

کاوش و استخراج برهم کنش‌های پروموتور/انهنسر بالقوه در ژنوم افراد سرطانی با استفاده از یک الگوریتم چندهدفه تکاملی

محمدجواد حسین پور^۱، حمید پروین^۲، صمد نجاتیان^{۳*}، وحیده رضایی^۴

• پذیرش مقاله: ۹۷/۳/۲۲

• دریافت مقاله: ۹۶/۷/۳۰

مقدمه: سرطان به عنوان یکی از شایع‌ترین انواع بیماری‌ها، سلامت بسیاری از انسان‌ها را تحت تأثیر قرار داده است. هدف اصلی در این مقاله، ارائه یک الگوریتم تکاملی چندهدفه می‌باشد. این الگوریتم، با استفاده از اطلاعات برهم کنش بین ژنومی در کروموزوم‌های افراد مبتلا به سرطان، ناحیه‌های بالقوه پروموتور/انهنسر را کشف و استخراج می‌کند. استخراج صحیح این ناحیه‌ها می‌تواند به علم پزشکی در تشخیص زودهنگام بیماری سرطان کمک کند.

روش: پژوهش حاضر از نوع کاربردی و توصیفی می‌باشد. در این پژوهش از مجموعه داده Hi-C شامل اطلاعات مربوط به برهم کنش‌های بین ژنومی، در سلول GM12878 استفاده شد. جهت کشف و استخراج پروموتور/انهنسرهای بالقوه از الگوریتمی تکاملی و چندهدفه استفاده شد. الگوریتم مذکور با استفاده از نرم‌افزار متلب پیاده‌سازی گردید. همچنین کارایی این الگوریتم نیز، با استفاده از دو معیار مورد ارزیابی قرار گرفت. معیار اول، تابع تناسبی است که میزان برهم کنش‌های بین ژنومی نسبت به طول نواحی ژنوم را محاسبه می‌کند و معیار دوم تعداد پروموتور/انهنسرهای بالقوه کشف شده می‌باشد.

نتایج: نتایج و مقایسه‌های انجام گرفته در این پژوهش، نشان از کارایی بالا و بهینه بودن روش پیشنهادی در کشف پروموتور/انهنسر با طول متغیر نسبت به روش HiC-Pro می‌باشد؛ بنابراین روش پیشنهادی می‌تواند پروموتور/انهنسرهای بالقوه‌ای را کشف کند که روش HiC-Pro قادر به کشف آن نیست.

نتیجه‌گیری: با توجه به نتایج به دست آمده، الگوریتم پیشنهادی قادر به کشف و استخراج بهینه پروموتور/انهنسرهای بالقوه با طول متغیر می‌باشد، که می‌تواند در تشخیص زودهنگام سرطان کمک شایانی به علم پزشکی کند.

کلید واژه‌ها: پروموتور، انهنسر، مجموعه داده Hi-C، الگوریتم شبیه‌سازی تبرید چند هدفه MOSA، روش HiC-Pro

ارجاع: حسین پورمحمدجواد، پروین حمید، نجاتیان صمد، رضایی وحیده، کاوش و استخراج برهم کنش‌های پروموتور/انهنسر بالقوه در ژنوم افراد سرطانی با استفاده از یک الگوریتم چندهدفه تکاملی. مجله انفورماتیک سلامت و زیست پزشکی ۱۳۹۷؛ ۵(۲): ۳۱۳-۳۰۴.

۱. دانشجوی دکتری دانشکده مهندسی کامپیوتر، واحد یاسوج، دانشگاه آزاد اسلامی، یاسوج، ایران

۲. دکتری مهندسی کامپیوتر- هوش مصنوعی استادیار، دانشکده مهندسی کامپیوتر، واحد نورآباد ممسنی، دانشگاه آزاد اسلامی، نورآباد ممسنی، ایران

۳. دکتری مهندسی برق، استادیار، دانشکده مهندسی برق، واحد یاسوج، دانشگاه آزاد اسلامی، یاسوج، ایران

۴. دکتری ریاضیات، استادیار، دانشکده ریاضیات، واحد یاسوج، دانشگاه آزاد اسلامی، یاسوج، ایران

۵. عضو باشگاه پژوهشگران جوان و نخبگان، واحد یاسوج، دانشگاه آزاد اسلامی، یاسوج، ایران

۶. عضو باشگاه پژوهشگران جوان و نخبگان، واحد نورآباد ممسنی، دانشگاه آزاد اسلامی، نورآباد ممسنی، ایران

* نویسنده مسئول: دانشکده مهندسی برق، واحد یاسوج، دانشگاه آزاد اسلامی، یاسوج، ایران.

