

بهبود تشخیص کمپلکس‌های پروتئینی مبتنی بر خوشه بندی دوگانه داده‌های بیان ژن

امیر لکی زاده*

گروه مهندسی کامپیوتر، دانشکده فنی و مهندسی، دانشگاه قم
پست الکترونیکی: lakizadeh@qom.ac.ir

چکیده

با توجه به نقش کمپلکس‌های پروتئینی در انجام بسیاری از کارکردهای سلولی موجودات زنده، کشف آن‌ها می‌تواند به درک عمیق‌تر سازوکارهای تنظیمی سلول، توصیف فرآیند تکامل در سیگنال‌های سلولی، پیش‌بینی عملکرد زیستی پروتئین‌های کشف‌شده و از همه مهم‌تر تحقق اهداف درمانی (تشخیص بیماری و طراحی دارو) منجر شود. رویکردهای محاسباتی ارائه‌شده تاکنون به منظور تشخیص کمپلکس‌های پروتئینی، به طور عمده بر خوشه‌بندی شبکه برهم‌کنش پروتئین-پروتئین تمرکز دارند و این در حالی است که شبکه‌های برهم‌کنش، علاوه بر این که نوفه‌دار می‌باشند، به‌تنهایی، فاقد سازوکار لازم برای در نظر گرفتن ماهیت پویای سلول در فرآیند تشخیص کمپلکس‌های پروتئینی می‌باشند.

در این مقاله، مدلی سه لایه مبتنی بر رویکرد خوشه‌بندی دوگانه برای تشخیص کمپلکس‌های پروتئینی از منابع داده‌ای مختلف ارائه شده است. لایه‌های اول، دوم و سوم مدل پیشنهادی، به ترتیب وظیفه پویاسازی، کاهش نوفه و تشخیص نهایی کمپلکس‌های پروتئینی را بر عهده دارند. مجموعه ارزیابی‌های مختلف، نشان می‌دهد که روش پیشنهادی توانسته است با استفاده از منابع داده‌ای

* نویسنده مسئول

مختلف در لایه‌های سه‌گانه، ضمن مدیریت بهتر چالش‌های اصلی مسئله مانند نوفه‌دار بودن منابع داده‌ای و ضرورت مدل‌سازی پویای فرآیند تشخیص به کمک داده‌های بیان ژن، دقت فرآیند تشخیص کمپلکس‌های پروتئینی را بر اساس سنجنده‌های مرتبط و در شرایط گوناگون، به میزان قابل توجهی بهبود دهد.

واژه‌های کلیدی: کمپلکس پروتئینی، شبکه برهم‌کنش پروتئین-پروتئین، داده‌های بیان ژن، خوشه‌بندی دوگانه، شباهت معنایی.

۱. مقدمه

اکثر پروتئین‌ها عملکردهای زیستی خود را در سلول زنده در قالب تشکیل کمپلکس‌ها انجام می‌دهند [۱]. در حقیقت بسیاری از پروتئین‌ها به‌عنوان اجزای کمپلکس‌های ضروری و پایدار وجود دارند. برای مثال مولکول حیاتی هموگلوبین، یک ترکیب پایدار از چهار پروتئین گلوبولین‌دار مشخص می‌باشد. همچنین بسیاری از آنزیم‌ها که نقش راه‌انداز را در بسیاری از واکنش‌های درون‌سلولی به عهده‌دارند، مجموعه‌ای از پروتئین‌ها می‌باشند که به‌طور هم‌زمان به یکدیگر متصل شده و ساختار موردنظر را تشکیل می‌دهند. علاوه بر آن برهم‌کنش‌های ناپایدار و

موقت مانند اثرات هورمون و یا اثرات یک سیگنال موقت، مستلزم تشکیل کمپلکس‌های پروتئینی می‌باشد [۲]. از لحاظ زیستی، کمپلکس‌های پروتئینی اجزاء مولکولی کلیدی برای انجام بسیاری از عملکردهای زیستی از قبیل رونوشت DNA، ترجمه mRNA، انتقال سیگنال سلولی، چرخه سلولی و غیره می‌باشد. برای مثال کمپلکس RNA پلیمراز در فرآیند رونویسی اطلاعات ژنتیکی DNA و تشکیل rRNA به‌منظور تولید پروتئین، نقش تعیین‌کننده‌ای دارد و یا کمپلکس ریز منفذ هسته، مسئول مبادله محافظت‌شده اجزاء سلولی میان هسته و سیتوپلاسم پیرامون آن می‌باشد و از انتقال مواد غیرمجاز به درون غشاء هسته جلوگیری می‌کند [۳].

به‌طورکلی با توجه به اهمیت و جایگاه کمپلکس‌های پروتئینی، کشف آن‌ها می‌تواند به درک عمیق‌تر سازوکارهای تنظیمی سلول، توصیف فرآیند تکامل در سیگنال‌های سلولی، پیش‌بینی عملکرد زیستی پروتئین‌های کشف شده و از همه مهم‌تر تحقق اهداف درمانی (تشخیص بیماری و طراحی دارو) منجر شود [۴].

روش‌های آزمایشگاهی برای تشخیص کمپلکس‌های پروتئینی دارای محدودیت‌های مختلفی می‌باشند. برای مثال روش TAP که به کمک طیف‌سنج جرمی انجام می‌شود، مستلزم چندین مرحله شستشو و تصفیه است که منجر به حذف نامطلوب‌جاذبه (ضعیف و ناپایدار) میان پروتئین‌ها در کمپلکس می‌شود. علاوه بر آن پروتئین‌هایی که به‌عنوان برچسب در این نوع آزمایش‌های تجربی استفاده می‌شوند، می‌توانند در یک اثر نامطلوب، مانع تشکیل کمپلکس‌های پروتئینی شوند. بنابراین، استفاده از روش‌های آزمایشگاهی برای تشخیص کمپلکس‌های پروتئینی بسیار پرهزینه و زمان‌بر می‌باشد و در نتیجه رویکردهای محاسباتی در تشخیص کمپلکس‌ها اهمیت فوق‌العاده‌ای می‌یابد [۵].

در سال‌های اخیر به‌کارگیری فناوری‌های پیشرفته آزمایشگاهی در تعیین برهم‌کنش میان پروتئین‌ها منجر به ایجاد حجم بسیار بالایی از داده‌های تجربی شده است که

کارکرد پروتئین‌ها را در شبکه‌های سلولی پیچیده منعکس می‌کند. شبکه‌های برهم‌کنش پروتئین-پروتئین منبع داده‌ای اصلی به منظور کشف کمپلکس‌های پروتئینی می‌باشد. انگیزه اصلی از آنجا ناشی شده است که پژوهشگران به تجربه دریافته‌اند که پروتئین‌ها در سلول و در قالب پیمانها و کمپلکس‌های پروتئینی مختلف وظایف خود را انجام می‌دهند و پروتئین‌هایی که با یکدیگر برهم‌کنش دارند، اغلب در فرآیندهای زیستی یکسانی شرکت می‌کنند و در نتیجه می‌توان در اغلب موارد پیمانها (کمپلکس‌های) پروتئینی را به یک عملکرد زیستی مشخص نگاشت کرد به‌طوری‌که پروتئین‌های یک پیمانها، قرابت عملکردی بیشتری با یکدیگر نسبت به پروتئین‌های یک پیمانها دیگر دارند [۶].

مسئله تشخیص کمپلکس‌های پروتئینی از درون شبکه‌های PPI از لحاظ محاسباتی به مسئله خوشه‌بندی این شبکه‌ها نگاشت می‌شود [۷]. انگیزه این کار نیز بر اساس مشاهدات پژوهشگران تجربی از تناظر میان یک کمپلکس واقعی با یک خوشه از پروتئین‌ها در یک شبکه PPI می‌باشد. در یک تقسیم‌بندی، چالش‌های پژوهش در تشخیص کمپلکس‌های پروتئینی را می‌توان به دو دسته کلی چالش‌های ناشی از مجموعه شبکه‌های PPI و چالش‌های روش‌های محاسباتی موجود، تقسیم‌بندی کرد [۸].

