



دانشکده مهندسی کامپیوتر

پیاده‌سازی و کاربست یک الگوریتم جاسازی گراف برای جاسازی شبکه‌ی اندرکنش ناهمگن دارو- پروتئینی برای پیش‌بینی دارو- هدف

پروژه کارشناسی مهندسی کامپیوتر گرایش مهندسی نرم افزار

فاطمه فتاحی

استاد راهنما

دکتر مینایی

شهریور ۱۳۹۸



تأییدیه‌ی هیأت داوران جلسه‌ی دفاع از پروژه

نام دانشکده: دانشکده مهندسی کامپیوتر

نام دانشجو: فاطمه فتاحی

عنوان پروژه: پیاده‌سازی و کاربست یک الگوریتم جاسازی گراف برای جاسازی شبکه‌ی اندرکنش ناهمگن

دارو- پروتئینی برای پیش‌بینی دارو- هدف

تاریخ دفاع: شهریور ۱۳۹۸

رشته: مهندسی کامپیوتر

گرایش: مهندسی نرم افزار

ردیف	سمت	نام و نام خانوادگی	مرتبه دانشگاهی	دانشگاه یا مؤسسه	امضا
۱	استاد راهنما				
۲	استاد مشاور				
۳	استاد مدعو خارجی				
۴	استاد خارجی				
۵	استاد مدعو داخلی				
۶	استاد مدعو داخلی				
۷	استاد مدعو داخلی				

تأییدیه‌ی صحت و اصالت نتایج

باسمه تعالی

اینجانب فاطمه فتاحی به شماره دانشجویی ۹۴۵۲۱۱۳۵ دانشجوی رشته مهندسی کامپیوتر مقطع تحصیلی کارشناسی تأیید می‌نمایم که کلیه‌ی نتایج این پروژه حاصل کار اینجانب و بدون هرگونه دخل و تصرف است و موارد نسخه‌برداری‌شده از آثار دیگران را با ذکر کامل مشخصات منبع ذکر کرده‌ام. در صورت اثبات خلاف مندرجات فوق، به تشخیص دانشگاه مطابق با ضوابط و مقررات حاکم (قانون حمایت از حقوق مؤلفان و مصنفان و قانون ترجمه و تکثیر کتب و نشریات و آثار صوتی، ضوابط و مقررات آموزشی، پژوهشی و انضباطی ...) با اینجانب رفتار خواهد شد و حق هرگونه اعتراض درخصوص احقاق حقوق مکتسب و تشخیص و تعیین تخلف و مجازات را از خویش سلب می‌نمایم. در ضمن، مسؤولیت هرگونه پاسخگویی به اشخاص اعم از حقیقی و حقوقی و مراجع ذیصلاح (اعم از اداری و قضایی) به عهده‌ی اینجانب خواهد بود و دانشگاه هیچ‌گونه مسؤولیتی در این خصوص نخواهد داشت.

نام و نام خانوادگی: فاطمه فتاحی

تاریخ و امضا:

مجوز بهره‌برداری از پایان‌نامه

بهره‌برداری از این پایان‌نامه در چهارچوب مقررات کتابخانه و با توجه به محدودیتی که توسط استاد راهنما به شرح زیر تعیین می‌شود، بلامانع است:

- ☐ بهره‌برداری از این پایان‌نامه برای همگان بلامانع است.
- ☐ بهره‌برداری از این پایان‌نامه با اخذ مجوز از استاد راهنما، بلامانع است.
- ☐ بهره‌برداری از این پایان‌نامه تا تاریخ ممنوع است.

استاد راهنما: دکتر مینایی

تاریخ:

امضا:

فاطمه فتاحی

شهریور ۱۳۹۸

چکیده

پیش بینی ارتباطات دارو-هدف به عنوان یک مرحله بنیادی در فرایند کشف دارو و استفاده مجدد از آن محسوب می شود. روش های آزمایشگاهی در این حوزه، علی رغم داشتن نتایج دقیق و سودمند بسیار هزینه بر و زمان گیر هستند. در چند دهه ی اخیر، روش های محاسباتی جدیدی به منظور حل این مسئله مطرح شده اند که علاوه بر داشتن نتایج بسیار سودمند، به طور قابل توجهی کم هزینه بوده و زمان کمتری را نسبت به روش های آزمایشگاهی در بر می گیرد. از جمله ی این روش ها که نتایج حاصل از آن بسیار قابل توجه ارزیابی شده است، پیش بینی روابط بر پایه ی شبکه^۱ در شبکه های زیست پزشکی می باشد.

امروزه، روش های مبتنی بر شبکه متعددی برای پیش بینی انواع روابط میان داروها و پروتئین ها پیشنهاد گردیده است. با وجود عملکرد بسیار ارزنده آن ها، این راهکارها تنها برای شبکه های همگن^۲ به خوبی کار می کنند. این روش ها تنها اطلاعات مربوط به گره^۳ های شبکه را در اختیار ما می گذارند و مفاهیم یال^۴ ها را در بر نمی گیرند. در این پروژه، ما با استفاده از روش ادج توک^۵، به جاسازی نودها در شبکه ناهمگن^۶ دارو-پروتئینی پرداخته ایم و سپس با به کارگیری الگوریتم ماشین بردار پشتیبان^۷، ارتباطات دارو-هدف را پیش بینی نموده ایم. نتایج به دست آمده حاصل از این روش نسبت به سایر روش های موجود مطرح شده در این حوزه به طور چشمگیری بالاتر است.

واژگان کلیدی: پیش بینی ارتباطات دارو-هدف، جاسازی نودها، اج توک، ماشین بردار پشتیبان، شبکه ناهمگن

¹Network

²Homogeneous

³Node

⁴Edge

⁵Edge2vec

⁶Heterogeneous

⁷Support Vector Machine

فهرست مطالب

چ	فهرست تصاویر
ح	فهرست جداول
خ	فهرست الگوریتم‌ها
د	فهرست علائم اختصاری
۱	فصل ۱: مقدمه
۲	۱-۱ مقدمه
۴	فصل ۲: مروری بر منابع
۵	۲-۱ مقدمه
۵	۲-۲ تعاریف، اصول و مبانی نظری
۵	۲-۲-۱ شبکه
۶	۲-۳ مروری بر ادبیات موضوع
۶	۲-۳-۱ جاسازی نودها در شبکه
۹	۲-۴ نتیجه گیری
۱۱	فصل ۳: روش تحقیق
۱۲	۳-۱ مقدمه
۱۲	۳-۲ تشریح کامل روش تحقیق