مقدمه

امروزه در دانش پزشکی میزان داده‌های مربوط به علائم بیماران سرطانی و روش‌های کمکی برای تشخیص این بیماری‌ها، بسیار وسیع و گسترده شده است، به طوری که معمولاً تحلیل و در نظر گرفتن کلیه عوامل دخیل، دشوار به نظر می‌رسد. در واقع سرطان یکی از شایع‌ترین انواع بیماری‌ها در انسان است که خیلی از انسان‌ها را در مراحل مختلف زندگی تحت تأثیر قرار می‌دهد. در سال‌های اخیر، میزان شیوع سرطان روند رو به رشدی داشته و بررسی داده‌ها نشان می‌دهد که میزان بقای بیماران تا پنج سال پس از تشخیص، ۸۸ درصد و ده سال پس از تشخیص ۸۰ درصد بوده است [۱]؛ بنابراین لازم است که این بیماری سریع تشخیص داده شود. یکی از روش‌های تشخیص زودهنگام سرطان کشف برهم‌کنش‌های غیر معمول موجود در ژنوم افراد بیمار و استخراج نواحی درگیر در این برهم‌کنش‌ها است. در واقع ژنوم‌های انسانی دارای دو عنصر پروموتور (Promoter) و انهنسر (Enhancer) هستند. حال، زمانی که یک ژن بخواهد بیان (Gene Expression) شود و فعالیتی را انجام دهد، این دو عنصر باید با یکدیگر برهم‌کنش (Interaction) داشته باشند. در این صورت اگر آن ژن، مربوط به تومور سرطانی در شخص بیمار باشد، باعث می‌شود که آن ژن در آن فرد سرطانی بیشتر فعالیت کرده و گسترش پیدا کند [۲]. حال تشخیص دادن این نواحی راه‌انداز/انهنسر مسئله مهمی است که به علم پزشکی در تشخیص زودهنگام سرطان کمک کرده و همچنین در درمان این بیماری می‌تواند مؤثر باشد. در این راستا، تکنولوژی Hi-C به جهت کشف برهم‌کنش‌های بیشتر در ژنوم به وجود آمد. این تکنولوژی برای اولین بار توسط آیدین و همکاران جهت شبیه‌سازی کردن همه برهم‌کنش‌های کروموزومی در یک آزمایش خاص مطرح گردید. در حقیقت تکنولوژی Hi-C به همراه ابزارهای آن برهم‌کنش‌های بین نقاط مختلف ژنوم را در یک ساختار سه بعدی ارائه می‌دهد [۳]. بدین معنی که از نظر ساختاری چه ناحیه‌ای از ژنوم با چه ناحیه‌ای از ژنوم دیگر برهم‌کنش دارد. این اطلاعات می‌توانند دانش زیادی را در زمینه ساختار برهم‌کنش ژنوم ارائه کند و بسیار گران قیمت نیز می‌باشند. از جمله ابزارهای Hi-C می‌توان به (Humer, HiCup, HiC-Pro, etc..) اشاره کرد، که در تحقیقات گذشته [۴-۷] به منظور کشف پروموتور/انهنسر و برهم‌کنش‌های موجود بین آن‌ها مورد استفاده قرار گرفته است، که در بین آن‌ها کارایی ابزار HiC-Pro از دیگر ابزارها بالاتر می‌باشد.

در سال‌های گذشته پژوهش‌های زیادی در این زمینه صورت گرفته است. Jin و همکاران، در مقاله‌ای نگاشتی از برهم‌کنش‌های سه بعدی کروماتین‌ها در سلول‌های انسانی ارائه کردند. در این پژوهش آن‌ها یک نگاهت از برهم-کنش‌های کروماتین ایجاد شده در فیروبلات‌های انسانی با استفاده از یک روش تحلیل C^۳ با ژنوم گسترده (Hi-C) به دست آوردند. همچنین بیشتر از یک میلیون برهم‌کنش‌های کروماتین با دامنه طویل در فرگمنت‌های 10-15kb تعیین شد و اصول کلی سازمان کروماتین را در انواع متفاوت ویژگی‌های ژنومی مشخص گردید. مشاهدات آن‌ها پیشنهاد می‌دهد که چشم انداز کروماتین سه بعدی (زمانی که در یک نوع سلولی خاص ایجاد شد) تقریباً ثابت است و می‌تواند بر انتخاب یا فعال‌سازی ژن‌های مقصد توسط یک فعال ساز رونویسی موجود در همه جا به یک روش خاص سلول تأثیرگذار باشد [۸]. Cabrerros و همکاران، الگوریتم تشخیص اجماعی بر روی داده‌های Hi-C جهت کشف اجماعی از نواحی در ساختار سه بعدی دی‌ان‌ای موش و انسان به کار بردند. الگوریتم پیشنهادی آن‌ها قادر به کشف تعداد متغیری از اجماع بود. همچنین این الگوریتم، می‌توانست اجماعی از مکان‌های دی‌ان‌ای که در دور و نزدیک از هم قرار دارند، را به صورت دنباله‌ای کشف کند [۹]. Mifsud و همکاران، در مقاله‌ای نگاهت تماس‌های پروموتور با طول دامنه بالا در سلول‌های انسانی با ضبط Hi-C با کیفیت بالا ارائه کردند. آن‌ها در این پژوهش در آزمایشگاه ژنتیک از ضبط Hi-C (Chi-C) یک آزمایش ژنوم اتخاذ شده) استفاده کردند تا تعاملاتی با دامنه طویل، تقریباً ۲۲۰۰۰ پروموتور را در ۲ نوع سلول‌های خون انسان مورد سنجش قرار دهند. این رویکرد می‌تواند توالی‌سازی عمیق ناحیه خاص را مشخص کرده و باعث شود تا نگاهت برهم‌کنش‌ها با کیفیت بالا بین مناطق انجام شود [۱۰]. Servant و همکاران، در مقاله‌ای ابزاری به نام HiC-Pro جهت کشف برهم‌کنش‌های بین ژنومی در کروموزوم‌های انسانی ارائه دادند. ابزار پیشنهادی آن‌ها می‌توانست برهم‌کنش‌های یک ژنوم خاص با دیگر ژنوم را به دست آورد. همچنین HiC-Pro قابلیت اجرا به صورت موازی نیز داشت، که منجر به سرعت بالای آن در آنالیز مجموعه داده‌های Hi-C شده بود [۵]. Fotuhi و همکاران، در مقاله‌ای یک الگوریتم خوشه‌بندی چند منظوره برای مجموعه داده‌های به دست آمده از پیکربندی کروموزوم جهت شناسایی مدل‌های برهم‌کنش‌های کروموزومی