شبکه‌های PPI به‌دست‌آمده از روش‌های آزمایشگاهی دارای مقدار زیادی خطای مثبت کاذب می‌باشند و در نتیجه شبکه PPI مورد استفاده خوشه‌بندی، دارای مقدار زیادی نوفه می‌باشد [۹]. با توجه به هزینه و زمان مورد نیاز در انجام روش‌های آزمایشگاهی، استفاده از روش‌های محاسباتی جهت ارتقای کیفیت شبکه‌های برهم‌کنش از طریق وزن‌دار کردن برهم‌کنش‌ها ضروری است. بدیهی است که افزایش اطمینان در منابع داده‌ای موجود به افزایش دقت روش‌های محاسباتی موجود در کاربردهای مختلف منجر می‌شود. متداول‌ترین راهکار جهت بهره‌گیری بیشتر از منابع زیستی و کاهش نوفه (برای مثال، تداخل موجود در شبکه‌های PPI تولیدشده توسط روش‌های آزمایشگاهی

طیف‌سنج جرمی، Yeast2Hybrid و in-silico، وزن‌دار کردن برهم‌کنش‌ها با استفاده از منابع داده‌ای متمایز و با قالب‌های مختلف (مانند پایگاه داده‌های دنباله‌ای، ساختاری، بیان ژن و انواع شبکه‌های با ساختارهای مختلف) می‌باشد [۱۰].

روش‌های تشخیص کمپلکس را می‌توان از دیدهای متفاوتی دسته‌بندی کرد. در یک دسته‌بندی کلان، می‌توان روش‌های تشخیص را برحسب توانایی یا عدم توانایی آن‌ها در مدل‌سازی شرایط پویای سلول در فرآیند تشخیص، به دو دسته روش‌های پویا و روش‌های ایستا تقسیم‌بندی نمود. روش‌های ایستای تشخیص کمپلکس را می‌توان برحسب استفاده و یا عدم استفاده آن‌ها از سایر منابع داده‌ای زیستی به دو دسته تقسیم‌بندی کرد. روش‌های مبتنی بر فقط خوشه‌بندی گراف و روش‌های مبتنی بر خوشه‌بندی گراف همراه با بعضی از دانش‌های اضافی زیست‌شناسی. دانش‌های زیست‌شناسی نیز در قالب قالب هسته-افزونه، شباهت‌های عملکردی، پایسته‌های تکاملی و برهم‌کنش‌های متقابل (یا انحصاری) دسته‌بندی می‌شوند. از مهم‌ترین روش‌هایی که فقط از خوشه‌بندی گراف استفاده می‌کنند می‌توان از روش‌های MCOD[1]، DPCLUS[11]، CFINDER[12]، SPICi[13]، CMC[8]، MCL[14]، و ClusterONE[15] را نام برد. از مهم‌ترین روش‌هایی نیز که از اطلاعات زیستی در کنار اطلاعات مربوط به پیکربندی شبکه برهم‌کنش در استخراج کمپلکس‌های پروتئینی استفاده می‌کنند می‌توان به [16] HSS[19]، RNSC[18]، CORE[17]، COACH[16] و [20] TINCD اشاره کرد. از مهم‌ترین روش‌های پویای تشخیص کمپلکس نیز می‌توان به روش‌های TS-[22]PCD، TSN-[21]PCD و OCD-PCD [23] اشاره کرد.

در یک دسته‌بندی دیگر، روش‌های تشخیص بر اساس رویکرد جستجوی کمپلکس‌ها، به چهار حالت مبتنی بر ادغام و توسعه خوشه‌ها، مبتنی بر افزاز شبکه مبتنی بر انطباق شبکه و روش‌های فرامکاشفه‌ای تقسیم می‌شوند.

در روش‌های مبتنی بر ادغام و توسعه، کمپلکس‌ها در یک رویکرد پایین به بالا، با شروع از تعدادی خوشه پایه (که می‌تواند یک مثلث و یا کلیک باشد) ایجاد می‌شوند. در روش‌های مبتنی بر افزاز شبکه، در یک رویکرد بالا به پایین، هر شبکه به تعدادی زیرشبکه تقسیم می‌شود و در روش‌های مبتنی بر انطباق شبکه نیز از طریق انطباق شبکه‌های چند ارگانسیم زیستی با یکدیگر و استخراج مناطق مشترک، به کشف کمپلکس‌های جدید پرداخته می‌شود. روش‌های مبتنی بر ادغام و توسعه خوشه‌ها نیز به نوبه خود به دو دسته جستجوی چگالی همسایگی محلی LD و مبتنی بر معیارهای آماری SM تقسیم‌بندی می‌شوند. در روش‌های LD، هدف پیدا کردن زیرگراف‌های چگال با یک جستجوی محلی در مجاورت هر خوشه اولیه می‌باشد و روش‌های SM، از معیارهای آماری برای ترکیب خوشه‌ها استفاده می‌کنند. روش‌های مبتنی بر افزاز نیز به دو دسته جستجوی محلی مبتنی بر هزینه CL و شبیه‌سازی جریان FS تقسیم‌بندی می‌شوند. در روش‌های CL به افزاز شبکه بر اساس بیشینه‌کردن یک تابع هزینه پرداخته می‌شود و در روش‌های FS از مفهوم گام زدن تصادفی برای افزاز شبکه استفاده می‌شود [۲۳، ۲۴].

در یک جمع‌بندی درباره نقاط ضعف و قوت رویکردهای جستجوی مختلف باید یادآور شد که روش‌های مبتنی بر ادغام و توسعه از این حیث که امکان کشف خوشه (کمپلکس)‌های همپوشان را دارا می‌باشند، نسبت به روش‌های مبتنی بر افزاز که تنها خوشه‌های مجزا را کشف می‌کنند، مزیت دارند. همچنین کارایی روش‌های مبتنی بر انطباق، با توجه به محدودیت و ناقص بودن دانش در دسترس از ارگانسیم‌های مختلف عملاً با محدودیت‌های بسیاری مواجه می‌باشد. روش‌های فرامکاشفه‌ای نیز در تشخیص کمپلکس‌های پروتئینی نتوانسته‌اند به مزیت قابل توجهی نسبت به سایر روش‌ها دست یابند. هرچند امروزه، روش‌های تشخیص متنوعی در رویکرد مبتنی بر ادغام و توسعه و در هر دو دسته LD و SM ارائه شده

است که در زمره بهترین روش‌های تشخیص می‌باشند، اما همچنان چالش اصلی در یک روش تشخیص، چگونگی استفاده از دیگر داده‌های زیستی به منظور مقابله با نوفه موجود در شبکه‌های PPI و در نظر گرفتن شرایط پویای سلول در تشخیص کمپلکس‌ها می‌باشد [۲۵].

۲. مفاهیم اولیه

در این بخش، مروری بر مفاهیم اولیه‌ای که برای درک بهتر روش ارائه شده ضروری است انجام می‌گیرد.

شبکه برهم‌کنش پروتئین - پروتئین

به‌طور کلی اکثر فرآیندهای درون‌سلولی مانند انتقال سیگنال‌های الکتروشیمیایی، نقل و انتقال مواد، تنظیم موازنه یونی درون هسته و غیره، مستلزم تشکیل پیوند (برهم‌کنش) میان محصولات ژنی (اعم از پروتئین‌ها، RNA و غیره) می‌باشد. در سال‌های اخیر توسعه روش‌های آزمایشگاهی پیشرفته با بازده بالا (از قبیل ژل الکتروفورز دو بعدی و طیف‌سنج جرمی) در تعیین برهم‌کنش میان پروتئین‌ها منجر به ایجاد حجم بسیار بالایی از داده‌های تجربی شده است که کارکرد پروتئین‌ها را در شبکه‌های سلولی پیچیده منعکس می‌کند. شبکه‌های برهم‌کنش پروتئین-پروتئین یکی از منابع داده‌ای اصلی در زیست‌شناسی سامانه‌ها و مهم‌ترین منبع داده‌ای به منظور کشف کمپلکس‌های پروتئینی می‌باشند. در این شبکه‌ها، رئوس گراف متناظر با پروتئین‌ها و یال‌ها نتیجه یک نوع برهم‌کنش میان آن‌ها می‌باشند.

