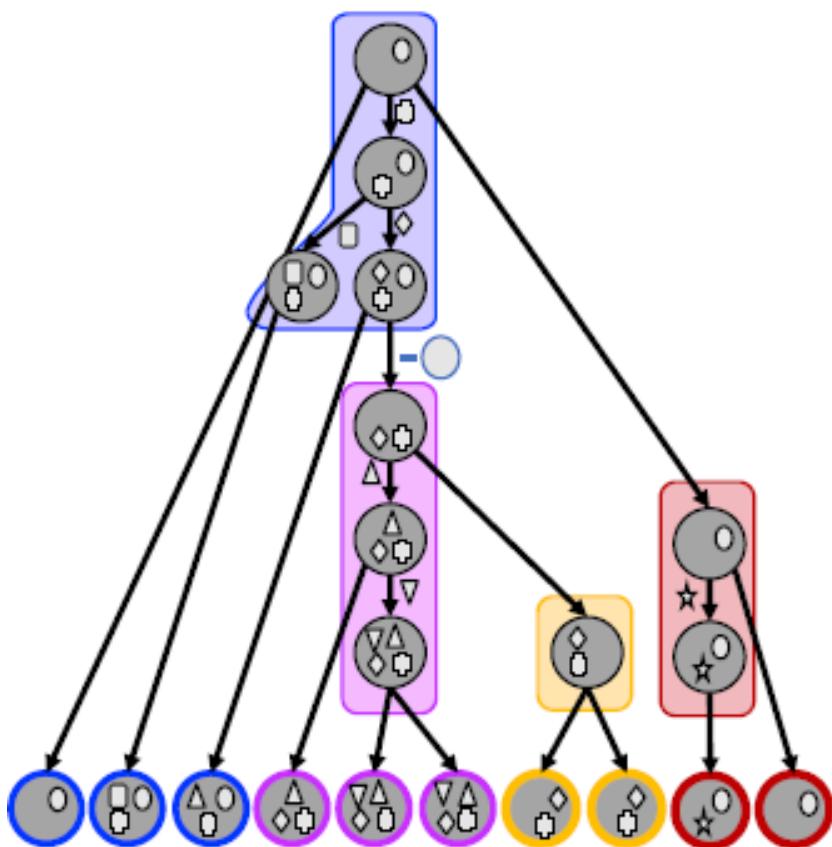
ساختن درخت اتصالات T' :

و در نهایت ماتریس جهش‌ها B^* با بیشینه درست‌نمایی:

الگوریتم اسکارلت باید مسئله بیشینه درست‌نمایی همراه با انتخاب بهترین حذف‌ها حل کند. این الگوریتم از طریق یافتن ماتریس جهش با بیشینه درست‌نمایی B^* انجام خواهد گرفت. در اینجا $L(T)$ مجموعه برگ‌های درخت T را بیان می‌کند.

$$y = xxxxxxxxxxxxxxxxxx$$

الگوریتم اسکارلت از ۲ قسمت اصلی زیر تشکیل شده است.



شکل ۱۸.۳: ساختن درخت اتصالات 'T'

- محاسبه وضعیت‌های جهش با بیشینه درست‌نمایی^{۲۸} R^* از زیردرخت‌ها.
- استنتاج هر زیردرخت به صورت مستقل با هدف بیشینه درست‌نمایی با شرط داشتن R^* .

در اینجا $I(T)$ مجموعه نودهای داخلی درخت T را بیان می‌کند.

$$y = xxxxxxxxxxxxxxxxxxxxxxx$$

مرحله اول:

$$y = xxxxxxxxxxxxxxxxxxxxxxx$$

در اینجا با فرض داشتن R راست نمایی محاسبه خواهد شد و R^* با شمارش حالت جهش‌های معتبر برای هر موقعیت مکانی جهش a محاسبه شده و سپس بیشینه راست نمایی بالا حساب خواهد شد.

مرحله دوم:

یافتن زیردرخت‌ها پایش شده:

²⁸Subtrees

Mutations

	□	□	△	☆	◊	▽	○
Cells	0	0	0	0	0	0	1
	0	1	0	0	0	0	1
	0	1	1	0	0	0	1
	0	1	1	0	1	0	0
	0	1	1	0	1	1	0
	0	1	1	0	1	1	0
	0	1	0	0	1	0	0
	1	1	0	0	1	0	0
	0	0	0	1	0	0	1
	0	0	0	1	0	0	1

شکل ۱۹.۳: ماتریس جهش‌ها B^* با پیشینه درست‌نمایی

تعریف ماتریس سه‌تایی (ternary) : مولفه‌های این ماتریس مقادیر ۰ و ۱ و ؟ می‌باشند.

$$y = xxxxxxxxxxxxxxxxxxxxxxx$$

و در نهایت حل معادله برنامه‌ریزی خطی عدد صحیح زیر:

$$y = xxxxxxxxxxxxxxxxxxxxxxx$$

منوط به شروط زیر:

با فرض اینکه M عدد ثابت و بزرگی است:

$$y = xxxxxxxxxxxxxxxxxxxxxxx$$

$$y = xxxxxxxxxxxxxxxxxxxxxxx F_{w,a,b}$$

$$\text{نقض } G_{w,a,b}$$

$$y = xxxxxxxxxxxxxxxxxxxxxxx F_{w,a,b} \text{lr}$$

$$\text{نقض } H_{w,a,b}$$

$$y = xxxxxxxxxxxxxxxxxxxxxxx F_{w,a,b} \text{lr}$$

به طور خلاصه:

جهش‌های سومایتک تومورها در تمام مقیاس‌های ژنومی، از دگرگونی‌های تک‌هسته‌ای (SNV) تا جهش‌های حذف و تغییر تعداد کپی (CNA) وجود دارد. تا به امروز، بیشتر روش‌های ساخت فیلوزنی توموری از داده‌های توالی‌یابی تک‌سلولی DNA فقط از دگرگونی‌های تک‌هسته‌ای استفاده می‌کردند. [۲۸، ۵۶، ۵۲] و جهش‌های حذف و تغییر تعداد کپی و در نتیجه اطلاعات مهم استنباط فیلوزنیک تومور را نادیده می‌گرفتند.

در این مقاله الگوریتم scarlet معرفی شد که در آن به طور همزمان از دگرگونی تک‌هسته‌ای (SNV) و جهش‌های حذف و تغییر تعداد کپی (CNA) از داده‌های توالی‌یابی تک‌سلولی برای استنباط فیلوزنی تومور استفاده شد. این الگوریتم، یک مدل تکاملی بر اساس در نظر گرفتن خطای ناشی از حذف جهش است که حذف جهش‌ها را محدود به مکان‌هایی می‌کند که شواهدی از حذف جهش‌های حذف و تغییر تعداد کپی موجود باشد. مدل‌های فیلوزنی با در نظر گرفتن حذف خطای ناشی از اطلاعات جهش‌های حذف و تغییر تعداد کپی که به آسانی در داده‌های دگرگونی تک‌هسته‌ای موجود است، نسب به مدل‌های دولو یا فرض مکان‌های بی‌نهایت، ابهام کمتری در استنباط درخت فیلوزنی دارند. اگر چه به صورت طبیعی در داده‌های دگرگونی تک‌هسته‌ای یک عدم قطعیت ذاتی در حضور یا عدم حضور جهش در سلول‌ها وجود دارد، اما کاهش میزان ابهام در استنباط فیلوزنی تومور منجر به افزایش دقت فیلوزنی استنباط شده است. در این مقاله نشان داده شد که فیلوزنی توموری استنباط شده برای بیماران مبتلا به سرطان روده از دقت و تکرارپذیری بیشتری برخوردار است و این الگوریتم در نهایت فیلوزنی‌هایی را استنباط کرد که در آن ^۳ حذف جهش رخ داده بود. البته این الگوریتم محدودیت‌های خاص خود را دارد. به عنوان مثال، این نوع پیاده‌سازی از الگوریتم اسکارلت مستلزم درخت حذف و تغییر تعداد کپی به عنوان ورودی و میزان درست‌نمایی هر یک از این درختان است. این رویکرد در موقعی که تعداد مشخصی از تغییرات تعداد کپی وجود دارد قابل اجراست اما هنگامی که داده‌های توالی‌یابی تک‌سلولی در مقیاس بزرگ انجام شود، به درختان زیادی از جهش‌های حذف و تغییر تعداد کپی نیاز خواهد بود.

۷.۳ الگوریتم [۹] :Deepphylo

همانطور که می‌دانیم، سرطان یک بیماری تکاملی است که با تجمع تدریجی جهش بدنی ^{۲۹} در سلول‌های تومور مشخص می‌شود. رمزگشایی از تاریخچه تکاملی یک تومور، یک چالش مهم در مطالعات سرطان است و می‌تواند از جنبه‌های مهم بالینی از جمله پیشرفت تومور ^{۳۰}، گشتش متاستاتیک ^{۳۱} و وجود زیرکلون‌های واگرا ^{۳۲}

²⁹Somatic mutation

³⁰Tumor progression

³¹Metastatic spread

³²Divergent subclones

در شاخه‌های مختلف درخت فیلوزنیک تومور درک بهتری از تومور در اختیار ما بگذارد. با توجه به اهمیت مسئله، تحولات سریعی در طراحی روش‌های محاسباتی اصولی برای استنباط فیلوزنی تومور وجود داشته است. بسیاری از این روش‌ها از داده‌های توالی یابی‌های آنبوه^{۳۳} استفاده می‌کنند که DNA میلیون‌ها سلول سرطانی و طبیعی با هم یک توالی را تشکیل می‌دهند. استنباط درخت فیلوزنی با استفاده از این نوع داده‌ها، معمولاً^{۳۴} بر مبنای دگرگونی‌های شناسایی شده^{۳۴} از بخش‌های مختلف سلول‌های سرطانی انجام می‌شود. به عنوان مثال: حذف و تغییر تک‌نوکلئوتیدها [۷۵، ۲۹، ۵۳، ۶۵، ۲۶]^{۴۵}، حذف و تغییر تعداد کپی [۸۵]^{۴۶}، دگرگونی‌های ساختاری^{۳۵} [۵۹، ۲۷].

اگرچه استنباط درخت فیلوزنی با استفاده از این نوع داده مقرن به صرفه است اما رزولوشن^{۳۶} پایین داده‌های توالی یابی‌های آنبوه یک فاکتور محدود کننده در مدل‌سازی تکامل تومور است. به طور خاص داده‌های توالی یابی‌های آنبوه ناشی از یک نمونه تومور به طور معمول یک توپولوژی خطی را به عنوان یک راه حل بهینه در تعیین درخت فیلوزنی تومور در نظر می‌گیرد. [۲۶]

با این حال، دانستن اینکه آیا تومور شامل زیرکلون‌های واگرایی است که از طریق شاخه‌های متمایزی از فیلوزنی تکامل می‌یابند، گام مهمی در جهت درک بهتر پیشرفت تومور و بهبود طرح درمانی است. تحولات اخیر تکنولوژی، محققان را قادر به انجام آزمایش‌های توالی یابی تک سلولی کرده است، جایی که DNA از یک سلول استخراج، تکثیر و توالی یابی می‌شود. توالی یابی تک سلولی، داده‌هایی با رزولوشن بالا برای مطالعه تکامل تومور با جزئیات زیاد را فراهم می‌کند، به عنوان مثال، امکان شناسایی توپولوژی شاخه‌ای با اطمینان بالا یا حل مشکل کلی استنباط کامل تاریخ تکامل تومور را فراهم می‌کند، حتی زمانی که تمام سلول‌های تک توالی که از یک نمونه بیوپستی^{۳۷} توموری استخراج شده باشد. روش‌های متعددی برای استنباط تاریخچه تکاملی تومور از طریق توالی یابی تک سلولی وجود دارد که از مهمترین آنها می‌توان به موارد زیر اشاره کرد:

- رویکردهای مبتنی بر آمار و احتمالات که از فرض مکان‌های بی‌نهایت استفاده می‌کنند. مثل الگوریتم OncoNEM^{۴۶} و IrSCITE^{۴۰}.

- رویکردهایی که از فرض مکان‌های بی‌نهایت استفاده نمی‌کنند و فرض را بر این می‌گذارند که تخطی‌های در شکل‌گیری درخت تکاملی فیلوزنی تا یک مقدار خطأ مشخص وجود دارد، مثل الگوریتم SiFit^{۸۷}.

به تازگی الگوریتم‌هایی مثل SPhyR که از یک رویکرد بهینه‌سازی ترکیبی مبتنی بر زوجیت دلو^{۳۸} استفاده

³³Bulk sequencing data

³⁴Detected variants

³⁵Structural variant

³⁶Resolution

³⁷Biopsy

³⁸Dollo parsimony

می‌کنند یا الگوریتم SiCloneFit که بهینه یافته الگوریتم SiFit می‌باشد، ارائه شده است. [۸۶، ۲۸] شایان ذکر است که روش‌های همچون PhISCS-BnB، که از روش‌های بهینه‌سازی بر مبنای شاخه-مرز^{۳۹} استفاده می‌کنند، و یا روش‌هایی مثل ScisTree، که بر مبنای اتصال اکتشافی همسایگی^{۴۰} عمل می‌کند، به منظور بهبود زمان محاسباتی استنباط درخت فیلوزنی تومور ارائه شده‌اند. [۸۳، ۶۲]

در حالتی که هم داده‌های توالی‌یابی‌های انبو و هم داده‌های توالی‌یابی تک سلولی موجود باشد می‌توان تقریب دقیق‌تری از درخت فیلوزنی تومور بدست آورد. [۵۴، ۵۱]

همانطور که در بالا خلاصه شد، روش‌های موجود برای بازسازی فیلوزنی تومور با استفاده از داده‌های توالی‌یابی تک سلولی محدودیت‌های مهمی دارند. اولاً^{۴۱}، بسیاری از این روش‌ها، فرض مکان‌های بی‌نهایت را به کار می‌گیرند (حتی در مواقعي که شرایطی برای خطای محدود^{۴۲} و افزایش همزمان جهش‌ها^{۴۳} در نظر گرفته شود) و سطح نویز یکنواختی را در نظر می‌گیرند (منفی کاذب و همچنین نرخ مثبت کاذب) هر دو این محدودیت‌ها، با پیشرفت درک ما از تکامل تومور و فناوری توالی‌یابی تک سلولی تغییر می‌کند. مهمتر از همه، هدف از این روش‌ها استنباط محتمل‌ترین درخت فیلوزنی توموری است و برای حذف نویز (به دلیل مثال، ترک آلل یا پوشش توالی کم^{۴۴}) از روش‌های همچون بیشینه درست‌نمایی^{۴۵} یا حداقل زوجیت^{۴۶} استفاده می‌کنند. به بیان دیگر این روش‌ها قصد دارند تا یک مساله پارامتری از مرتبه^{۱۱} را حل کنند ولی بدليل عدم مقیاس‌بندی داده‌های توالی‌یابی تک سلولی به مرتبه‌های بزرگتر، در حل دقیق این مساله ناتوان هستند. حتی وقتی هدف این است که به جای بازسازی کامل درخت فیلوزنی تومور، فقط ویژگی‌های اساسی توپولوژی فیلوزنی تومور را استنباط کنیم، این روش‌ها نمی‌توانند به راحتی داده‌های توالی‌یابی تک سلولی شامل چند صد جهش و سلول را کنترل کنند. در نتیجه، تکنیک‌های سریع برای استنباط ویژگی‌های کلیدی فیلوزنی تومور، به عنوان مثال، مواردی که می‌توانند توپولوژی‌های شاخه‌ای را از هم تفکیک کنند، به ویژه برای مجموعه داده‌های توالی‌یابی تک سلولی با سطح نویز بالا از محبوبیت بیشتری برخوردار هستند. به همین منظور، بهتر است در ابتدا به این سوال پاسخ داده شود که آیا حذف نویز برای ساخت فیلوزنی کامل لازم است یا خیر. سرانجام، هر یک از ابزارهای موجود به تلاش انسانی زیادی در طراحی و اجرای الگوریتمی نیاز داشته است، زیرا هر پیشرفت تکنولوژیکی در تولید داده‌ها، توسعه روش‌های کاملاً جدید را ضروری می‌کند. بنابراین داشتن یک رویکرد محاسباتی کلی که بتواند با تغییر منطقی تکنیکی سازگار شود، صرفاً از طریق آموخت آن با داده‌های جدید، بدون نیاز به مدل‌سازی صریح مشخصات نویز، بسیار مطلوب است.

