



دانشگاه تهران
پردیس دانشکده های فنی
دانشکده مهندسی کامپیوتر

دسته بندی بسته در شبکه های نرم افزار محور

پایان نامه جهت دریافت درجه کارشناسی ارشد
در رشته مهندسی کامپیوتر گرایش معماری

نگارش:

محمدرضا پیروزی

استاد راهنما:

دکتر ناصر یزدانی

بهمن ماه ۱۳۹۵

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه تهران
پردیس دانشکده های فنی
دانشکده مهندسی کامپیوتر



دسته بندی بسته در شبکه های نرم افزار محور

پایان نامه جهت دریافت درجه کارشناسی ارشد
در رشته مهندسی کامپیوتر گرایش معماری

نگارش:

محمدرضا پیروزی

استاد راهنما:

دکتر ناصر یزدانی

بهمن ماه ۱۳۹۵



دانشگاه تهران
پردیس دانشکده‌های فنی
دانشکده مهندسی برق و کامپیوتر

گواهی دفاع از پایان‌نامه کارشناسی ارشد

هیأت داوران پایان‌نامه کارشناسی ارشد آقای / خانم به شماره
دانشجویی در رشته گرایش را در تاریخ
..... با عنوان

به عدد به حروف

با نمره نهایی	
---------------	--

و درجه ارزیابی کرد.

--

ردیف	مشخصات هیأت داوران	نام و نام خانوادگی	مرتبه دانشگاهی	دانشگاه یا مؤسسه	امضاء
۱	استاد راهنما استاد راهنمای دوم (حسب مورد)				
۲	استاد مشاور				
۳	استاد داور داخلی				
۴	استاد مدعو				
۵	نماینده کمیته تحصیلات تکمیلی دانشکده / گروه				

نام و نام خانوادگی معاون تحصیلات تکمیلی و نام و نام خانوادگی معاون آموزشی و تحصیلات تکمیلی
پژوهشی دانشکده / گروه: پردیس دانشکده‌های فنی:
تاریخ و امضاء: تاریخ و امضاء:

تعهدنامه اصالت اثر

باسمه تعالی

اینجانب مسعود هاشمیان تأیید می‌کنم که مطالب مندرج در این پایان‌نامه حاصل کار پژوهشی اینجانب است و به دستاوردهای پژوهشی دیگران که در این نوشته از آنها استفاده شده است مطابق مقررات ارجاع گردیده است. این پایان‌نامه قبلاً برای احراز هیچ مدرک هم سطح یا بالاتر ارائه نشده است.

کلیه حقوق مادی و معنوی این اثر متعلق به دانشکده فنی دانشگاه تهران می‌باشد.

نام و نام خانوادگی دانشجو :

امضای دانشجو :

تقدیم به پدر و مادرم و همسر عزیزم

چکیده

در شبکه‌های نرم‌افزار محور، یک جریان بر اساس تعداد دلخواهی از فیلدها در هر لایه از سرآیند بسته، قابل تعریف است. به عنوان مثال در نسخه ۱,۳,۱ پروتکل OpenFlow از ۱۵ فیلد برای تعریف یک جریان استفاده می‌شود. این امر سبب شده تا اولاً عرض مدخل در جدول جریان افزایش یابد، ثانیاً کنترل ریزتر بر روی ترافیک شبکه باعث افزایش تعداد مدخل‌ها در جدول جریان شده است. در نتیجه طراحی یک راه‌حل دسته‌بندی بسته در محیط‌هایی که به کارایی بالا نیاز دارند، به یک چالش بزرگ تبدیل گشته است. برای حل این مسئله در سوئیچ‌های نرم‌افزاری سعی می‌شود تا حد امکان اندازه فضای نمونه برای مطابقت دادن یک چندتایی با قوانین، کوچک و کوچک‌تر شود. در الگوریتم MC-SBC یک ساختار مبتنی بر درخت تصمیم دو سطحی برای جدول‌های جستجو مطرح شده، و در آن بیت‌های موثر جهت افراز مجموعه قوانین و ایجاد درخت، با استفاده از یک روش آماری مشخص می‌گردند. در این پایان‌نامه از الگوریتم ژنتیک برای پیدا کردن مجموعه بیت‌های موثر جهت افراز مجموعه قوانین استفاده شده است. نتایج گزارش شده نشان می‌دهد مجموعه قوانین، با استفاده از این الگوریتم به صورت بهتری افراز می‌شوند، به طوری که هم نسبت تکرار قوانین در آن کاهش می‌یابد، و هم بیشینه تعداد قوانین در برگ‌های درخت کمتر می‌شود. بهبود کیفیت افراز مجموعه قوانین باعث می‌شود هنگامی که اندازه مجموعه قوانین بزرگ است، سرعت عمل دسته‌بندی بسته‌ها در نمونه پیاده‌سازی شده بر روی پردازنده گرافیکی تا حدود ۱۰ درصد بهبود یابد.

