



دانشگاه صنعتی اصفهان

دانشکده مهندسی برق و کامپیو تر

پیدا کردن موظیف در شبکه های وزن دار بر هم کنش پروتئین - پروتئین

پایان نامه کارشناسی ارشد مهندسی کامپیو تر - هوش مصنوعی

مهسا ایمانی

استاد راهنما

دکتر نیلوفر قیصری

دکتر مهدی صادقی

بِنَامِ خَدا



دانشگاه صنعتی اصفهان

دانشکده برق و کامپیوتر

پیدا کردن موظیف در شبکه‌های وزن‌دار برهم‌کنش پروتئین-پروتئین

پایان‌نامه کارشناسی ارشد مهندسی کامپیوتر- هوش مصنوعی

مهرسا ایمانی

استاد راهنما

دکتر نیلوفر قیصری

دکتر مهدی صادقی

نمی دانم چطور با کلمات سپاس خدایی را بگویم که دست مهربان یاری دهندهش را در تمام مراحل زندگی ام احساس کردم و حتی قبل از آن که قدمی برای انجام کاری بردارم او صدقه برای کمک به من برمی داشت. خدا را سپاس می گویم سپاسی که بر سپاس های دیگر برتری داشته باشد مانند برتری که پروردگار نسبت به آفریدگان دارد.

از اساتید راهنمایم، سرکار خانم دکتر قیصری و جناب آقای دکتر صادقی که علاوه بر اینکه استاد راهنمای راهگشای من بودند از نظر اخلاقی هم استاد و الگوی من بودند کمال تشکر را دارم. از جناب آقای دکتر موسوی که افتخار شاگردی ایشان را داشتم و ایشان همیشه چه از نظر شیوه تدریس و چه شیوه بخورد با دانشجو و اخلاقی استاد و الگوی من بودند و سرکار خانم دکتر عمومی که داوری این پایان نامه را بر عهده گرفتند، به واسطه پیشنهادهای بسیار پربارشان برای بهبود کار تشکر می نمایم.

بر خود لازم می دانم از مادرم که اگر موقعيتی در زندگی ام کسب کردم آن را مديون فداکاری ها و دلسوزی هایشان هستم، خواهر و برادرانم برای همدلی ها و همراهی هایشان و پدرم که استقلال فکری و عملی را به من آموخت، خالصانه تشکر و قدردانی کنم.

از سرکار خانم نکویی و سرکار خانم مظاہری، پرسنل مهربان و خوش خلق دفتر تحصیلات تکمیلی دانشکده و جناب آقای دکتر سید محمود مدرس هاشمی رئیس تحصیلات تکمیلی دانشکده برق و کامپیوتر نیز بابت زحماتشان در طول دوره تشکر و قدردانی می کنم.

از تمامی دوستان عزیزم در آزمایشگاه هوش مصنوعی و همکلاسی های خوبیم که حضورشان همواره دلگرم کننده بود، صمیمانه تشکر می کنم. بر خود لازم می بینم که از عصمت پاکیزه عزیزم که همکاری اش در کارهای علمی مشترکمان باعث افتخار و نشاط بود و نیز برای یک دنیا مرام و دوستی، از مژگان قربانی و هدی نوری خواهران و مشاوران عزیزم که همیشه از من حمایت می کردند، از اعظم عموزادی نازین، دوست و خواهri که همیشه در کم می کرد ، از پروین رزاقی عزیز که لحظات زندگی ام در خوابگاه را شیرین کرد و از جیمبول برای لبخند زیبای الهام بخشش، خالصانه تشکر کنم.

کلیهی حقوق مادی مترتب بر نتایج مطالعات،
ابتكارات و نوآوری‌های ناشی از تحقیق موضوع
این پایان‌نامه (رساله) متعلق به دانشگاه صنعتی
اصفهان است.

تعدیم به بال و پر زندگیم،

مادر ناز نینم

که وجودش سهل عشق و دوست داشتن است،

به پاس همراهی های خالصانه، فدا کاری ها، دلواپسی ها و آرزو های بلند شان برای من.

فهرست مطالب

<u>صفحه</u>	<u>عنوان</u>
.....	فهرست مطالب
.....	هشت
.....	چکیده
.....	فصل اول: مقدمه
۱.....	۱-۱ معرفی
۲.....	۲-۱ پروتومیکس
۴.....	۳-۱ علم شبکه
۵.....	۱-۳-۱ شبکه‌های برهمکنش پروتئین-پروتئین
۵.....	۲-۳-۱ موتیف‌های شبکه
۶.....	۴-۱ تبیین مساله پایان نامه
۸.....	۵-۱ اهداف و انگیزه پایان نامه
۸.....	۶-۱ روش پیشنهادی
۹.....	۷-۱ دستاوردها
۱۰.....	۸-۱ روند ارائه‌ی مطالب
.....	فصل دوم: پیش زمینه زیست شناختی
۱۱.....	۱-۲ مقدمه
۱۱.....	۲-۲ پروتئین و کار کرد آن
۱۲.....	۳-۲ پروتئین و بیماری‌ها
۱۲.....	۱-۳-۲ پروتئینهای دفاعی از قبیل آنتی‌بادیهای سیستم ایمنی
۱۳.....	۴-۲ برهمکنش پروتئین-پروتئین (ب.پ.پ)
۱۳.....	۱-۴-۲ انواع برهمکنش‌ها
۱۴.....	۵-۲ شبکه برهمکنش پروتئین-پروتئین
۱۶.....	۶-۲ آزمایش‌های زیستی برای کشف برهمکنش پروتئین-پروتئین
۱۶.....	۷-۲ در دسترس بدون داده‌های برهمکنش پروتئین-پروتئین
۱۷.....	۸-۲ نتیجه‌گیری
.....	فصل سوم: یکپارچه سازی باتک‌های اطلاعاتی
۱۸.....	۱-۳ مقدمه
۱۸.....	۲-۳ بررسی اهمیت یکپارچه سازی داده‌ها در حوزه‌ی داده‌های زیست شناختی
۲۰.....	۳-۳ انگیزه‌های یکپارچه سازی پایگاه‌های زیست‌شناختی
۲۱.....	۴-۳ مزایای ترکیب اطلاعات منابع اطلاعاتی
۲۱.....	۱-۴-۳ افرونگی
۲۲.....	۲-۴-۳ تکمیل شدن
۲۲.....	۳-۴-۳ کاهش زمان
۲۲.....	۴-۴-۳ کاهش هزینه