طرح مسئله

اگر مجموعه داده را به صورت n قلم داده در نظر گرفته شود. $D = \{g_1 \rightarrow g_{1 \leq i \leq n}, g_2 \rightarrow g_{1 \leq i \leq n}, \dots, g_n \rightarrow g_{1 \leq i \leq n}\}$ در اینجا g نشان دهنده یک ناحیه از ژنوم موجود در یک کروموزوم است. حال هدف مسئله یافتن اقلام پرتکراری مانند $P \rightarrow E$ که P و E به شکل زیر باشد:

$$P = \{g_{i_1}, g_{i_1+1}, \dots, g_{i_2}\}, i_1 < i_2$$

$$E = \{g_{i_3}, g_{i_3+1}, \dots, g_{i_4}\}, i_3 < i_4$$

در اینجا g_{i_1} تا g_{i_2} نواحی بهم پیوسته از ژنوم یک کروموزوم می باشد که مشخص کننده یک پروموتور بالقوه و g_{i_3} تا g_{i_4} نیز نواحی بهم پیوسته از ژنوم یک کروموزوم است که مشخص کننده یک انهنسر بالقوه می باشد. مجموعه داده مورد استفاده در این تحقیق مجموعه داده Hi-C می باشد، که این مجموعه داده برهم کنش های بین نقاط مختلف ژنوم را در افراد سرطانی نشان می دهد. اگر مجموعه داده Hi-C را برابر با D در نظر گرفته شود. مسئله مورد نظر یافتن رابطه برهم کنش $P \rightarrow E$ می باشد، که هر یک از نواحی P با تمامی نواحی E برهم کنش داشته باشند. که در این رابطه P پروموتور و E انهنسر در نظر گرفته می شود.

کشف پروموتور/انهنسر چندهدفه

در این بخش روش پیشنهادی به منظور کاوش و استخراج پروموتور/انهنسر در ژنوم افراد سرطانی توضیح داده می شود. رویکرد پیشنهادی بر اساس الگوریتم شبیه سازی تبرید چندهدفه می باشد، که در این بخش هر یک از اهداف مورد استفاده در الگوریتم پیشنهاد شده آورده شد. به منظور تشریح هر یک از اهداف فرض می شود که $P \rightarrow E$ مطابق با آنچه در بخش قبل توضیح داده شد، یک جواب مسئله می باشد.

هدف اول: ضریب نسبی (Ratio)

ضریب نسبی، برابر با جمع تمامی برهم کنش های هر قلم داده در مجموعه P با هر قلم داده در مجموعه E تقسیم بر جمع تمامی برهم کنش های موجود در مجموعه P با دیگر قلم داده های موجود در مجموعه D می باشد. مقدار ضریب اطمینان در بازه $[0, 1]$ می باشد و هر چقدر این مقدار بیشتر باشد ضریب نسبی بهتری را ارائه می دهد.

ارائه کردند. در این مقاله روش Arboretum-Hi-C (روش خوشه بندی طیفی چند منظوره) برای شناسایی جنبه های مشترک و در زمینه خاصی از ساختار ژنوم ارائه گردید. در مقایسه با روش خوشه بندی استاندارد، روش Arboretum-Hi-C الگوهای بیولوژیکی را که برای نگهداری سازگارتر بودند، به دست آورد. همچنین خوشه هایی که این الگوریتم خوشه بندی ارائه داد بهتر از روش های استاندارد بود [۱۱]. Charalampos و همکاران، ابزار HiC-bench جهت آنالیز مجموعه داده های Hi-C ارائه کردند. این ابزار، علاوه بر کشف پروموتور/انهنسر، برهم کنش های بین آن ها را نیز استخراج می کند. همچنین ابزار پیشنهادی قادر به اجرا به صورت موازی بر روی چند سیستم، مصورسازی نتایج (visualize)، کشف برهم کنش های دورن و بین کروموزومی و اجرا در سریع ترین زمان ممکن می باشد [۱۲].

اما مشکل اصلی موجود در مقالات مورد بررسی این است که، فناوری Hi-C و ابزارهای آن تنها می توانند فرگمت های با طول ثابت (Fix-bin size) از ژنوم جهت تشکیل پروموتور و انهنسر کشف و استخراج کرده و برهم کنش های بین آن ها را شمارش کنند. در این صورت خروجی این ابزار و روش های ارائه شده قبلی، که با فرگمت های با طول ثابت کار می کنند، نمی توانند بهینه باشد. در این صورت ارائه یک راه حل که بتواند پروموتور/انهنسر با طول متغیر را کشف و استخراج کرده و برهم کنش های بین آن ها را شمارش کنند، حس می شود. در این تحقیق یک الگوریتم چندهدفه مبتنی بر شبیه سازی تبرید ارائه شد که بتواند فرگمت های با طول متغیر از ژنوم که پروموتور و انهنسر بالقوه (potential promoter/enhancer) می باشند را کشف و استخراج کرده و برهم کنش های بین آن ها را شمارش کند.

روش

پژوهش حاضر از نوع کاربردی و توصیفی می باشد. در این پژوهش، با استفاده از الگوریتمی چندهدفه مبتنی بر شبیه سازی تبرید نواحی با طول متغیر پروموتور/انهنسر را ژنوم بیماران سرطانی کشف و استخراج شد و تمامی عملیات الگوریتم پیشنهادی بر روی مجموعه داده Hi-C، که اطلاعات مربوط به برهم کنش های بین ژنومی در سلول GM12878 انسانی می باشد، انجام گرفت.

$$Ratio(P \rightarrow E) = \frac{\sum_{m=i_1}^{i_2} \sum_{p=i_3}^{i_4} (g_m - g_p)}{\sum_{m=i_1}^{i_2} \sum_{o=1}^n (g_m - g_o)} \quad (1)$$

هدف دوم: ضریب پیوستگی

پیوستگی CL(Continuous Limitation) آن مجموعه.

ضریب پیوستگی، برابر با حاصل ضرب طول تمامی قلم داده‌های موجود در مجموعه موردنظر (P یا E) در محدودیت

$$Sequence(P) = Length(P) \times CL(P) \quad (2)$$

تابع تناسب برابر است با حاصل تقسیم ضریب نسبی بر مجموع ضریب پیوستگی مجموعه‌های P و E. در حقیقت این معیار ضریب برهم‌کنش‌های بین مجموعه اقلام P با مجموعه اقلام E را محاسبه می‌کند. مقدار تابع تناسب عددی بین ۱- تا ۱ است که هر چه به ۱ نزدیک‌تر باشد مقدار تابع تناسب بهتر خواهد بود و نشان‌دهنده وجود برهم‌کنش بیشتر بین هر دو نواحی P و E می‌باشد.