۱۳ ۳-۲-۱ الگوریتم ادج توک

۱۴ ۳-۲-۲ روش ماشین بردار پشتیبان

۱۵ فصل ۴: نتایج و تفسیر آنها

۱۶ ۴-۱ مقدمه

۱۶ ۴-۲ محتوا

۱۶ ۴-۲-۱ مجموعه داده ها

۱۶ ۴-۲-۲ ارزیابی نتایج

۲۰ فصل ۵: جمع بندی و پیشنهادها

۲۱ ۵-۱ مقدمه

۲۱ ۵-۲ جمع بندی

۲۱ ۵-۳ نوآوری

۲۱ ۵-۴ پیشنهادها

۲۳ مراجع

۲۵ واژه‌نامه فارسی به انگلیسی

۲۷ واژه‌نامه انگلیسی به فارسی

فهرست تصاویر

- ۲-۱ جاسازی نودها در شبکه ۶
- ۲-۲ استراتژی پیاده روی تصادفی ۸
- ۲-۳ چگونگی تعریف متابت برای یک شبکه زیست پزشکی ۱۰
- ۳-۱ مراحل روش پیشنهادی ۱۳
- ۴-۱ نمودارهای مربوط به تنظیم پارامترهای الگوریتم ادج توک ۱۷
- ۴-۲ مقادیر بازخوانی، دقت، امتیاز اف ۱ و خطای همینگ به ازای مقادیر مختلف دو پارامتر
پی و آر ۱۸
- ۴-۳ مقایسه میزان عملکرد میان روش های موجود با روش پیشنهادی ۱۹

فهرست جداول

فهرست الگوریتم‌ها

فهرست علائم اختصاری

فصل ۱

مقدمه

پیش بینی ارتباطات دارو-هدف پروتئینی به عنوان یک مرحله ضروری و مهم در حوزه پیدایش دارو به شمار می رود. در حقیقت این امر موجب ظهور دارو ها و هم چنین اهداف جدید برای دارو های موجود می شود. در چند دهه ی اخیر، روش های آزمایشگاهی بسیار متعددی در این موضوع مطرح گردیده است. اما اکثر این روش ها بسیار هزینه بر بوده و از لحاظ زمانی نیز وقت به خصوصی را می گیرد. آمار و ارقام نشان دهنده این است که برای تولید یک داروی جدید به طور متوسط، نیاز به ۱۰ تا ۱۵ سال زمان و حدود ۸۰۰ میلیون دلار سرمایه است [۳].

با در نظر گرفتن این مسئله، روش های محاسباتی می توانند جایگزین مناسبی برای آن ها باشند. این روش ها علاوه بر پشت سر گذاشتن مشکلات ذکر شده، دقت و عملکرد بسیار خوبی نیز ارائه می کنند. امروزه با افزایش تنوع و ناهمگونی دارو ها و روابط میان شان با سایر عناصر، استفاده از این روش های محاسباتی بسیار ارزشمند تلقی می شود. در چند سال اخیر، تلاش های بسیاری به منظور تشخیص این چنین روابطی با استفاده از این روش ها انجام شده است. تقریباً در تمام این روش ها فرض بر آن است که دارو های با ویژگی های مشابه، هدف های مشابهی نیز دارند.

از آن جا که دانش زیست پزشکی مدام در حال رشد و پیچیدگی است و ما با تعداد بسیار زیادی مدل و روابط بینشان روبه رو هستیم، باید از یک مدل داده ای که بتواند به خوبی آن ها را توصیف نماید استفاده نماییم که شبکه نام دارد. در این مدل ما با انواع گره ها و یال ها که هر کدام اطلاعات مخصوص به خود را نگاه می دارند در ارتباطیم.

تقریباً اغلب روش های مطرح شده برای پیش بینی ارتباطات دارو-هدف برای شبکه های همگن و مدل های دو بخشی^۱ کاربرد دارند. به تازگی نیز روش هایی برای شبکه های ناهمگن نیز مطرح گردیده است که در این گونه روش ها علاوه بر در نظر گرفتن اطلاعات نودها، مفاهیم مربوط به یال ها را نیز در اختیار داریم. در این پروژه، روشی برای پیش بینی ارتباطات دارو-هدف و گرفتن دقت و بازخوانی بسیار بالا بر روی یک شبکه متشکل از انواع متنوعی از روابط و موجودیت ها ارائه شده است. هم چنین به تنظیم پارامترهای موجود در الگوریتم ادج توک نیز پرداخته ایم و بالاترین نتایج را ذکر نموده ایم.

در بخش های آتی، ابتدا به بررسی کارهای پیشین در این حوزه و سپس به ارائه دقیق روش ارائه شده می

¹bipartite models

پردازیم. در انتها نتایج را بررسی و به جمع بندی و ارائه ی پیشنهادات می پردازیم.

فصل ۲

مروری بر منابع

۲-۱-۱ مقدمه

در این فصل ابتدا به بررسی مفاهیم شبکه و اصول آن و سپس، به بررسی موضوع جاسازی نود ها در شبکه و انواع روش های تعیین شباهت نود می پردازیم. در بخش انواع روش های تعیین شباهت نود، در حقیقت همان روش های پیشین برای پیش بینی ارتباطات دارو-هدف که مبتنی بر شبکه هستند پرداخته خواهد شد.

۲-۲ تعاریف، اصول و مبانی نظری

در این بخش ابتدا به بررسی شبکه و کاربرد آن و هم چنین انواع آن می پردازیم.

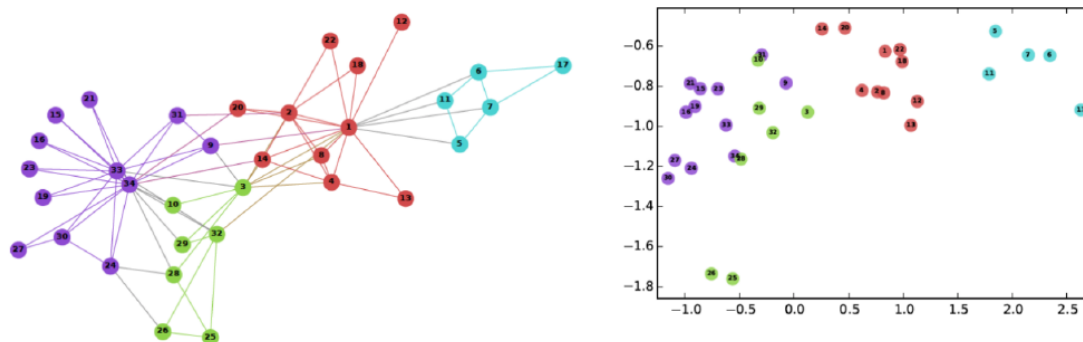
۲-۲-۱ شبکه

^۱ در بسیاری از مسائل از جمله سیستم های اجتماعی و اقتصادی برای درک و توصیف رفتار و روابط موجود در سیستم از مدل داده ای گراف برای ساده سازی شهود استفاده می گردد. در حقیقت، شبکه ها یک زبان عمومی برای توصیف و مدل سازی سیستم های پیچیده هستند. در علم زیست پزشکی نیز به دلیل پیچیدگی و تنوع روابط و موجودیت ها از این نوع ساختار داده به منظور درک بهتر استفاده می شود.