کمپلکس پروتئینی

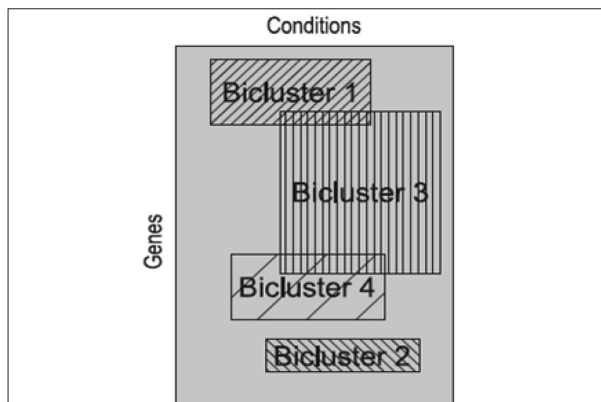
یک کمپلکس پروتئینی به منزله یک ماشین مولکولی است که شامل تعدادی پروتئین می‌باشد که در یک بازه زمانی و مکانی مشخص به یکدیگر متصل می‌شوند. اکثر پروتئین‌ها عملکردهای زیستی خود را در سلول زنده در قالب تشکیل کمپلکس‌ها انجام می‌دهند. در حقیقت بسیاری از پروتئین‌ها به عنوان اجزای کمپلکس‌های ضروری و

پایدار وجود دارند. برای مثال مولکول حیاتی هموگلوبین، یک ترکیب پایدار از چهار پروتئین گلوبول دار مشخص می‌باشد. همچنین بسیاری از آنزیم‌ها که نقش راه‌انداز را در بسیاری از واکنش‌های درون سلول به عهده‌دارند، مجموعه‌ای از پروتئین‌ها می‌باشند که به‌طور هم‌زمان به یکدیگر متصل شده و یک آنزیم را شکل می‌دهند. علاوه بر آن برهم‌کنش‌های ناپایدار و موقت مانند اثرات هورمون و یا اثرات یک سیگنال موقت، مستلزم تشکیل کمپلکس‌های پروتئینی می‌باشد. از لحاظ زیستی، کمپلکس‌های پروتئینی اجزاء کلیدی مولکولی برای انجام بسیاری از عملکردهای زیستی از قبیل رونوشت DNA، ترجمه mRNA، انتقال سیگنال سلولی، چرخه سلولی و غیره می‌باشند.

بیان ژن و ماتریس آن

ریزآرایه دی‌ان‌ای^۱، یک سطح جامد مشبک متشکل از هزاران خانه می‌باشد، به‌طوری‌که هر خانه در سطح شبکه، متناظر با یک ژن می‌باشد. به منظور انجام یک آزمایش ریزآرایه، دو نمونه سالم (یا کنترل) و نمونه هدف (که معمولاً یک mRNA از قبیل عامل رونویسی یا یک عامل بیماری و یا هر محصول ژنی مورد مطالعه می‌باشد، تهیه می‌شود. سپس، به هر یک از دو نمونه، یک رنگ به‌عنوان برچسب نظیر می‌شود. در مرحله بعدی، با مخلوط کردن محتوای دو نمونه، مقدار متناسبی از آن به صورت ترکیبات مایع بر روی سطح ریزآرایه ریخته می‌شود و در نهایت برحسب شدت و ضعف پیوند (با توجه به تمایل ترکیبات نوکلئوتیدی به اتصال به مکمل خود، این پیوند برقرار می‌شود) میان نمونه با هر یک از ژن‌های سطح ریزآرایه، به ازای هر ژن در هر آزمایش و به کمک روش‌های پردازش تصویر یک کمیّت عددی به دست می‌آید که بیانگر مقدار بیان ژن در آزمایش مربوطه است. به‌طور خلاصه می‌توان گفت که با یک آزمایش ریزآرایه امکان ثبت مقدار بیان هزاران ژن فراهم می‌شود. در هنگام استفاده از ریزآرایه، معمولاً تعدادی آزمایش متوالی به‌صورت یک سری از فواصل

1- DNA Micro array



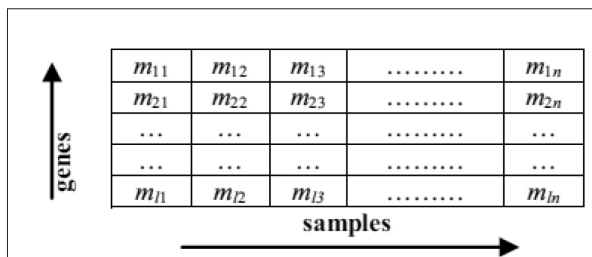
شکل ۲: خوشه‌بندی دوگانه ماتریس بیان ژن [۲۷]

حاوی تعدادی نمونه داده است که در زیرمجموعه‌ای از خصیصه‌ها دارای بیشترین شباهت به یکدیگر می‌باشند. شکل ۲ خوشه‌بندی دوگانه ماتریس بیان ژن، جهت استخراج مجموعه ژن‌های با الگوی بیان به هم‌وابسته را نشان می‌دهد.

گراف هستان‌شناسی ژن^۵

توسعه فناوری منجر به تولید داده‌ها و دانش زیادی در هریک از زمینه‌های پژوهشی شده است و این مسئله ضرورت استانداردسازی و سازمان‌دهی این حجم بالای اطلاعات را به‌منظور تسهیل دسترسی و استخراج دانش نهفته در آنها اجتناب‌ناپذیر ساخته است. با استفاده از گراف مفهومی هستان‌شناسی، مدل‌سازی دانش کشف‌شده در قالب تعدادی مفاهیم (ترم اطلاعاتی) و رابطه معنایی میان آنها امکان‌پذیر است. در یک گراف هستان‌شناسی، به طور عمده، دو رابطه is-a و part-of میان ترم‌های اطلاعاتی وجود دارد. رابطه is-a، به معنای یک نمونه خاص از یک مفهوم (مانند درخت سیب، یک نمونه از درخت است) و رابطه part-of، به معنای بخشی از یک مفهوم دیگر است (مانند منقار و پر در عقاب). یک گراف هستان‌شناسی ژن یا به اختصار GO، یک ساختار سلسله‌مراتبی درخت‌واره از ترم‌های زیستی است که رابطه is-a و part-of این ترم‌ها را نشان می‌دهد. در یک نگاه دقیق‌تر، گراف هستان‌شناسی ژن از سه گراف هستان‌شناسی بدون دور متناظر با مؤلفه

5- Gene Ontology Graph (GO graph)



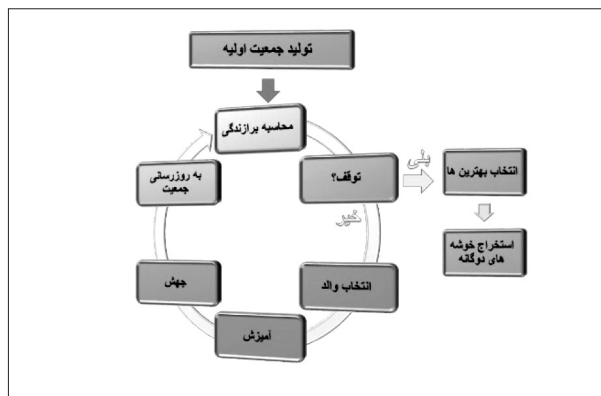
شکل ۱: ماتریس بیان ژن [۲۶]

زمانی یکسان درون چرخه سلولی و یا تعدادی آزمایش در شرایط آزمایشگاهی متفاوت (برای مثال بافت‌های مختلف و یا آزمایش در حضور محصولات ژنتیکی مختلف) انجام می‌گیرد. خروجی n آزمایش یک ریزآرایه شامل l ژن، یک ماتریس $M_{(l \times n)}$ به نام ماتریس بیان ژن می‌باشد که در آن m_{ij} مقدار بیان ژن i ام را در آزمایش j ام نشان می‌دهد (شکل ۱). به سطر متناظر با مجموعه مقادیر بیان یک ژن در ماتریس، پروفایل بیان آن ژن می‌گویند و ژن‌های دارای پروفایل بیان شبیه به هم را هم-بیان^۲ و ژن‌های دارای عملکرد یکسان در یک یا تعدادی از شرایط را هم-تنظیم^۳ می‌گویند.

خوشه‌بندی دوگانه^۴

فرض کنید یک مجموعه داده شامل n نمونه داده باشد که هر نمونه m خصیصه دارد. در الگوریتم‌های خوشه‌بندی، محاسبه فاصله دو نمونه داده، فاصله بین تمامی m خصیصه در نظر گرفته می‌شود. هر عنصر داده‌ای را در تعیین فاصله عنصر داده‌ای با بقیه عناصر لحاظ می‌کند و این در حالی است که دو عنصر داده‌ای می‌توانند فقط در زیرمجموعه‌ای از خصیصه‌ها به یکدیگر شبیه باشند. خوشه‌بندی دوگانه که به آن خوشه‌بندی دوبعدی نیز می‌گویند، یکی از روش‌های یادگیری ماشین است که عمل خوشه‌بندی را به طور هم‌زمان بر روی سطرها (نمونه‌ها) و ستون‌های (خصیصه‌های) ماتریس داده انجام می‌دهد. خروجی یک الگوریتم خوشه‌بندی دوگانه، تعدادی خوشه دوگانه از نمونه‌ها می‌باشد به طوری که هر خوشه دوگانه،