محدویت پیوستگی اقلام داده یک متغیر نوع منطقی (Boolean) می‌باشد، که در صورت پیوسته بودن اقلام مجموعه مورد نظر مقدار ۱ و در غیر این صورت مقدار ۱- را می‌گیرد. پیوسته بودن، بدین معنی است که مجموعه‌ای از ژنوم که تشکیل پروموتور یا انهنسر را داده‌اند به صورت ناحیه‌ای پیوسته از کروموزوم مربوطه باشند.

تابع تناسب (Fitness)

$$Fitness(P \rightarrow E) = \left(\frac{Ratio(P \rightarrow E)}{(Sequence(P) + Sequence(E))} \right) \quad (3)$$

تکرار شد. در نزدیکی جواب بهینه حاصله موجود به اندازه ۱۰ همسایه، جمعیت جدید را به وجود می‌آید، سپس توسط مقدار تابع تناسب هر یک از جمعیت تولیدی جدید را با بهترین جواب حاصله مقایسه شد. اگر پاسخ جدید از بهترین پاسخ یافته شده بهتر بود، پاسخ جدید را به عنوان بهترین پاسخ جدید در نظر گرفته شد؛ اما اگر بهتر نبود آن را با درصد احتمالی به صورت پاسخ مشروط در نظر گرفته می‌شود. سپس بهترین پاسخ یافته شده را به‌روزرسانی کرده و کاهش دما نیز به عنوان آخرین مرحله انجام می‌گیرد. حال این روند به اندازه تعداد تکرار ادامه پیدا می‌کند و در هر تکرار پاسخ بهینه ذخیره می‌شوند. در هر تکرار پاسخ بهینه، بهترین پروموتور/انهنسر را نشان می‌دهد، که می‌تواند مؤثر در برهم‌کنش‌های بین ژنومی در بیماران سرطانی باشند.

ساختار یک عنصر از جمعیت ورودی الگوریتم پیشنهادی
ساختار هر فرد (Individual) از جمعیت ورودی الگوریتم پیشنهادی به صورت زیر می‌باشد. در ادامه هر یک از قسمت‌های این ساختار توضیح داده می‌شود.

الگوریتم پیشنهادی

تعداد جمعیت اولیه در الگوریتم پیشنهادی برابر با ۲۵، بیشینه تکرار برابر با ۱۰۰، تعداد زیر تکرار برابر با ۲۰ و تعداد همسایگان قابل کشف در نزدیکی جواب بهینه ۱۰ می‌باشد. ورودی الگوریتم، مجموعه داده Hi-C حاصل از برهم‌کنش موجود بین ژنوم و تعداد کل فرگمنت‌های مورد استفاده در مجموعه داده می‌باشد.

در مرحله پیش پردازش ساده‌سازی مجموعه داده Hi-C انجام می‌گیرد. ساده‌سازی به‌منظور حذف تمامی برهم‌کنش‌های بین کروموزومی اجرا می‌شود. در این صورت مجموعه داده موردنظر شامل تمامی برهم‌کنش‌های درون کروموزومی می‌باشد. بعد از مرحله پیش‌پردازش، جمعیت اولیه به صورت تصادفی ساخته می‌شود، سپس این جمعیت توسط تابع تناسب مطابق با فرمول ۳ ارزیابی می‌شوند. پس از ارزیابی، فردی از جمعیت ورودی که نشان‌دهنده پروموتور/انهنسر با بهترین مقدار تابع تناسب است به‌عنوان بهترین پاسخ فعلی در نظر گرفته می‌شود. سپس دمای اولیه T_0 نیز تنظیم می‌شود که در این جا $T_0 = 10$ و میزان کاهش دما نیز $\alpha = 0.99$ در نظر گرفته شد. حال به اندازه بیشینه تکرار، مراحل زیر را

طول انهنسر	انهنسر	طول پروموتور	پروموتور
------------	--------	--------------	----------

- راه انداز: مشخص کننده ابتدای ژنوم تشکیل دهنده پروموتور می باشد.
- طول راه انداز: مشخص کننده طول پروموتور می باشد، که به صورت تصادفی در نظر گرفته می شود.
- انهنسر: مشخص کننده ابتدای ژنوم تشکیل دهنده انهنسر می باشد.
- طول انهنسر: مشخص کننده طول انهنسر می باشد، که به صورت تصادفی در نظر گرفته می شود.

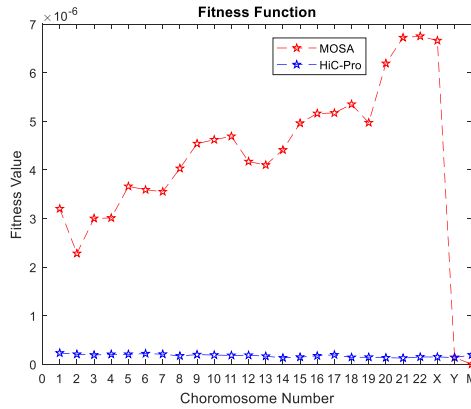
مجموعه داده

در این تحقیق از مجموعه داده Hi-C استفاده شد. مجموعه داده مورد بررسی، حاوی ۲۵۰ میلیون رکورد مربوط به اطلاعات برهم کنش‌های بین ژنومی در سلول GM12878 افراد مبتلا به سرطان می باشد. سلول GM12878 یک سلول لیمفوبلاستوئیدی (lymphoblastoid) تولید شده از خون انسان با استفاده از تبدیل EBD (Epstein-Barr virus) می باشد. همچنین مجموعه داده مورد نظر در مطالعات متنوع ژنتیکی مانند [۱۰،۱۳] استفاده شده است.