انواع شبکه

همان طور که می دانیم، ما در بسیاری از شبکه ها تنها با یک نوع نود و رابطه سر و کار داریم مثل رابطه پزشک با بیمار؛ در واقع، به این شبکه ها که در آن ما کاری به نوع یال نداریم شبکه های همگن^۲ می نامیم. اما دسته ی دیگری از شبکه ها وجود دارند که ما در دنیای امروز اغلب با آن ها سروکار داریم. در این نوع شبکه ها، انواع متنوعی از روابط وجود دارد که ما آن ها را در نظر می گیریم. این دسته را شبکه های ناهمگن^۳ می خوانند.

^۱Network^۲Homogeneous^۳Heterogeneous



شکل ۲-۱: جاسازی نودها در شبکه

۲-۳. مروری بر ادبیات موضوع

در این بخش ابتدا به بررسی مفهوم جاسازی نودها در گراف و اهمیت آن می پردازیم. سپس به روش های جاسازی نودها در گراف که تاکنون مطرح شده است، اشاره خواهیم کرد.

۲-۳-۱ جاسازی نودها در شبکه

^۴ همان طور که در بخش قبل توضیح داده شد، در شبکه های زیست پزشکی به دلیل پیچیدگی بسیار زیاد روابط و انواع زیاد نودها نمی توان اطلاعات مهم و مورد نیاز را از آن استخراج نمود. به همین منظور، روشی تحت عنوان جاسازی نودها در گراف مطرح گردیده است که نودها را در یک فضای با ابعاد پایین تر تعبیه می کند؛ به گونه ای که اطلاعات مربوط به آن ها از دست نرود. به این ترتیب گره های مشابه در گراف، در فضای جاسازی نیز نزدیک به هم هستند. هدف رمزگذاری گره ها به این گونه است که شباهت در فضای تعبیه شده، شباهت در شبکه ی اصلی را نشان می دهد. در شکل ۲-۳-۱ یک شبکه مربوط به یک مجموعه داده خاص، در یک فضای جاسازی شده قرار گرفته است.

در این حوزه، روش های بسیاری برای تعیین شباهت نودها در گراف ایجاد شده است. به طور کلی این روش ها به دسته زیر تقسیم بندی می شوند:

^۴Node embedding

- شباهت مبتنی بر وابستگی^۵
 - شباهت چند گامی^۶
 - روش های پیاده روی تصادفی^۷
- روش های پیاده روی تصادفی به دلیل بهینه بودن از لحاظ پیچیدگی زمانی و سرعت بیشتر آن ها نسبت به دو روش دیگر از توجه بسیاری برخوردارند و دو روش دیگر به ندرت به کار گرفته می شوند. در روش های پیاده روی تصادفی، احتمال ملاقات یک گره دلخواه در یک پیاده روی تصادفی که از گره مشخصی با استفاده از یک استراتژی^۸ مشخص آغاز می شود، محاسبه می گردد. مزیت های این روش به صورت زیر می باشد:
- انعطاف پذیری^۹: علاوه بر اطلاعات همسایه های محلی شامل اطلاعات همسایه های مرتبه های بالاتر نیز می شود.
 - کارایی: به دلیل در نظر نگرفتن همه ی نود ها در فرآیند آموزش^{۱۰}، پیچیدگی کمتری دارد.

استراتژی های پیاده روی تصادفی

همان گونه که در قسمت قبل مطرح شد، روش های پیاده روی تصادفی با استفاده از یک سری استراتژی های مشخص، به تعیین شباهت نود در گراف می پردازند. در ادامه تعدادی از استراتژی های مطرح شده را توضیح خواهیم داد.

یکی از ساده ترین ایده های مطرح شده توسط پروزی و همکارانش در [۱۱]، استراتژی دیپ واک^{۱۱} بود که از هر نود، پیاده روی های تصادفی با طول ثابت و بدون هیچ تأثیری^{۱۲} را اجرا می کنیم. کمی بعد از آن در سال ۲۰۱۶، گورور و لسکوک، این استراتژی را با دو پارامتر پی^{۱۳} به معنای پارامتر بازگشت

⁵Adjacency-based Similarity

⁶Multi-hop Similarity

⁷Random Walk Approaches

⁸Strategy

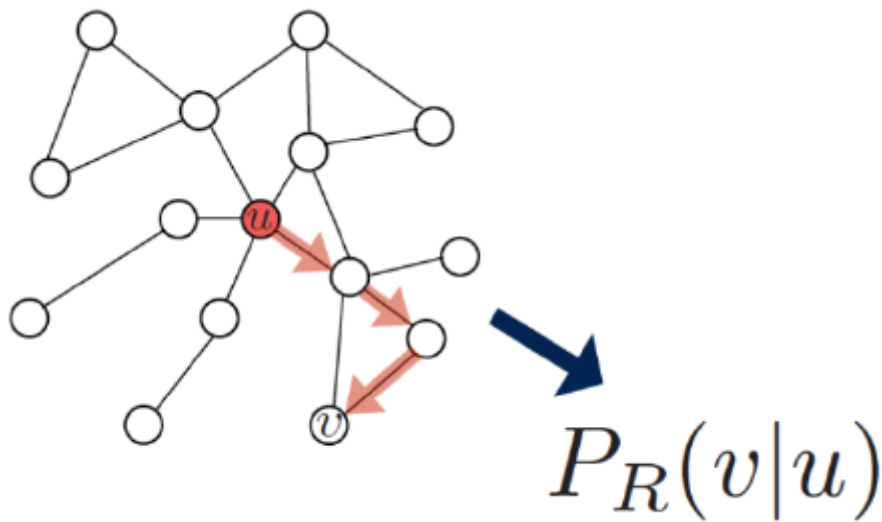
⁹Flexibility

¹⁰Training

¹¹DeepWalk

¹²Unbiased

¹³p



شکل ۲-۲: استراتژی پیاده روی تصادفی

^{۱۴} و ^{۱۵} کیو به معنای دور شدن ^{۱۶} تحت تأثیر قرار دادند و روشی تحت عنوان نود تو وک ^{۱۷} نتایج با دقت بالاتری را کسب کردند [۶]. در واقع، این دو پارامتر باعث یک توازن ^{۱۸} بین دید های محلی و جهانی شبکه می شود.