۵-۳	پایگاه داده STRING	۲۳
۶-۳	نتیجه‌گیری	۲۳
فصل چهارم: بررسی کارهای مرتبط		
۱-۴	مقدمه	۲۵
۲-۴	اصطلاحات مورد نیاز گراف/ شبکه	۲۵
۱-۲-۴	زیرگراف القابی و زیرگراف غیر القابی	۲۶
۱-۲-۴	یکریختی گراف	۲۶
۲-۲-۴	گراف تصادفی	۲۶
۳-۴	آمار و خصوصیات شبکه	۲۷
۱-۳-۴	شبکه‌های بی مقیاس	۲۷
۴-۴	موتیف‌های شبکه	۲۹
۱-۴-۴	تعریف موتیف	۳۰
۲-۴-۴	فراوانی موتیف‌ها	۳۰
۵-۴	موتیف‌ها در شبکه‌های بدون وزن	۳۱
۱-۵-۴	تصادفی کردن شبکه	۳۲
۲-۵-۴	روش‌های فعلی پیدا کردن موتیف بدون وزن	۳۲
۳-۵-۴	موانع فعلی	۳۴
۴-۵-۴	ابزارهای شناخته شده برای پیدا کردن موتیف	۳۵
۵-۵-۴	مقایسه ابزارهای شناخته شده برای پیدا کردن موتیف	۴۰
۶-۴	موتیف‌ها در شبکه‌های وزن دار	۴۱
۱-۶-۴	مقدمه	۴۱
۲-۶-۴	روش‌های فعلی پیدا کردن موتیف در شبکه‌های وزن دار	۴۱
۷-۴	نتیجه‌گیری	۴۵
فصل پنجم: روش پیشنهادی		
۱-۵	مقدمه	۴۶
۲-۵	الگوریتم پیشنهادی	۴۹
۱-۲-۵	یافتن همهٔ زیرگراف‌ها	۴۹
۲-۲-۵	گروه بندی بردارهای ویژگی بدست آمده از زیرگراف‌ها	۵۳
۳-۲-۵	تولید شبکه‌های تصادفی	۵۶
۴-۲-۵	تعیین موتیف بودن	۵۶
۳-۵	نتیجه‌گیری	۵۶
فصل ششم: نتایج تجربی		
۱-۶	مقدمه	۵۸
۱-۱-۶	نتایج بر روی شبکه ورودی تصادفی	۵۹
۲-۱-۶	نتایج بر روی شبکه ورودی واقعی	۶۰
۲-۶	معایب و برتری‌های الگوریتم پیشنهادی	۶۳
۳-۶	بحث و نتیجه‌گیری	۶۴
فصل هفتم: جمع‌بندی و ارائه پیشنهادات		