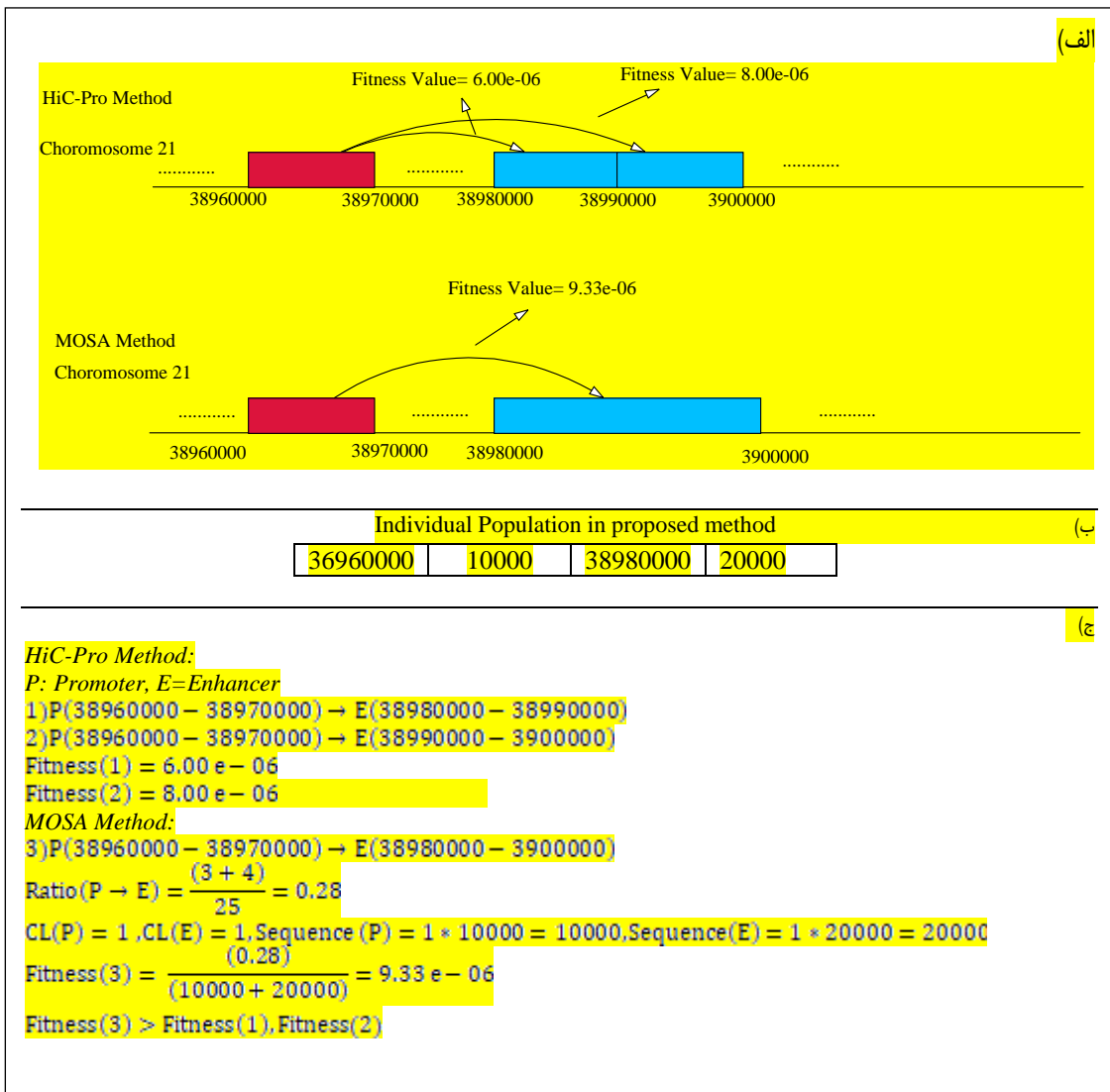
شکل ۱، میانگین تابع تناسب مربوط به روش پیشنهادی و روش HiC-Pro در همه کروموزوم‌ها بر روی مجموعه داده مورد نظر به تصویر کشیده شده است. مطابق این شکل روش پیشنهادی در همه کروموزوم‌ها بهتر از روش HiC-Pro می باشد. در ادامه، مطابق شکل ۲، روند محاسبه تابع ارزیابی با یک مثال در هر دو روش پیشنهادی و HiC-Pro توضیح داده شده است.

همانطور که شکل ۲(الف) نشان می دهد، ناحیه ۳۸۹۶۰۰۰۰ تا ۳۸۹۷۰۰۰۰ از کروموزوم ۲۱ را به عنوان یک پروموتور بالقوه تشخیص می دهد؛ اما ناحیه ای که به عنوان انهنسر در نظر گرفته است، ناحیه ۳۸۹۸۰۰۰۰ تا ۳۸۹۰۰۰۰۰ از کروموزوم ۲۱

می باشد. شکل ۲(ب)، یک فرد از جمعیت ورودی الگوریتم پیشنهادی را نشان می دهد. در اینجا مقدار اول ابتدای راه انداز، مقدار دوم طول راه انداز، مقدار سوم ابتدای انهنسر و مقدار چهارم طول انهنسر را بیان می کند. مطابق این ورودی، الگوریتم پیشنهادی، وجود برهم کنش بین پروموتور و انهنسر را بر روی مجموعه داده Hi-C بررسی کرده و مطابق آن تابع تناسب مربوط به آن را محاسبه می کند (در شکل ۲(ج) روند محاسبه تابع تناسب به تصویر کشیده شده است). حال اگر مقدار تابع تناسب مربوط به برهم کنش بین نواحی پروموتور/انهنسر جدید از مقادیر تابع تناسب هر کدام از برهم کنش‌های قبلی موجود در مجموعه داده بهتر بود، ناحیه جدید که حاصل ادغام دو ناحیه قبلی است را به عنوان انهنسر جدید در نظر می گیرد و برهم کنش بین این دو ناحیه را کشف و استخراج می کند. این نشان می دهد که، روش پیشنهادی در کشف نواحی پروموتور/انهنسر بالقوه توانایی بهتری نسبت به روش HiC-Pro دارد. چون روش HiC-Pro، ناحیه‌های ۳۸۹۶۰۰۰۰ تا ۳۸۹۷۰۰۰۰ را به عنوان پروموتور و ناحیه‌های ۳۸۹۸۰۰۰۰ تا ۳۸۹۹۰۰۰۰ و ۳۹۰۰۰۰۰ را به عنوان دو انهنسر که هر یک با پروموتور دارای برهم کنش جداگانه هستند در نظر گرفته است و مقدار تابع تناسب حاصل از هر برهم کنش، به ترتیب $6.00E-06$ و $8.00E-06$ می باشد؛ اما روش پیشنهادی، ناحیه ۳۸۹۶۰۰۰۰ تا ۳۸۹۷۰۰۰۰ را به عنوان پروموتور و ناحیه ۳۸۹۸۰۰۰۰ تا ۳۹۰۰۰۰۰ را به عنوان انهنسر در نظر می گیرد، که مقدار تابع تناسب حاصل از برهم کنش بین آن‌ها $9.33E-06$ می باشد و از هر یک از دو مقدار تناسب ثبت شده در روش HiC-Pro بهتر است و این نشان می دهد که ناحیه‌های مذکور برهم کنش شدیدتری نسبت به ناحیه‌های کشف شده در روش HiC-Pro دارند؛ بنابراین ناحیه پروموتور/انهنسر کشف شده در روش پیشنهادی دقیق تر از نواحی کشف شده در روش HiC-Pro است.