در این میان روش های دیگری نیز بر همین اساس ایجاد شده اند. از جمله این روش ها می توان به [۱] اشاره نمود که بر مبنای وزن های یادگیرنده به وجود آمده است.

استراتژی دیگری که در همین دوره بعد از نود تو وک کانون توجه پژوهشگران قرار گرفت، الگوریتم لاین ^{۱۹} بود [۱۲]. این روش، یک روش بهینه سازی جایگزین نود تو وک محسوب می شود که براساس احتمالات پیاده روی های تصادفی تک گام ^{۲۰} و دو گام طرح ریزی شده است. این استراتژی از آن جهت بسیار مورد توجه قرار گرفت که بسیار مقیاس پذیر بوده و برای شبکه های بی وزن ^{۲۱}، وزن دار ^{۲۲}، بی جهت ^{۲۳} و جهت دار ^{۲۴}

¹⁴Return parameter

¹⁵q

¹⁶Walk away parameter

¹⁷Node2vec

¹⁸Trade off

¹⁹LINE

²⁰One-hop

²¹Unweighted

²²Weighted

²³Undirected

²⁴Directed

به خوبی کار می کند. تاکنون روش های متعددی مطرح شده است؛ اما به نظر می رسد این روش عملکرد و نتایج بسیار دقیق تری نسبت به سایر روش های موجود دارد.

مشکلی که در همه ی این روش ها به چشم می خورد این است که تنها برای شبکه های همگن به خوبی عمل می کنند. در حالی که شبکه های زیست پزشکی غالباً به صورت ناهمگن می باشند و لذا این گونه روش ها نمی توانند به خوبی عملیات جاسازی نودها در شبکه را نشان دهند.

برای حل این مشکلات، روش دیگری تحت عنوان متاپت توک^{۲۵} معرفی شد که علاوه بر شبکه های همگن برای شبکه های ناهمگن به خوبی کار می کند [۴]. مبنای این روش، تعریف یک سری متاپت^{۲۶} می باشد که بر مبنای آن ها در داخل گراف عمل پیاده روی تصادفی را انجام می دهد. این استراتژی، همبستگی بین انواع مختلف نود را در نظر می گیرد. علی رغم کارایی این روش برای شبکه های ناهمگن معایبی نیز دارد که به شرح زیر می باشد:

- برای تعریف هر متاپت نیاز به دانش در آن حوزه می باشیم. برای مثال باید بدانیم که هر پروتئین با چه عناصری می تواند ارتباط داشته باشد و زنجیره مان را طبق آن تعریف نماییم.
- تنها از یک متاپت در هر لحظه برای تولید یک پیاده روی تصادفی استفاده میکند.
- مفاهیم مربوط به یال ها را در نظر نمی گیرد.

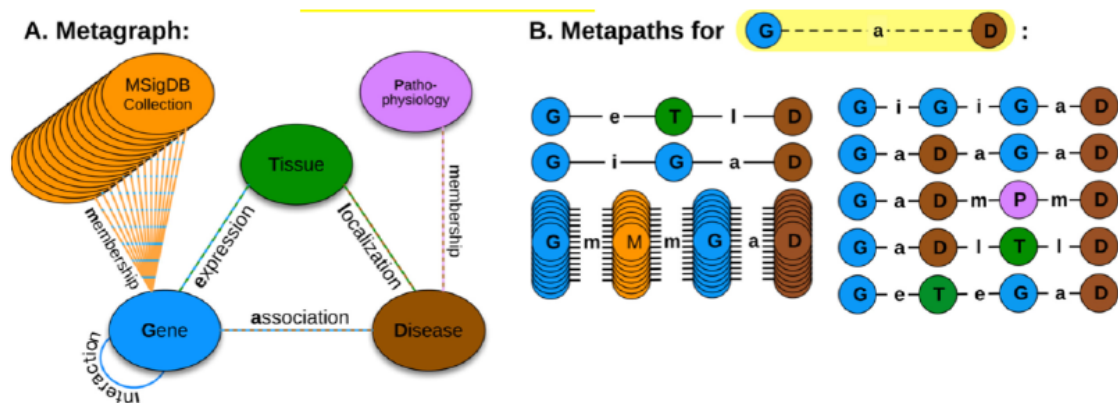
همان گونه که در شکل ۲-۳-۱ مشاهده می کنید تمام ارتباطات ممکن بین ژن و بیماری را به عنوان متاپت تعریف نموده است.

۲-۴ نتیجه گیری

در این بخش مفاهیم و مطالعات پیشین مربوط به حوزه مدنظر بررسی شدند. همان گونه که مطالعه شد، تقریباً تمامی روش ها و استراتژی های مربوط به عمل پیاده روی تصادفی در شبکه به منظور جاسازی نود ها در فضایی با ابعاد پایین تر محدود به شبکه های همگن می شدند و روش های جدید نیز معایب زیادی دارند که به آن ها اشاره شد. در نتیجه، نیاز به یک استراتژی جدید برای تعیین شباهت نود ها در گراف احساس می شود. در مطالعه پیش رو سعی در استفاده از یک روش جدید است.

²⁵ Metapath2vec

²⁶ Metapath



شکل ۲-۳: چگونگی تعریف متاپت برای یک شبکه زیست پزشکی

فصل ۳

روش تحقیق

۳-۱ مقدمه

در این بخش به بررسی روش ارائه شده در این پروژه برای پیش بینی ارتباطات دارو-هدف در شبکه ناهمگن می پردازیم.