³⁹Branch-bound⁴⁰Joining-based heuristic⁴¹Limited loss⁴²concordant gain of mutations⁴³Low sequence coverage⁴⁴Maximum-likelihood⁴⁵Maximum parsimony

رفع این محدودیت‌ها از طریق رویکرد یادگیری ماشینی یا رویکردهای "داده محور" امکان پذیر است که مجموعه‌ای کلی از توابع را در نظر گرفته و تابعی را در نهایت انتخاب می‌کند که برآورد بهتری از مجموعه داده‌های آموزشی (دادگان واقعی یا شبیه‌سازی شده) باشد. چنین رویکردی نه تنها می‌تواند از عدم دقت در مدل‌سازی مشخصات نویز بکاهد بلکه الگوهای اساسی ضمنی را در داده‌ها یا مسئله را برای توسعه اهداف واقع بینانه‌تر شناسایی می‌کند. پیشرفت‌های اخیر در یادگیری عمیق تعمیم قابل توجهی از فرمول‌بندی‌ها را برای حل بسیاری از مشکلات نشان داده است. [۵۰، ۶۹، ۲۴]

این امکان وجود دارد که یک معماری یادگیری عمیق، زمانی که بتواند در تعداد کافی مجموعه داده آموزش را دیده باشد، بتواند در استباط خواص متمایز از فیلوزنی‌های تومور موفق شود. در سالهای اخیر، بسیاری از برنامه‌های محاسباتی، رویکرد الگوریتمی خود را به رویکردهای داده محور تغییر داده‌اند. مانند رمزگشایی متن دست نوشته برای شناسایی رقم [۱۶] و پردازش زبان طبیعی. [۲۴]

مسانی که در بایولوژی ساختار یافته، فرمول‌سازی هدفمند یا کمی‌سازی آنها مشکل است (مانند استباط ساختار سه بعدی توالی پروتئینی) از روش‌های مبنی بر یادگیری عمیق بشرطین استفاده را در جهت حل مسائل خواهند کرد. [۶۸]

با این حال این مقاله، اولین مقاله استباط درخت فیلوزنی تومور مبنی بر رویکردهای داده محور است. در این مقاله، اولین روش‌های بازسازی فیلوزنی تومور مبنی بر داده را برای رفع محدودیت‌های استراتژی‌های موجود ارائه شده است. نویسنده‌گان این مقاله از داده‌های توالی‌یابی تک سلولی در کنار شبکه‌های عصبی عمیق و یادگیری تقویتی برای استباط ویژگی‌های تپولوژیکی فیلوزنی تومور و همچنین محتمل‌ترین سابقه تکاملی تومور استفاده شده‌است. برای رسیدن به این هدف، چندین چالش وجود داشت:

۱. شبکه عصبی در حالت ایده‌آل باید طوری طراحی شود که بتواند تعداد متفاوتی از سلول‌ها و جهش‌ها را کنترل کند. متناظر باً، برای مدل‌هایی با ورودی‌هایی با اندازه ثابت، بهتر است که از دانش خود در زمینه تهیه داده استفاده شود تا داده‌ها به روشی تهیه شود تا موفقیت در پیش‌بینی‌ها را تسهیل کند.

۲. با توجه به استفاده از شبکه‌های عصبی، برای آموزش مناسب به تعداد زیادی نمونه نیاز است. متاسفانه، تعداد مجموعه داده‌های توالی‌یابی تک سلولی تومور در دسترس عموم برای آموزش مدل‌های یادگیری عمیق به اندازه کافی زیاد نیست. بنابراین، نیاز به تولید تعداد زیادی مجموعه داده شبیه‌سازی شده داده‌های توالی‌یابی تک سلولی وجود دارد.

۳. نویز و خطاهای موجود در داده‌های توالی‌یابی تک سلولی پیچیدگی بیشتری را به این مسئله می‌افزاید و چارچوب پیشنهادی یادگیری عمیق باید از نظر تحمل نویز ارزیابی شود.

۴. معماری انتخاب شده مستلزم نوع خاصی از نظارت است که ما باید قادر به تامین آن باشیم.

به منظور کاهش یا حذف نویز در ورودی "ماتریس ژنوتیپ" استخراج شده از داده‌های توالی‌یابی تک‌سلولی، می‌توان نظارت را به صورت مجموعه داده‌ای از ورودی‌های نویزدار به همراه با ورودی‌های بدون نویز ارائه داد. یک نظارت جایگزین و ارزان‌تر توسط مکانیزم بازخورد^{۴۶} است که تعیین می‌کند که آیا یک خروجی از شبکه عصبی با موفقیت بدون نویز شده است یا خیر. گزینه سوم توسط یکتابع هزینه ارائه می‌شود که به طور غیرمستقیم کمک به نظارت بر فرایند یادگیری تقویتی می‌کند.

در این مقاله با الهام از رویکردهای جدید یادگیری عمیق برای مسائل گوناگون مانند "الگوریتم گرادیان سیاست تقویتی" برای مساله فروشنده دوره گرد^{۴۷} [۶۷] NeuroSAT [۸۲] برای مساله رضایتمندی با استفاده از نظارت تکیتی، یک چارچوب محاسباتی ایجاد شد تا همه چالش‌های فوق را به شرح زیر با موفقیت حل کند:

۱. یک رویکرد مبتنی بر یادگیری تقویتی به منظور آموزش مدلی جهت از بین بردن نویز داده‌ها بدون نیاز به استاندارد مرجع^{۴۸} به کار گرفته شد. تابع هزینه استفاده شده در این مدل یکتابع هزینه خاص برای رفع مساله از بین بردن نویز بود.

$$y = xxxxxxxxxxxxxxxxxxxxxxxx$$

که در آن X ماتریس خروجی ناشی از ورودی A' است.

۲. داده‌های ماتریس ورودی، که از مجموعه دادگان نویزی توالی‌یابی تک‌سلولی استخراج شده، در کنار نرخ نویز و موقعیت مکانی به عنوان ورودی به شبکه داده شده است. این رویکرد در مجموعه دادگانی با سایز متفاوت همچنان کارآمد است و مستقل از جایگایی در سطر و ستون ماتریس ورودی است.

۳. یک مرحله پیش‌پردازش دیتا، به منظور به کارگیری دانش حاصل از تجربه در نظر گفته شده است تا هر گونه عملکردی را که می‌تواند پیش‌بینی مدل را بهبود بخشد، بر روی داده‌ها اعمال گردد.

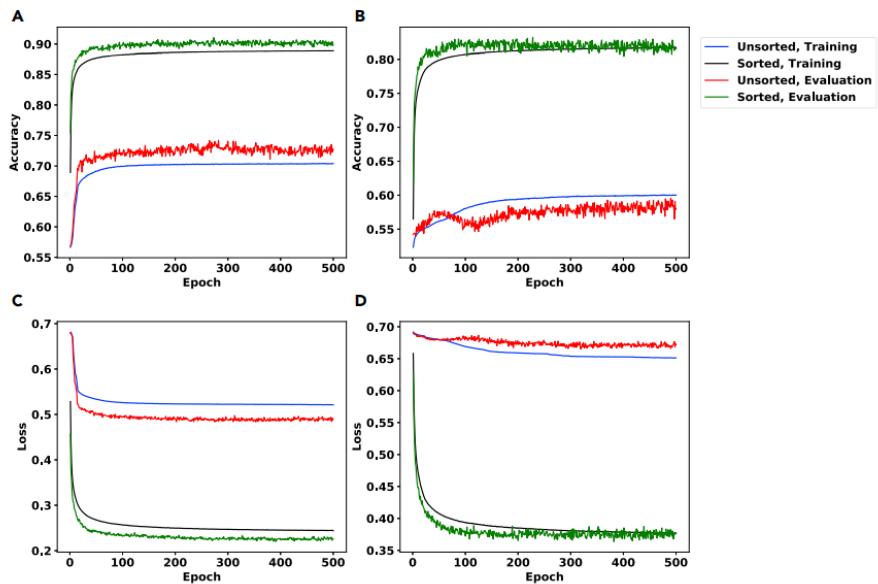
۴. داده‌های شبیه‌سازی شده ورودی مدل از طریق یک چاچوربی که راستی آزمایی شده است، توسعه یافته است.

در نمودار ۲۰.۳، میزان دقیقت علمکرد شبکه در حذف نویز داده ورودی و تاثیر مرحله پیش‌پردازش بر خروجی الگوریتم را مشاهده می‌کنید. همچنین تاثیر میزان نرخ نویزی بودن داده‌ها در خروجی شبکه قابل توجه است.

⁴⁶feedback

⁴⁷Gold standard

تصاویر A و C میزان دقت شبکه در حذف نویز داده‌ای را نشان می‌دهد که با نرخ نویزهای $\alpha = 0.02$ و $\beta = 0.1$ نمونه‌برداری شده‌اند اما تصاویر B و D میزان دقت شبکه در حذف نویز داده‌ای را نشان می‌دهد که با نرخ‌های کاذب مثبت $\alpha = 0.00004$ و کاذب منفی $\beta = 0.002$ نمونه‌برداری شده است.



شکل ۲۰.۳: عکس‌های مقایسه این تأثیر مرحله پیش‌پردازش دیتا در دقت خروجی مدل در حذف نویز از دیتا را مشاهده می‌کنید که میزان دقت حذف نویز بهبود قابل قبولی داشته است.

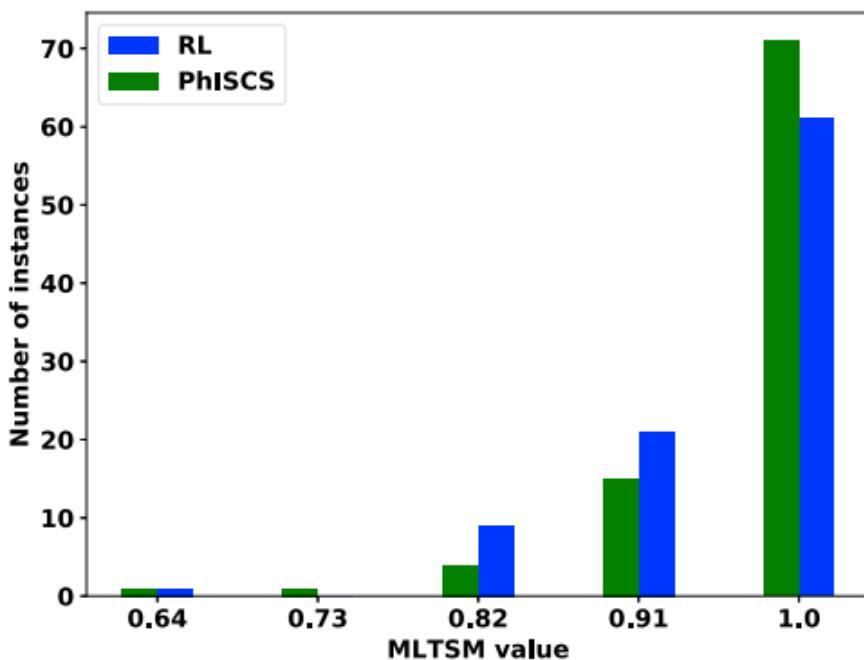
همچنین در جدول ۱.۲ تأثیر مرحله پیش‌پردازش دیتا در دقت خروجی مدل در حذف نویز از دیتا را مشاهده می‌کنید که میزان دقت حذف نویز بهبود قابل قبولی داشته است.

Input MAtrix Size	A	B	Unsorted Acc.	Sorted Acc.
10*10	0.002	0.1	72	90
10*10	$4 * 10^{-4}$	0.02	60	81
25*25	$3.2 * 10^{-4}$	0.016	50	77
25*25	$6.4 * 10^{-4}$	0.0032	52	65

fffffffffffff 3.1:

در نهایت مقایسه بین عملکرد الگوریتم پیشنهادی در این مقاله و الگوریتم PhISCS با استفاده از معیار شباهت MLTSM84 انجام شد که نتیجه این مقایسه در شکل ۲۱.۳ آمده است. همانطور که در شکل

مشهود است عملکرد الگوریتم پیشنهادی در میزان شباهت‌های مشابه، تعداد استنباطهای بیشتری از فیلوزنی تومور را شامل می‌شود.



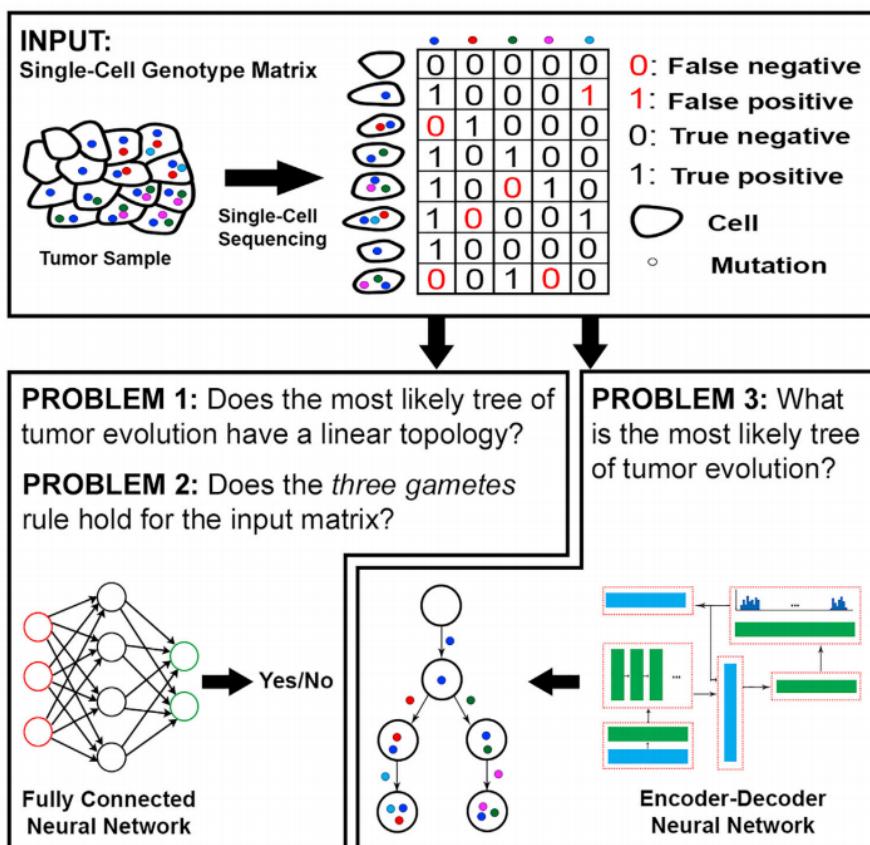
شکل ۲۱.۳: عرضه شده در شکل ۲۲.۳ آمده است:

۱۰.۷.۳ جمع‌بندی

وجود ناهمگنی‌های درون توموری باعث ناکارآمدی درمان‌های دارویی تومور می‌شود زیرا که هر یک از این روش‌های درمانی به طور موثر فقط بر روی تعداد محدودی از کلون‌های توموری اثر می‌گذاردند و همه این زیرنواحی را تحت تاثیر قرار نمی‌دهند. با مطالعه بر روی تومورهای مختلف این امکان حاصل می‌شود تا الگوهای درون توموری بهتر شناخته شوند و درمان داریی تومور کارآمدتر و بهینه‌تر از گذشته صورت پذیرد. مطالعه بر روی داده‌های توالی‌یابی تک‌سلولی یکی از زمینه‌های تحقیقاتی است که می‌تواند منجر به افزایش دانش از نحوه شکل‌گیری و تکامل تومور شود. پیدا کردن سیر زمانی تکامل تومور با استفاده از داده‌های توالی‌یابی تک‌سلولی چالشی است که اخیراً مورد توجه قرار گرفته است که در جدول زیر خلاصه‌ای از روش‌هایی که در این فصل مورد بررسی قرار گرفته‌اند، مشاهده می‌شود.