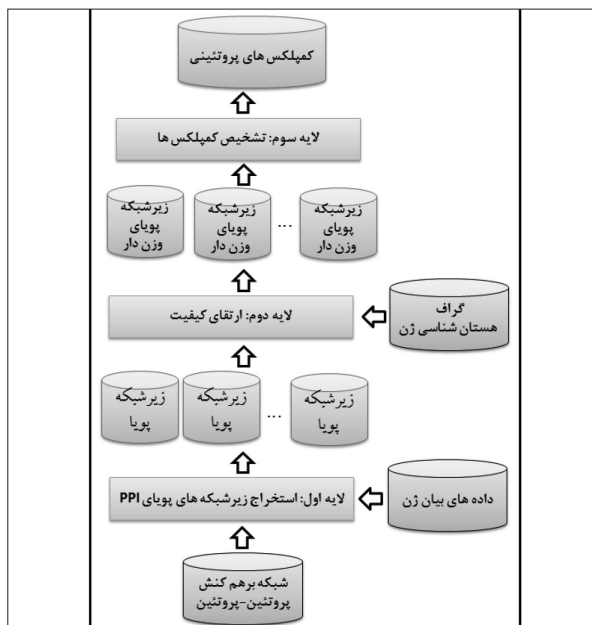
2- Co-express Genes
3- Co-regulated Genes
4- BiClustering



شکل ۴: مراحل اصلی الگوریتم GA-GCT

می باشد. نحوه کار بدین صورت است که ابتدا، توسط الگوریتم ژنتیک پیشنهادی GA-DCT، به خوشه بندی دوگانه داده های بیان ژن پرداخته می شود. در این الگوریتم (شکل ۴) به منظور استخراج خوشه های دوگانه از ماتریس بیان ژن استفاده می شود. به طور خلاصه، در مقایسه با سایر روش های خوشه بندی دوگانه، این الگوریتم دارای دو ویژگی متمایز می باشد: (۱) مبتنی بر روش فرامکاشفای الگوریتم ژنتیک و در نتیجه جستجوی خوب در فضای راه حل های مسئله (۲) مبتنی بر گسسته سازی ماتریس بیان ژن و در نتیجه مقاوم به نوفه موجود در این داده ها. نحوه کار بدین صورت است که ماتریس بیان ژن M با m سطر و n ستون داده شده است هر عنصر m_{ij} از ماتریس M ، مقدار بیان i امین ژن را در j امین نقطه زمانی نشان می دهد. یک خوشه دوگانه $B(i, j)$ یک زیر ماتریس از M می باشد. چنین خوشه دوگانه ای، به کمک یک رشته دودویی (به عنوان کروموزوم) شامل دو قسمت نشان داده می شود. قسمت اول کروموزوم به طول m ، متناظر با نمایه های ژن های موجود در خوشه دوگانه و قسمت دوم به طول n ، متناظر با نمایه های نقاط زمانی موجود در خوشه دوگانه می باشد. بنابراین، هر خوشه دوگانه به وسیله یک بردار دودویی با طول ثابت $n+m$ که برحسب سطرها و ستون های موجود در خوشه، تعدادی از بیت های آن در دو قسمت دارای مقدار ۱ و مابقی صفر می باشد (شکل ۵).

پس از حذف خوشه های کم اهمیت، متناظر با هر خوشه دوگانه، یک زیرشبکه PPI از شبکه اولیه PPI استخراج



شکل ۳: مدل سه لایه پیشنهادی

سلولی یا به اختصار CC، عملکرد مولکولی یا به اختصار MF و فرآیند زیستی یا به اختصار BP تشکیل شده است. هر ترم زیستی در GO، شامل تعدادی محصول ژنی از قبیل پروتئین است. برای مثال، یک ترم زیستی می تواند بیانگر یک فرآیند زیستی باشد که از چندین محصول ژنی تشکیل شده است. از این حیث، می توان گفت که هر ترم، تعدادی ژن را تفسیر می کند. این تفسیر، می تواند مستقیم و یا غیرمستقیم و از طریق رابطه سلسله مراتبی درخت باشد و معمولاً پایه ای برای محاسبه مقدار اطلاعات فراهم شده توسط ترم می باشد.

۳. مواد و روش ها

شکل ۳ مدل پیشنهادی مقاله را در قالب سه لایه نشان می دهد. در ادامه، به معرفی مختصر هر یک از لایه های مدل پرداخته می شود. با توجه به این که در نظر گرفتن پویایی سلول در تشخیص کمپلکس های پروتئینی، یکی از چالش های اصلی مسئله تشخیص می باشد. لایه اول، با استفاده از داده های بیان ژن، راهکاری را به این منظور پیشنهاد می دهد. ورودی لایه اول، شبکه PPI و یک ماتریس بیان ژن می باشد و خروجی این لایه، تعدادی زیرشبکه PPI

1	2	3	4	...	m	1	2	3	...	n
0	1	1	0	...	1	1	0	1	...	0
Genes						Time-points				

شکل ۵: بیان یک خوشه دوگانه در یک کروموزوم

یک سنجنده شباهت معنایی ترکیبی (مبتنی بر گره و مبتنی بر یال) دو مرحله‌ای به نام TSSS ارائه شده است که توانسته است نسبت به سنجنده‌های متعارف معنایی به دقت بهتری دست یابد.

تابع $SSW(x,y)$ وزن متناظر با برهم‌کنش میان پروتئین‌های x و y را با استفاده از مقادیر شباهت مبتنی بر پیکربندی و شباهت معنایی این دو پروتئین محاسبه می‌کند. براساس رابطه (۱)، $SSW(x,y)$ ترکیبی وزن‌دار از دو تابع $TSSS(x,y)$ و $Sim_{top}(x,y)$ می‌باشد که α وزن هر یک از این دو نوع شباهت در محاسبه SSW می‌باشد.

$$SSW(x,y) = \alpha \times Sim_{top}(x,y) + (1-\alpha) \times TSSS(x,y) \quad (1)$$

در این رابطه، $Sim_{top}(x,y)$ می‌تواند هر یک از توابع محاسبه مقدار شباهت مبتنی بر پیکربندی شبکه PPI باشد. بدیهی است که با انتخاب مقدار مناسب برای α ، می‌توان تأثیر هر یک از دو مقدار شباهت مبتنی بر پیکربندی شبکه و گراف هستان‌شناسی را کنترل کرد. تابع $TSSS(x,y)$ مقدار شباهت معنایی پروتئین‌های x و y را در دو مرحله محاسبه می‌کند. در مرحله اول ماتریس شباهت‌های ترم-ترم SS را برای تمامی ترم‌های موجود در گراف GO محاسبه می‌کند و در مرحله دوم، با استفاده از ماتریس SS شباهت‌های معنایی میان پروتئین‌ها محاسبه می‌کند. در مرحله اول، شباهت معنایی میان دو ترم x و y در GO بر اساس رابطه (۲) محاسبه می‌شود که در آن C_x و C_y دو ضریب شباهت معنایی مبتنی بر ترم می‌باشند. رابطه $C_t(x,y)$ تابعی از اجداد مشترک دو ترم x و y می‌باشد و بر اساس (۳) محاسبه می‌شود. منظور از A_x و A_y مجموعه اجداد ترم‌های x و y و منظور از $IIC(a)$ نسخه معکوس شده $IC(a)$ ، تابع محتوای اطلاعات^۱ بر اساس رابطه (۴) برای ترم a می‌باشد.