۳-۲ تشریح کامل روش تحقیق

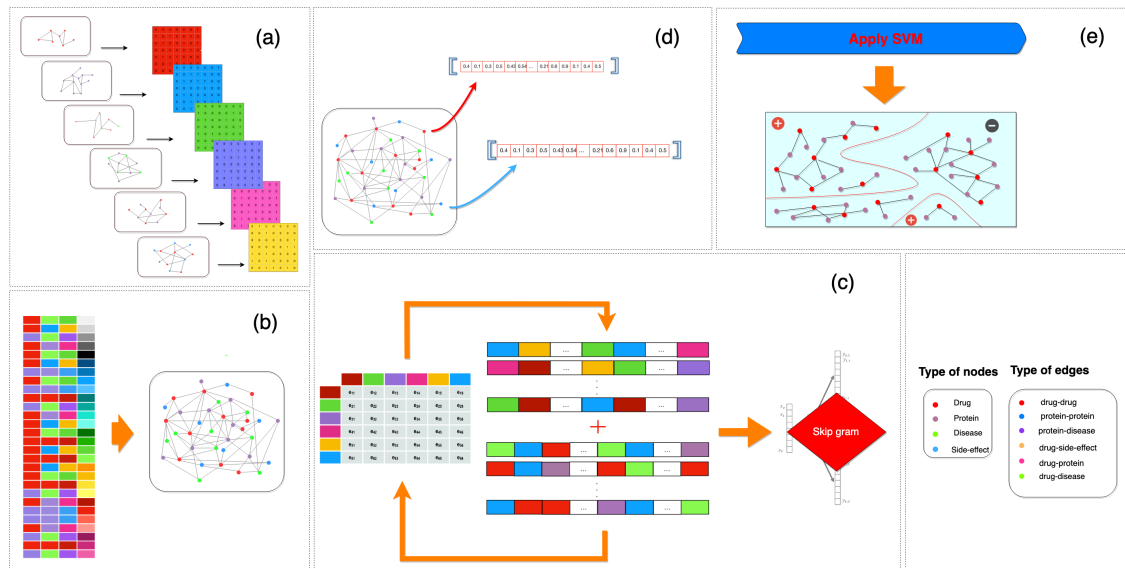
به طور کلی، روش پیشنهادی ما از چند گام که در ادامه به آن اشاره خواهد شد، تشکیل شده است. ابتدا، ما پایگاه داده ای مطابق با مقاله [۱۰] با استفاده از چهار پایگاه داده ایجاد نمودیم. این پایگاه داده که جزئیات آن در فصل بعد آمده است، از چهار نوع نود و شش نوع یال تشکیل یافته است. به عبارت دیگر شش نوع ماتریس^۱ که روابط مختلف بین نود ها را نشان می دهد داریم. عدد صفر در خانه های آرایه به منزله عدم وجود رابطه بین دو راس و نبود یال و عدد یک حاکی از وجود یک رابطه یا یال بین آن دو است.

در مرحله بعد که در شکل ۳-۲ به خوبی واضح می باشد، برای ساختن شبکه مربوطه مان تمام سطر هایی که نشان دهنده یک ارتباط میان دو نود می باشد را استخراج کرده و شبکه مورد نظرمون را می سازیم.

حال نوبت به آن رسیده است که با استفاده از یکی از روش های جاسازی نود در گراف، شبکه جاسازی شده را به دست بیاوریم. همان گونه که در فصل پیشین اشاره شد، روش های بررسی شده هر کدام مزایا و معایب خود را داشتند. ما برای این پروژه، از روش جدیدی تحت عنوان ادج توک که در مقاله [۵] به خوبی شرح داده شده است استفاده نمودیم که با عملکرد دقیق تر و بهتری به جاسازی نود ها در شبکه پردازیم.

پس از انجام عمل جاسازی نودها در شبکه، هر نود با استفاده از یک بردار نماینده^۲ با ابعاد^۳ ۱۲۸ بازنمایی می شود. به منظور پیش بینی و طبقه بندی یال ها باید به اندازه تمام یال های با برجسب^۴ مثبت^۵، یال هایی با برجسب منفی^۶ تولید نماییم. منظور از برجسب منفی، یال هایی ست که در حقیقت در گراف وجود ندارند یا ارتباطی میان آن دو نود مشخص وجود ندارد. به این نوع داده ها، نمونه های منفی می گویند. حال در گام آخر، با استفاده از یک طبقه بند به پیش بینی ارتباطات دارو-هدف می پردازیم. ما در این آزمایش از طبقه بند

^۱Matrix^۲representative Vector^۳Dimension^۴Label^۵Positive^۶Negative



شکل ۳-۱: مراحل روش پیشنهادی

ماشین بردار پشتیبان استفاده نمودیم.

در ادامه به بررسی دقیق الگوریتم ادج توک و روش ماشین بردار پشتیبان می پردازیم.

۳-۲-۱ الگوریتم ادج توک

در شبکه های ناهمگن ما با انواع نودها و روابط سروکار داریم. ژنگ جاو و همکارانش با استفاده از فرایند پیاده روی تصادفی، یک ماتریس انتقال نوع یال براساس وزن های انتقال میان انواع یال ها ایجاد نمودند. او با استفاده از فرمول ۳-۱، احتمال نودهای همسایه یک نوع نود مشخص را به حداکثر می رساند:

$$\arg \max_{\theta, M} \prod_{v \in V} \prod_{c \in N(v)} p(c | v; \theta; M) \quad (3-1)$$

در شبکه $G(V, E)$ ، V به مجموعه ی نودها اشاره می کند و E مجموعه ی یال هارا مشخص می نماید. هم چنین $N(V)$ یک گروه از همسایگان نود V می باشد. θ به پارامتری شدن جاسازی نودها اشاره می کند. ماتریس M نیز همان ماتریس انتقال نوع یال می باشد.

به طور کلی، الگوریتم ادج تو وک مقادیر ماتریس M را ابتدا به ۱ مقدار دهی می کند. سپس با استفاده از چهارچوب^۷ انتظار حداکثر^۸، این ماتریس آموزش می بیند. چهارچوب انتظار حداکثر از دو مرحله انتظار و حداکثر سازی تشکیل یافته است. در مرحله انتظار، با استفاده از فرمول^{۳-۲} همبستگی^۹ میان انواع یال ها را بهینه می کند:

$$M(e_i, e_j) = \text{Sigmoid}\left(\frac{E[(\vec{v}_i - \mu(\vec{v}_i))(\vec{v}_j - \mu(\vec{v}_j))]}{\sigma(\vec{v}_i)\sigma(\vec{v}_j)}\right) \quad (2-3)$$