شکل ۱: مقایسه روش پیشنهادی با روش HiC-Pro بر روی مجموعه داده



شکل ۲: الف) یک نمونه کشف شده از پروموتور/انهنسر بالقوه با طول متغیر توسط روش پیشنهادی در کروموزوم ۲۱ سلول GM12878 در مقایسه با روش HiC-Pro که با فرگمنت‌های با طول ثابت کار می‌کند. ب) یک نمونه جمعیت ورودی الگوریتم جهت کشف پروموتور/انهنسر بالقوه شکل الف. ج) محاسبه تابع ارزیابی مربوط به پروموتور/انهنسر بالقوه کشف شده در روش پیشنهادی و روش HiC-Pro.

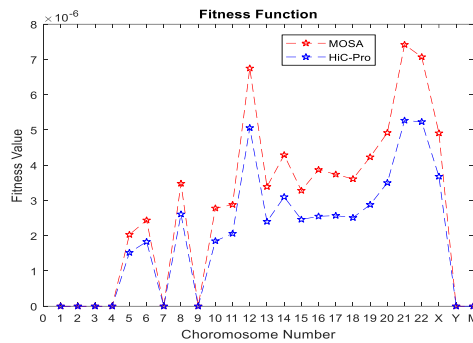
الگوریتم پیشنهادی نسبت به میانگین مقدار تابع تناسب در روش HiC-Pro در همان کروموزوم نیز بیان شد. مطابق با جدول ۱، روش پیشنهادی تعداد ۱۵۱ برهم کنش که روش HiC-Pro قادر به کشف آنها نبوده را کشف کرده است. برهم کنش‌های کشف شده حاصل از وجود برهم-کنش‌های بین نواحی پروموتور/انهنسر بالقوه با طول متغیر می‌باشد که روش HiC-Pro قادر به کشف آنها نیست.

در ادامه روش پیشنهادی با روش HiC-Pro، مورد استفاده در [۴-۷]، مقایسه شده است. روش HiC-Pro از سایر روش‌ها، مانند HiClib، Homer، HiCdat، HiCUP و... در آنالیز مجموعه داده HiC قدرتمندتر می‌باشد. جدول ۱، تمامی برهم کنش‌های کشف شده توسط روش پیشنهادی که روش HiC-Pro قادر به کشف آنها نبوده است را به تفکیک کروموزوم نشان می‌دهد. همچنین میانگین مقدار تابع تناسب مربوط به هر کروموزوم حاصل از برهم کنش‌های کشف شده در

شماره کروموزوم	تعداد برهم‌کنش روش پیشنهادی	میانگین مقدار تابع تناسب روش پیشنهادی	میانگین مقدار تابع تناسب روش HiC-Pro
۱	۰	0.00E+00	0
۲	۰	0.00E+00	0
۳	۰	0.00E+00	0
۴	۰	0.00E+00	0
۵	۱	2.03E-06	1.52E-06
۶	۱	2.44E-06	1.83E-06
۷	۰	0.00E+00	0
۸	۳	3.48E-06	2.61E-06
۹	۰	0.00E+00	0
۱۰	۱	2.78E-06	1.85E-06
۱۱	۴	2.88E-06	2.06E-06
۱۲	۲	6.75E-06	5.06E-06
۱۳	۴	3.39E-06	2.40E-06
۱۴	۴	4.29E-06	3.10E-06
۱۵	۵	3.28E-06	2.46E-06
۱۶	۱۴	3.87E-06	2.55E-06
۱۷	۱۲	3.74E-06	2.57E-06
۱۸	۱۵	3.61E-06	2.51E-06
۱۹	۱۰	4.23E-06	2.88E-06
۲۰	۱۷	4.92E-06	3.50E-06
۲۱	۱۸	7.42E-06	5.27E-06
۲۲	۷	7.07E-06	5.23E-06
X	۳۳	4.91E-06	3.68E-06
Y	۰	0.00E+00	0.00E+00
M	۰	0.00E+00	0.00E+00

مطابق شکل ۳، همان‌طور که مشاهده می‌شود، مقادیر تابع تناسب برهم کنش‌های استخراج شده توسط روش پیشنهادی در اکثر کروموزوم‌ها بهتر از روش HiC-Pro می‌باشد. این بدین معنی است که، برهم کنش‌های استخراج شده توسط الگوریتم پیشنهادی در این نواحی از هر کروموزوم که بیان کننده پروموتور/انهنسر بالقوه می‌باشد، در مقایسه با روش HiC-Pro بهینه خواهد بود.

همان‌طور که مشاهده می‌کنید بیشترین تعداد برهم کنش‌های کشف شده مربوط به کروموزوم X، بعد از آن مربوط به کروموزوم‌های ۲۱، ۲۰، ۱۸، ۱۶ و کمترین آن مربوط به کروموزوم‌های ۱ تا ۹ در مجموعه داده مورد مطالعه می‌باشد. شکل ۳، میانگین مقدار تابع تناسب روش پیشنهادی در کشف برهم کنش‌های نواحی‌های پروموتور/انهنسر بالقوه به تفکیک هر کروموزوم در مقایسه با روش HiC-Pro در همان نواحی را نشان می‌دهد.