واژه‌های کلیدی: شبکه‌های نرم‌افزار محور، دسته‌بندی بسته، سوئیچ OpenFlow، پردازنده گرافیکی،

الگوریتم ژنتیک

فهرست مطالب

فصل ۱: مقدمه	۲
۱-۱- تعریف مسئله	۴
۱-۲- روش انجام پژوهش	۵
۱-۳- ساختار پایان نامه	۶
فصل ۲: زمینه‌های تحقیق	۷
۲-۱- شبکه‌های همگون و ناهمگون	۷
۲-۲- پیش‌بینی پیوند	۸
۲-۳- فرض جهان باز و جهان بسته	۹
۲-۴- گرادیان نزولی تصادفی (SGD)	۱۰
فصل ۳: پیشینه پژوهش	۱۱
۳-۱- روش‌های مبتنی بر ویژگی‌های پنهان	۱۲
۳-۲- تقسیم‌بندی داده‌های پایگاه دانش	۱۴
۳-۳- مراحل آموزش روش‌های مبتنی بر ویژگی‌های پنهان	۱۶
۳-۳-۱- مرحله‌ی آموزش	۱۷
۳-۳-۲- مرحله‌ی ارزیابی	۱۷
۳-۳-۳- مرحله‌ی آزمون	۱۹
۳-۴- بررسی روش‌های موجود	۱۹
۳-۴-۱- روش‌های خطی	۱۹
۳-۴-۲- روش‌های دوخطی - رسکال	۲۰
۳-۴-۳- روش‌های ادراک چندلایه‌ای	۲۶

۳۱	3-4-4 شبکه‌های عصبی تنسور.....
۳۳	۳-۴-۵- روش‌های فاصله‌ی پنهان.....
۳۵	3-4-6- مدل TransE.....
۳۸	3-4-7- مدل TransH.....
۴۱	فصل ۴: قوانین انجمنی.....
۴۱	۴-۱- معیارهای اطمینان.....
۴۲	۴-۱-۱- معیار پشتیبانی قانون:.....
۴۳	۴-۱-۲- معیار پوشش سر:.....
۴۳	۴-۱-۳- معیار اطمینان استاندارد:.....
۴۷	۴-۱-۴- معیار اطمینان با فرض نیمه‌کامل.....
۴۸	۴-۲- الگوریتم‌ها.....
۵۱	۴-۳- قوانین هدف.....
۵۲	۴-۳-۱- قانون R-subsumption.....
۵۱	۴-۳-۲- قانون R-equivalence.....
۵۳	۴-۳-۳- قانون 2-hope translation.....
۵۳	۴-۳-۴- قانون Triangle alignment.....
۵۴	۴-۳-۵- قانون Specific R-subsumption.....
۵۵	۴-۴- جمع‌بندی.....
۵۶	فصل ۵: روش پیشنهادی.....
۵۶	۵-۱- چالش‌ها.....
۵۶	۵-۲- عملکرد روش‌های موجود.....

۵۷	۵-۲-۱- روش رسکال:
۵۷	۵-۲-۲- روش NTN:
۵۸	۵-2-3- روش Structured Embedding
۵۹	۵-۲-۴- روش TransE:
۵۹	۵-۲-۵- روش TransH:
۶۰	۵-۳- نقاط قوت و ضعف روش‌های موجود:
۶۷	۵-۴- بررسی عمل‌کرد روش‌ها در یک نگاه:
۷۰	۵-۵- استفاده از قوانین انجمنی برای بهبود نتایج:
۷۱	۵-۵-۱- قانون R-subsumtion
۷۱	۵-5-2- قانون R-equivalence
۷۱	۵-۵-۳- قانون تعدی (2-hope translation)
۷۳	۵-5-4- قانون Triangle alignment
۷۳	۵-5-5- قانون Specefic R-subsumtion
۷۴	فصل ۶: ارزیابی
۷۴	۶-۱- مقدمه:
۷۴	۶-۲- آزمایش‌ها:
۸۲	۶-۳- جمع‌بندی:
۸۶	فصل ۷: نتیجه‌گیری
۸۶	۷-۱- نتیجه‌گیری:
۸۶	۷-۲- کارهای آینده:
۸۹	فصل ۸: مراجع

فصل ۹: واژه‌نامه انگلیسی به فارسی ۹۲

فصل ۱۰: واژه‌نامه فارسی به انگلیسی ۹۷

فهرست اشکال

No table of figures entries found.