۷۱ پیشنهادات	۱-۷
۷۲ مراجع	

چکیده

شبکه‌های زیست‌شناسی، که عموماً شبکه‌های پیچیده و وسیع هستند، حاوی اطلاعات مهمی می‌باشند. امروزه نظر به پیشرفت قابل ملاحظه آنها، عملیاتی بر روی این گونه شبکه‌ها برای استخراج اطلاعاتی که هر کدام از آنها دربردارند اعمال می‌شود. از میان مهم‌ترین آنها پیدا کردن موتیف‌های شبکه است. موتیف به صورت الگوهایی از ارتباطات تعریف می‌شود که در شبکه با فراوانی بیشتری نسبت به حالت تصادفی مشاهده می‌شود. تحقیقات انجام شده نشان می‌دهند که چنین ساختارهایی می‌توانند از لحاظ فعالیت در شبکه دارای اهمیت باشند. تاکنون، الگوریتم‌های مختلفی برای پیدا کردن زیرگراف‌های با فراوانی بالا در شبکه‌های بدون وزن ارائه شده که دو محدودیت اساسی دارند. اول آنکه تقریباً تمام این روش‌ها فقط قادر به پیدا کردن موتیف در شبکه‌های بدون وزن هستند. با توجه به رشد شبکه‌ها در طول دهه اخیر، بررسی خصوصیات دیگری نظیر ناهمگونی یال‌ها که ورای خصوصیات توبولوژیکی آنها است، نیز حائز اهمیت می‌باشد. بررسی شبکه‌های وزن دار به جای شبکه‌های بی وزن اجازه می‌دهد ناهمگونی یال‌ها در شبکه‌ها نیز مورد بررسی قرار گیرد. یکی از دلایل اصلی فقدان چنین مطالعاتی آن است که در شمارش زیرگراف‌ها فقط زیرگراف‌های القایی در نظر گرفته شده است. محدودیت دیگر روش‌های موجود آن است که در نظر گرفتن وزن یال‌ها چالش محاسباتی جدیدی را ایجاد می‌کند. در نظر گرفتن زیرگراف‌های غیر القایی یک موتیف شبکه در شبکه‌های برهم‌کنش پروتئین، مسئله‌ای چالش برانگیز و کاملاً مطلوب است زیرا این شبکه‌ها تا کامل شدن و عاری از خطای بودن فاصله بسیار دارند. برهم‌کنش‌های گزارش شده توسط این شبکه‌ها شامل برهمکنش‌هایی می‌شوند که به طور اشتباه تشخیص داده شدند و نیز بسیاری از برهمکنش‌ها تشخیص داده نشده‌اند. بنابراین وقوع یک موتیف شبکه خاص در یک شبکه ممکن است شامل یال‌های اضافه در وقوعش در شبکه دیگری باشد و برعکس. یکی از دلایلی که تحلیل الگوهای گراف‌ها را دشوارتر می‌کند این است که آنها ذاتاً به شکل بردار نیستند و به سادگی نیز قابل تبدیل به بردار نیستند. در این پایان نامه روش جدیدی برای حل دو مشکل مطرح شده در روش‌های قبل پیشنهاد داده شده است که برخلاف روش‌های قبلی، که هنگام نسبت دادن زیرگراف‌های یافت شده به گروه‌های مشابه به بررسی یک ریختی زیرگراف‌ها (انطباق دقیق) می‌پردازند، از تطبیق غیر دقیق گراف استفاده می‌کند. تطبیق غیر دقیق زیرگراف‌ها در شبکه‌های برهمکنش موجود که دارای میزان زیادی خطای هستند کاملاً مطلوب است. روش پیشنهادی موتیف‌هایی وزن دار را معرفی می‌کند که وزن هر یال نشان دهنده امید حضور آن یال در موتیف اجماع حاصل است. برای این منظور با استفاده از چندجمله‌ای‌های متقارن ویژگی‌های طیفی زیرگراف‌ها بدست آورده می‌شود و سپس به انطباق و تحلیل بردارهای ویژگی حاصل با استفاده از خوش بندی به جای طبقه بندی زیرگراف‌ها می‌پردازد. نتایج تجربی نشان می‌دهد دخیل کردن وزن‌ها می‌تواند سبب پیدا کردن موتیف و ضدموتیف‌هایی متفاوت از آنچه توسط روش‌های موجود بدست می‌آید، شود. این طور به نظر می‌رسد که ذات ارتباط بین وزن‌ها کلی نیست و وابسته به عملکرد شبکه است.

کلمات کلیدی:، شبکه‌های وزن دار برهم‌کنش پروتئین، موتیف وزن دار، انطباق غیر دقیق گراف، زیرگراف غیر القایی، ویژگی‌های طیفی گراف، خوش بندی

فصل اول

فصل اول: مقدمه

۱ + معرفی

پیدا کردن موتیف های شبکه، یک از مسائل مطرح در علم شبکه است و تاکنون روش های مختلفی برای آن ارائه شده است. در این پایان نامه قصد داریم به پیدا کردن موتیف های شبکه های وزن دار برهمنکش پروتئین-پروتئین پردازیم. این شبکه ها یکی از موضوعات اصلی در پروتئومیکس و زیست سیستم هستند. بدلیل روش های تولید این برهمنکش ها این شبکه ها دارای خطای بالا هستند و تا کامل شدن فاصله بسیار دارند. این امر پیدا کردن موتیف های شبکه را بسیار پیچیده تر می کند.

۲ پروتئومیکس

توجه به جزئیات به تنها یی نمی تواند پاسخگوی سوالات بشر و جهان پیچیده باشد همچنان که باید جزئیات را دید باید نگاهی نیز به تصویر کلی حیات داشت در چند سال گذشته دانشمندان متوجه شده اند برای مشاهده تصویری بزرگ از هستی باید علاوه بر کشف ذرات دیدی کلی به جهان داشت. به دانش نگرش سیستمیک به حیات و پدیده های مرتبط زیستی، زیست سیستم گفته می شود.

یکی از کاربردهای اصلی این رشته که امروزه هزینه ها و تمرکز زیادی را به خود اختصاص داده است، استفاده از سیستم بیولوژی در مدل سازی اثر دارو ها در سلول های هدف و مسیر های بیوشیمیایی در گیر است.

دانشمندان امیدوارند که در آینده نزدیک با مدل سازی های دقیق تر بتوانند دارو های اختصاصی تر با اثرات جانبی کمتر و بهینه را طراحی کنند.

موضوعات اصلی در زیست سیستم شامل آنالیز شبکه های مختلف است که مهمترین این شبکه ها عبارتند از [۲]:

- شبکه های ژنی
- شبکه های برهمکنش پروتئین-پروتئین (ب.پ.پ)^۱
- شبکه های متابولیسمی
- شبکه های سیگنالی

برای مدل سازی هر کدام از این شبکه ها با توجه به ویژگی هایی که هر کدام از آنها دارند از روش ها و الگوریتم های خاصی استفاده می شود ولی آن چه مهم است این است که در نهایت هدف زیست سیستم رسیدن به یک مدل واحد و کشف ارتباطات بین این شبکه ها و نقش آنها در حیات سلول ، ارتباط بین سلول ها و شرایط محیطی و در نهایت پاسخ موجود به تغییرات محیطی است.