در این فرمول، e_i و e_j همان نوع یال ها هستند و اشاره به ماتریس انتقال نوع لبه به روز شده می کند. همچنین، v_i و v_j بردارهایی هستند که e_i و e_j را بازنمایی می کنند. مقدار $E[\cdot]$ همان مقدار انتظار می باشد. در مرحله حداکثر سازی، آخرین نسخه به روز شده از ماتریس انتقال نوع یال را دریافت کرده و پیاده روی های تصادفی را بر مبنای این ماتریس انجام می دهد. احتمال پیاده روی های تصادفی با استفاده از فرمول زیر تعیین می شود:

$$p(n | v; u; M) = \frac{\omega_{vn} \cdot M_{T(u,v)T(v,n)} \cdot \alpha_{pq}(n, u)}{\sum_{k \in N(v)} \omega_{vk} \cdot M_{T(u,v)T(v,k)} \cdot \alpha_{pq}(k, u)} \quad (3-3)$$

۳-۲-۲ روش ماشین بردار پشتیبان

بردارهای پشتیبان به زبان ساده، مجموعه ای از نقاط در فضای n بعدی داده ها هستند که مرز دسته ها را مشخص می کنند و مرزبندی و دسته بندی داده ها براساس آنها انجام می شود و با جابجایی یکی از آنها، خروجی دسته بندی ممکن است تغییر کند. ماشین بردار پشتیبان، یک دسته بند یا مرزی است که با معیار قرار دادن بردارهای پشتیبان، بهترین دسته بندی و تفکیک بین داده ها را برای ما مشخص می کند. در این روش فقط داده های قرار گرفته در بردارهای پشتیبان مبنای یادگیری ماشین و ساخت مدل قرار می گیرند و این الگوریتم به سایر نقاط داده حساس نیست و هدف آن هم یافتن بهترین مرز در بین داده هاست به گونه ای که بیشترین فاصله ممکن را از تمام دسته ها (بردارهای پشتیبان آنها) داشته باشد.

⁷Framework

⁸Expectation Maximization

⁹Correlation

فصل ۴

نتایج و تفسیر آنها

۴-۱ مقدمه

در این بخش نتایج حاصل از روش پیشنهادی مطرح شده در فصل سوم را بر روی مجموعه داده ای که ساخته ایم، بررسی می کنیم.

۴-۲ محتوا

۴-۲-۱ مجموعه داده ها

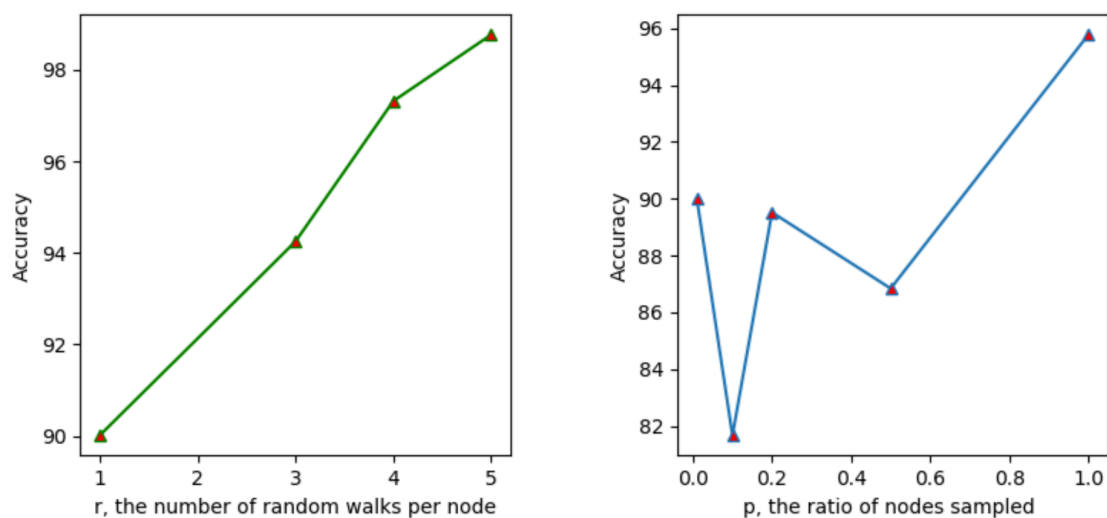
ما مجموعه داده هایمان را طبق مقاله [۱۰] که ترکیبی از چند پایگاه داده ی مرتبط با دارو است، ایجاد نمودیم. این مجموعه داده از چهار نوع نود که شامل دارو، پروتئین، اثر جانبی و بیماری و انواع یال تشکیل شده است. یال ها نیز انواع مختلفی مثل دارو-دارو، دارو-پروتئین، دارو-اثر جانبی، دارو-بیماری، پروتئین-پروتئین و پروتئین-بیماری دارند. این مجموعه داده حدود ۱۱۴۸۸ نود و ۱۸۹۵۴۳۳ یال دارد. نود های دارو و روابط دارو-هدف از مجموعه داده دراگ بانک^۱ [۸] برگرفته شده است. هم چنین نود های پروتئینی و روابط پروتئین-پروتئین از مجموعه داده اچ پی آردی^۲ [۷] گرفته شده است. علاوه بر این، نود های اثرات جانبی و روابط دارو-اثر جانبی از پایگاه داده سایدرا^۳ [۹] است. در نهایت، نود های بیماری، یال های دارو-بیماری و پروتئین-بیماری از پایگاه داده [۲] استفاده شده است.

۴-۲-۲ ارزیابی نتایج

چندین آزمایش به منظور عملکرد روش پیشنهادی انجام گرفته است. ما از اعتبار سنجی متقابل کی فولد^۴ به منظور افزایش اطمینان از ارزیابی مان استفاده کرده ایم. در این مطالعه، ما روش ماشین بردار پشتیبان با تابع هسته^۵ آربی اف^۶ را برای رسیدن به بالاترین دقت به کار گرفته ایم.