شکل ۳: میانگین مقدار تابع تناسب روش پیشنهادی و روش HiC-Pro به تفکیک هر کروموزوم

بحث و نتیجه گیری

امروزه، علم ژنتیک ثابت کرده است که، نواحی پرموتر/انهنسر در کروموزوم‌های انسانی که با یکدیگر برهم‌کنش دارند، دارای طول ثابت (Fix-bin size) نمی‌باشند، بلکه طول آن‌ها می‌تواند متغیر باشد [۱۴]. همچنین تا به امروز، تحقیقات انجام شده [۵۶، ۱۵] قادر به کشف برهم‌کنش‌های بین نواحی ثابت ژنوم می‌باشند، که ممکن است ناحیه‌های کشف شده (پرموتر/انهنسرهای بالقوه) دقیق نبوده و برهم‌کنش‌های شمارش شده بین آن‌ها صحیح نباشد.

در این پژوهش، یک الگوریتم فراابتکاری چند هدفه مبتنی بر شبیه‌سازی تبرید جهت حل مشکل فوق ارائه شد. الگوریتم پیشنهادی قادر به کشف پرموتر/انهنسرهای بالقوه با طول متغیر و شمارش برهم‌کنش‌های بین آن‌ها می‌باشد. نتایج و مقایسه‌های انجام گرفته نشان از کارایی بالا و بهتر بودن روش پیشنهادی در کشف پرموتر/انهنسر با طول متغیر نسبت به روش HiC-Pro است. در واقع روش HiC-Pro یکی از ابزارهای قدرتمند آنالیز داده‌های HiC می‌باشد، که قادر به کشف پرموتر/انهنسر با طول متغیر نمی‌باشد.

روش پیشنهادی، قادر به کشف تعداد ۱۵۱ برهم‌کنش در تمامی کروموزوم‌ها می‌باشد، که روش HiC-Pro نمی‌تواند آن را کشف کند. برهم‌کنش‌های کشف شده حاصل از وجود برهم‌کنش‌های بین نواحی پرموتر/انهنسر بالقوه با طول متغیر می‌باشد که روش HiC-Pro قادر به کشف آن‌ها نیست. همچنین میانگین مقدار تابع تناسب روش پیشنهادی در کشف برهم‌کنش‌های نواحی‌های پرموتر/انهنسر بالقوه به تفکیک هر کروموزوم در مقایسه با روش HiC-Pro در همان نواحی را نشان می‌دهد. همان‌طور که مشاهده شد، مقادیر تابع تناسب برهم‌کنش‌های استخراج شده توسط روش پیشنهادی در اکثر کروموزوم‌ها بهتر از روش HiC-Pro می‌باشد. این بدین معنی

است که، برهم‌کنش‌های استخراج شده توسط الگوریتم پیشنهادی در این نواحی از هر کروموزوم که بیان‌کننده پرموتر/انهنسر بالقوه می‌باشد در مقایسه با روش HiC-Pro بهینه خواهد بود.

نتایج به‌دست آمده نشان می‌دهند که، روش پیشنهادی نسبت به روش HiC-Pro موفق بوده است. روش HiC-Pro، یکی از ابزارهای قدرتمند در زمینه آنالیز داده‌های Hi-C می‌باشد [۵]. کارایی این ابزار در آنالیز مجموعه داده‌های Hi-C از دیگر ابزارهای معرفی شده [۴-۷، ۹-۱۲] که در بخش مقدمه ذکر گردید، بالاتر می‌باشد. همچنین روش پیشنهادی قادر به کشف و استخراج نواحی با طول متغیر از ژنوم است که این نواحی به صورت بالقوه می‌توانند پرموتر/انهنسر باشند. علاوه بر این، روش پیشنهادی توانایی شمارش تمامی برهم‌کنش‌های بین این نواحی را دارد. تشخیص دادن ناحیه‌های پرموتر/انهنسر که با یکدیگر برهم‌کنش داشته باشند، مسئله مهمی است که به علم پزشکی در تشخیص زودهنگام سرطان کمک کرده و همچنین در درمان این بیماری می‌تواند مؤثر باشد؛ بنابراین پیشنهاد می‌شود، روش ارائه شده که قادر به استخراج صحیح این نواحی می‌باشد، بتواند جهت بهبود تشخیص بیماری سرطان به کار گرفته شود.

تشکر و قدردانی

مقاله مورد نظر مستخرج از پایان‌نامه دکترای تحت عنوان «کشف و استخراج برهم‌کنش‌های پرموتر/انهنسر و اجماعی از آن‌ها در ژنوم افراد سالم و ژنوم افراد سرطانی با استفاده از الگوریتم‌های فراابتکاری جدید» می‌باشد و از تمامی زحمات دکتر حمید علی‌نژاد رکنی عضو هیأت علمی دانشکده پزشکی دانشگاه ولز جنوبی سیدنی استرالیا به‌خاطر در اختیار گذاشتن

تضاد منافع

بدین وسیله نویسندگان تصریح می‌نمایند که هیچ‌گونه تضاد منافی در خصوص پژوهش حاضر وجود ندارد.

مجموعه داده Hi-C مربوط به سلول GM12878 نیز تشکر و قدردانی می‌گردد.