جدول ۲.۳: Comparison

Method	Dataset	Algorithm	Output	Evaluation Method	Limitation
Kim & Simon approach	Thrombocytopenia Essential (TE)	Minimal spanning tree of the Edmonds' algorithm	Phylogenetic tree	Leave one out cross validation	high computational time and excluding uncertainty dataset error
BitPhylogeny	JAK2 negative myeloproliferative	TSSB, MCMC	Evolutionary clonal tree	V measure comparison with K-Centroids and Hierarchical Clustering	High computational time, infinite sites assumption and homozigous differentiation
SCITE	JAK2 negative myeloproliferative, clear cell and renal cell carcinoma, estrogen-receptor positive breast cancer	Maximum bayesian MCMC likelihood,	Phylogenetic tree	Better performance in real dataset in comparison with biphylogeny algorithm	Infinite sites assumption
ONCONEM	Muscle invasive bladder transitional cell carcinoma	Neighbor joining, MCMC	Phylogenetic tree, evolutionary clonal tree	Score function extracted from nested model	infinite sites assumption, homozigous heterozygous differentiation
SASC	Muscle invasive bladder transitional cell carcinoma, Thrombocytopenia Essential (TE) Robust approach based on scDNA-seq data from a metastatic colorectal cancer patient	Simulated annealing	Phylogenetic tree	Better performance in real dataset in comparison with SCITE algorithm	Limited mutation assumption
SCARLET	scDNA-seq data from a metastatic colorectal cancer patient	Loss-Supported phylogenetic Model	Phylogenetic tree	Mutation matrix error and pairwise ancestral relationship error	Mutation loss due to the dollo assumption
DeepPhylo	Acute Lymphoblastic Leukemia, TNBC dataset	Critic-actor reinforcement learning	Phylogenetic tree	Accuracy, maximum likelihood	Fixed input dimension, lack of empirical experiment



شکل ۲۲.۳: عرضه مسئله‌های پیشین

فصل ۴

روش پیشنهادی

۱.۴ مقدمه

پس از آشنایی با روش‌های پیشین که برای حل مسئله مشابه مورد استفاده قرار گرفته‌اند، حال می‌توانیم به معرفی و تشریح روش‌های پیشنهادی خود برای حل مسئله پیش رو بپردازیم. در این فصل ابتدا داده‌های ورودی مسئله را همراه با فرضیات در نظر گرفته شده بیان می‌کنیم و پس از آن دو روش پیشنهادی متفاوت را بیان خواهیم نمود. در روش اول که به رویکردهای پیشین نزدیک‌تر است با تغییری از جنس روش‌های نوین در مراحل میانی به یک روش جدید می‌رسیم که به علت افزایش سرعت همگرایی می‌توان فرض و داده‌های جدیدی را از طریق حذف و تغییر تعداد کپی به آن افزود و پاسخ گرفت. اما روش دوم کاملاً متفاوت بوده و با رویکردی جدید در حوزه یادگیری ماشین همراه است که به کمک یادگیری تقویتی به حل مسئله مورد نظر می‌پردازد.

۲.۴ معرفی دادگان ورودی

قبل از وارد شدن به بخش روش‌های پیشنهادی نیاز است تا دادگان ورودی را مشخص و معرفی نماییم. دادگان ورودی در این پایان‌نامه همگی به صورت فایل‌های خام اسکی^۱ هستند که حاوی اطلاعات جهش‌های ماتریس ژن-سلول (SNV) و اطلاعات مربوط به حذف و تغییر تعداد کپی هستند.

¹Ascii

در ادامه جدول ۱.۴ را برای معرفی اندیس‌های بکار گرفته شده در روابط مربوط به روش پیشنهادی اول معرفی می‌نماییم.

۳.۴ روش پیشنهادی برای مدیریت داده‌های از دست رفته

در این بخش به معرفی روش پیشنهادی پرداخته خواهد شد. در ابتدا به دلیل وجود داده‌های از دست رفته در پایگاه داده‌های مورد استفاده به بررسی و رویکرد حل این مشکل خواهیم پرداخت و در ادامه پس از معرفی روش‌های پیشنهادی برای آن و هر کدام از آن‌ها را به طور مفصل شرح خواهیم داد. همان‌گونه که در داده‌های حقیقی مشاهده شد در پایگاه داده‌های حقیقی ما با اطلاعات از دست رفته مواجه هستیم و به همین دلیل نیز سعی کردیم تا در پایگاه داده مجازی تولید شده نیز به مشابه داده‌های حقیقی، شامل اطلاعات از دست رفته باشد. در این بخش به رویکرد روش محاسبه استاتیک برای مدیریت این داده‌های از دست رفته می‌پردازیم و در بخش بعد به معرفی روشی برای بدست آوردن درخت فیلوزنی پرداخته خواهد شد. همان‌گونه که در ادامه بررسی خواهد شد، این اطلاعات از دست رفته در پایگاه داده‌های مختلف نرخ‌های متفاوتی دارد که تاثیر این تغییرات نیز در روشی پیشنهادی بررسی خواهد شد.

۱.۳.۴ روش محاسبه استاتیک

در این روش قصد داریم تا به یکباره بتوانیم مقادیر مناسب برای داده‌هایی که از دست رفته‌اند را تخمین بزنیم. در این روش باید توجه شود که ما لزوماً به دنبال جایگذاری مقدار از دست رفته با مقدار درست واقعی نیستیم. اگرچه چنین بیانی در نگاه اول ممکن است تعجب‌آور باشد اما با دقت بیشتر متوجه خواهیم شد که ما در آینده برای خطاهای موجود در پایگاه داده مدل‌سازی‌های محدودی داریم. مدل‌هایی که بهترین آن‌ها نیز ممکن است با واقعیت نویز افزوده شده به دادگان متفاوت باشد. در نتیجه اگر مطمئن بودیم که تمام داده‌هایی که موجود می‌باشند بدون خطا هستند در آن صورت ما نیز به دنبال یافتن جایگذاری با مقدار واقعی بودیم اما در حال حاضر که درصدی از داده‌های در دسترس خود همراه با خطای می‌باشند، ما به دنبال جایگذاری‌ای هستیم که بتواند در مجموع با مدل‌سازی خطای که در نظر می‌گیریم بیشترین سازگاری را داشته باشد کما اینکه ممکن است در حقیقت جایگذاری اشتباهی انجام داده باشیم. حال با توجه به توضیحی که بیان شد به تشریح این روش می‌پردازیم.

با توجه به فرض مدل مکان‌های بی‌نهایت می‌دانیم که جهش‌های اتفاق افتاده در والد در تمامی نسل‌های

آنده باقی خواهد ماند. بنابرین اگر تمامی جهش‌های نمونه (سلول) a در نمونه‌ای دیگر مانند b قرار داشته باشد، بنابرین می‌توان نتیجه گرفت که a یکی از اجداد b خواهد بود. همین فرضیه هسته اصلی روش پیشنهادی درنظر گرفته شده را تشکیل می‌دهد. بنابرین اگر جهش i در سلول a از دست رفته است، با توجه به اینکه آن جهش در سلول b چه وضعیتی دارد می‌توان تصمیم‌گیری کرد. اگر $= b(i)$ باشد، در این صورت (i) حتماً باید باشد و گرنه فرض اولیه مدل مکان‌های بینهایت نقض خواهد شد. اما اگر $= 1(b)$ باشد، آنگاه نتیجه خاصی نمی‌توان گرفت و باید به دنبال نمونه والد a یعنی نمونه d باشیم. حال اگر $= 1(d(i))$ باشد، آنگاه (i) حتماً باید باشد. اما اگر $= d(i)$ بود آنگاه انتخاب هر مقداری برای (i) تقریباً آزاد خواهد بود زیرا با فرض اولیه تناقضی ندارد و اینکه ساختار فیلوزنی را تغییر نمی‌دهد. اما از آنجایی که خود داده‌های در دسترس شامل خطایی باشند و هر نمونه‌ای که حاوی اطلاعات از دست رفته است لزوماً یک نواده یا یک والد ندارد، مجموعه‌ای از سلول‌های فرزند یا والد خواهند بود که متناسب با پارمترهای خطایی که در نظر می‌گیریم و فاصله‌زنی ای که دارند می‌توانند در تصمیم‌گیری تاثیرگذار باشند. صورت دقیق‌تر توضیحات داده شده را می‌توان به صورت فرمولی که در ادامه آمده است به نمایش درآورد.

در ابتدا تابعی به نام $F_s(D_{ij})$ تعریف می‌کیم که به نوعی با توجه به ارزشی که به سلول‌های نواده شده از سلول j می‌دهد سعی دارد تا اطمینان \circ بودن داده از دست رفته D_{ij} را بیان کند. برای محاسبه این تابع می‌دانیم که ابتدا سلول‌های مختلف با توجه به احتمال نواده بودنشان باید رتبه‌بندی شوند و وزن بگیرند. پس از آن هر سلول متناسب با ارزش تاثیرگذاری خود می‌تواند در مورد جایگاه جهش i برای سلول j نظر دهد.

$$F_s(D_{ij}) = \sum_{n \in \mathcal{N}} (1 - D_{mj}) \prod_{m=1}^M W(D_{mn}, D_{mj}) \quad (1.4)$$

در فرمول ۱.۴ مجموعه \mathcal{N} برابر با مجموعه سلول‌های متمایز از هم است. زیرا که در بسیاری از پایگاه‌داده‌ها از یک نمونه سلول ممکن است چندین نمونه وجود داشته باشد که وجود آن‌ها باعث بایس در محاسبات ما خواهد شد. همچنین تابع $W_s(c, p)$ به ارزش‌دهی جهش c در برابر p به عنوان نواده بودن می‌پردازد که در فرمول ۲.۴

تعریف شده است.

$$W(c, p) = \begin{cases} 1 & \text{if } c = 1, p = 1 \\ 1 - \xi & \text{if } c = 1, p = 0 \\ 0 & \text{if } c = 0, p = 1 \\ 1 & \text{if } c = 0, p = 0 \end{cases} \quad (2.4)$$

مقدار یک عددی بین (۰، ۱) است که پارامتری در جهت میزان ارزش دهی به نوادگان با فواصل مختلف می باشد.

هرچه این عدد بزرگتر باشد به معنی کم ارزش تر شدن نوادگان با فواصل بیشتر است و بلافاصله.

به همین صورت برای اولاد سلول j نیز می توان مشابه حالت قبل عمل کرد که روابط آن به صورت فرمول ۳.۴ خواهد شد.

$$F_a(D_{ij}) = \sum_{n \in \mathcal{N}} D_{mj} \prod_{m=1}^M W(D_{mj}, D_{mn}) \quad (3.4)$$

حال دو نکته در استفاده از روابط بالا باقی خواهد ماند.

نکته اول وجود داده های دیگر از دست رفته در محاسبه توابع است که به دو صورت می توان با آنها برخورد نمود.

رویکرد اول این است که در آنجایگاه ژنی از محاسبه آن خود داری شود و رویکرد دوم استفاده از از مقدار ۵٪

یا فراوانی نسبی آن جهش در محاسبات است که ما رویکرد اول را در این گزارش استفاده خواهیم کرد.

نکته دوم وجود خطأ در داده هاست. برای مدیریت این مشکل می توان با مدل سازی خطأ که به صورت فرمول ۴.۴ بیان می شود، برخورد کرد.

$$\begin{aligned} P(D_{ij} = 1 | E_{ij} = 0) &= \alpha, & P(D_{ij} = 0 | E_{ij} = 0) &= 1 - \alpha \\ P(D_{ij} = 0 | E_{ij} = 1) &= \beta, & P(D_{ij} = 1 | E_{ij} = 1) &= 1 - \beta \end{aligned} \quad (4.4)$$

پس از تعریف مدل سازی خطأ می توان روابط قبلی را مجددا به صورتی که در ادامه آمده است بازنویسی کرد.

$$W_e(c, p) = \sum_{i,j \in \{0,1\}} P(c|E_c = i)P(p|E_p = j)W(i, j) \quad (5.4)$$

که در این صورت توابع F_p و F_a نیز به صورت زیر همراه با مدل‌سازی خطاب بازتعریف خواهند شد.

$$\begin{aligned}\hat{F}_s(D_{ij}) &= \sum_{n \in \mathcal{N}} [1 - D_{mj}(1 - \alpha)] \prod_{m=1}^M W_e(D_{mn}, D_{mj}) \\ \hat{F}_a(D_{ij}) &= \sum_{n \in \mathcal{N}} D_{mj}(1 - \beta) \prod_{m=1}^M W_e(D_{mj}, D_{mn})\end{aligned}\quad (6.4)$$

حال پس از محاسبه مقادیر \hat{F}_s و \hat{F}_a می‌توان در مورد داده نامعلوم D_{ij} به صورت فرمول ۷.۴ تصمیم گرفت.

$$D_{ij} = \begin{cases} 0 & \text{if } \hat{F}_s \geq \hat{F}_a \\ 1 & \text{if } \hat{F}_s < \hat{F}_a \end{cases} \quad (7.4)$$

همچنین با کمی دقت در فرمول‌بندی انجام شده اگر برای تمام j, i ‌های ماتریس D این مقادیر توابع \hat{F} محاسبه شوند، خود می‌توانند معیاری برای ارزیابی پایگاهداده در دسترس و احتمال درستی فرض مدل مکان‌های بی‌نهایت باشند.

۱.۱.۳.۴ تصادفی

پر کردن کاملاً تصادفی میس‌ها

جدول ۱.۴: اندیس‌های به کار رفته در روابط روش پیشنهادی اول

ماتریس داده نویزی در دسترس که مقادیر ۰ و ۱ در آن قرار دارد	D
ماتریس داده حقیقی بدون نویز که به دنبال آن هستیم	E
درخت فیلوزنی جهش‌ها	T
بردار انتصابات	σ
ماتریس متناظر درخت	X_T
تعداد سلول‌های نمونه	N
تعداد جهش‌ها	M
نرخ خطای مثبت کاذب	α
نرخ خطای منفی کاذب	β

جدول ۲.۴: پارامترهای مدل ریاضی

زمان خدمت‌دهی به بیمار در مرحله k ام	t_{ik}
زمان فاری خدمت‌دهی به بیمار در محله k ام	\tilde{t}_{ik}
مقدار بدینانه (حداکثر) برای زمان خدمت‌دهی به بیمار در مرحله k ام	t_{ik}^p
محتمل‌ترین مقدار برای زمان خدمت‌دهی به بیمار در مرحله k ام	t_{ik}^m
مقدار خوشبینانه (حداقل) برای زمان خدمت‌دهی به بیمار در مرحله k ام	t_{ik}^o

جدول ۳.۴: متغیرهای مدل ریاضی

متغیر صفر-یک تخصیص بیمار به تخت/اتاق عمل	X_{ild_k}
زمان شروع خدمت‌دهی به بیمار	S_{ild_k}
متغیر صفر-یک توالی بیماران	Y_{ijkl_k}
متغیر صفر-یک تخصیص جراح به بیمار	V_{ni}

۴.۴ روش پیشنهادی اول (درخت‌بازی)

۱.۴.۴ پیش‌پردازش

قبل از شروع باید بر روی داده‌ها یک پیش‌پردازش اعمال کنیم که وابسته به سیاست درنظر گفته شده می‌تواند باعث تغییر در پاسخ نهایی نیز شود. به این منظور داده‌هایی که miss شده‌اند با روش‌های زیر می‌تواند برای ورود به مرحله بعد تخمین زده شود.

فصل ۵

نتایج تجربی

۱.۵ پایگاه داده‌های ورودی

قبل از اینکه وارد روش پیشنهادی شویم به تشریح وردی‌های مسئله و داده‌هایی که مورد استفاده قرار خواهیم داد می‌پردازیم. داده‌های ورودی برابر ماتریس $D_{m \times n}$ می‌باشد که بعد اول M برابر با زن‌ها و بعد دوم N برابر سلول‌های نمونه‌برداری شده می‌باشد. در هر خانه $d_{i,j}$ یک بردار داده قرار دارد که حاوی اطلاعات زن j در سلول i می‌باشد.

۱.۱.۵ پایگاه داده مصنوعی^۱

با توجه به این نکته که از درخت فیلوژنی حقیقی^۲ داده‌های حقیقی موجود اطلاعی نداریم، به سراغ ساخت پایگاه داده مصنوعی می‌رویم. با استفاده از این پایگاه داده مصنوعی می‌توانیم در مورد روش‌هایی که در ادامه بیان خواهیم کرد یک معیار ارزیابی نسبتاً مناسبی داشته باشیم و تا حدودی از مشکلات روش‌های پیشنهادی آکاه شویم و به تصحیح آن بپردازیم. برای ساخت پایگاه داده مصنوعی که همان ماتریس ورودی $D_{m \times n}$ می‌باشد، از دو روش مختلف با دو فرض مختلف استفاده خواهیم کرد که در ادامه به تشریح هر کدام خواهیم پرداخت. برای ایجاد پایگاه داده در این حالت ابتدا درختی تصادفی با پارامترهای n ،^۳ ایجاد می‌کنیم که n تعداد زن‌ها (جهش‌ها) بوده و ζ عددی در بازه $(0, \infty)$ است که یک پارامتر کنترلی است که وظیفه اش کنترل کلی تعداد

¹Synthetic Dataset

²Ground-truth Phylogeny Tree

نسلهای مختلف را از یک جمعیت در درخت فیلوزنی می‌باشد. حال برای تولید پایگاه داده مصنوعی به ترتیب سه گام زیر باید انجام شود.