$$SS(x,y) = C_t(x,y) \times C_g(x,y) \quad (2)$$

$$C_t(x,y) = \frac{2 \times IIC(a)}{SV(x) + SV(y)} \quad (3)$$

$$IIC(a) = 1 - IC(a) \quad (4)$$

می‌شود به طوری که هر زیرشبکه، به شامل برهم‌کنش‌های بین پروتئینی ایجاد شده در بازه‌ای از زمان می‌باشد و در نتیجه می‌توان آن را به عنوان یک زیرشبکه پویا نام‌گذاری کرد. مجموعه زیرشبکه‌های پویای PPI استخراج شده، ورودی لایه دوم مدل می‌باشد. یکی از چالش‌های مهم مسئله، «وجود نوفه در شبکه PPI می‌باشد [۲۸] و از طرفی منابع اطلاعاتی دیگری مانند گراف هستان‌شناسی ژن، وجود دارد که استفاده از آن‌ها در کنار شبکه‌های PPI می‌تواند به کاهش نوفه و افزایش دقت این شبکه‌ها منجر شود. هدف اصلی لایه دوم، ارتقای کیفیت شبکه‌های PPI موجود از طریق وزن‌دار کردن برهم‌کنش‌های بین پروتئینی می‌باشد. منظور از وزن نیز یک کمیت عددی بین صفر و یک می‌باشد که بیانگر قوت و یا ضعف اتصال میان دو پروتئین می‌باشد. در این لایه، هر زیرشبکه از مجموعه زیرشبکه‌های پویای PPI به دست آمده در لایه قبلی، به صورت جداگانه وزن‌دار می‌شود. این عمل باعث می‌شود که مبنای پویایی به دست آمده در لایه قبل، در این لایه نیز لحاظ گردد و در واقع، عمل وزن‌دار کردن زیرشبکه‌های پویای PPI به صورت محلی انجام می‌گیرد. وزن محاسبه شده به ازای هر زوج پروتئین در دو سر یک یال (برهم‌کنش) با روش پیشنهادی SSW تعیین می‌شود که ترکیبی از دو مقدار شباهت دو پروتئین متناسب با دو منبع داده‌ای مختلف می‌باشد. این منابع داده‌ای عبارتند از شبکه PPI و گراف هستان‌شناسی ژن. شبکه PPI، اطلاعات پیکربندی (برهم‌کنش‌های فیزیکی) لازم برای اندازه‌گیری شباهت را منعکس می‌کند و گراف هستان‌شناسی ژن، دربرگیرنده رابطه و جایگاه عملکردی پروتئین‌ها در کل ژنوم و در یک ساختار سلسله مراتبی می‌باشد. در این لایه، برای محاسبه شباهت مبتنی بر گراف هستان‌شناسی ژن،

CAMWI Method
Input: G(V, E): A weighted PPI network, α, β : seed and set thresholds
Output: Complexes: the detected protein complexes
1-Seeds \leftarrow SeedSelection(G, α)
2-Cores \leftarrow CoreFinding(G, Seeds)
3-InitComplexes \leftarrow CoreGrowing(G, Cores, β)
4-Complexes \leftarrow Filtering(InitComplexes)
شکل ۶: مراحل مختلف روش CAMWI

زیاد و نزدیک به هم برای اکثر زوج‌های پروتئینی جلوگیری شود.

لایه سوم، به هدف اصلی مدل یعنی تشخیص کمپلکس‌های پروتئینی از درون مجموعه زیرشبکه‌های پویای وزن‌دار PPI به دست آمده از لایه دوم می‌پردازد. روش ارائه شده با نام CAMWI، کمپلکس‌ها را در یک قالب هسته-افزونه تشخیص می‌دهد و ذاتاً برای شبکه‌های وزن‌دار ارائه شده است هرچند که استفاده از آن در شبکه‌های بدون وزن نیز برتری آن را نسبت به سایر روش‌ها نشان می‌دهد. چهار قدم اساسی این روش عبارت است از (۱) انتخاب دانه‌های اولیه به کمک مفهوم ضریب خوشه‌بندی وزن‌دار (۲) شناسایی هسته پیرامون هر دانه با استفاده از مفهوم چگالی وزن‌دار در یک رویه حریصانه، به طوری که هر هسته، دارای چگالی وزن‌دار بیشینه می‌باشد (۳) توسعه هر هسته کمپلکس با اضافه کردن پروتئین‌های افزونه و تشکیل کمپلکس (۴) پالایش و حذف کمپلکس‌های تکراری و یا بسیار شبیه به هم. شکل ۶ الگوریتم CAMWI را نشان می‌دهد.

۴. نتایج

در این بخش، نتایج ارزیابی مدل پیشنهادی که با تجمیع لایه‌های سه‌گانه به دست می‌آید، ارائه می‌شود. مدل پیشنهادی، در لایه اول، به کمک روش پیشنهادی GA-DCT، به پویاسازی شبکه PPI می‌پردازد، به گونه‌ای که لایه اول، با دریافت شبکه PPI ورودی، زیرشبکه‌های پویای PPI را بر مبنای داده‌های بیان ژن تولید می‌کند. سپس، مجموعه

مقدار $SV(a)$ ، نیز مقدار معنایی ترم a ، متناظر با جمع معنایی مقادیر اطلاعات معکوس شده در رابطه (۵) می‌باشد.

پارامتر W_t نیز با استفاده از رابطه (۶) محاسبه می‌شود.

$$SV(a) = \sum_{t \in Aa} W_t \times IIC(t) \quad (5)$$

$$W_t = \text{Min}(\text{distance } t \text{ to GO root}) / \text{Min}(\text{length of path from GO root to a leaf node cross } t) \quad (6)$$

مطابق روابط اخیر، ضریب $C_g(x,y)$ در رابطه (۳)، ترکیبی از رویکرد مبتنی بر گره و مبتنی بر یال می‌باشد. مقدار $C_g(x,y)$ نیز، به عنوان دومین ضریب شباهت معنایی در یک رویکرد مبتنی بر مجموعه، بر اساس رابطه (۷) محاسبه می‌شود که D_x و D_y مجموعه نواده‌های ترم‌های x و y می‌باشد. بنابراین، روش TSS در مرحله اول، ترکیبی از رویکردهای مبتنی بر گره، مبتنی بر یال و مبتنی بر مجموعه را برای اندازه‌گیری شباهت معنایی استفاده می‌کند. یادآوری می‌گردد این مرحله فقط یک بار پس از دریافت گراف هستان‌شناسی ژن انجام می‌شود.

$$C_g(x,y) = |D_x \cap D_y| / |D_x \cup D_y| \quad (7)$$

در مرحله دوم، شباهت معنایی میان پروتئین‌های a و b $TSS_t(a,b)$ ، بر اساس رابطه (۸) و با استفاده از ماتریس SS حاصل از مرحله اول محاسبه می‌شود.

$$TSS_t(a,b) = ss_t(a,b) \times g(a,b) \quad (8)$$

$$t \in \{MAX, AVG, BMA\}$$

$sst(a,b)$ ، با رویکردی مبتنی بر زوج، شباهت میان پروتئین‌های a و b را برحسب شباهت‌های ترم-ترم در مجموعه ترم‌های تفسیرکننده آن‌ها و به کمک یکی از قوانین MAX, AVG و یا BMA محاسبه می‌کند. $g(a,b)$ نیز یک سنجنده مبتنی بر گراف است که با استفاده از رابطه (۹) محاسبه می‌شود.

$$g(a,b) = \frac{\sum_{t \in f_\theta(N_a \cap N_b)} IC(t)}{\sum_{t \in f_\theta(N_a \cup N_b)} IC(t)} \quad (9)$$

که N_a و N_b مجموعه ترم‌های تفسیر پروتئین‌های a و b و f_θ ، یک تابع پالایه می‌باشد که برحسب یک مقدار آستانه θ ، ترم‌های با عمق کم (نزدیک به ریشه درخت GO) را از مجموعه از ترم‌ها حذف می‌کند تا از تولید مقادیر شباهت

جدول ۱: منابع داده ای مورد استفاده

توصیف	نوع	مجموعه داده
تعداد ۵۹۷۴۸ برهم کنش میان ۵۶۴۰ پروتئین	شبکه برهم کنش	BioGrid [۲۹]
سطح بیان ۳۵۵۲ ژن در ۱۲ آزمایش	ماتریس بیان ژن	YMC [۳۰]

جدول ۲: مشخصات مجموعه کمپلکس های محک مورد استفاده در ارزیابی روش

مجموعه محک	تعداد کمپلکس	تعداد پروتئین	اندازه کمپلکس			
			3 >	10 - 3	11-25	25 <
CYC2008 [31]	408	1627	172	204	27	5

ابتدا، برهم کنش های بین پروتئینی به روش SSW (لایه دوم) و زن دار می شود، سپس عمل تشخیص به کمک روش CAMWI انجام می گیرد.

- لایه اول، سپس لایه سوم، با نام DCT+CAMWI: این روش به نام Bi-CAMWI نام گذاری شد. در این روش، ابتدا، با استفاده از روش پیشنهادی GA-DCT، زیر شبکه های پویا، از شبکه PPI اولیه، استخراج می شود و سپس عمل تشخیص به کمک روش CAMWI بر روی هر زیر شبکه انجام می شود. - لایه دوم، لایه اول و سپس لایه سوم، با نام SSW+DCT+CAMWI: در این روش، در ابتدا، شبکه PPI به کمک روش SSW به صورت سراسری و زن دار می شود و سپس با استفاده از روش پیشنهادی GA-DCT، زیر شبکه های پویای زن دار، از شبکه PPI استخراج می شود و سپس عمل تشخیص به کمک روش CAMWI بر روی هر زیر شبکه انجام می شود.

- لایه اول، لایه دوم و سپس لایه سوم، با نام DCT+SSW+CAMWI: این روش، مدل سه لایه پیشنهادی مقاله می باشد.