ژن ها با عملکرد پروتئین هایی که می سازند نقش خود را اعمال می کنند. بنابراین لازمه شناخت کامل ژن ها ، شناخت کامل پروتئین ها است. در واقع هم اکنون مطالعه پروتئین ها با توجه به اهمیت آنها و پیشرفت های سریع ابزاری و تکنیکی در این رشته، جایگاه خاصی یافته است. با تکمیل طرح ژنوم انسان مشخص شد که مکانیسم ملکولی رفتار سلول ها در شرایط مختلف را نمی توان از روی توالی ژن های آنها پیشگویی کرد. رفتار سلولی و تمام فعالیت هایی که در سلول انجام می شود به عهده پروتئین های بیان شده به وسیله ژن است. در واقع برای ارتباط ژنوم با رفتار سلول ها باید پروتئین های سلول ها را شناخت [۱ و ۲].

بنابراین برای شناسایی مکانیسم های مولکولی رفتار سلولی و واکنش های زیستی، لازم است پروتئین هایی که در یک سلول بیان می شود، تغییرات آنها در شرایط مختلف، عملکرد آنها و همچنین برهمکنش های بین پروتئین های مختلف در یک سلول، بررسی شود. به مجموعه این بررسیها، نقشه برداری پروتئوم یا پروتئومیک، گفته می شود. امروزه پروتئومیک و ژنومیک دو فعالیت برجسته ای هستند که بیوانفورماتیک دانان به آن مشغول هستند [۱ و ۲].

ژنومیک شامل تجزیه و تحلیل داده ها و اطلاعات ژنتیکی ژنوم موجودات است. پروتئومیک به آنالیز پروتئین های یک موجود زنده گفته می شود. سنجش مقایسه ای پروتئین ها به صورت جامع در یک مقیاس وسیع موضوع علم پروتئوم یا پروتئومیکس می باشد. از دیدگاه بیوشیمیست ها پروتئومیکس، مطالعه بیش از یک پروتئین

^۱ Protein-protein interaction (PPI)

در یک زمان خاص است. اولین مطالعات بر روی پروتئین در سال ۱۳۵۴ با ابداع ژل دو بعدی آغاز شد. علم پروتومیکس بسیار گسترده است. شناخت پروتئین‌ها، بررسی کمی آنها در سلول‌ها، بافت‌ها و مایعات زیست شناختی، سنجش تغییرات در بیان پروتئین‌ها در سلول‌های بیمار در مقابل سلول‌های طبیعی، توصیف تغییرات پس از ترجمه، مطالعه برهمکنش‌های پروتئین-پروتئین، تعیین موقعیت، شناسایی عملکرد سلولی در سطح پروتئین‌ها، شناسایی ژن‌های ناشناخته به کمک پروتئین‌ها، بسیاری از کاربردها و جوانب دیگر آن می‌باشد. از مهمترین اهداف تحقیقات پروتومیکس نیز شرح و توصیف مکانیسم‌های مولکولی دخیل در فرایندهای سلولی، ویژگی شبکه‌های پیچیده پروتئینی و اختلال در آنها، کشف بیومارکرهای پروتئینی برای آشکار سازی و تشخیص بیماری‌ها و شناخت اهدافی برای طراحی درمانهای داروئی می‌باشد. در مقابل ژنوم که نسبتاً ثابت و بی تغییر است، پروتومیکس دینامیک و متغیر است. بنابراین، می‌توان پروتومیکس را علم "پس از ژنوم" نامید، چرا که پروتومیکس نقش قابل توجهی در عبور از ژنومیکس به کاربردهای سودمند بالینی به ویژه در عرصه‌های تشخیص و پیشگیری دارد [۲۱].

به دلیل اهمیت پروتومیک، دهه آینده در بین بیولوژیستهای مولکولی به دهه پروتومیک مشهور شده است. حجم مقالاتی که در چند سال اخیر در مورد پروتومیک انتشار یافته است در این سال‌ها با شتاب بسیار زیاد، در حال افزایش است.

۱ ۳ علم شبکه

علم شبکه، مطالعه‌ی سیستم‌های پیچیده (تشکیل شده از بخش‌های در حال تعامل) از دید تئوری گراف است. ریشه‌ی علم شبکه به مطالعه‌ی شبکه‌های اجتماعی در دهه‌ی ۱۳۳۰ باز می‌گردد. علم شبکه بینش وسیعی در مورد چگونگی عملکرد اجزای مختلف جهان ما، از سلول تا اینترنت و اقتصاد فراهم کرده است. در دهه اخیر به دلیل رشد روز افزون میزان دسترسی به داده‌های شبکه نظیر شبکه جهانی اینترنت²، شبکه‌های برهمکنش پروتئین، شبکه‌های غذا و نیز میزان دسترسی به کامپیوترهای شخصی قدرتمند برای پردازش داده، علم شبکه مورد توجه زیادی قرار گرفته است. امروزه علم شبکه تقریباً مورد توجه محققان تمامی رشته‌های علمی و صنعتی است.