- ایجاد یک درخت فیلوزنی تصادفی
 - تبدیل درخت فیلوزنی به ماتریس اطلاعات سلول-ژن (E)
 - اضافه کردن نویز به ماتریس E و تبدیل آن به ماتریس نویزی D
- در ادامه هر بخش به صورت جداگانه به تفضیل شرح داده خواهد شد.

۱.۱.۱.۵ ساخت درخت تصادفی

برای ساخت درخت تصادفی از دو روش مختلف استفاده شده است که هرکدام جداگانه توضیح داده شده است.

روش اول: با استفاده از درخت تصادفی دودویی ژنولوژی^۳

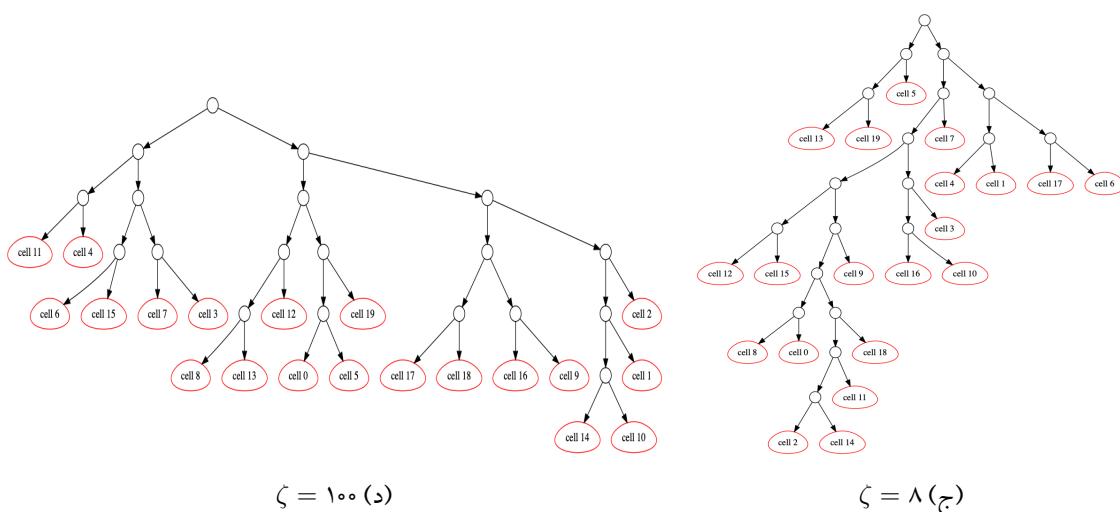
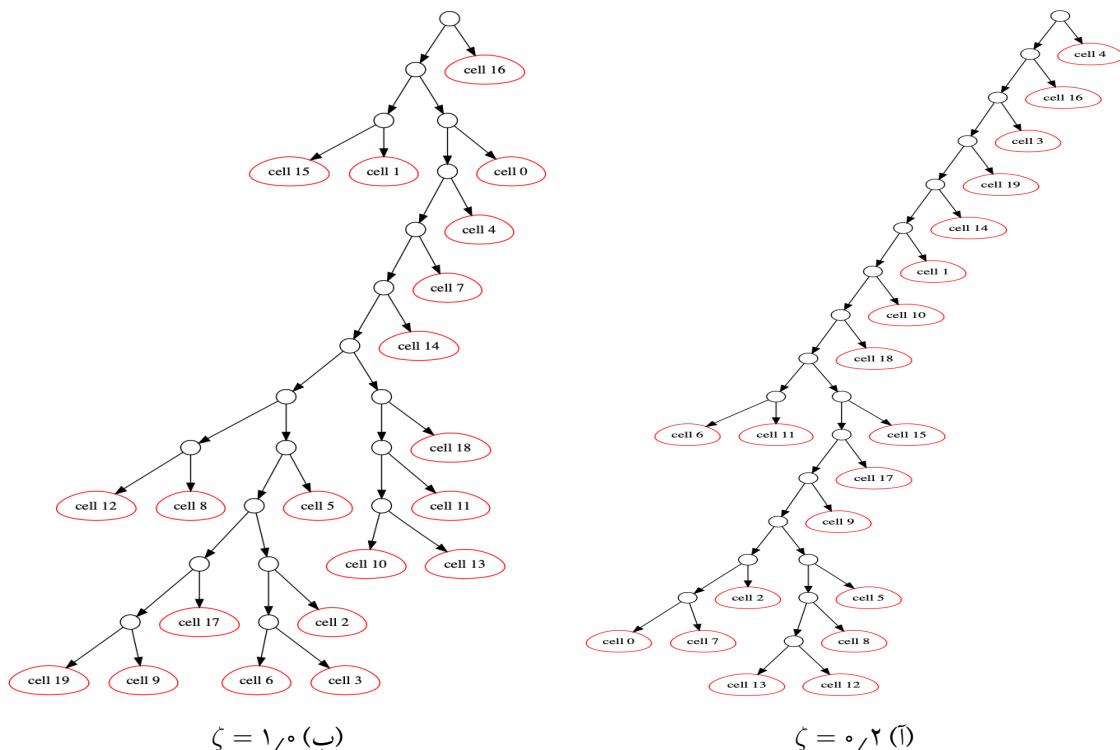
در این روش همان‌گونه که از نام آن مشخص است با استفاده از درخت تصادفی دودویی ژنولوژی به ساخت ماتریس داده ورودی مسله می‌پردازیم که برای ساخت این دادگان از فرض‌های که در ادامه آمده است استفاده خواهیم کرد.

در مرحله اول که ساخت درخت است به این صورت عمل می‌کنیم که به تعداد n گونه (سلول) در نظر می‌گیریم. سپس به ترتیب مراحل زیر را انجام می‌دهیم تا به درخت تصادفی مورد نظر برسیم.

- به هر کدام از n گونه متمایز در ابتدا وزن $1 = w_i$ را اختصاص می‌دهیم که متناسب با احتمال انتخاب هر گونه در مراحل بعدی خواهد بود.
- برای هر گونه i تابع جرم احتمال را در ادامه به صورت $F_i = \frac{w_i}{\sum_{i=1}^n w_i}$ در نظر می‌گیریم
- با استفاده از F دو گونه متمایز v, u را انتخاب می‌کنیم و به هم متصل می‌کنیم
- به جای دو گونه v, u یک گونه جدید uv با وزن $w_{uv} = \frac{w_u + w_v}{\sqrt{\zeta}}$ را قرار می‌دهیم.
- تعداد گونه‌ها یک واحد کم شده است. بررسی می‌کنیم اگر تعداد گونه‌های باقی‌مانده از ۲ کمتر باشد درخت تصادفی ساخته شده است و پایان کار است. در غیر این صورت به مرحله اول بازمی‌گردیم.

³Random Binary Genealogical Tree

پارامتر ζ به گونه‌ای کنترل کننده میزان ناپایداری در طی نسل‌ها می‌باشد. بطوریکه نمونه‌ای از نتایج مقادیر مختلف آن برای $n = 20$ در شکل ۱.۵ آورده شده است. پس از ساخت درخت تصادفی به سراغ مرحله بعد یعنی تبدیل درخت به ماتریس ژن-سلول E می‌رویم.



شکل ۱.۵: درخت فیلوزنی تصادفی تولید شده برای $n = 20$ و ζ ‌های مختلف

در ادامه با توجه به اینکه تعداد دلخواه جهش‌ها چه عددی بوده است یکی از گام‌های زیر را برمی‌داریم.

- اگر تعداد جهش‌ها $N > M$ بوده باشد در آن صورت به صورت تصادفی به تعداد دفعات اختلاف یکی از انشعاب‌ها در درخت را به صورت تصادفی انتخاب کرده و آن جهش اضافه شده را تا تمامی نوادگان پیش خواهیم برد.
- اگر تعداد جهش‌ها $N < M$ بوده باشد آنگاه مجدداً به اندازه تعداد اختلاف انشعاب‌هایی را انتخاب کرده و این بار جهش در آن انشعاب را تا تمامی نوادگان حذف می‌کنیم.

به این ترتیب تمامی سلول‌ها را با تعداد جهش‌های انتخابی خواهیم داشت. در نهایت برای اخیرین تغییر در جهش‌ها می‌توان یک گام دیگر برداشت که آن تولید یه عدد تصادفی کوچکتر از $\frac{M}{2}$ است که به آن تعداد می‌توان جهش‌های موجود را از انشعابی برداشت و بر روی انشعابی دیگر قرار داد. با این کار ممکن است تعداد جهش‌ها در انشعاب‌های مختلف تغییر کند و چه بسا به مدل‌های واقعی نزدیکتر شود که البته در این پایان‌نامه از گام آخر صرف نظر کرده‌ایم.

حال کار ما با پخش تصادفی جهش‌ها در پایگاه‌داده مجازی پایان یافته است. تا به اینجا ما در فرض خود از هر نمونه جمعیت مختلف یک سلول داشته‌ایم. اما در بعضی مواقع در پایگاه داده‌های واقعی ممکن است از یک جمعیت بیش از یک نمونه وجود داشته باشد که البته این امر لزوماً درست نیست به این دلیل که بعد از افزوده شدن نویز به داده‌ها ممکن است برخی سلول‌ها جهش‌هایشان مسابه هم شود. اما به هر حال اگر چنین چیزی را بخواهیم که داشته باشیم با انتخاب تصادفی برخی سلول‌ها (برگ‌ها) در درخت و کپی کردن آن‌ها می‌توان به چنین مقصودی رسید.

روش دوم: با استفاده از درخت تصادفی جهش‌های ^۴ ثانی

این روش نیز تا حدود زیادی مشابه روش قبل است با این تفاوت که در اینجا به جای اینکه درخت تصادفی را با توجه سلول‌ها از پایین به بالا بسازیم، ابتدا یک درخت تصادفی بدون در نظر گرفتن سلول‌ها ایجاد می‌کنیم و سپس به تخصیص جهش‌ها به آن می‌پردازیم و در نهایت برای آخرین مرحله به تعداد دلخواه سلول را به درخت اضافه کرده و درخت را تکمیل می‌کنیم. در گام اول به تعداد $1 + M$ نود در نظر می‌گیریم. مشابه حالت قبل با طی مراحلی بکه در ادامه آمده است به ساختار یک درخت تصادفی می‌رسیم.

- به هر کدام از m نود متمایز در ابتدا وزن $w_i = 1$ را اختصاص می‌دهیم که متناسب با روند حرکتی تومور به سمت آن جهش‌ها در مراحل بعدی خواهد بود.

⁴Random Mutation History Tree

- برای هر نود v تابع جرم احتمال را در ادامه به صورت $F_v = \sum_{i=1}^{w_v} \frac{w_i}{w_v}$ بیان می‌شود در نظر می‌گیریم.
- با استفاده از F دو نود متمایز u, v را انتخاب می‌کنیم و به هم متصل می‌کنیم.
- به جای دو گونه u, v یک نود جدید uv با وزن $w_{uv} = \sqrt{\zeta} \cdot \frac{w_u + w_v}{\zeta}$ را قرار می‌دهیم.
- تعداد نودها یک واحد کم شده است. بررسی می‌کنیم اگر تعداد نودهای باقیمانده از ۲ کمتر باشد به مرحله بعد می‌رویم و در غیر این صورت به مرحله اول بازمی‌گردیم.
- در این مرحله تمامی برگ‌های درخت ساخته شده را حذف می‌کنیم و تنها باقیمانده را به عنوان درخت تصادفی جهش‌ها در نظر می‌گیریم.

پس از به پایان رسیدن مراحلی که بیان شد درخت تصادفی آماده است و حال نوبت به تخصیص دادن خود ژن‌ها به هرکدام از این نودهای درخت است. برای این منظور به هرکدام از M نود یک ژن را به صورت تصادفی تخصیص می‌دهیم. پس از آن برای نهایی سازی درخت جهش‌ها از پارامتر دلخواه $A = [\gamma * (M - 1) / \gamma]$ استفاده می‌کنیم که γ عددی بین $(0, 1)$ است و A تعداد یال‌هایی است که در درخت باید برداشته شود و دو نود آن با یکدیگر ادغام شود. این کار باعث می‌شود تا در درخت جهش‌ها در برخی نودها به جای یک جهش چند جهش داشته باشیم که بتواند به مدل داده‌های واقعی نزدیکتر باشد.

پس از تکمیل درخت جهش‌ها نوبت قرار دادن نمونه‌هایی بر روی آن است. به همین منظور با فرض اینکه $N \geq M$ است. به تعداد M تا از سلول‌ها را به هر کدام از نودهای درخت جهش به عنوان برگ‌های جدید اضافه می‌کنیم و برای $N - m$ سلول باقیمانده همین کار را این‌بار به صورت تصادفی انجام می‌دهیم. در نهایت درخت تصادفی جهش‌ها ساخته شده است که نمونه‌ای از آن را در شکل ۲.۵ قابل مشاهده است.

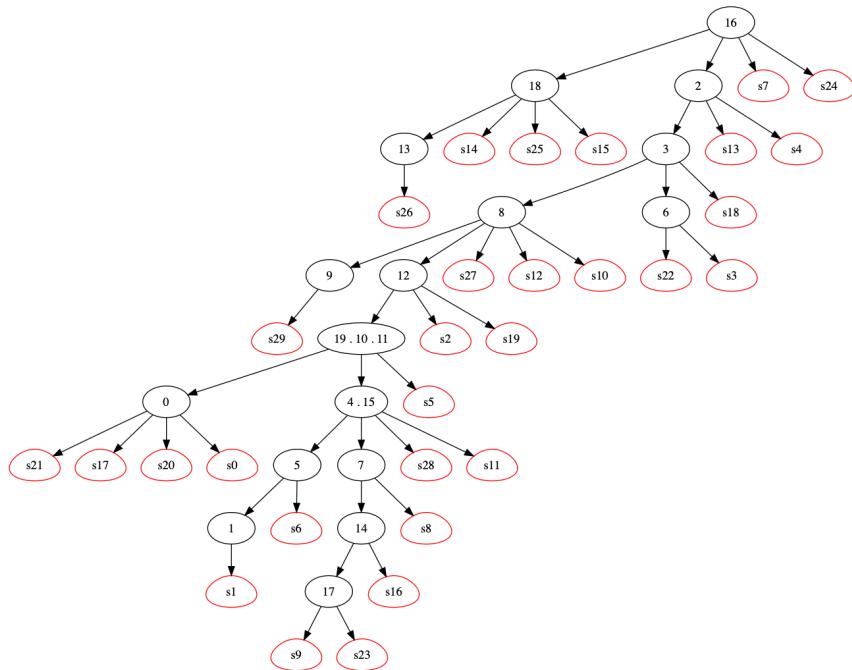
۲.۱.۱.۵ تبدیل درخت به ماتریس ژن-سلول

با داشتن درخت (تولید شده با هرکدام از روش‌ها تفاوتی ندارد) در ادامه از فرض‌های مختلف در تولید ماتریس E می‌توان استفاده کرد.

فرض مدل مکان‌های بی‌نهایت^۵

در این حالت فرض می‌کنیم که هر جهش اتفاق افتاده در درخت فیلوزنی در تمامی نسل‌های پس از آن باقی می‌ماند و هیچ‌گاه از بین نمی‌رود. در چنین حالتی درخت حاصل از این روش درختی یکتا بوده که به نام درخت

⁵Infinite Site Models



شکل ۲.۵: درخت جهش تصادفی با پارامترهای $N = ۳۰, M = ۲۰, \zeta = ۱, \gamma = ۰.۱۵$

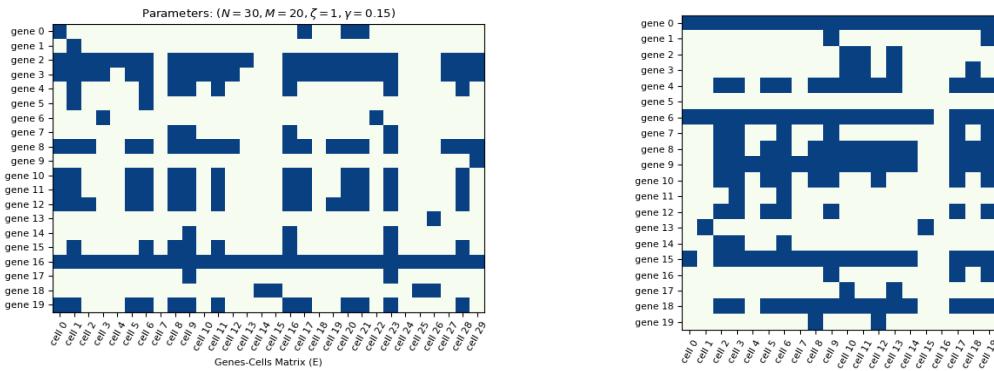
فیلوژنی کامل^۶ شناخته می‌شود.