دقت در نتایج به دست آمده نتایج زیر را نشان می دهد: ۱. مقایسه میان دو روش CAMWI و SSW+CAMWI نشان می دهد که وزن دار کردن شبکه PPI، از این حیث که باعث می شود زیر برنامه انتخاب دانه ها به صورت هوشمندانه تر، زیر مجموعه کوچکتری از دانه های اولیه را انتخاب کند؛ منجر به افزایش مقدار سنجنده صحت می شود. از طرف دیگر بدیهی است که کاهش تعداد کمپلکس های تشخیصی به کاهش مقدار سنجنده بازخوانی منجر شده است، هر چند که در کل مقدار سنجنده $f_measure$ افزایش

زیر شبکه های پویای PPI، جهت ارتقای کیفیت در اختیار لایه دوم قرار می گیرد. لایه دوم، با استفاده از روش پیشنهادی SSW، به وزن دار کردن مستقل زیر شبکه های پویا می پردازد و مجموعه زیر شبکه های پویای وزن دار PPI در اختیار لایه سوم که وظیفه اصلی تشخیص کمپلکس ها را بر عهده دارد، قرار می گیرد. لایه سوم، با استفاده از روش پیشنهادی CAM-WI، کمپلکس های پروتئینی را در هر زیر شبکه پویای وزن دار تشخیص می دهد. مجموعه کمپلکس های به دست آمده به ازای تمامی زیر شبکه ها، پس از پالایش و حذف کمپلکس های تکراری و یا بسیار شبیه به هم، مجموعه کمپلکس های نهایی تشخیص داده شده توسط مدل پیشنهادی می باشد.

به منظور ارزیابی روش محاسباتی ارائه شده از سه نوع مجموعه داده مختلف استفاده شده است: یک شبکه برهم کنش پروتئین-پروتئین، یک مجموعه داده محک از کمپلکس های پروتئینی کشف شده به روش آزمایشگاهی و یک مجموعه داده بیان ژن. جدول ۱ و ۲ ویژگی منابع داده ای مورد استفاده در این ارزیابی را نشان می دهد.

از آنجا که در تحلیل های ارائه شده و در مقایسه با روش های گوناگون، با توجه به شرایط و دستگاه مقایسه، از مجموعه داده های متفاوتی استفاده شده است، در این بخش، در یک دستگاه ارزیابی یکسان، هر یک از مراحل توسعه مدل به عنوان یک روش نام گذاری می شود و با حالت کامل مدل پیشنهادی (شامل سه لایه) مقایسه شد. نتایج مربوطه در شکل ۲ نشان داده شده است. در این مقایسه، از روش های زیر استفاده شده است:

- فقط لایه سوم، با نام CAMWI.

- لایه دوم، سپس لایه سوم، با نام SSW+CAMWI: در

یافته است. این نتایج، تأییدی است بر وجود مقادیر خطای بالای مثبت کاذب در شبکه PPI، به گونه‌ای که کاهش نوفه موجود در شبکه PPI از طریق وزن‌دار کردن آن، توانسته است از تولید جواب‌های اشتباه فراوان توسط روش جلوگیری نموده، منجر به افزایش قابل ملاحظه‌ای در مقدار سنجنده صحت شود. این مسئله در بهبود مقدار سنجنده MMR نیز مشاهده می‌شود که بیانگر افزایش میزان انطباق کمپلکس‌های پیش‌بینی شده با کمپلکس‌های واقعی پس از ارتقای کیفیت شبکه PPI با روش SSW می‌باشد.

۲. مقایسه میان دو روش CAMWI و DCT+CAMWI نشان می‌دهد که در نظرگرفتن شرایط پویای سلول، منجر به بازیابی مجموعه‌ای جدید از کمپلکس‌های پروتئینی شده است و به همین دلیل نرخ سنجنده بازخوانی در DCT+CAMWI افزایش قابل ملاحظه‌ای نسبت به CAMWI پیدا کرده است.

۳. مقایسه میان دو روش SSW+CAMWI و DCT+CAMWI نشان می‌دهد که تعدادی از کمپلکس‌های پروتئینی وجود دارند که جز از طریق در نظرگرفتن شرایط پویای سلول قابل بازیابی نمی‌باشند؛ به همین دلیل نرخ سنجنده بازخوانی در DCT+CAMWI افزایش قابل ملاحظه‌ای نسبت به SSW+CAMWI پیدا کرده است؛ از طرف دیگر به علت وزن‌دار بودن برهم‌کنش‌ها در SSW+CAMWI با استدلالی مشابه حالت ۱، نرخ صحت در آن بیشتر می‌باشد.

۴. در مقایسه مدل کامل DCT+SSW+CAMWI با روش دوم یعنی SSW+CAMWI، بدیهی است که در نظرگرفتن پویایی از این حیث که توانایی تشخیص کمپلکس‌های پویا را امکان‌پذیر می‌سازد، منجر به افزایش نرخ بازخوانی شده است و از طرف دیگر به علت افزایش تعداد کمپلکس‌های تشخیصی (به علت افزایش تعداد شبکه‌های ورودی) نرخ صحت کاهش یافته است.

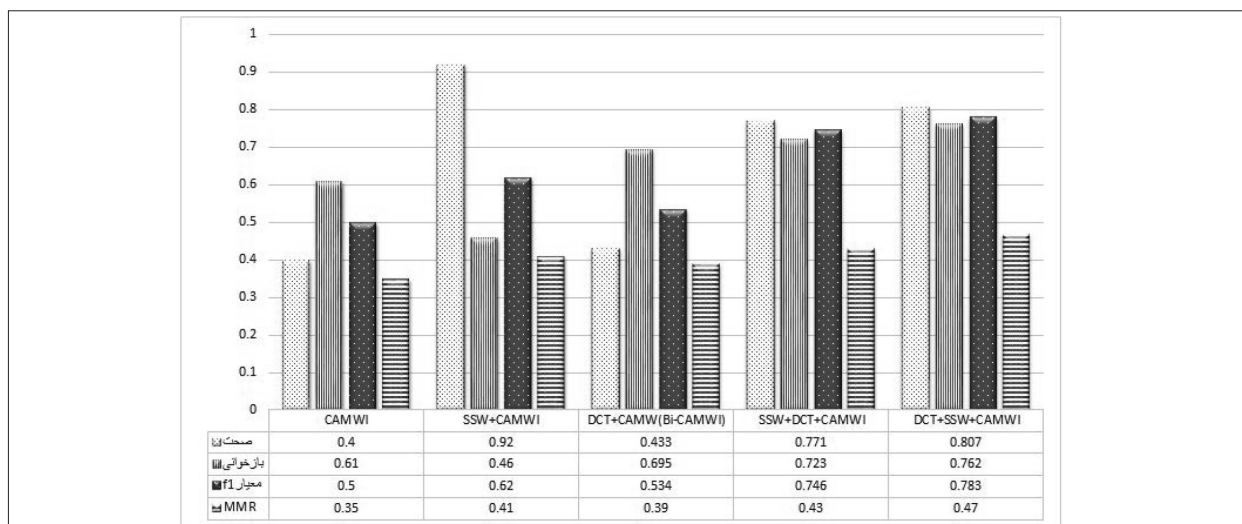
۵. در مقایسه مدل کامل با روش سوم یعنی DCT+CAMWI، بدیهی است که وزن‌دار کردن برهم‌کنش‌ها در مدل کامل، توانسته است نرخ صحت را افزایش دهد.

همچنین، این بار برعکس مقایسه روش اول و دوم، از آنجا که عمل وزن‌دار کردن به صورت محلی (برای هر زیرشبکه، به صورت مستقل) انجام گرفته است، مقدار بازخوانی نیز افزایش یافته است. در واقع، وزن‌دار کردن محلی، به بازخوانی بهتر کمپلکس‌ها منجر شده است. این مسئله در بهبود مقدار سنجنده MMR نیز مشاهده می‌شود که بیانگر افزایش میزان انطباق کمپلکس‌های پیش‌بینی شده با کمپلکس‌های واقعی در مدل کامل می‌باشد.

۶. مقایسه مدل چهارم SSW+DCT+CAMWI با مدل کامل DCT+SSW+CAMWI نیز نشان می‌دهد که وزن‌دار کردن محلی در مدل کامل، منجر به افزایش مقدار سنجنده صحت شده است، این مسئله نشان دهنده تاثیر منفی برهم‌کنش‌های بین پروتئینی نوفه‌دار در حالت وزن‌دار کردن سراسری است. همچنین، افزایش کیفیت اتصالات زیرشبکه پویا در حالت مدل کامل، به بازخوانی بهتر کمپلکس‌ها نیز منجر شده است. این مسئله در بهبود مقدار سنجنده MMR نیز مشاهده می‌شود که بیانگر افزایش میزان انطباق کمپلکس‌های پیش‌بینی شده با کمپلکس‌های واقعی در مدل کامل می‌باشد.