رفتار سلولی و تمام فعالیت‌هایی که در سلول انجام می‌شود بر عهده پروتئین‌ها است. همه پروتئین‌ها با هم برهمنکش دارند و تقریباً می‌توان گفت که همه پروتئین‌ها اثر خود را با همکاری پروتئین‌های دیگر در سلول اعمال می‌کنند و هیچ پروتئینی نیست که در یاخته به تنها یی عمل کند. اینکه کدام پروتئین در کدام زمان، و در کدام

² WWW

سلول ایجاد شود، رویکردی است که براساس مکانیسم های تنظیم سلولی و مبتنی بر مقتضیات و نیازهای (کوتاه مدت، میان مدت و درازمدت) سلول شکل می گیرد. این رخدادها، به نوبه خود، نیز ویژگی های متفاوت سلول های یک موجود منفرد را به روشنی تبیین می کند. لذا برای بررسی مکانیسم عمل پروتئین ها در سلول از شبکه های برهmekش پروتئین-پروتئین که نحوه‌ی تعامل پروتئین ها را نشان می دهد، استفاده می شود.

۱ ۴ + شبکه های برهmekش پروتئین-پروتئین

فرآیندهای سلولی به وسیله تعامل پروتئین ها با یکدیگر، با محیط و با DNA تعریف می شوند. فهمیدن این-که چطور پروتئین ها با یکدیگر در سطح شیمیایی تعامل می کنند، گام اولیه برای مدل کردن فرآیندهای سلولی است. فهم این موضوع می تواند مفاهیم مهمی برای درک سیستم های بزرگتر آشکار شده به وسیله شبکه های برهmekش پروتئین-پروتئین (ب.پ.پ) و شبکه های تنظیمی ژن^۳ را بیان کند. در شبکه های برهmekش پروتئین، گره ها پروتئین ها و یال های غیر جهت دار، برهم کنش های خروجی پروتئین-پروتئین را نشان می دهد. در شبکه های تنظیمی ژن، گره ها ژن ها و پروتئین ها، و یال های جهت دار، تولید یک پروتئین به وسیله ی یک ژن یا تنظیم یک ژن به وسیله ی یک پروتئین را نشان می دهد.

در چند سال گذشته شبکه های ب.پ.پ، کاربردهایی نظری پیش بینی برهم کنش های دامنه-دامنه^۴ و پیش-بینی برهم کنش های پروتئین-پیتید^۵ را براساس موتیف های شبکه یافته اند. از وابسته سازی یک شبکه ب.پ.پ با یک شبکه پروتئین-DNA می توان برای تفسیر ژن هایی که قبلا طبقه بندی نشده اند، استفاده کرد.

۱ ۴ + موتیف های شبکه

پیمانه ای بودن^۶ روشهای استاندارد در طراحی سیستم ها و مهندسی می باشد. ساختار پیمانه ای امکان استفاده مجدد از زیربخش های متداول را فراهم می کند. مهندسان به منظور کنترل و مدیریت پیچیدگی سیستم های بزرگتر از سازمان دهی سلسله مراتبی استفاده می کنند. علاوه بر این در علم شبکه، این ویژگی در شبکه هایی که به طور طبیعی در حال رشد و تکامل هستند، نیز شایع است. مطالعه ی چنین زیر شبکه هایی که به طور طبیعی رخداده اند، منجر به ایجاد بینشی در مورد نقش مجموعه ای از گره های یک شبکه در پردازش اطلاعات شده است [۳].

³ Gene Regularly Networks

⁴ Domain-Domain interactions

⁵ De Novo Protein-Peptide Interactions

⁶ Modularity

موتیف‌های شبکه برای در ک پیمانه ای بودن ساختار کلی شبکه‌ها مفید هستند. موتیف‌ها برای اولین بار توسط میلو و همکارانش معرفی شدند [۳] و اهمیت آنها به عنوان واحد پردازش اطلاعات از نظر تئوری مدل شده است [۴ و ۵]. همچنین موتیف‌ها به صورت تجربی مورد آزمایش و بررسی قرار گرفته‌اند [۹-۶].

موتیف‌های شبکه به صورت زیر شبکه‌هایی که مهم و غیرتصادفی هستند، تعریف می‌شوند [۳]. موتیف شبکه به طور معمول از تکرار زیاد آن در شبکه واقعی در مقایسه با نسخه‌های تصادفی همان شبکه، تعیین می‌شوند [۴ و ۳]. به همین ترتیب، ضد موتیف از تکرار کم آن در شبکه اصلی در مقایسه با نسخه‌های تصادفی همان شبکه تعیین می‌شوند.

مطالعات اخیر در زمینه هم‌تراز کردن^۷ شبکه‌های زیست‌شناسی، زیر شبکه‌های حفظ شده را آشکار کرده است [۱۰].

پیدا کردن همه‌ی نمونه‌های یک گراف معین به عنوان زیر گراف یک شبکه از جهات زیر لازم است:

- تعیین این نکته که آیا یک گراف معین (شاید به صورت تجربی تعیین شده باشد) یک موتیف معنادار هست.
- بررسی خوشبندی موتیف به منظور تعیین این که چطور موتیف‌ها می‌توانند بخش‌هایی از ساختارهایی بزرگ‌تر باشند و تعیین این نکته که معناداری موتیف چقدر به دقت شبکه وابسته است.
- پیدا کردن موتیف‌های مبتنی بر فراوانی با پیدا کردن موتیف‌های مبتنی بر بقا به روش‌های خاص (هم‌تراز کردن شبکه) ترکیب شود.