علاوه براین، ما به تنظیم پارامترهای الگوریتم ادج توک به منظور یافتن بهترین نتیجه و میزان تغییرات

^۱ DrugBank^۲ HPRD^۳ SIDER^۴ k-Fold cross validation^۵ Kernel Function^۶ RBF



شکل ۴-۱: نمودارهای مربوط به تنظیم پارامترهای الگوریتم ادج تو وک

پرداختیم. این روش شامل چهار پارامتر اساسی آر^۷، دابلویو^۸، پی^۹ و این^{۱۰}. آر به تعداد پیاده روی های تصادفی روی هر نود اشاره می کند که مقدار پیش فرض آن برابر با یک می باشد. پارامتر دابلویو به طول هر پیاده روی اشاره کرده که معادل با ۵۰ می باشد. علاوه بر آن، پارامتر پی همان نسبت نود های نمونه گذاری شده است که مقدار پیش فرض آن ۰.۱۰ است و این تعداد تکرارها برای آموزش ماتریس انتقال نوع یال است که برابر ۱۰ می باشد. ما در این آزمایش مقدار دو متغیر آر و پی را تغییر داده و دو پارامتر دیگر را ثابت نگه داشتیم. شکل ۴-۲-۲ این موضوع را به خوبی نشان می دهد. با افزایش آر، به درصد قابل توجهی می رسمیم که این به دلیل تاثیر روابط یال های بیشتر در طول هر نود می باشد. اما با افزایش مقدار پی، دقت ما به طور نامعینی می باشد و گاهی زیاد و گاهی کم می شود. لذا از تغییر مقدار پیش فرض آن صرف نظر کرده و آن را ثابت نگاه می داریم.

معیار های ارزیابی زیادی برای اندازه گیری عملکرد مطالعه مان وجود دارد. ما چهار معیار ارزیابی امتیاز

⁷ r
⁸ W
⁹ p
¹⁰ N

		Recall	Precision	F1 score	Hamming loss
p	0.1	0.81	0.81	0.81	0.18
	0.2	0.89	0.89	0.89	0.10
	0.5	0.86	0.86	0.86	0.13
	1	0.95	0.95	0.95	0.04
r	3	0.94	0.94	0.94	0.05
	4	0.97	0.97	0.97	0.02
	5	0.98	0.98	0.98	0.01

شکل ۴-۲: مقادیر بازخوانی، دقت، امتیاز اف ۱ و خطای همینگ به ازای مقادیر مختلف دو پارامتر پی و آر

اف ۱^{۱۱}، بازخوانی^{۱۲}، دقت^{۱۳} و فاصله همینگ^{۱۴} را برای هریک از مقادیر پارامترهای پی و آر را محاسبه کرده و در شکل ۴-۲-۲ یادداشت کرده ایم

در شکل ۴-۲-۲، ما میزان امتیاز آیو آر اُسی^{۱۵} را با روش هایی که تاکنون موردتوجه قرار گرفته اند، مورد مقایسه قرار دادیم. همان گونه که مشخص است، این روش نسبت به روش های موجود به دقت و عملکرد بسیاری دست یافته است.

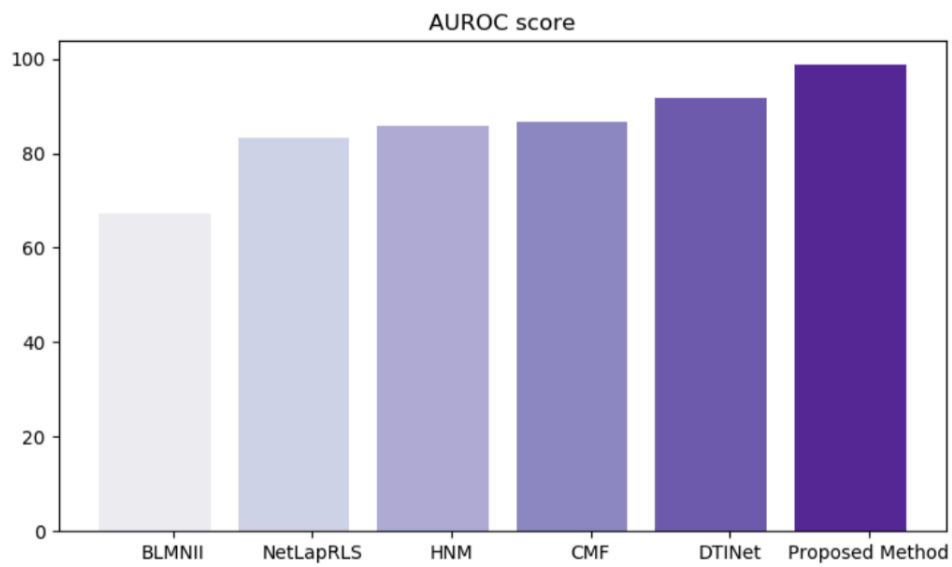
¹¹F1-score

¹²Recall

¹³Precision

¹⁴Hamming loss

¹⁵AUROC Score



شکل ۴-۳: مقایسه میزان عملکرد میان روش های موجود با روش پیشنهادی

فصل ۵

جمع بندی و پیشنهادها

۵-۱ مقدمه

پیش بینی ارتباطات دارو-هدف از موضوعات مطرح در حوزه کشف دارو و استفاده مجدد از آن محسوب می شود، که با الگوریتم ادج توک به این کار پرداختیم. در این الگوریتم، علاوه بر مفاهیم نود ها اطلاعات مربوط به یال ها نیز حفظ می شود. در این بخش به معرفی کلی مطالعه ی انجام شده و جمع بندی می پردازیم. در انتها نیز درباره ی نوآوری و پیشنهادات آتی برای این مطالعه بحث می کنیم.

۵-۲ جمع بندی

در این مطالعه، سعی در پیش بینی و تشخیص ارتباطات و لینک های میان دارو و هدف کردیم. برای این کار، ما ابتدا مجموعه داده ای از روابط مختلف میان عناصر را طبق مجموعه داده های معتبر به وجود آورده و سپس الگوریتم ادج توک که به تازگی ارائه شده است را به منظور جاسازی نود ها به کار بردیم. در انتها با استفاده از روش ماشین بردار پشتیبان به پیش بینی این ارتباطات پرداختیم. هم چنین به منظور رسیدن به دقت بالاتر، به تنظیم پارامتر های این دو الگوریتم پرداختیم. روشی که در این پژوهش پیشنهاد شده است نسبت به روش های موجود از دقت بالاتری برخوردار می باشد.

۵-۳ نوآوری

روشی که در این پژوهش از آن برای پیش بینی ارتباطات دارو-هدف استفاده شده است، تاکنون مورد توجه قرار نگرفته است. هم چنین این روش نسبت به روش هایی که تاکنون ارائه شده است، میزان عملکرد و دقت قابل توجهی را ایفا می کند.

۵-۴ پیشنهادها

برای کارهای آتی می توان روابط بیشتری هم چون ارتباطات پروتئین-پروتئین^۱ و یا دارو-دارو^۲ را با به کار گیری این روش پیش بینی نمود. هم چنین، اعمال روش پیشنهادی بر روی مجموعه داده های دیگر می تواند

^۱Protein-Protein interaction(PPI)

^۲Drug-Drug interaction(DDI)

به پیدایش و کشف روابط بیشتر دارو-هدف شود.