در ادامه، در یک ارزیابی متفاوت، روش‌های تشخیص قبلی که از ترکیبات مختلف لایه‌های سه‌گانه مدل پیشنهادی ساخته شده‌اند، از حیث تعداد کمپلکس‌های تشخیص داده شده در بازه‌های مختلف برحسب اندازه کمپلکس، مقایسه شدند. نتایج مربوطه در جدول ۲ نشان داده شده است. این نتایج، به طور خاص، اثر وزن‌دار کردن شبکه و کاهش نوفه آن به کمک روش SSW را در کاهش تعداد کمپلکس‌های پیش‌بینی شده با اندازه کوچکتر از ۳ (کمپلکس‌های نوفه‌دار) نشان می‌دهد.

با توجه به ماهیت تصادفی مدل سه لایه پیشنهادی، نتایج آن از یک اجرا به اجرای دیگر متفاوت است. بنابراین، به منظور اندازه‌گیری بازه تغییرات نتایج، مدل پیشنهادی در بهترین مقادیر پارامتر به ازای شبکه برهم‌کنش BioGrid و داده‌های بیان ژن YMC و مجموعه محک CYC2008،



شکل ۷: نتایج بررسی مدل پیشنهادی

جدول ۳: مقادیر میانگین و انحراف معیار سنجنده‌های مختلف در ۱۰ اجرای یکسان مدل پیشنهادی با بهترین مقادیر پارامتر

وارینانس	میانگین	سنجنده
۰.۰۰۱۲	۰.۷۸	صحت
۰.۰۰۱۴	۰.۷۳	بازخوانی
۰.۰۰۱۳۵	۰.۷۶	f ₁

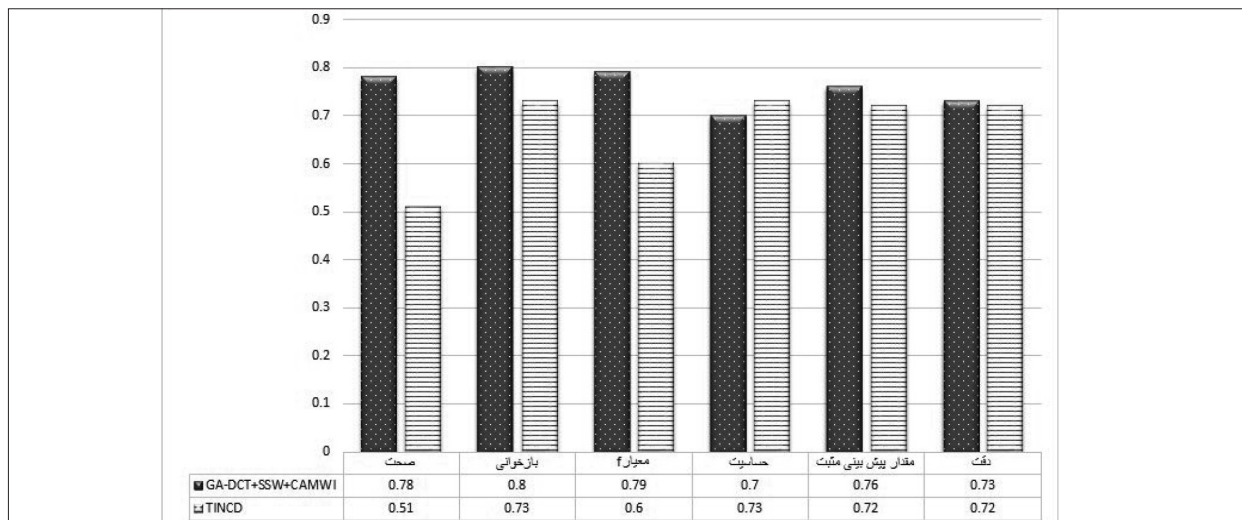
نشان داده شده است. مطابق نتایج به دست آمده، با وجود این که روش TINCD از نتایج ۱۱ روش دیگر تشخیص کمپلکس از روی شبکه PPI و ۵ روش تشخیص کمپلکس‌ها از داده‌های آزمایشگاهی TAP استفاده می‌کند، اما هم‌چنان روش پیشنهادی مقاله در اکثر سنجنده‌های ارزیابی به نتایج بهتری دست یافته است. به طور خاص، یکی از نقاط ضعف روش TINCD، تعداد بسیار زیاد کمپلکس‌های گزارش شده توسط آن است که منجر به کاهش مقدار سنجنده صحت شده است [۳۲]. نتایج فوق، تأیید دیگری بر ضرورت پویاسازی فرآیند تشخیص و کاهش نوفه شبکه PPI است و توانایی مدل پیشنهادی را نشان می‌دهد که توانسته است با وجود عدم استفاده از داده‌های TAP، به دقت بیشتری دست یابد. این نتایج، به ازای شبکه PPI ورودی DIP و مجموعه محک CYC2008 به دست آمده است. هم‌چنین، یادآوری می‌گردد که نتایج روش TINCD، بر مبنای مجموعه کمپلکس‌های گزارش شده توسط آن‌ها به

جدول ۲: مقایسه ترکیبات مختلف لایه‌های سه‌گانه مدل پیشنهادی بر حسب تعداد کمپلکس‌های پیش‌بینی شده در اندازه‌های مختلف

روش	تعداد کمپلکس‌های پیش‌بینی شده			
	کمتر از ۳	۳-۱۰	۱۱-۲۵	بیشتر از ۲۵
CAMWI	۱۰۵	۲۱۱	۱۱	۱
SSW+CAMWI	۵۱	۵۷	۵	۲
DCT+CAMWI	۹۶	۱۷۱	۶	۱
SSW+DCT+CAMWI	۷۵	۶۸	۴	۳
DCT+SSW+CAMWI	۶۷	۷۲	۳	۲

به تعداد ۱۰ بار اجرا شد و مقادیر سنجنده‌های صحت، بازخوانی و سنجنده اندازه‌گیری شد. محاسبه میانگین و انحراف معیار هر یک از سنجنده‌ها در جدول ۳ نشان داده شده است. این نتایج نشان می‌دهد که برای دستیابی به یک دقت قابل قبول برای روش، اجرای حداقل ۱۰ تکرار از الگوریتم ضروری است.

در ادامه ارزیابی‌های این بخش، به مقایسه مدل سه لایه پیشنهادی با روش TINCD [۲۰] پرداخته شد. روش TINCD، از این حیث که یکی از جامع‌ترین روش‌هایی است که تاکنون به منظور تشخیص کمپلکس‌های پروتئینی ارائه شده است، یک معیار ارزیابی مناسب می‌باشد. نتایج ارزیابی به ازای سنجنده‌های صحت، بازخوانی، f₁، حساسیت، مقدار پیش‌بینی‌شده مثبت و دقت، در شکل ۳



شکل ۸-۷: مقایسه مدل پیشنهادی و روش TINCD

ازای همین مجموعه داده‌ها می‌باشد.

۵. نتیجه‌گیری

با توجه به نقش کمپلکس‌های پروتئینی در انجام بسیاری از کارکردهای سلولی موجودات زنده، کشف آن‌ها می‌تواند به درک بهتر فرآیندهای سلولی و توسعه کاربردهای مبتنی بر مهندسی زیستی منجر شود. پیشرفت ابزارها و فناوری‌های آزمایشگاهی در دهه گذشته منجر به تولید و انباشت داده‌های زیادی از کارکرد اجزاء سلول شده است. این امر، در کنار محدودیت‌های روش‌های آزمایشگاهی، ضرورت مدل‌سازی سیستمی از طریق بررسی و شناخت سلول در سطح سیستمی و در قالب برهم‌کنش میان اجزاء را اجتناب‌ناپذیر ساخته است. رویکردهای محاسباتی ارائه‌شده تاکنون به منظور تشخیص کمپلکس‌های پروتئینی، به‌طور عمده بر خوشه‌بندی شبکه PPI تمرکز دارند. چالش‌های عمده مسئله تشخیص کمپلکس‌های پروتئینی عبارت است از (۱) ضرورت در نظر گرفتن شرایط پویای سلول در فرآیند تشخیص، (۲) مقدار زیاد نوفه از نوع مثبت کاذب در شبکه‌های PPI موجود و (۳) استفاده بیشینه از دانش زیستی موجود در فرآیند تشخیص. در این مقاله مدلی سه لایه مبتنی بر رویکرد زیست‌شناسی سامانه‌ها برای تشخیص کمپلکس‌های

پروتئینی از شبکه‌های PPI ارائه شده است. طراحی لایه‌های سه‌گانه به‌گونه‌ای است که هر لایه، وظیفه مدیریت و کاهش یکی از چالش‌های سه‌گانه مسئله تشخیص را بر عهده دارد. نوآوری‌های اصلی مدل پیشنهادی عبارت‌اند از:

۱- پیشنهاد رویکرد مبتنی بر خوشه‌بندی دوگانه ماتریس بیان ژن، به منظور در نظر گرفتن ویژگی پویای برهم‌کنش‌های بین پروتئینی در تشخیص کمپلکس‌های پروتئینی.