در این قسمت باید با استفاده از درخت تصادفی تولید بتوانیم ماتریس جهش‌ها را برای سلول‌های مختلف با فرض مکان‌های بی‌نهایت بدست آوریم. در ابتدا ماتریس E را به ابعاد $M \times N$ ایجاد می‌کنیم و برای هر درایه j, i در آن که i شماره جهش و j شماره سلول است به صورت فرمولی که در ادامه آمده است مقداردهی می‌کنیم.

$$E_{i,j} = \begin{cases} 1 & \text{if mutation } i \text{ is an ancestor of cell } j \\ 0 & \text{o.w} \end{cases} \quad (1.5)$$

به این ترتیب با فرض مدل مکان‌های بی‌نهایت ماتریس بدون خطای E را داریم که برای تصاویر دو روش درخت مرحله قبل در شکل ۳.۵ بدست آمده‌اند.

⁶Perfect Phylogeny Tree



(ب) ماتریس درخت شکل ۲.۵

(ا) ماتریس درخت شکل ۱.۵ ب.

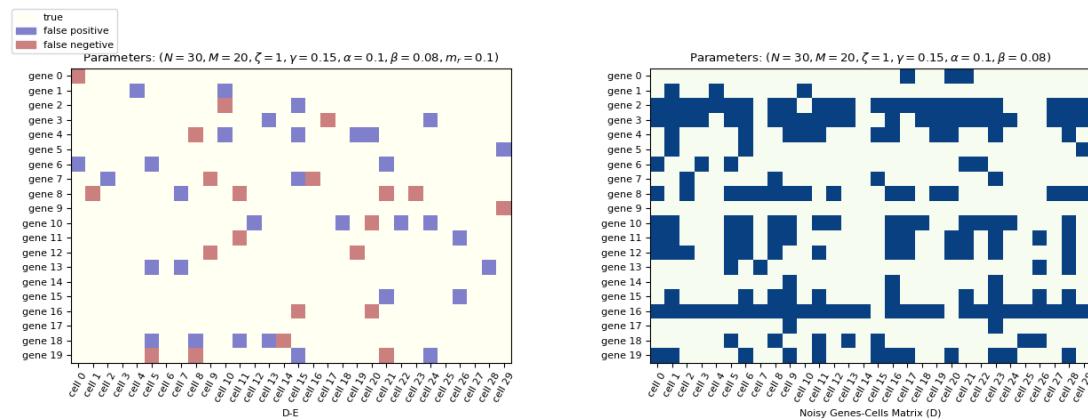
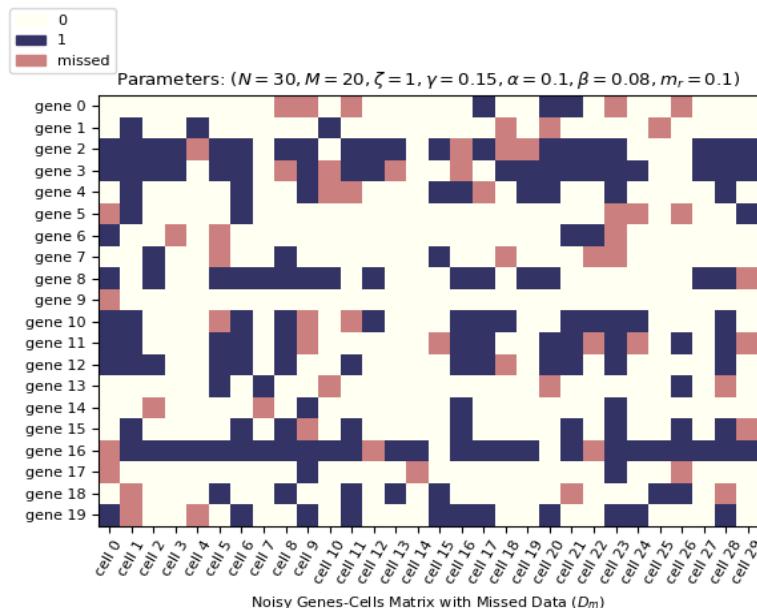
شکل ۳.۵: ماتریس‌های ژن-سلول (E) بدست آمده از درخت‌های تصادفی ساخته شده

۳.۱.۱.۵ اضافه کردن نویز به ماتریس ژن-جهش

برای قسمت نهایی آمده سازی پایگاه داده مجازی نیاز است تا به ماتریس E با پارامتر $\Theta = (\alpha, \beta, m_r)$ نویز اضافه کنیم و آن را به ماتریس D تبدیل کنیم که $\alpha = P(D_{ij} | E_{ij} = 0)$ و $\beta = P(D_{ij} | E_{ij} = 1)$ است و همچنین $m_r \in (0, 1)$ که نرخ داده‌های از دست رفته را مشخص می‌کند.

برای این منظور به ازای تمامی درایه‌های E هر بار یک عدد تصادفی با توزیع یکنواخت بین $(0, 1]$ بوجود می‌آوریم و اگر عدد تولید شده کوچکتر از α بود آنگاه ان درایه در ماتریس D را برابر با ۰ قرار می‌دهیم. به همین ترتیب مجدداً این بار برای درایه‌های E این کار را تکرار می‌کنیم و اگر عدد تصادفی تولید شده کوچکتر از β شد، درایه متناظر را در ماتریس D برابر با ۱ قرار می‌دهیم.

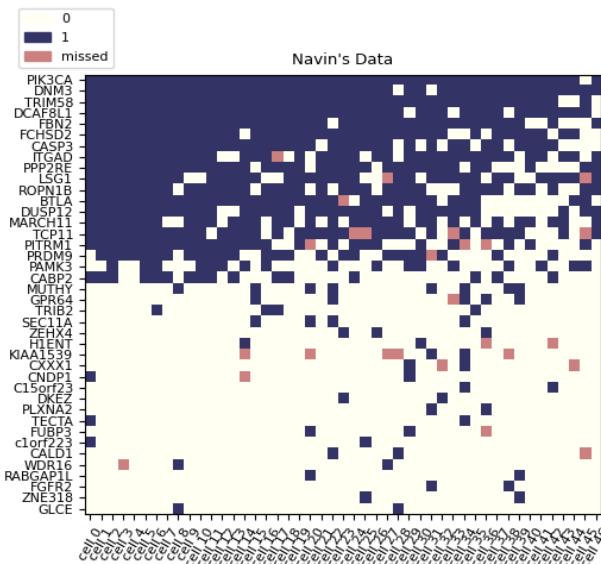
پس از اتمام کار نوبت به اضافه کردن داده‌های از دست رفته است. برای این منظور با نرخ m_r بعضی از درایه‌های ماتریس D را برابر با ۲ قرار می‌دهیم که به منزله در دسترس نبودن اطلاعات است. نام ماتریس نهایی را که شامل داده‌های از دست رفته است D_m می‌گزاریم. در ادامه تصاویر اضافه شدن نویز به ماتریس شکل ۳.۵ ب در شکل ۴.۵ آمده است.

(ب) نویزی اضافه شده با پارامترهای $\alpha = 0.1, \beta = 0.08$ (آ) ماتریس نویزی با $\alpha = 0.1, \beta = 0.08$ (ج) ماتریس نویزی به همراه دادهای از دست رفته با پارامترهای $\alpha = 0.1, \beta = 0.08, m_r = 0.1$

شکل ۴.۵: ماتریس‌های ژن-سلول همراه با نویز و داده‌های از دست رفته شکل ۳.۵ ب که برای ورودی مسله آماده شده است.

۲.۱.۵ پایگاه داده حقیقی^۷

به عنوان پایگاه داده حقیقی از پایگاه داده استفاده شده در مقاله SCITE به عنوان پایگاه داده حقیقی اصلی استفاده خواهیم کرد که ماتریس داده ورودی آن به صورت شکل ۵.۵ می‌باشد. همچنین پایگاه داده حقیقی Xu



شکل ۵.۵: داده‌های حقیقی Navin در مقاله SCITE

نیز که در مقاله SCITE مورد استفاده قرار گرفته است در شکل ۶.۵ آمده است.

⁷Real Dataset



شکل ۶.۵: داده‌های حقیقی Xu در مقاله SCITE

۲.۵ روش پیشنهادی بدست آوردن درخت فیلوزنی

پس از تخمین داده‌های از دست رفته، در این بخش به معرفی روش پیشنهادی برای یافتن درخت فیلوزنی می‌پردازیم. در روش‌های گذشته که رویکرد آن در ادامه بیان شده است به موفقیت نرسیدیم و این‌بار در نظر داریم تا با استفاده از یک درخت در ساختار شبکه ژن‌ها بتوانیم به یک درخت فیلوزنی مناسب دست یابیم.

- استفاده از شبکه ژن‌ها

۷ استفاده از یک گراف ابتدایی و سپس تغییر و هرس کردنش تا رسیدن به درخت جهش‌ها

۳ استفاده از یک درخت نمونه نابهینه و تغییر اتصالات تا رسیدن به درخت بهینه جهش‌ها

- استفاده از شبکه سلول‌ها

۷ استفاده از یک گراف سلول‌ها و بهینه‌کردن ارتباطات بین آن‌ها و سپس تبدیل آن به درخت فیلوزنی

در رویکرد اول ما از شبکه‌های ژنی استفاده خواهیم نمود. این شبکه‌ها نودهایی معادل با یک ژن متمایز را در نظر می‌گیرند. در گذشته با استفاده از شبکه‌ای کامل با وزن‌های متفاوت که بر حسب اطلاعات ورودی به الگوریتم تعیین می‌شد، متساقانه به موفقیت خاصی نرسیدیم. همچنین مشابه همین رویکرد را در ساختار شبکه‌های سلولی دنبال کردیم که مجدداً پیشرفت قابل ملاحظه‌ای حاصل نشد. به همین جهت این‌بار در این گزارش با تغییری اساسی به دنبال یافتن روشی مناسب برای استنتاج درخت فیلوزنی می‌باشیم.

۱.۲.۵ استفاده از شبکه ژن‌ها برای یافتن درخت فیلوژنی

در این رویکرد با استفاده از شبکه‌ای که نودهایی معادل ژن‌ها داشته باشد سعی داریم تا به درخت فیلوژنی

بهینه برسیم.

۱.۱.۲.۵ استفاده از یک درخت نمونه نابهینه و تغییر اتصالات تا رسیدن به درخت فیلوژنی بهینه

در این روش قصد داریم تا با شروع از یک درخت نمونه که در ابتدا به صورت تصادفی از اتصال ژن‌ها بوجود آمده است، به بهینه‌ترین درخت ممکن برسیم. این روش به صورت تکرارواره با تغییر اتصالات درخت سعی در بدست آوردن درختی مطلوب‌تر دارد که شرایط و روابط تاثیرگزار در آن به تفضیل شرح داده خواهد شد. در واقع این روش پیشنهادی یک جست‌وجوی حریصانه می‌باشد که طی شرایطی می‌توان انتظار داشت که به پاسخ بهینه دست یافته شود. این روش به نام روش زنجیره مارکو مونت-کارلو^۸ شناخته می‌شود که در بسیاری از مقالات مرتبط نیز مورد استفاده قرار گرفته شده است.

برای شروع یک درخت تصادفی T را با نودهایی معادل ژن‌های پایگاه داده ورودی در نظر می‌گیریم که در گام اول به صورت تصادفی ساخته شده است. در گام‌های بعدی یک نود n_1 را از درخت T به صورت تصادفی انتخاب می‌کنیم. سپس زیردرخت با ریشه این نود را از درخت کم می‌کنیم. حال در درخت باقی‌مانده یک نود دیگر n_2 را به صورت تصادفی انتخاب می‌کنیم و آن زیر درخت قبلی با ریشه n_1 را به n_2 متصل می‌کنیم و درخت جدید را T_n نام‌گذاری می‌کنیم. پس از آن با احتمال،

$$P = \min \left(1, \frac{Eng(T)}{Eng(T_n)} \right) \quad (2.5)$$

درخت جدید بدست آمده T_n را به عنوان نتیجه این گام می‌پذیریم و در غیر این صورت درخت این گام نیز همان درخت سابق T باقی خواهد ماند. در رابطه ۲.۵،تابع Eng برای یک درخت در واقع انرژی آن درخت را محاسبه می‌کند و ما به دنبال پایدارترین درخت هستیم که کمترین انرژی را داشته باشد. تعریف این تابع برای یک درخت به این صورت است که با توجه به نمونه‌هایی که در دادگان ورودی D قرار دارد و اینکه کدام ژن بالاتر یا پایین‌تر از دیگر ژن‌ها قرار دارد به درخت یک نمره انرژی منصوب می‌کند که به صورت فرمول ۳.۵ بیان می‌شود.

$$Eng(T) = ||E - \hat{E}|| \quad (3.5)$$

⁸Markov Chain Monte Carlo

که در اینجا \hat{E} ماتریس تخمین زده شده روش پیشنهادی با توجه به درخت نهایی بدست آمده خواهد بود. در واقع ماتریس E همان ماتریس صحیح بدون خطای خطا است که جهش‌های مختلف را به ازای سلول‌های مختلف مشخص می‌کند. هنگامی که ما درخت ساخته شده فرضی T را داشته باشیم می‌توانیم در دو گام به \hat{E} برسیم. توجه به این نکته ضروری است که اگر در واقعیت فرض ما که همان مکان بی‌نهایت بود کامل برقرار باشد و E را داشته باشیم، حتماً باید بتوانیم به درختی با $= 0$ $Eng(T)$ دست یابیم. اما از آنجایی که ما D را به عنوانی از تخمین E داریم بنابرین محاسبه خطای واقعی خواهد بود نه خود آن که در اصل به صورت فرمول ۴.۵ می‌شود.

$$Eng(T) \approx Err(T) = ||D - \hat{E}|| \quad (4.5)$$

می‌دانیم که هر نود از این درخت T یک مکان برای اتصال سلولی می‌تواند باشد که در این صورت معنی آن اینگونه خواهد بود که سلول ضمیمه شده به آن نود تمام جهش‌های والد خود را داشته است. بنابرین در گام اول نیاز است تا هر مکان از درخت مشخص شود که چه نمونه‌هایی می‌تواند تولید نماید. این اطلاع توسط ماتریس A مشخص می‌شود که به صورت زیر از روی درخت ساخته خواهد شد.

$$A_{i,j} = \begin{cases} 1 & \text{اگر } j = i \text{ یا جهش } i \text{ والد جهش } j \text{ باشد} \\ 0 & \text{در غیر این صورت} \end{cases} \quad (5.5)$$

حال در گام دوم کافی می‌توانیم با توجه به یک معیار بهترین انتخاب را برای ضمیمه کردن سلول‌ها (نمونه‌های) موجود به درخت داشته باشیم. به همین جهت در ماتریس D که هر ستون آن برابر با نمایش یک سلول است، می‌تواند با هر ستون از ماتریس A مقایسه شود و بهترین ستونی که از A انتخاب شود برابر با جایگاه مناسب ضمیمه شدن نمونه با مقداری خطای درخت T است. حال با توجه به اینکه فرض مکان‌های بی‌نهایت را داشتیم ماتریس E را به صورت زیر می‌سازیم.

$$\hat{E}_{i,j} = A_{i,\sigma_j} \quad (6.5)$$

که σ_i برابر با بهترین نود (زن) برای اتصال نمونه بردار d_j است که بهترین جایگاه به صورت فرمول زیر انتخاب

می شود.

$$\sigma_j = \arg \max_{x \in [1 \rightarrow M]} \sum_{i=1}^M \left[A_{i,x} D_{i,j} (1 - \beta) + (1 - A_{i,x}) (1 - D_{i,j}) (1 - \alpha) + A_{i,x} (1 - D_{i,j}) \beta + (1 - A_{i,x}) D_{i,j} \alpha \right] \quad (7.5)$$

حال با داشتن ماتریس \hat{E} می توان خطای درخت را محاسبه نمود و با هدایت mcmc طبق فرمول ۸.۵ به درخت بهینه T_{op} رسید.