مراجع

- [1] Abu-El-Haija, S., Perozzi, B., Al-Rfou, R., and Alemi, A. A. Watch your step: Learning node embeddings via graph attention. in *Advances in Neural Information Processing Systems* (2018), pp. 9180–9190.
- [2] Davis, A. P., Murphy, C. G., Johnson, R., Lay, J. M., Lennon-Hopkins, K., Saraceni-Richards, C., Sciaky, D., King, B. L., Rosenstein, M. C., Wiegers, T. C., et al. The comparative toxicogenomics database: update 2013. *Nucleic acids research* 41, D1 (2012), D1104–D1114.
- [3] DiMasi, J. A. New drug development in the united states from 1963 to 1999. *Clinical Pharmacology & Therapeutics* 69, 5 (2001), 286–296.
- [4] Fu, G., Ding, Y., Seal, A., Chen, B., Sun, Y., and Bolton, E. Predicting drug target interactions using meta-path-based semantic network analysis. *BMC bioinformatics* 17, 1 (2016), 160.
- [5] Gao, Z., Fu, G., Ouyang, C., Tsutsui, S., Liu, X., Yang, J., Gessner, C., Foote, B., Wild, D., Ding, Y., and Yu, Q. edge2vec: Representation learning using edge semantics for biomedical knowledge discovery. *BMC Bioinformatics* 20, 1 (Jun 2019), 306.
- [6] Grover, A., and Leskovec, J. node2vec: Scalable feature learning for networks. in *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining* (2016), ACM, pp. 855–864.
- [7] Keshava Prasad, T., Goel, R., Kandasamy, K., Keerthikumar, S., Kumar, S., Mathivanan, S., Telikicherla, D., Raju, R., Shafreen, B., Venugopal, A., et al. Human protein reference database—2009 update. *Nucleic acids research* 37, suppl_1 (2008), D767–D772.
- [8] Knox, C., Law, V., Jewison, T., Liu, P., Ly, S., Frolkis, A., Pon, A., Banco, K., Mak, C., Neveu, V., et al. Drugbank 3.0: a comprehensive resource for ‘omics’ research on drugs. *Nucleic acids research* 39, suppl_1 (2010), D1035–D1041.

- [9] Kuhn, M., Campillos, M., Letunic, I., Jensen, L. J., and Bork, P. A side effect resource to capture phenotypic effects of drugs. *Molecular systems biology* 6, 1 (2010).
- [10] Luo, Y., Zhao, X., Zhou, J., Yang, J., Zhang, Y., Kuang, W., Peng, J., Chen, L., and Zeng, J. A network integration approach for drug-target interaction prediction and computational drug repositioning from heterogeneous information. *Nature communications* 8, 1 (2017), 573.
- [11] Perozzi, B., Al-Rfou, R., and Skiena, S. Deepwalk: Online learning of social representations. in *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining* (2014), ACM, pp. 701–710.
- [12] Tang, J., Qu, M., Wang, M., Zhang, M., Yan, J., and Mei, Q. Line: Large-scale information network embedding. in *Proceedings of the 24th international conference on world wide web* (2015), International World Wide Web Conferences Steering Committee, pp. 1067–1077.

واژه‌نامه فارسی به انگلیسی

Flexibility	انعطاف پذیری
Training	آموزش
Unbiased	بدون تاثیر
Return parameter	پارامتر بازگشت
Walk away parameter	پارامتر دور شدن
Dimension	ابعاد
Framework	چارچوب
Correlation	همبستگی
Cross validation	اعتبارسنجی متقابل
Kernel Function	تابع هسته
Network based	مبتنی بر شبکه
Support vector machine	ماشین بردار پشتیبانی
Node embedding	جاسازی نودها
Bipartite Models	مدل های دو بخشی
Drug-target interactions	ارتباطات دارو-هدف
Network	شبکه
Trade off	توازن
Unweighted	بی وزن
Directed	جهت دار
Representative Vector	بردار نماینده
Homogeneous	همگن
Heterogeneous	ناهمگن
Adjacency-based Similarity	شباهت مبتنی بر وابستگی
Multi-hop Similarity	شباهت چند گامی
Random Walk Approaches	روش های پیاده روی تصادفی
Strategy	استراتژی

واژه‌نامه انگلیسی به فارسی

Node embedding	جاسازی نودها
Bipartite Models	مدل های دو بخشی
Drug-target interactions	ارتباطات دارو-هدف
Network	شبکه
Homogeneous	همگن
Heterogeneous	ناهمگن
Adjacency-based Similarity	شباهت مبتنی بر وابستگی
Multi-hop Similarity	شباهت چند گامی
Random Walk Approaches	روش های پیاده روی تصادفی
Strategy	استراتژی
Flexibility	انعطاف پذیری
Training	آموزش
Unbiased	بدون تاثیر
Return parameter	پارامتر بازگشت
Walk away parameter	پارامتر دور شدن
Trade off	توازن
Unweighted	بی وزن
Directed	جهت دار
Representative Vector	بردار نماینده
Dimension	ابعاد
Framework	چارچوب
Correlation	همبستگی
Cross validation	اعتبارسنجی متقابل
Kernel Function	تابع هسته
Network based	مبتنی بر شبکه
Support vector machine	ماشین بردار پشتیبانی

Abstract:

Drug-target protein interaction prediction is considered as a fundamental step in the process of drug discovery and re-purposing. The experimental methods in this domain, despite having accurate and beneficial results, are time-consuming. In the last few years, not only have these computational methods been overcome to these problems, but the results of them have been extremely efficient. One type of these techniques is based on network prediction on biomedical networks. Today, different network-based approaches have been suggested for identifying relations between drugs and targets which work on homogeneous networks well. But these approaches do not consider edge semantics in the network. In this project, we employ the edge2vec algorithm for node embedding on the heterogeneous network and then using the support vector machine method to predict drug-target interactions. Our method has obtained a significant accuracy than the existing methods.

Keywords: Drug-Target interaction prediction, Node embedding, Edge2vec, Support Vector Machine, Heterogeneous network



Iran University of Science and Technology
Computer Engineering Department

Drug-Target Interaction Prediction by Embedded Vectors of Nodes in Heterogeneous Interactions Network

Bachelor of Science Thesis in Computer Engineering

By:

Fatemeh Fattahi

Supervisor:

Dr. Minaei

September 2019