فهرست جداول

جدول (۲-۱) تغییرات نسخه‌های اصلی OpenFlow **Error! Bookmark not defined.**

جدول (۲-۲) سرویس‌های ارائه شده توسط *ISP1* برای مشتریان. **Error! Bookmark not defined.**

جدول (۲-۳) جریان‌های مربوط به بسته‌های ورودی از پورت *X*. **Error! Bookmark not defined.**

جدول (۲-۴) فیلدهای سرآیند بسته که در OpenFlow پشتیبانی می‌شوند. **Error! Bookmark not defined.**

جدول (۲-۵) نسبت تعداد مقادیر یکتا به تعداد کل قوانین در مجموعه قوانین **Error! Bookmark not defined.**

جدول (۳-۱) نتایج به دست آمده در [۱۲] **Error! Bookmark not defined.**

جدول (۳-۲) گذردهی بخش برخط در MC-SBC **Error! Bookmark not defined.**

جدول (۴-۱) گذردهی بخش برخط در MC-SBC **Error! Bookmark not defined.**

جدول (۴-۲) کمینه، میانگین، و بیشینه تعداد قوانین در برگ‌ها برای مجموعه قوانین ACL. **Error! Bookmark not defined.**

جدول (۴-۳) کمینه، میانگین، و بیشینه تعداد قوانین در برگ‌ها برای مجموعه قوانینFW **Error! Bookmark not defined.**

جدول (۴-۴) کمینه، میانگین، و بیشینه تعداد قوانین در برگ‌ها برای مجموعه قوانین ...IPC **Error! Bookmark not defined.**

جدول (۴-۵) مشخصات پردازنده گرافیکی استفاده شده **Error! Bookmark not defined.**

فصل ۱: مقدمه

در سال های اخیر شاهد رشد بسیار زیادی در شبکه های اجتماعی بوده ایم و مطالعات زیادی روی این شبکه ها انجام گرفته است. داده های شبکه های اجتماعی یکی از ابزار محبوب برای مدل کردن رابطه و رفتار افراد و جامعه یا گروهی که در آن عضو هستن بشمار می رود. این داده ها معمولا به صورت گرافی نمایش داده می شود که در آن گره ها افراد و لبه ها روابط بین این افراد می باشد.

همچنین یادگیری ماشین مدت زیادی است که در علوم کامپیوتر جایگاه خود را پیدا کرده و به عنوان ابزار قدرتمندی برای کمک به انسان در زمینه های مختلف محسوب می شود و ماشین را بیش از پیش در خدمت انسان در آورده است.

در روش های اولیه یادگیری ماشین، عمدتا از داده ها و متن های خام در زمینه ی یادگیری استفاده می شد. اخیرا از طرف برخی شرکت ها و موسسات بزرگ، همچون گوگل، ای بی ام، مایکروسافت و ... پایگاه های دانشی معرفی شده است که انجام راه کارهای مختلف یادگیری ماشین را ساده تر و کاربردی تر کرده است. در این پایگاه های دانش اطلاعات مورد نیاز برای عملیات های مختلف یادگیری به صورت منظم و نیمه منظم موجود است و دغدغه ی نرمال کردن و رفع خطا و استخراج حقایق را به مقدار زیادی کم کرده است. این پایگاه های دانش عمدتا به صورت یک شبکه از موجودیت ها و روابط بین آن ها که می توان آن را به صورت یک گراف داده نمایش داده به این صورت که گره ها موجودیت ها و یال های بین گره ها نشانگر روابط بین آن ها باشند، که این یال ها می توانند از یک نوع باشند یعنی گراف نشانگر یه شبکه تک-رابطه ای باشد (همگون) یا هر یال با یک برچسب، نشانگر نوع رابطه باشد و شبکه نشانگر یک گراف چند رابطه ای (ناهمگون) باشد.