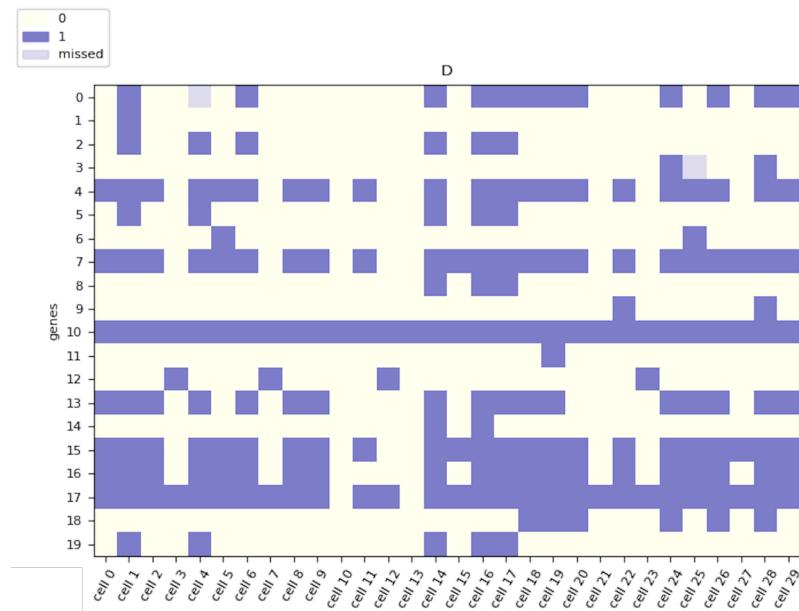
$$T_{op} = \min_{T \in \text{All possible } T} \left(\|D - \hat{E}_T\| \right) \quad (8.5)$$

۳.۵ نتایج تجربی

در این بخش به نتایج بدست آمده برای روش پیشنهادی می‌پردازیم و برای هر دو داده مصنوعی و حقیقی نتایج بدست آمده را تحلیل خواهیم نمود.

۱.۳.۵ نتایج بر روی پایگاه داده مصنوعی

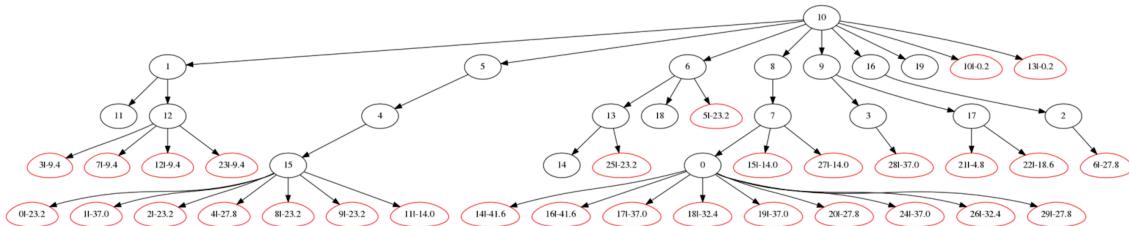
همان‌گونه که در بخش دوم توضیح داده شد با توجه به سختی دسترسی به پایگاه داده‌های حقیقی و اینکه در آن‌ها نیز حقیقت داده‌ها (E) وجود ندارد تصمیم به ایجاد پایگاه داده‌ای مصنوعی گرفته شد که با کمک آن بتوان ارزیابی مناسبی از روش پیشنهادی و میزان کارایی و مقاومت روش را نسبت به تغییر پارامترها سنجید. فرض کنید ماتریس ورودی شکل ۷.۵ را در اختیار داریم و میخواهیم بهترین درخت فیلوزنی را برای آن بیابیم.



شکل ۷.۵: نمونه‌ای تصادفی از ماتریس ورودی D

حال یک درخت تصادفی به صورت شکل ۸.۵ می‌سازیم. در درخت شکل ۸.۵ نمونه‌ها (سلول‌ها) با رنگ قرمز به درخت متصل شده‌اند که البته این ضمیمه بهترین ضمیمه ممکن است و میزان انرژی (خطای) هر ضمیمه نیز در کادر قرمز رنگ سلول‌ها به صورتی عددی منفی نوشته شده است. پس از این مرحله اگر ۳۰۰۰ گام MCMC را اجرا نماییم می‌توانیم نتیجه حاصله را در شکل ۹.۵ مشاهده کیم. در

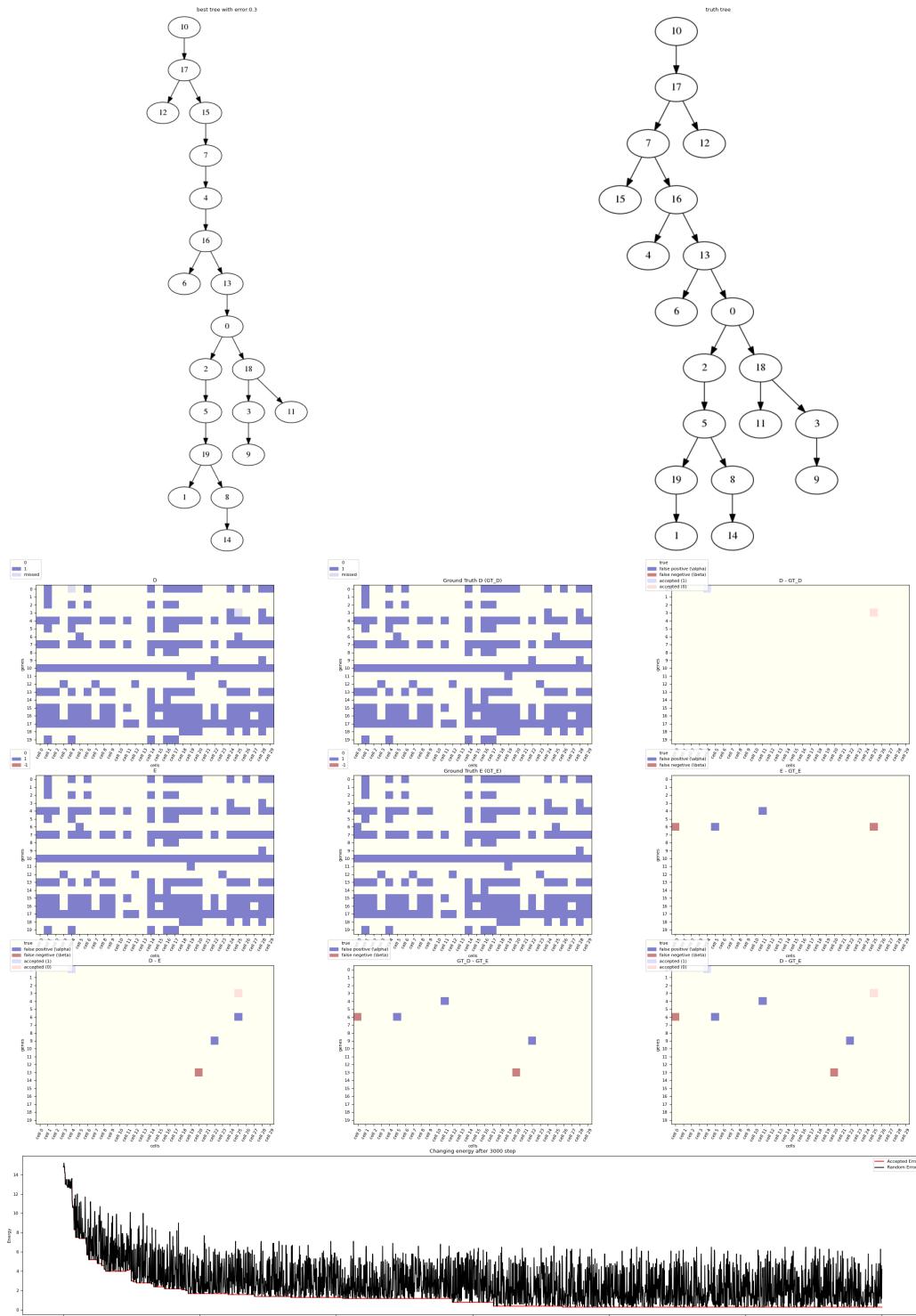
best tree with error:14.800000000000004



شکل ۷.۵: درخت تصادفی ایجاد شده به عنوان درخت اولیه شکل ۷.۵

این شکل دو درخت وجود دارد که درخت سمت راستی درخت حقیقی است که به دنبال آن بودیم و درخت سمت چپ بهترین درخت یافته شده است. همچنین در پایین شکل، ۹ ماتریس مشاهده می‌شود که ماتریس‌ها سمت راست و پایین به نوعی بیان‌کننده میزان خطای میان خطای میان ماتریس سمت چپ بالا می‌باشند. در بالای هر ماتریس نام آن نوشته شده است و در نهایت در انتهای تصویر نیز روند کاهش خطای تلاش‌های MCMC در گام‌های مختلف قابل مشاهده است. فقط نکته‌ای که وجود دارد این است که خطای نوشته شده در تصاویر برابر $1/\sqrt{0}$ مقیاس نوشته شده است.

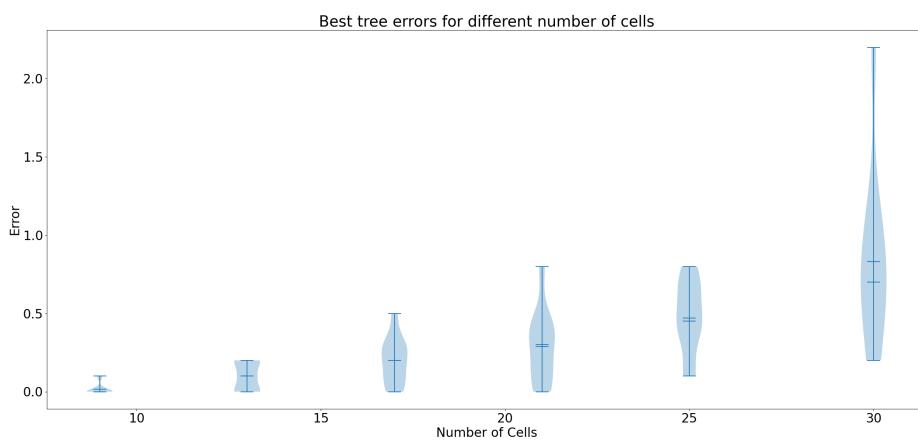
همان‌گونه که مشخص است در ماتریس D دوداده از دست رفته وجود دارد که یکی از آن‌ها در حقیقت جهش یافته و دیگر خیر. اگر ما در محاسبات خود این دوداده را در محاسبه خطای نظر نگیریم و با تغییر ۳ داده دیگر می‌توانیم به ماتریس \hat{E} (در شکل به نام E نوشته شده است) بررسیم که معادل بهترین درخت بدست آمده است. که این یعنی ماتریس D ما با ۵ تغییر بدست ما رسیده است. حال اگر حقیقت داده‌ها و درخت اصلی را مشاهده کنیم می‌بینیم که در آنجا نیز ۵ خطای وارد شده است که ۲ تای آن‌ها را درست کشف شده است. بنابرین الگوریتم بدون اطلاع از حقیقت توانسته با حداقل ۵ خطای یک درخت فیلوزنی مناسب دست بیابد که در ساختار نیز شباهات بسیار زیادی به حقیقت دارد. بنابرین روش پیشنهادی توانسته درخت فیلوزنی را با صحت $\frac{9916}{20*30} = 0.9916$ بازسازی کند که عددی قابل قبول می‌باشد.



شکل ۹.۵: نتیجه اجرای روش پیشنهادی برای ماتریس شکل ۷.۵

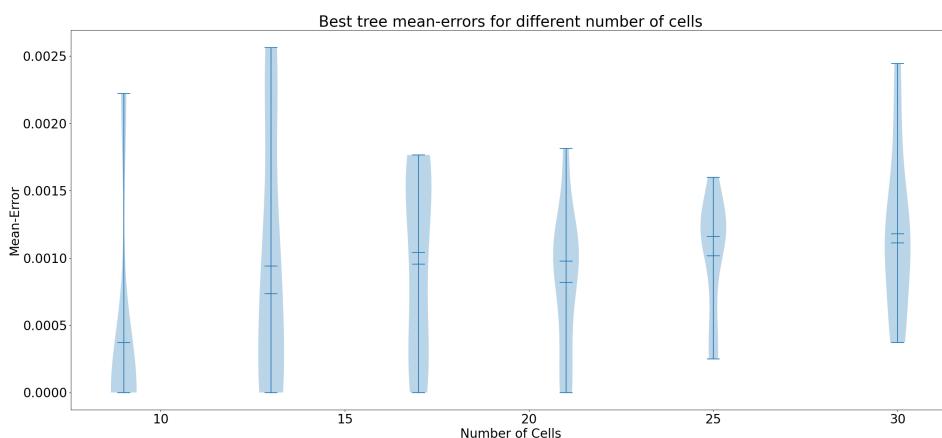
اما برای بررسی مناسب‌تر تعدادی تست را به ازای M و N ‌های مختلف اجرا نمودیم که به صورت خلاصه نتایج حاصل از آن در ادامه قابل مشاهده است.

در شکل ۱۰.۵ مقدار خطای درخت بهینه یافته شده قابل مشاهده است که نشان می‌دهد هر چه تعداد نمونه‌ها افزایش پیدا می‌کند و اندازه ماتریس ورودی بزرگ‌تر می‌شود، مقدار خطا نیز افزایش می‌یابد. در این اجرا تعداد جهش‌ها نیز عددی بین تعداد نمونه‌ها و نصف تعداد نمونه‌ها بوده است. حال برای اینکه متوجه شویم آیا این



شکل ۱۰.۵: نتیجه اجرای روش پیشنهادی برای تعداد نمونه‌های مختلف

افزایش خطا بخاطر ضعف روش پیشنهادی است یا ماهیت داده‌های ورودی میزان خطا را در هر اجرا بر تعداد خانه‌های ماتریس D تقسیم می‌کنیم که در آن صورت به نمودار شکل ۱۱.۵ می‌رسیم. در این نمودار جدید



شکل ۱۱.۵: نتیجه اجرای روش پیشنهادی برای تعداد نمونه‌های مختلف

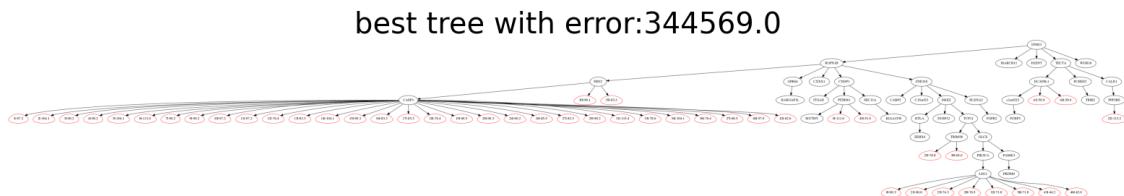
مشخص می‌شود که با افزایش اندازه ماتریس ورودی روش پیشنهادی سعی می‌کند تا خط را به ازای هر داده کنترل کند که نشان از کارآمدی روش پیشنهادی می‌باشد.

۲.۳.۵ نتایج بر روی داده‌های حقیقی

در این قسمت به ارائه گزارش و نتایج حاصل از روش‌های پیشنهادی با استفاده از داده‌های حقیقی برای بدست آوردن درخت فیلوزنی خواهیم پرداخت.