References

1. Jemal A, Bray F, Center MM, Ferlay J, Ward E, Forman D. Global cancer statistics. *CA Cancer J Clin* 2011;61(2):69-90.
2. Delen D, Walker G, Kadam A. Predicting breast cancer survivability: a comparison of three data mining methods. *Artif Intell Med* 2005;34(2):113-27.
3. Imakaev M, Fudenberg G, McCord RP, Naumova N, Goloborodko A, Lajoie BR, et al. Iterative correction of Hi-C data reveals hallmarks of chromosome organization. *Nature methods*. *Nat Methods* 2012;9(10):999-1003.
4. Hwang YC, Zheng Q, Gregory BD, Wang LS. High-throughput identification of long-range regulatory elements and their target promoters in the human genome. *Nucleic Acids Res* 2013;41(9):4835-46.
5. Servant N, Varoquaux N, Lajoie BR, Viara E, Chen CJ, Vert JP, et al. HiC-Pro: an optimized and flexible pipeline for Hi-C data processing. *Genome Biol* 2015;16:259.
6. Hwang YC, Lin CF, Valladares O, Malamon J, Kuksa PP, Zheng Q, et al. HIPPIE: a high-throughput identification pipeline for promoter interacting enhancer elements. *Bioinformatics* 2015;31(8):1290-2.
7. Wingett S, Ewels P, Furlan-Magaril M, Nagano T, Schoenfelder S, Fraser P, et al. HiCUP: pipeline for mapping and processing Hi-C data. *F1000Res* 2015;4:1310.
8. Jin F, Li Y, Dixon JR, Selvaraj S, Ye Z, Lee AY, et al. A high-resolution map of the three-dimensional chromatin interactome in human cells. *Nature* 2013;503(7475):290-4.
9. Cabrerros I, Abbe E, Tsirig A. Detecting community structures in Hi-C genomic data. *Annual Conference on Information Science and Systems (CISS)*; 2016 Mar 16-18; Princeton, NJ, USA: IEEE; 2016. p. 1-17.
10. Mifsud B, Tavares-Cadete F, Young AN, Sugar R. Mapping long-range promoter contacts in human cells with high-resolution capture Hi-C. *Nat Genet* 2015;47(6):598-606.
11. Fotuhi Siahpirani A, Ay F, Roy S. A multi-task graph-clustering approach for chromosome conformation capture data sets identifies conserved modules of chromosomal interactions. *Genome Biology* 2016; 17:114.
12. Lazaris C, Kelly S, Ntziachristos P, Aifantis I, Tsirigos A. HiC-bench: comprehensive and reproducible Hi-C data analysis designed for parameter exploration and benchmarking. *BMC Genomics* 2017;18(1):22.
13. Arvey A, Agius P, Noble WS, Leslie CS. Sequence and chromatin determinants of cell-type-specific transcription factor binding. *Genome Res* 2012; 22(9):1723-34.
14. Whalen S, Truty RM, Pollard KS. Enhancer-promoter interactions are encoded by complex genomic signatures on looping chromatin. *Nat Genet* 2016; 48(5): 488-96.
15. Schmid MW, Grob S, Grossniklaus U. HiCdat: a fast and easy-to-use Hi-C data analysis tool. *BMC Bioinformatics* 2015;16:277.

Detection and Extraction of Potential Promoter/Enhancer Interactions in Genome of Cancer Patients using an Evolutionary Multi-Objective Algorithm

Hosseinpoor Mohammadjavad¹, Parvin Hamid^{2,6}, Nejatian Samad^{3,5*}, Rezaei Vahideh^{4,5}

• Received: 22 Oct, 2017

• Accepted: 12 Jun, 2018

Introduction: Cancer, as one of the most common diseases, has influenced the health of many people. The main aim of this study was to present a multi-objective evolutionary algorithm. The algorithm is capable of detecting and extracting potentially promoter/enhancer areas in the chromosomes of the affected people using the information concerning inter-genomic interactions. The correct extraction of these areas can help early diagnosis of cancer.

Methods: In this applied and descriptive research, Hi-C data set including information on inter-genomic interactions in the GM12878 cell was used. Multi-objective evolutionary algorithm was used in order to discover and extract potential promoter /enhancer interactions. The mentioned algorithm was implemented using MATLAB software. Furthermore, the efficiency of this algorithm was evaluated using two criteria. The first criterion is a proportional function that calculates the magnitude of inter-genomic interactions relative to the length of the genome regions; and the second criterion is the number of discovered potential promoters/enhancers.

Results: The results and comparisons showed higher efficiency and optimality of the suggested method in discovering promoter/Enhancer interactions with variable length in comparison to HiC-Pro method. Therefore, the suggested method is able to discover the potential promoter/ enhancer interactions that cannot be discovered by HiC-Pro method.

Conclusion: The suggested algorithm is able to optimally discover and extract potential promoter/enhancer with variable length. This is a great help in medical science for early diagnosis of cancer.

Keywords: Promoter, Enhancer, Hi-C Dataset, Multi-Objective Simulation Annealing Algorithm (MOSA), HiC- Pro Method

• **Citation:** Hosseinpoor MJ, Parvin H, Nejatian S, Rezaie V. Detection and Extraction of Potential Promoter/Enhancer Interactions in Genome of Cancer Patients using an Evolutionary Multi-Objective Algorithm. *Journal of Health and Biomedical Informatics* 2018; 5(2): 304-313.

1. Ph.D. student in computer engineering, Computer Engineering Dept., Yasooj Branch, Islamic Azad University, Yasooj, Iran
2. Ph.D. in Computer Engineering - Artificial Intelligence, Assistant Professor, Computer Engineering Dept., Nourabad Mamasani Branch, Islamic Azad University, Nourabad Mamasani, Iran
3. Ph.D. in Electrical Engineering, Assistant Professor, Electrical Engineering Dept., Yasooj Branch, Islamic Azad University, Yasooj, Iran
4. Ph.D. in Mathematics, Assistant Professor, Mathematics Dept., Yasooj Branch, Islamic Azad University, Yasooj, Iran
5. Member of Young Researchers and Elite Club, Yasooj Branch, Islamic Azad University, Yasooj, Iran.
6. Member of Young Researchers and Elite Club, Nourabad Mamasani Branch, Islamic Azad University, Nourabad Mamasani, Iran.

*Correspondence: Electrical Engineering Dept., Yasooj Branch, Islamic Azad University, Yasooj, Iran.

• Tel: 09173414008

• Email: nejatian@iauyasooj.ac.ir