یکی از مسائل بنیادی و اساسی در یادگیری ماشین روی شبکه های اجتماعی، پیش بینی پیوند

در شبکه های همگون و ناهمگون است به این معنی که از اطلاعات موجود در گراف دانش استفاده کرده و وجود یا عدم وجود یک یال را پیش بینی کرد، یا به عبارت دیگر مساله ی پیش بینی پیوند این است که نمایی از یک شبکه به ما داده می شود و ما مایل هستیم که بدانیم در آینده نزدیک، احتمالا چه تراکنش هایی میان اعضاء فعلی شبکه روی می دهد و یا اینکه کدام یک از تراکنش های موجود را از دست می دهیم.

این راهکار در زمینه های مختلف یادگیری ماشین مورد استفاده قرار می گیرد و کاربرد گسترده ای در زندگی انسان پیدا کرده است. برای مثال از این راهکار در سیستم های توصیه گر در فروشگاه های اینترنتی، سیستم های تشخیص پزشکی، جواب گویی به سوال و ... استفاده می شود. اکثر مطالعات انجام شده در این زمینه روی شبکه های تک رابطه ای بوده است. به این معنا که روابط بین موجودیت ها از یک نوع است و این روابط به صورت دوتایی های مرتب مورد استفاده قرار می گیرند. برای مثال اگر در یک شبکه اجتماعی رابطه را دوستی بین افراد در نظر بگیریم یال های گراف شبکه به صورت «الف، ب» خواهد بود به این معنی که شخص الف با شخص ب رابطه دوستی دارد.

پایگاه های دانشی که اخیرا معرفی شده اند عمدتا داده ها را به صورت داده های چند رابطه ای ذخیره می کنند و اطلاعات بیشتری از یک رابطه دوتایی بلی یا خیر به ما می دهند. منظور از داده های چند رابطه ای گراف جهت داری است متشکل از موجودیت ها و روابط بین آنها که بصورت «مبدا h، رابطه r، مقصد t» نمایش داده می شود، به این معنی که یک رابطه r بین موجودیت های h و t وجود دارد. برای مثال سه تایی «تهران، واقع در، ایران» این اطلاع را به ما می دهد که استان تهران داخل کشور ایران قرار دارد. در این نوع پایگاه دانش هم انواع مختلف موجودیت وجود دارد و هم انواع مختلف رابطه بین موجودیت ها. پایگاه های دانشی مانند Graph Google Knowledge، FreeBase و WordNet وجود دارند که شامل تعداد زیادی نمونه چندرابطه ای می باشند و تعداد زیادی موجودیت و روابط بین آنها را می توان در آنها یافت و از آن برای یادگیری مدل استفاده کرد.

[۵].

۱-۱- تعریف مسئله

روش های مختلفی برای حل مساله ی پیش بینی پیوند در پایگاه های دانش ناهمگون ارائه شده است که از رویکردهای مختلفی سعی به حل این مساله می کنند از جمله روش های آماری، روش های ویژگی های پنهان، روش های ویژگی های گراف و ... تمرکز ما در این مقاله بر روی روش های پیش بینی پیوند مبتنی بر ویژگی های پنهان خواهد بود.

در روش های مبتنی بر ویژگی های پنهان، با استفاده از ویژگی هایی که در موجودیت ها و روابط بین آن ها وجود دارد سعی می شود میزان ارتباط بین موجودیت ها را تشخیص دهیم و به این صورت وجود یک پیوند را تایید یا رد کنیم. برای مثال اگر دو شخص با هم همکار هستند، به احتمال زیادی ویژگی های مشترکی دارند، مثلاً هر دو اهل یک شهر هستند، هر دو در یک رشته ی دانشگاهی تحصیل کرده اند، خصوصیات اخلاقی یکسانی دارند و ... در رابطه ی همکار بودن به هیچ یک از این ویژگی ها به طور مستقیم اشاره نشده است و این ویژگی ها به طور ضمنی در این موجودیت ها قرار دارند که با استفاده از آن می توانیم وجود یا عدم وجود رابطه ی همکار بودن را حدس بزنیم. پس هر موجودیت می تواند تعداد زیادی ویژگی پنهان داشته باشد که رابطه ها به این ویژگی ها وزن می دهند، مثلاً در رابطه ی همکاری احتمال اینکه رشته ی تحصیلی دو شخص در همکار شدن آن ها تاثیر گذار باشد بیشتر از ویژگی رنگ پوست دو شخص است، پس وزن ویژگی رشته ی تحصیلی در این رابطه بیشتر از ویژگی رنگ پوست است.