۱.۲.۳.۵ نتایج بهینه‌سازی درخت ثالث

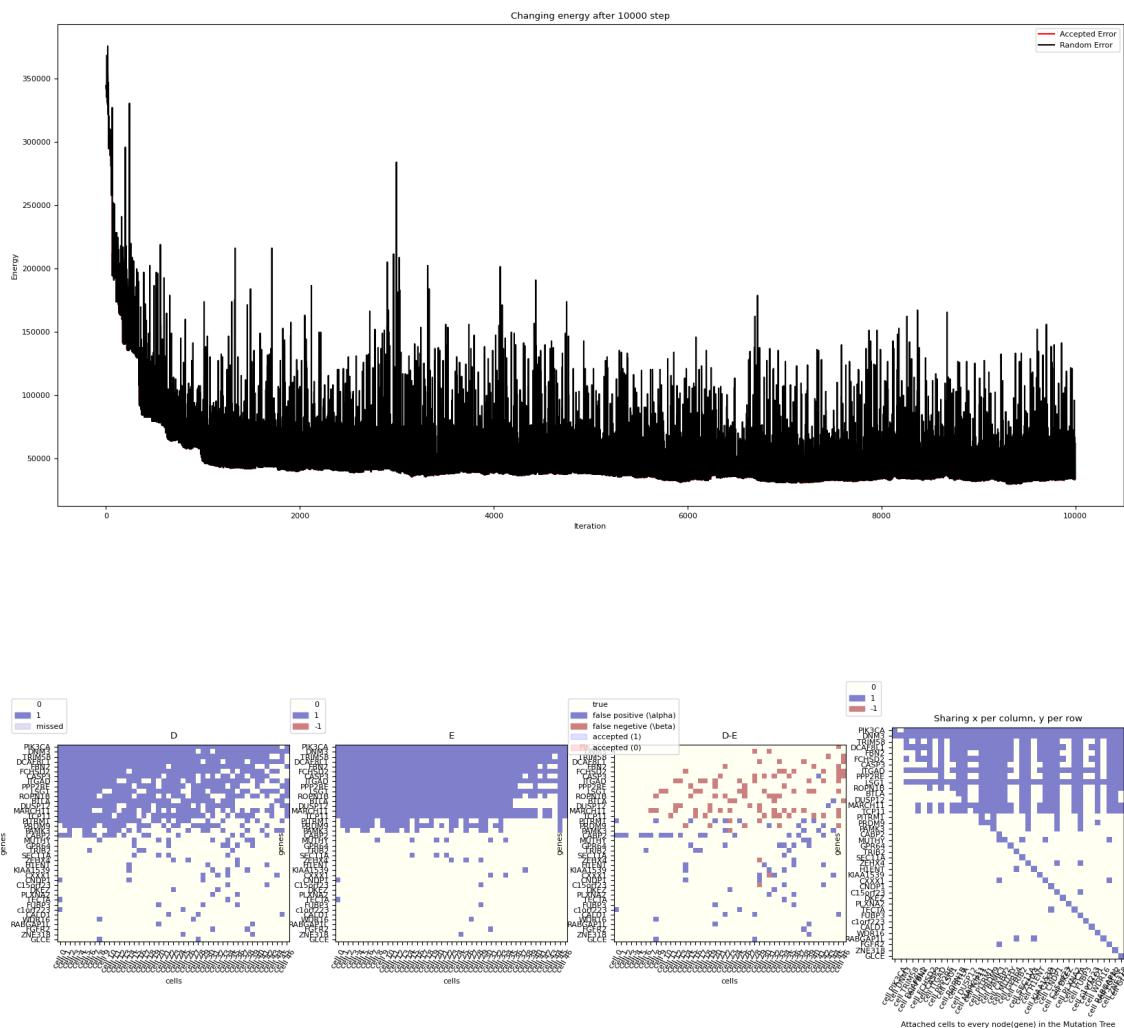
همان‌طور که در فصل قبل بیان شد، یکی از روش‌های بدست آوردن درخت فیلوزنی استفاده از یک درخت تصادفی نابهینه ثالثی بود که طی تکرار گام‌هایی سعی در تغییر اتصالات و یافتن درخت بهینه داشت که بتواند روند صحیح تغییرات ثالثی را در تومور مورد نظر نمایش دهد. نتیجه بدست آمده بر روی پایگاه داده حقیقی Navin به شرح زیر می‌باشد که عکس ۱۲.۵ درخت تصادفی اولیه الگوریتم را نشان می‌دهد که انرژی آن نیز بالای تصویر نوشته شده است.



شکل ۱۲.۵: درخت تصادفی اولیه

تصویر ۱۳.۵ نیز نمودار تغییر انرژی را طی گام‌های مختلف نمایش پیشنهادی مشخص می‌کند. در نهایت تصویر بهترین درخت یافته شده به همراه انرژی آن.

در نهایت برای مقایسه نیز تصویر درخت حاصله در مقاله اصلی SCITE را در شکل ۱۵.۵ نمایش داده شده است. همان‌طور که مشخص است مقدار انرژی بدست آمده برای خروجی الگوریتم پیشنهادی بهتر (کمتر) از انرژی درخت SCITE می‌باشد که دلیل بر بهینه‌تر بودن درخت روش پیشنهادی ارائه شده در این گزارش است.



شکل ۱۳.۵: نمودار تغییر انرژی در طی گام‌های مختلف

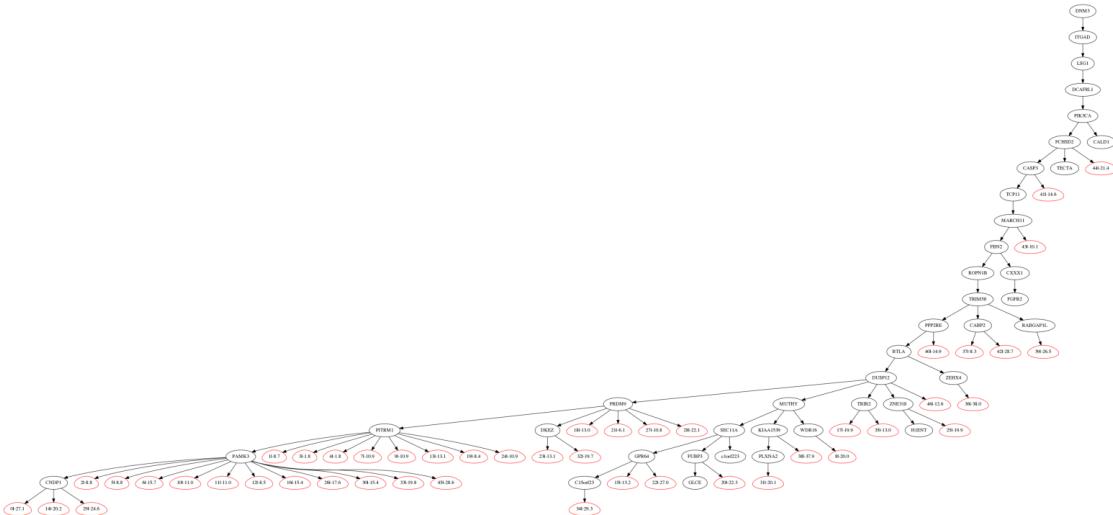
۴.۵ گام‌های آتی

در ادامه برای تکمیل روش پیشنهادی در دو قسمت نیاز به بهبود وجود دارد.
قسمت اول مربوط به درخت اولیه است و قسمت دیگر مربوط به سرعت MCMC می‌باشد.

۱.۴.۵ بهبود در ساخت درخت اولیه

در حال حاضر ما درخت اولیه را به صورت تصادفی انتخاب می‌کنیم که می‌توان در این مرحله درخت اولیه را با استفاده از مفروضات مدل مکان‌های بینهایت و با توجه به ماتریس ورودی بهبود بخشید. این کار باعث

best tree with error:29929.0



شکل ۱۴.۵: بهترین درخت یافته شده و خروجی الگوریتم برای مقاله SCITE

می شود تا شروع الگوریتم از نقطه بهتری باشد که در این صورت هم گام های لازم برای رسیدن به درخت بهینه می تواند کمتر شود و هم اینکه احتمال قرار گرفتن در نقاط اکسٹرمم نسبی را کاهش می دهیم.

۲.۴.۵ افزایش سرعت همگرایی MCMC

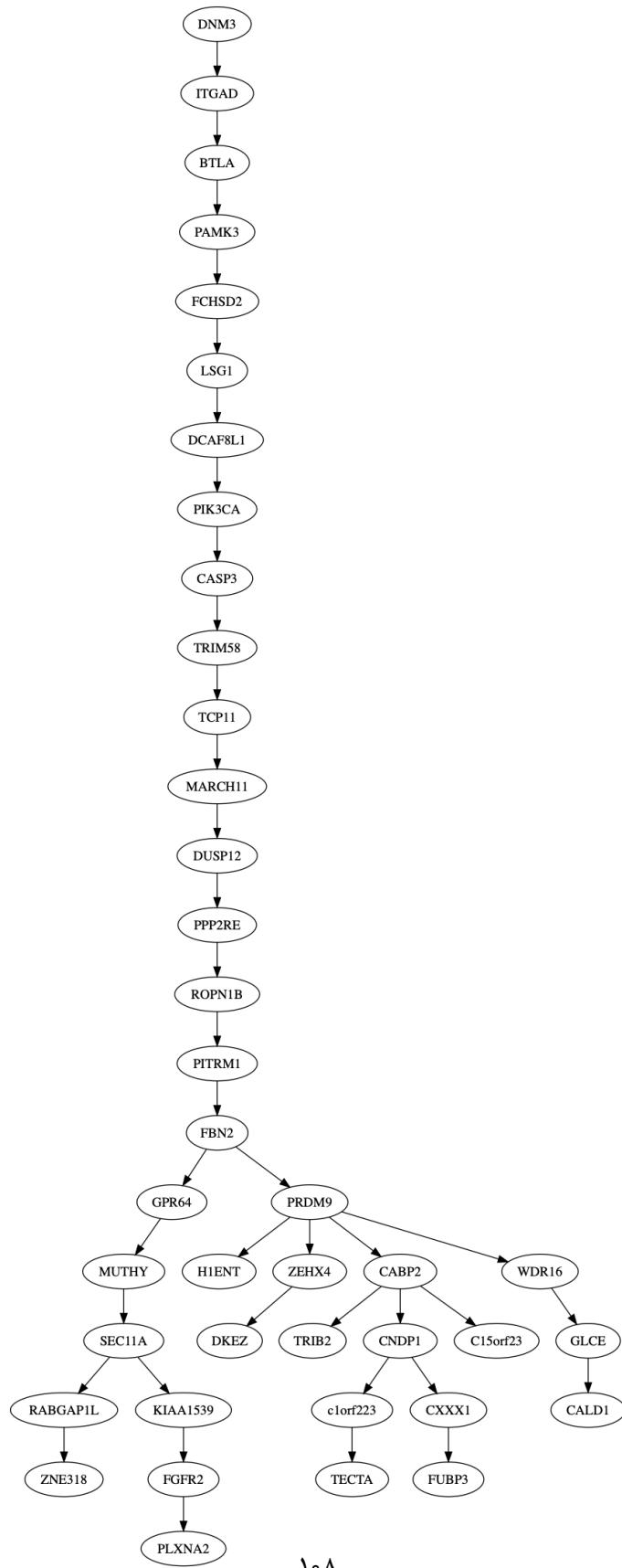
در این بخش نیاز است تا در دو قسمت روش پیشنهادی بهبود پاید.

۱.۲.۴.۵ تنوع در گام‌ها با استراتژی معقول

برای این بخش کاری که باید انجام شود این است که بتوان برای افزایش سرعت همگرایی از روش‌های مختلف در گام‌ها استفاده کرد. برای مثال در حال حاضر می‌توان از سه روش مختلف در هر گام استفاده نمود. روش اول تعویض دو نود در درخت می‌باشد. روش دوم جدایی یک زیر درخت و اتصال آن به محلی دیگر می‌باشد و در نهایت روش سوم تعویض دو زیر درخت با یکدیگر می‌باشد. با انتخاب یک استراتژی مناسب بین هر کدام از این روش‌ها در گام‌های مختلف احتمالاً بتوان سرعت همگرایی را افزایش داد.

۲.۲.۴.۵ قرار دادن احتمال وزن دار به ازای هر انتخاب

در حال حاضر ما در هر کدام از روش‌های مختلف که در بخش قبل برای گام‌های MCMC بیان کردیم، انتخاب نودها را به صورت کاملاً یکنواخت انجام می‌دهیم. در صورتی که احتمالاً بتوان با تعریف فرمولی مناسب این احتمال انتخاب بین نودهای مختلف در درخت را از حالت یکنواخت خارج کرد و در نتیجه مجدداً سرعت همگرایی الگوریتم را افزایش داد.



فصل ٦

بحث و نتیجه‌گیری

مراجع

- [1] Nci dictionary of cancer terms: somatic mutation definition. <https://www.cancer.gov/publications/dictionaries/cancer-terms/def/somatic-mutation?redirect=true>.
- [2] Ii neoplasms. 19 June 2014.
- [3] Cancer - activity 1 - glossary. page page 4 of 5, 2008.
- [4] Abrams, Gerald. Neoplasia i. 23 January 2012.
- [5] Akselrod-Ballin, Ayelet, Karlinsky, Leonid, Hazan, Alon, Bakalo, Ran, Horesh, Ami Ben, Shoshan, Yoel, and Barkan, Ella. Deep learning for automatic detection of abnormal findings in breast mammography. In *Deep learning in medical image analysis and multimodal learning for clinical decision support*, pages 321–329. Springer, 2017.
- [6] Alberts, Bruce, Johnson, Alexander, Lewis, Julian, Raff, Martin, Roberts, Keith, and Walter, Peter. Molecular biology of the cell 4th edition. New York: Garland Science, 1463, 2002.
- [7] Anderson, Kristina, Lutz, Christoph, Van Delft, Frederik W, Bateman, Caroline M, Guo, Yanping, Colman, Susan M, Kempski, Helena, Moorman, Anthony V, Titley, Ian, Swansbury, John, et al. Genetic variegation of clonal architecture and propagating cells in leukaemia. *Nature*, 469(7330):356–361, 2011.
- [8] Andor, Noemi, Graham, Trevor A, Jansen, Marnix, Xia, Li C, Aktipis, C Athena, Petritsch, Claudia, Ji, Hanlee P, and Maley, Carlo C. Pan-cancer analysis of the extent and consequences of intratumor heterogeneity. *Nature medicine*, 22(1):105–113, 2016.
- [9] Azer, Erfan Sadeqi, Ebrahimabadi, Mohammad Haghiri, Malikić, Salem, Khardon, Roni, and Sahinalp, S Cenk. Tumor phylogeny topology inference via deep learning. *Iscience*, 23(11):101655, 2020.

- [10] Beerenwinkel, Niko, Schwarz, Roland F, Gerstung, Moritz, and Markowetz, Florian. Cancer evolution: mathematical models and computational inference. *Systematic biology*, 64(1):e1–e25, 2015.
- [11] Behjati, Sam, Huch, Meritxell, van Boxtel, Ruben, Karthaus, Wouter, Wedge, David C, Tamuri, Asif U, Martincorena, Iñigo, Petljak, Mia, Alexandrov, Ludmil B, Gundem, Gunes, et al. Genome sequencing of normal cells reveals developmental lineages and mutational processes. *Nature*, 513(7518):422–425, 2014.
- [12] Birbrair, Alexander, Zhang, Tan, Wang, Zhong-Min, Messi, Maria Laura, Olson, John D, Mintz, Akiva, and Delbono, Osvaldo. Type-2 pericytes participate in normal and tumoral angiogenesis. *American Journal of Physiology-Cell Physiology*, 307(1):C25–C38, 2014.
- [13] Bishop, Christopher M. Pattern recognition. *Machine learning*, 128(9), 2006.
- [14] Burrell, Rebecca A and Swanton, Charles. Tumour heterogeneity and the evolution of poly-clonal drug resistance. *Molecular oncology*, 8(6):1095–1111, 2014.
- [15] Chen, Rui, Mias, George I, Li-Pook-Than, Jennifer, Jiang, Lihua, Lam, Hugo YK, Chen, Rong, Miriami, Elana, Karczewski, Konrad J, Hariharan, Manoj, Dewey, Frederick E, et al. Personal omics profiling reveals dynamic molecular and medical phenotypes. *Cell*, 148(6):1293–1307, 2012.
- [16] Ciregan, Dan, Meier, Ueli, and Schmidhuber, Jürgen. Multi-column deep neural networks for image classification. In *2012 IEEE conference on computer vision and pattern recognition*, pages 3642–3649. IEEE, 2012.
- [17] Cooper, Geoffrey M. *Elements of human cancer*. Jones & Bartlett Learning, 1992.
- [18] Davis, Alexander and Navin, Nicholas E. Computing tumor trees from single cells. *Genome biology*, 17(1):1–4, 2016.
- [19] de Visser, J Arjan GM and Rozen, Daniel E. Clonal interference and the periodic selection of new beneficial mutations in escherichia coli. *Genetics*, 172(4):2093–2100, 2006.
- [20] Dean, Frank B, Nelson, John R, Giesler, Theresa L, and Lasken, Roger S. Rapid amplification of plasmid and phage dna using phi29 dna polymerase and multiply-primed rolling circle amplification. *Genome research*, 11(6):1095–1099, 2001.
- [21] Demichelis, R, Retsky, MW, Hrushesky, WJM, Baum, M, and Gukas, ID. The effects of surgery on tumor growth: a century of investigations. *Annals of oncology*, 19(11):1821–1828, 2008.