۲- ارائه روش خوشه‌بندی دوگانه ماتریس بیان ژن، مبتنی بر الگوریتم ژنتیک با تابع برازندگی کارا با قابلیت کشف الگوهای پیچیده (و سازگار با مسئله تشخیص کمپلکس‌های پروتئینی) و مقاوم به نوفه موجود در داده‌های بیان ژن.

۳- ارائه روش وزن‌دار کردن شبکه PPI از طریق ترکیب وزن‌دار شباهت معنایی مبتنی بر گراف هستان‌شناسی ژن و شباهت مبتنی بر پیکربندی شبکه PPI و همچنین، ارائه سنجنده TSSS به منظور محاسبه شباهت معنایی دو پروتئین از روی گراف هستان‌شناسی ژن.

۴- پیدا کردن هسته کمپلکس با استفاده از مفاهیم چگالی وزن‌دار و ضریب خوشه‌بندی وزن‌دار، در یک رویه حریصانه، به‌طوری‌که هسته هر کمپلکس دارای چگالی وزن‌دار بیشینه می‌باشد.

11. Altaf-Ul-Amin, M., et al., Development and implementation of an algorithm for detection of protein complexes in large interaction networks. *BMC bioinformatics*, 2006. 7(1): p. 207.

12. Adamcsek, B., et al., CFinder: locating cliques and overlapping modules in biological networks. *Bioinformatics*, 2006. 22(8): p. 1021-1023.

13. Jiang, P. and M. Singh, SPICi: a fast clustering algorithm for large biological networks. *Bioinformatics*, 2010. 26(8): p. 1105-1111.

14. Enright, A.J., S. Van Dongen, and C.A. Ouzounis, An efficient algorithm for large-scale detection of protein families. *Nucleic acids research*, 2002. 30(7): p. 1575-1584.

15. Nepusz, T., H. Yu, and A. Paccanaro, Detecting overlapping protein complexes in protein-protein interaction networks. *Nature methods*, 2012. 9(5): p. 471-472.

16. Wu, M., et al., A core-attachment based method to detect protein complexes in PPI networks. *BMC bioinformatics*, 2009. 10(1): p. 169.

17. Leung, H.C., et al., Predicting protein complexes from PPI data: a core-attachment approach. *Journal of Computational Biology*, 2009. 16: p. 133-144.

18. King, A.D., N. Pržulj, and I. Jurisica, Protein complex prediction via cost-based clustering. *Bioinformatics*, 2004. 20(17): p. 3013-3020.

19. Sharan, R., et al., Identification of protein complexes by comparative analysis of yeast and bacterial protein interaction data. *Journal of Computational Biology*, 2005. 12(6): p. 835-846.

20. Ou-Yang, L., et al., A two-layer integration framework for protein complex detection. *BMC bioinformatics*, 2016. 17(1): p. 100.

21. Li, M., et al., Towards the identification of protein complexes and functional modules by integrating PPI network and gene expression data. *BMC bioinformatics*, 2012. 13(1): p. 109.

22. Ou-Yang, L., et al., Detecting temporal protein complexes from dynamic protein-protein interaction networks. *BMC bioinformatics*, 2014. 15(1): p. 335.

23. Lakizadeh, A., S. Jalili, and S.-A. Marashi, PCD-GED: Protein complex detection considering PPI Dynamics based on time series gene expression data. *Journal of theoretical biology*, 2015. 378: p. 31-38.

24. Lakizadeh, A., S. Jalili, and S. Marashi, CAMWI: Detecting protein complexes using weighted clustering coefficient and weighted density. *Computational biology and chemistry*, 2015. 58: p. 231-240.

25. Lakizadeh, A. and S. Jalili, BiCAMWI: a genetic-based biclustering algorithm for detecting dynamic protein complexes. *PloS one*, 2016. 11(7).

26. Wang, J., et al., Recent advances in clustering methods for protein interaction networks. *BMC genomics*, 2010. 11(S3): p. S10.

27. Flores, J.L., et al., A new measure for gene expression biclustering based on non-parametric correlation. *Computer methods and programs in biomedicine*, 2013. 112(3): p. 367-397.

28. Yu, Y. and Z. Zheng, Protein Complex Identification Based on Weighted PPI Network with Multi-Source Information. *Journal of theoretical biology*, 2019.

29. Chatr-Aryamontri, A., et al., The BioGRID interaction database: 2013 update. *Nucleic acids research*, 2012. 41(D1): p. D816-D823.

30. Madeira, S.C. and A.L. Oliveira, A polynomial time biclustering algorithm for finding approximate expression patterns in gene expression time series. *Algorithms for Molecular Biology*, 2009. 4(1): p. 8.

31. Pu, S., et al., Up-to-date catalogues of yeast protein complexes. *Nucleic acids research*, 2009. 37(3): p. 825-831.

32. Beer, L.A. and D.W. Speicher, Protein detection in gels using fixation. *Current protocols in protein science*, 2018. 91(1): p. 10.5. 1-10.5. 20.

با توجه به گستردگی پژوهش‌های انجام گرفته در لایه‌های سه‌گانه مدل پیشنهادی، محورهای پژوهشی بسیار و متنوعی وجود دارد که توسعه هر یک می‌تواند به افزایش کارایی مدل پیشنهادی منجر شود. در ادامه به تعدادی از محورهای پژوهشی اشاره می‌شود. برای مثال همچنان، ضرورت ارائه یک روش کارای خوشه‌بندی گراف که با ماهیت داده‌های برهم‌کنش سازگاری بیشتری داشته باشد، وجود دارد. برای مثال، ارائه راهکارهای جدید و یا توسعه راهکارهای موجود از طریق مفاهیمی مانند خوشه‌بندی فازی گراف و یا خوشه‌بندی مبتنی بر شار عبوری از شبکه می‌تواند مفید باشد. همچنین، با توجه به نوفه بسیار زیاد شبکه‌های PPI، استفاده و توسعه روش‌های خوشه‌بندی مقاوم به نوفه یک رهیافت پژوهشی می‌باشد. همچنین، در مورد ارتقای شبکه PPI، همچنان، استخراج میزان شباهت با استفاده از منابع داده‌ای دیگر مانند داده‌های بیان ژن و یا شبکه‌های متابولیت و ... در ترکیب با سایر سنجنده‌های شباهت، می‌تواند به ارتقای کارتر شبکه PPI منجر شود.

منابع

1. Bader, G.D. and C.W. Hogue, An automated method for finding molecular complexes in large protein interaction networks. *BMC bioinformatics*, 2003. 4(1): p. 2.

2. Collins, S.R., et al., Toward a comprehensive atlas of the physical interactome of *Saccharomyces cerevisiae*. *Molecular & Cellular Proteomics*, 2007. 6(3): p. 439-450.

3. Chin, C.H., et al., A hub-attachment based method to detect functional modules from confidence-scored protein interactions and expression profiles. *BMC Bioinformatics*, 2010. 11: p. S25.

4. Gavin, A.C., et al., Proteome survey reveals modularity of the yeast cell machinery. *Nature*, 2006. 440: p. 631-636.

5. Jain, S. and G.D. Bader, An improved method for scoring protein-protein interactions using semantic similarity within the gene ontology. *BMC bioinformatics*, 2010. 11(1): p. 562.

6. Rašti, S. and C. Vogiatzis, A survey of computational methods in protein-protein interaction networks. *Annals of Operations Research*, 2019. 276(1-2): p. 35-87.

7. Zahiri, J., et al., Protein complex prediction: A survey. *Genomics*, 2019.

8. Liu, G., L. Wong, and H.N. Chua, Complex discovery from weighted PPI networks. *Bioinformatics*, 2009. 25(15): p. 1891-1897.

9. Wang, J., et al., Protein complex detection algorithm based on multiple topological characteristics in PPI networks. *Information Sciences*, 2019. 489: p. 78-92.

10. Lo, K., et al., Integrating external biological knowledge in the construction of regulatory networks from time-series expression data. *BMC systems biology*, 2012. 6(1): p. 101.