۱۴ تبیین مساله پایان نامه

امروزه شبکه‌های بزرگ، مانند شبکه‌های زیست‌شناسی که هر کدام شامل هزاران گره می‌باشند، بسیار مورد توجه محققان علم پرتونومیکس قرار گرفته‌اند. یکی از تلاش‌هایی که در این زمینه مورد پی‌گیری قرار گرفته است، پیدا کردن موتیف در این شبکه‌ها از جمله شبکه‌های برهم‌کنش پروتئین-پروتئین است. موتیف، به صورت الگوهایی از زیر گراف‌ها تعریف می‌شود که در شبکه با فراوانی بیشتری نسبت به حالت تصادفی مشاهده می‌شود. به نظر می‌رسد این چنین ساختارهایی از لحاظ فعالیت در شبکه دارای اهمیت باشند. الگوریتم‌های کنونی برای پیدا کردن موتیف‌های با اندازه‌ی بزرگ با مشکل مواجه شده‌اند. از طرفی داده‌های موجود برای شبکه‌های برهم‌کنش پروتئین-پروتئین خالی از خطای نیستند و هیچ یک از روش‌های پیش‌بینی و استخراج برهم‌کنش‌ها در شبکه، به

⁷ Graph Alignment

شبکه واقعی همگرا نمی‌شود. شبکه‌های گزارش شده توسط پایگاه داده‌های مختلف برای یک ارگانیسم دارای هم پوشانی بسیار کمی هستند. بنابراین این شبکه‌ها خالی از خطا نیستند و نیاز به منابع اطلاعاتی دیگر برای یکپارچه سازی داده^۸ و وزن‌دار کردن شبکه‌های برهمکنش پروتئین وجود دارد. به همین دلیل امروزه علاقه به استفاده از شبکه‌های برهم‌کنش وزن‌دار پروتئین به جای شبکه‌های بی‌وزن رو به افزایش است. در این شبکه‌ها وزن هر یال نشان‌دهنده میزان اعتبار آن است.

تاکنون، الگوریتم‌های مختلفی برای پیدا کردن زیر‌گراف‌های با فراوانی بالا در مجموعه شبکه‌ها یا در یک گراف ارائه شده است. این الگوریتم‌ها یا تنها به شمارش انواع خاصی از زیر‌گراف‌ها می‌پردازند و یا شرطی بر روی اندازهٔ زیر‌گراف دارند. در تکنیک‌های فعلی بدلیل محدودیت اندازه، موتیف‌های با حداکثر ۸ یا ۹ گره را پیدا می‌کنند. می‌توان این تکنیک‌ها را دو دسته تقسیم کرد، دسته اول به شمارش کامل زیر‌گراف‌ها پرداخته و دسته دوم از نمونه‌گیری زیر‌گراف استفاده می‌کند.

الگوریتم‌های موجود دو محدودیت اساسی دارند. اول آنکه تقریباً تمام این روش‌ها فقط قادر به پیدا کردن موتیف در شبکه‌های بدون وزن هستند. با توجه به رشد شبکه‌ها در طول دهه گذشته، بررسی خصوصیات دیگری نظری ناهمگونی یال‌ها که ورای خصوصیات توپولوژیکی آنها است، نیز حائز اهمیت می‌باشد. بررسی شبکه‌های وزن‌دار به جای شبکه‌های بی‌وزن اجازه می‌دهد ناهمگونی یال‌ها در شبکه‌ها نیز مورد بررسی قرار گیرد. یکی از دلایل اصلی فقدان چنین مطالعاتی بر روی شبکه‌های وزن‌دار آن است که در نظر گرفتن وزن یال‌ها چالش محاسباتی جدیدی را ایجاد می‌کند.

محدودیت دیگر اکثر روش‌های موجود آن است که در شمارش زیر‌گراف‌ها فقط زیر‌گراف‌های القابی در نظر گرفته شده است. در روش‌های پیشنهاد شده جدیدتر برای شبکه‌هایی با خطای بالا در نظر گرفتن زیر‌گراف‌های غیر القابی علاوه بر زیر‌گراف‌های القابی پیشنهاد شده است. در نظر گرفتن زیر‌گراف‌های غیر القابی یک موتیف شبکه در شبکه‌های برهمکنش پروتئین، مسئله‌ای چالش برانگیز و کاملاً مطلوب است زیرا این شبکه‌ها تا کامل شدن و عاری از خطا بودن فاصله بسیار دارند. برهمکنش‌های گزارش شده توسط این شبکه‌ها شامل برهمکنش‌هایی می‌شوند که به طور اشتباه تشخیص داده شدند و نیز بسیاری از برهمکنش‌ها تشخیص داده نشده اند. بنابراین وقوع یک موتیف شبکه خاص در یک شبکه ممکن است شامل یال‌های اضافه در وقوعش در شبکه دیگری باشد و بر عکس. همچنین یکی از دلایلی که تحلیل الگوهای گراف‌ها را دشوارتر می‌کند این است که آنها ذاتاً به شکل بردار نیستند و به سادگی نیز قابل تبدیل به بردار نیستند.