- [22] Dentro, Stefan C, Leshchiner, Ignaty, Haase, Kerstin, Tarabichi, Maxime, Wintersinger, Jeff, Deshwar, Amit G, Yu, Kaixian, Rubanova, Yulia, Macintyre, Geoff, Vázquez-García, Ignacio, et al. Portraits of genetic intra-tumour heterogeneity and subclonal selection across cancer types. *BioRxiv*, page 312041, 2018.
- [23] Deshwar, Amit G, Vembu, Shankar, Yung, Christina K, Jang, Gun Ho, Stein, Lincoln, and Morris, Quaid. Phylogenetic reconstruction of subclonal composition and evolution from whole-genome sequencing of tumors. *Genome biology*, 16(1):1–20, 2015.
- [24] Devlin, Jacob, Chang, Ming-Wei, Lee, Kenton, and Toutanova, Kristina. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [25] Dhungel, Neeraj, Carneiro, Gustavo, and Bradley, Andrew P. Fully automated classification of mammograms using deep residual neural networks. In *2017 IEEE 14th International Symposium on Biomedical Imaging (ISBI 2017)*, pages 310–314. IEEE, 2017.
- [26] Donmez, Nilgun, Malikic, Salem, Wyatt, Alexander W, Gleave, Martin E, Collins, Colin C, and Sahinalp, S Cenk. Clonality inference from single tumor samples using low coverage sequence data. In *International Conference on Research in Computational Molecular Biology*, pages 83–94. Springer, 2016.
- [27] Eaton, Jesse, Wang, Jingyi, and Schwartz, Russell. Deconvolution and phylogeny inference of structural variations in tumor genomic samples. *Bioinformatics*, 34(13):i357–i365, 2018.
- [28] El-Kebir, Mohammed. Sphyr: tumor phylogeny estimation from single-cell sequencing data under loss and error. *Bioinformatics*, 34(17):i671–i679, 2018.
- [29] El-Kebir, Mohammed, Oesper, Layla, Acheson-Field, Hannah, and Raphael, Benjamin J. Reconstruction of clonal trees and tumor composition from multi-sample sequencing data. *Bioinformatics*, 31(12):i62–i70, 2015.
- [30] Fearon, Eric R and Vogelstein, Bert. A genetic model for colorectal tumorigenesis. *cell*, 61(5):759–767, 1990.
- [31] Fedele, Clare, Tothill, Richard W, and McArthur, Grant A. Navigating the challenge of tumor heterogeneity in cancer therapy. *Cancer discovery*, 4(2):146–148, 2014.
- [32] Fisher, Rosie, Pusztai, Lazos, and Swanton, C. Cancer heterogeneity: implications for targeted therapeutics. *British journal of cancer*, 108(3):479–485, 2013.

- [33] Friedl, Peter and Wolf, Katarina. Plasticity of cell migration: a multiscale tuning model. *Journal of Cell Biology*, 188(1):11–19, 2010.
- [34] Fukushima, Kunihiko. Neocognitron. *Scholarpedia*, 2(1):1717, 2007.
- [35] Gelman, Andrew, Shirley, Kenneth, et al. Inference from simulations and monitoring convergence. *Handbook of markov chain monte carlo*, 6:163–174, 2011.
- [36] Gerlinger, Marco, Rowan, Andrew J, Horswell, Stuart, Larkin, James, Endesfelder, David, Gronroos, Eva, Martinez, Pierre, Matthews, Nicholas, Stewart, Aengus, Tarpey, Patrick, et al. Intratumor heterogeneity and branched evolution revealed by multiregion sequencing. *N Engl j Med*, 366:883–892, 2012.
- [37] Greaves, Mel and Maley, Carlo C. Clonal evolution in cancer. *Nature*, 481(7381):306–313, 2012.
- [38] Halford, S, Rowan, A, Sawyer, E, Talbot, I, and Tomlinson, Ian. O6-methylguanine methyltransferase in colorectal cancers: detection of mutations, loss of expression, and weak association with g: C> a: T transitions. *Gut*, 54(6):797–802, 2005.
- [39] Hanahan, Douglas and Weinberg, Robert A. The hallmarks of cancer. *cell*, 100(1):57–70, 2000.
- [40] Hanahan, Douglas and Weinberg, Robert A. Hallmarks of cancer: the next generation. *cell*, 144(5):646–674, 2011.
- [41] Handa, Osamu, Naito, Yuji, and Yoshikawa, Toshikazu. Redox biology and gastric carcinogenesis: the role of helicobacter pylori. *Redox Report*, 16(1):1–7, 2011.
- [42] Hastings, W Keith. Monte carlo sampling methods using markov chains and their applications. 1970.
- [43] Hou, Yong, Song, Luting, Zhu, Ping, Zhang, Bo, Tao, Ye, Xu, Xun, Li, Fuqiang, Wu, Kui, Liang, Jie, Shao, Di, et al. Single-cell exome sequencing and monoclonal evolution of a jak2-negative myeloproliferative neoplasm. *Cell*, 148(5):873–885, 2012.
- [44] Hugo, Honor, Ackland, M Leigh, Blick, Tony, Lawrence, Mitchell G, Clements, Judith A, Williams, Elizabeth D, and Thompson, Erik W. Epithelial—mesenchymal and mesenchymal—epithelial transitions in carcinoma progression. *Journal of cellular physiology*, 213(2):374–383, 2007.

- [45] Husić, Edin, Li, Xinyue, Hujdurović, Ademir, Mehine, Miika, Rizzi, Romeo, Mäkinen, Veli, Milanič, Martin, and Tomescu, Alexandru I. Mipup: minimum perfect unmixed phylogenies for multi-sampled tumors via branchings and ilp. *Bioinformatics*, 35(5):769–777, 2019.
- [46] Jahn, Katharina, Kuipers, Jack, and Beerenwinkel, Niko. Tree inference for single-cell data. *Genome biology*, 17(1):1–17, 2016.
- [47] Kim, Kyung In and Simon, Richard. Using single cell sequencing data to model the evolutionary history of a tumor. *BMC bioinformatics*, 15(1):1–13, 2014.
- [48] LeCun, Yann, Bengio, Yoshua, and Hinton, Geoffrey. Deep learning. *nature*, 521(7553):436–444, 2015.
- [49] Lee, Kyung-Hwa, Lee, Ji-Shin, Nam, Jong-Hee, Choi, Chan, Lee, Min-Cheol, Park, Chang-Soo, Juhng, Sang-Woo, and Lee, Jae-Hyuk. Promoter methylation status of hmlh1, hmsh2, and mgmt genes in colorectal cancer associated with adenoma–carcinoma sequence. *Langenbeck's archives of surgery*, 396(7):1017–1026, 2011.
- [50] Liu, Yinhan, Ott, Myle, Goyal, Naman, Du, Jingfei, Joshi, Mandar, Chen, Danqi, Levy, Omer, Lewis, Mike, Zettlemoyer, Luke, and Stoyanov, Veselin. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- [51] Malikic, Salem, Jahn, Katharina, Kuipers, Jack, Sahinalp, S Cenk, and Beerenwinkel, Niko. Integrative inference of subclonal tumour evolution from single-cell and bulk sequencing data. *Nature communications*, 10(1):1–12, 2019.
- [52] Malikic, Salem, McPherson, Andrew W, Donmez, Nilgun, and Sahinalp, Cenk S. Clonality inference in multiple tumor samples using phylogeny. *Bioinformatics*, 31(9):1349–1356, 2015.
- [53] Malikic, Salem, Mehrabadi, Farid Rashidi, Azer, Erfan Sadeqi, Ebrahimabadi, Mohammad Haghiri, and Sahinalp, S Cenk. Studying the history of tumor evolution from single-cell sequencing data by exploring the space of binary matrices. *bioRxiv*, 2020.
- [54] Malikic, Salem, Mehrabadi, Farid Rashidi, Ciccolella, Simone, Rahman, Md Khaledur, Rickerts, Camir, Haghshenas, Ehsan, Seidman, Daniel, Hach, Faraz, Hajirasouliha, Iman, and Sahinalp, S Cenk. Phiscs: a combinatorial approach for subperfect tumor phylogeny reconstruction via integrative use of single-cell and bulk sequencing data. *Genome research*, 29(11):1860–1877, 2019.
- [55] McGranahan, Nicholas and Swanton, Charles. Clonal heterogeneity and tumor evolution: past, present, and the future. *Cell*, 168(4):613–628, 2017.

- [56] McPherson, Andrew, Roth, Andrew, Laks, Emma, Masud, Tehmina, Bashashati, Ali, Zhang, Allen W, Ha, Gavin, Biele, Justina, Yap, Damian, Wan, Adrian, et al. Divergent modes of clonal spread and intraperitoneal mixing in high-grade serous ovarian cancer. *Nature genetics*, 48(7):758, 2016.
- [57] Nik-Zainal, Serena, Van Loo, Peter, Wedge, David C, Alexandrov, Ludmil B, Greenman, Christopher D, Lau, King Wai, Raine, Keiran, Jones, David, Marshall, John, Ramakrishna, Manasa, et al. The life history of 21 breast cancers. *Cell*, 149(5):994–1007, 2012.
- [58] Nowell, Peter C. The clonal evolution of tumor cell populations. *Science*, 194(4260):23–28, 1976.
- [59] Ricketts, Camir, Seidman, Daniel, Popic, Victoria, Hormozdiari, Fereydoun, Batzoglou, Serafim, and Hajirasouliha, Iman. Meltos: multi-sample tumor phylogeny reconstruction for structural variants. *Bioinformatics*, 36(4):1082–1090, 2020.
- [60] Ross, Edith M and Markowetz, Florian. Onconem: inferring tumor evolution from single-cell sequencing data. *Genome biology*, 17(1):1–14, 2016.
- [61] Sabeh, Farideh, Shimizu-Hirota, Ryoko, and Weiss, Stephen J. Protease-dependent versus-independent cancer cell invasion programs: three-dimensional amoeboid movement revisited. *Journal of Cell Biology*, 185(1):11–19, 2009.
- [62] Sadeqi Azer, Erfan, Rashidi Mehrabadi, Farid, Malikić, Salem, Li, Xuan Cindy, Bartok, Osnat, Litchfield, Kevin, Levy, Ronen, Samuels, Yardena, Schäffer, Alejandro A, Gertz, E Michael, et al. Phiscs-bnb: a fast branch and bound algorithm for the perfect tumor phylogeny reconstruction problem. *Bioinformatics*, 36(Supplement_1):i169–i176, 2020.
- [63] Sakr, WA, Haas, GP, Cassin, BF, Pontes, JE, and Crissman, JD. The frequency of carcinoma and intraepithelial neoplasia of the prostate in young male patients. *The Journal of urology*, 150(2):379–385, 1993.
- [64] Salehi, Sohrab, Steif, Adi, Roth, Andrew, Aparicio, Samuel, Bouchard-Côté, Alexandre, and Shah, Sohrab P. ddclone: joint statistical inference of clonal populations from single cell and bulk tumour sequencing data. *Genome biology*, 18(1):1–18, 2017.
- [65] Satas, Gryte and Raphael, Benjamin J. Tumor phylogeny inference using tree-constrained importance sampling. *Bioinformatics*, 33(14):i152–i160, 2017.
- [66] Satas, Gryte, Zaccaria, Simone, Mon, Geoffrey, and Raphael, Benjamin J. Scarlet: Single-cell tumor phylogeny inference with copy-number constrained mutation losses. *Cell Systems*, 10(4):323–332, 2020.

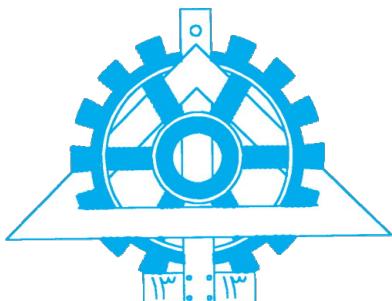
- [67] Selsam, Daniel, Lamm, Matthew, Bünz, Benedikt, Liang, Percy, de Moura, Leonardo, and Dill, David L. Learning a sat solver from single-bit supervision. *arXiv preprint arXiv:1802.03685*, 2018.
- [68] Senior, Andrew W, Evans, Richard, Jumper, John, Kirkpatrick, James, Sifre, Laurent, Green, Tim, Qin, Chongli, Žídek, Augustin, Nelson, Alexander WR, Bridgland, Alex, et al. Improved protein structure prediction using potentials from deep learning. *Nature*, 577(7792):706–710, 2020.
- [69] Silver, David, Schrittwieser, Julian, Simonyan, Karen, Antonoglou, Ioannis, Huang, Aja, Guez, Arthur, Hubert, Thomas, Baker, Lucas, Lai, Matthew, Bolton, Adrian, et al. Mastering the game of go without human knowledge. *nature*, 550(7676):354–359, 2017.
- [70] Singer, Jochen, Kuipers, Jack, Jahn, Katharina, and Beerenwinkel, Niko. Single-cell mutation identification via phylogenetic inference. *Nature communications*, 9(1):1–8, 2018.
- [71] Sokal, Alan. Monte carlo methods in statistical mechanics: foundations and new algorithms. In *Functional integration*, pages 131–192. Springer, 1997.
- [72] Srivastava, Nitish, Hinton, Geoffrey, Krizhevsky, Alex, Sutskever, Ilya, and Salakhutdinov, Ruslan. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014.
- [73] Stewart, BWKP and Wild, CP. World cancer report 2014. health, 2017.
- [74] Stratton, Michael R, Campbell, Peter J, and Futreal, P Andrew. The cancer genome. *Nature*, 458(7239):719–724, 2009.
- [75] Strino, Francesco, Parisi, Fabio, Micsinai, Mariann, and Kluger, Yuval. Trap: a tree approach for fingerprinting subclonal tumor composition. *Nucleic acids research*, 41(17):e165–e165, 2013.
- [76] Sun, Xiao-xiao and Yu, Qiang. Intra-tumor heterogeneity of cancer cells and its implications for cancer treatment. *Acta Pharmacologica Sinica*, 36(10):1219–1227, 2015.
- [77] Sutherland, NS. Outlines of a theory of visual pattern recognition in animals and man. *Proceedings of the Royal Society of London. Series B. Biological Sciences*, 171(1024):297–317, 1968.
- [78] Talbot, Simon J and Crawford, Dorothy H. Viruses and tumours—an update. *European Journal of Cancer*, 40(13):1998–2005, 2004.

- [79] Truninger, Kaspar, Menigatti, Mirco, Luz, Judith, Russell, Anna, Haider, Ritva, Gebbers, Jan-Olaf, Bannwart, Fridolin, Yurtsever, Hueseyin, Neuweiler, Joerg, Riehle, Hans-Martin, et al. Immunohistochemical analysis reveals high frequency of pms2 defects in colorectal cancer. *Gastroenterology*, 128(5):1160–1171, 2005.
- [80] Vander Heiden, Matthew G, Cantley, Lewis C, and Thompson, Craig B. Understanding the warburg effect: the metabolic requirements of cell proliferation. *science*, 324(5930):1029–1033, 2009.
- [81] Waclaw, Bartłomiej, Bozic, Ivana, Pittman, Meredith E, Hruban, Ralph H, Vogelstein, Bert, and Nowak, Martin A. A spatial model predicts that dispersal and cell turnover limit intratumour heterogeneity. *Nature*, 525(7568):261–264, 2015.
- [82] Williams, Ronald J. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.
- [83] Wu, Yufeng. Accurate and efficient cell lineage tree inference from noisy single cell data: the maximum likelihood perfect phylogeny approach. *Bioinformatics*, 36(3):742–750, 2020.
- [84] Yuan, Ke, Sakoparnig, Thomas, Markowetz, Florian, and Beerenwinkel, Niko. Bitphylogen: a probabilistic framework for reconstructing intra-tumor phylogenies. *Genome biology*, 16(1):1–16, 2015.
- [85] Zaccaria, Simone, El-Kebir, Mohammed, Klau, Gunnar W, and Raphael, Benjamin J. The copy-number tree mixture deconvolution problem and applications to multi-sample bulk sequencing tumor data. In *International Conference on Research in Computational Molecular Biology*, pages 318–335. Springer, 2017.
- [86] Zafar, Hamim, Navin, Nicholas, Chen, Ken, and Nakhleh, Luay. Syclonefit: Bayesian inference of population structure, genotype, and phylogeny of tumor clones from single-cell genome sequencing data. *Genome research*, 29(11):1847–1859, 2019.
- [87] Zafar, Hamim, Tzen, Anthony, Navin, Nicholas, Chen, Ken, and Nakhleh, Luay. Sifit: inferring tumor trees from single-cell sequencing data under finite-sites models. *Genome biology*, 18(1):1–20, 2017.
- [88] Zhu, Aizhi, Lee, Daniel, and Shim, Hyunsuk. Metabolic positron emission tomography imaging in cancer detection and therapy response. In *Seminars in oncology*, volume 38, pages 55–69. Elsevier, 2011.

Abstract

This thesis studies on writing projects, theses and dissertations using tehran-thesis class. It ...

Keywords SNV, CNV, Phylogenetic, Tree, Q-learning, Deep learning



University of Tehran
College of Engineering
**Faculty of New Science and
Technology
Network**



Inference of Phylogenetic Tree for Inter Tumor using Single Cell Mutations and CNV

A Thesis submitted to the Graduate Studies Office
In partial fulfillment of the requirements for
The degree of Master of Science
in Information Technology - Network Science

By:

Afshin Bozorgpour

Supervisors:

Dr. Saman Haratizadeh and Dr. Abolfazl Motahari

Jul 2021