⁸ Data integration

۱۵ اهداف و انتیزه پایان نامه

هدف اصلی این تحقیق وارد کردن وزن یال‌ها در مطالعه‌ی موتیف‌های شبکه‌های برهمکنش پروتئین است زیرا این شبکه‌ها هنوز با شبکه‌ای عاری از خط‌فاصله‌ی بسیار دارند. در این شبکه‌ها وزن هر یال نشان‌دهنده میزان اعتبار آن است. بدین ترتیب ناهمگونی یال‌ها در شبکه‌ها نیز مورد بررسی قرار گیرد. از آنجا که شبکه‌ها در طول دهه گذشته رشد کرده‌اند، باید به بررسی خصوصیات دیگری و رای خصوصیات توپولوژیکی آن‌ها پرداخت. یکی از این مطالعات بررسی شبکه‌های وزن دار به جای شبکه‌های بی‌وزن می‌باشد که اجزا می‌دهد ناهمگونی یال‌ها در شبکه‌ها نیز مورد بررسی قرار گیرد. اگر یک ارتباط کلی بین وزن‌ها و توپولوژی باشد، در نظر گرفتن وزن‌ها در بررسی سیستم‌های شبکه‌های اطلاعات جدیدی فراهم نمی‌کند. در دسترس بودن شبکه‌های وزن دار ما را قادر می‌سازد به بررسی ارتباط بین توپولوژی برهمکنش‌ها با قدرت این برهمکنش‌ها به پردازیم. این‌طور به نظر می‌رسد که ذات این ارتباط کلی نیست و مرتبط با عملکرد شبکه است. در نتیجه دخیل کردن وزن‌ها در بررسی موتیف‌ها می‌تواند در کمک ما را از سیستم‌های شبکه‌ای بهبود دهد.

از سوی دیگر فقط تعداد محدودی مطالعه در زمینه شبکه‌های زیستی وجود دارد که زیرگراف‌های غیر القایی را در پیدا کردن موتیف‌ها در نظر گفته است. دلیل آن این است که تعداد بسیار بیشتری زیرگراف غیر القایی یک ریخت با یک توپولوژی خاص وجود دارد و بنابراین شمارش زیرگراف‌های غیر القایی شبکه بسیار دشوار‌تر است. در نظر گرفتن زیرگراف‌های غیر القایی یک موتیف شبکه نه تنها چالش برانگیز است که کاملاً براساس شبکه‌های برهمکنش فعلی تا کامل شدن و عاری از خط‌فاصله بسیار دارند کاملاً مطلوب است.

۱۶ روش پیشنهادی

یکی از موانع برسر پیدا کردن موتیف‌های شبکه، مسئله ایزومورفیسم گراف است که از نظر محاسباتی پیچیده و در طبقه مسائل NP-complete قرار می‌گیرد. مسئله تطبیق گراف به دو دسته اصلی تطبیق دقیق گراف^۹ و تطبیق غیر دقیق گراف^{۱۰} تقسیم می‌شود. مسئله تطبیق دو گراف به صورت دقیق همان مسئله ایزومورفیسم می‌باشد که برای بدست آوردن زیرگراف‌های القایی موتیف‌ها از آن استفاده می‌شود. به عبارت دیگر در انطباق دو زیرگراف تعداد گره‌ها و یال‌ها و همچنین طریقه اتصال گره‌ها به یکدیگر باید یکسان باشد. استفاده از انطباق

⁹ exact graph matching

¹⁰ inexact graph matching

دقیق برای این شبکه‌های وزن دار بدلیل ناهمگونی بالای یال‌ها امکان‌پذیر نیست. در تطبیق غیر دقیق دو زیر‌گراف، تعداد گره‌ها و یال‌ها و همچنین طریقه اتصال گره‌ها به یکدیگر می‌تواند متفاوت باشد و هدف یافتن تطبیقی است که این اختلاف در آن حداقل باشد. تطبیق غیر دقیق گراف در طبقه مسائل NP-complete قرار می‌گیرد. ازین رو الگوریتم‌های زیادی برای ارائه تقریبی از راه حل مطلوب پیشنهاد شده اند [۱۱].

برای رفع دو محدودیت موجود در شمارش موتیف‌های غیر القایی در شبکه‌های وزن دار در این پایان نامه قصد داریم برخلاف روش‌های قبلی، به جای استفاده از تطبیق دقیق (ایزومورفیسم) از تطبیق غیر دقیق گراف استفاده کنیم. به عبارت دیگر در انطباق دو زیر‌گراف تعداد گره‌ها و یال‌ها و همچنین طریقه اتصال گره‌ها به یکدیگر می‌تواند متفاوت باشد. برای این منظور با تبدیل زیر‌گراف‌ها به بردارهای ویژگی با استفاده از ویژگی‌های های طیفی گراف به انطباق و تحلیل آنها می‌پردازیم. تطبیق غیر دقیق زیر‌گراف‌ها براساس شبکه‌های برهمکنش موجود که دارای میزان زیادی خطأ هستند کاملاً مطلوب است. در الگوریتم پیشنهادی مفهوم موتیف بسط داده می‌شود تا وزن یال‌ها که نشان دهنده‌ی قدرت برهمکنش (در مورد شبکه‌های برهمکنش پروتئین‌ها نشان دهنده‌ی قابلیت اعتماد برهمکنش‌ها) میان دو عنصر سازنده است در نظر گرفته شود.

یکی از دلایلی که تحلیل الگوهای گراف‌ها را دشوارتر می‌کند، آن است که آنها نه در ذات به شکل بردار هستند و نه به سادگی به بردار تبدیل می‌شوند. از این‌رو الگوریتم پیشنهادی برای تطبیق گراف‌ها، یک بردار از ویژگی‌های طیفی آنها می‌سازد. ویژگی‌های طیفی که ساخته می‌شود با تغییر جایگاه گره‌ها در ماتریس مجاورت تغییر نمی‌کنند و از تمام ماتریس‌طیفی بهره می‌برد. روش مورد استفاده از مقادیر چندجمله‌ای‌های متقارن^{۱۱} برای ساخت بردار ویژگی گراف استفاده می‌کند. سپس به خوشه‌بندی بردارهای ویژگی بدست آمده با استفاده از روش خوشه‌بندی جدید K-means++ می‌پردازد. سپس برای تعیین معناداری موتیف‌های اجماع بدست آمده به جای استفاده از تعریف Z-score ساده که در الگوریتم‌های پیشین استفاده می‌شد از تعمیم Z-score که براساس وزن یال‌های زیر‌گراف است استفاده می‌شود.

۱.۴ دستاوردها

روش پیشنهادی، اولین روش جامع در زمینه مطالعه‌ی موتیف‌های شبکه‌های وزن دار است که محدودیتی بر روی نوع زیر‌گراف‌های مورد بررسی ندارد و همزمان قادر به حل هر دو مشکل موجود در زمینه پیدا کردن

¹¹ symmetric polynomials

موتیف‌ها، یعنی دخیل کردن وزن یالها و درنظر گرفتن زیرگراف‌های غیرالقایی که در شبکه‌ای با خطای بالابسیار مورد نیاز است، می‌باشد.

این روش برای اولین بار از مفهوم انطباق غیر دقیق در انطباق زیرگراف‌ها به جای انطباق دقیق استفاده می‌کند. روش پیشنهاد شده برای انطباق غیر دقیق قابلیت اعمال به بقیه مسائل مطرح در زمینه شبکه‌های وزن‌دار مانند یافتن ترکیبات پروتئینی را نیز دارد.

در روش‌های قبلی، موتیف یافت شده یا وجود دارند یا ندارند. اما روش پیشنهادی مفهوم موتیف را از حالت وجود داشتن یا نداشتن به طیفی از وجود داشتن بسط می‌دهد. در موتیف ارائه شده وزن هر یال نشان‌دهنده‌ی امید حضور آن یال در موتیف است.

۱۸ روند ارائه‌ی مطالب

از آنجاکه مطالعه برهمکنش‌های پروتئین-پروتئین یکی از موضوعات مهم در زیست‌سیستم و علم پروتئومیکس است، در فصل دوم مفاهیم و روش‌های کنونی گردآوری این برهمکنش‌ها بیان می‌شود. در فصل سوم، بدلیل آنکه بخش بزرگی از داده‌های کنونی برهمکنش‌های پروتئین-پروتئین، دارای خطا می‌باشند و در این پایان نامه از شبکه‌های یکپارچه شده وزن‌دار استفاده می‌شود، مفهوم یکپارچه سازی بانک‌های اطلاعاتی و اهمیت آن در حوزه‌ی داده‌های زیست‌شناختی بیان می‌شود. سپس در فصل چهارم شبکه‌ها، مفهوم و اهمیت موتیف‌ها در شبکه‌ها و کارهای مرتبط در زمینه یافتن موتیف‌ها در شبکه‌ها از جمله در شبکه‌های برهمکنش‌های پروتئین-پروتئین بیان می‌شود. در فصل پنجم الگوریتم پیشنهادی و در فصل ششم نتایج تجربی آورده می‌شود. سپس از مطالب گفته شده در فصل آخر نتیجه گیری و پیشنهاداتی ارائه می‌شود.

فصل دوم

پیش زمینه زیست شناختی

در این فصل اطلاعات پایه در مورد پروتئین، برهم کنش‌های پروتئین-پروتئین و شبکه‌های برهم کنش پروتئین ارائه می‌شود. برهم کنش پروتئین-پروتئین به ارتباط مولکول‌های پروتئین با یکدیگر برمی‌گردد. این ارتباطات از چند جهت مانند نواحی تحقیقاتی کلی مثل بیوشیمی، بیوفیزیک، فرآیندهای زیست‌شناختی خاص، مسیرهایی از قبیل مسیر انتقال سیگنال^{۱۲} و مطالعات در سطح سیستم شبکه‌های ارگانیزم حائز اهمیت است.

۴ پروتئین و کارکرد آن

پروتئین‌ها مواد آلی بزرگ و یکی از انواع درشت‌ملکول‌های زیستی هستند که از زیرواحدهایی به نام اسید آمینه ساخته شده‌اند. اسیدهای آمینه مثل یک زنجیر خطی توسط پیوند پیتیدی میان گروه‌های کربوکسیل و آمین مجاور به یکدیگر متصل می‌شوند تا یک پلی پیتید را به وجود بیاورند. ترتیب اسیدهای آمینه در یک پروتئین توسط ژن مشخص می‌شود. به جز در بعضی از ارگانیسم‌ها، کد ژنتیک ۲۰ تا اسید آمینه استاندارد را معرفی می‌کند، پروتئین‌ها معمولاً به یکدیگر می‌پیوندند تا یک وظیفه‌ای را با یکدیگر انجام دهند که این خود باعث استوار شدن پروتئین می‌شود. چون ترتیب‌های نامحدودی در توالی و طول زنجیره اسید آمینه‌ها در تولید پروتئین‌ها وجود دارد، از این رو انواع بی‌شماری از پروتئین‌ها نیز می‌توانند وجود داشته باشند.

¹² Signal Transduction Pathway