روش های زیادی برای حل مساله ی پیش بینی پیوند مبتنی بر ویژگی های پنهان ارائه شده است. همه ی این روش ها از یک روال ثابتی برای حل مساله استفاده می کنند و هر کدام با نوآوری هایی که داشته اند بهبودهایی در نتایج بدست آمده حاصل کرده اند. این روش ها در بخش + به طور کامل معرفی خواهند شد.

در این پژوهش قصد داریم که با استفاده از قوانین انجمنی موجود در پایگاه های دانش به دو هدف برسیم،

۱. بهبود نتایج در روش های معرفی شده ی موجود: در ادامه پنج روش از روش ها معروف

مبتنی بر ویژگی های پنهان که به حل مساله ی پیش بینی پیوند پرداخته اند را انتخاب

می کنیم و در بخش + نشان می دهیم که استفاده از قوانین انجمنی در فرایند آموزش

این پنج روش باعث بهبود در نتایج این روش ها خواهد شد.

۲. ارائه ی چارچوبی برای سنجش نقاط ضعف و قوت روش های ارائه شده: همچنین با دسته بندی کردن قوانین انجمنی در بخش + نشان می دهیم که با اعمال جداگانه ی دسته های مختلف قوانین انجمنی می توانیم میزان قدرت و ضعف روش ها را در قانون های مختلف بررسی کنیم که با برطرف کردن نقاط ضعف می توان روش های قوی تری ارائه کرد.

۱-۲- روش انجام پژوهش

برای دستیابی به اهدافی که در بخش قبل مطرح شد، همانطور که اشاره شد از قوانین انجمنی استفاده خواهیم کرد. در روش های معرفی شده فقط از روابط موجود در پایگاه دانش در فرایند آموزش استفاده شده است. مثلاً اگر در پایگاه دانش داشته باشیم که «علی، پدر، حسن» به این معنی که علی پدر حسن است می توانیم به این نتیجه برسیم که «حسن، فرزند، علی» و از آن در فرایند آموزش استفاده کنیم اما در این روش ها این روابط مورد استفاده قرار نگرفته اند. در این پژوهش قصد داریم که اینگونه قوانین را به کمک روش هایی که در بخش + توضیح داده می شوند از پایگاه دانش استخراج کنیم و به کمک معیارهای کیفیت که در بخش + معرفی می شوند قوانینی که کیفیت مناسب دارند را انتخاب کنیم و در فرایند آموزش از آنها استفاده کنیم، نحوه ی استفاده از این قوانین در مدل های موجود در بخش + توضیح داده شده است.

۱-۳- کاربردهای پژوهش

مسائل پیش بینی پیوند در زمینه های زیادی قابل بکارگیری هستند و همین امر باعث شده است که در سال های گذشته بسیار مورد توجه و تحقیق قرار بگیرند. به عنوان مثال در ادامه تعدادی از این کاربردها را مطرح خواهیم کرد:

- پیش بینی پیوندهای احتمالی در شبکه های اجتماعی، برای مثال در شبکه های اجتماعی بین کاربران و مطالب ثبت شده، نظرات و ... بررسی شود که روابطی مانند دوستی، پسندیدن و

- نپسندیدن، روابط فامیلی و... وجود دارد یا خیر.
 - استفاده یه عنوان سیستم های توصیه گر، برای مثال کاربران و کالاها یا اشیاء موجودیت ها هستند و خریدن، امتیاز دادن، بررسی کردن و... رابطه ها هستند که می توان از بررسی این روابط و موجودیت ها اطلاعاتی نظیر کاربر کالای ۱۲ را خواهد خرید یا خیر یا اینکه کالای ۱ به تعداد بالا فروش خواهد رفت یا خیر بدست آورد.
 - استفاده در سیستم های تشخیص پزشکی، برای مثال شبکه ی بین بیماران، بیماری ها، داروها و... را در نظر بگیرید، با بررسی دقیق این شبکه می توان علل و درمان های بیماری ها را به کمک ماشین بدست آورد.
- علاوه بر کاربردهایی که در بالا برای پیش بینی پیوند اشاره شد، از روش ارائه شده در این پژوهش می توان برای کامل تر کردن پایگاه های داده نیز استفاده کرد، به این صورت که روابطی که از قوانین انجمنی استخراج شده از پایگاه دانش بدست می آید و کیفیت لازم را دارد را نیز به پایگاه دانش اضافه کنیم و آن را کامل تر کنیم.

۴-۱- ساختار پایان نامه

ادامه این پایان نامه بدین صورت سازمان دهی شده است: در فصل.

فصل ۲: زمینه‌های تحقیق

در این فصل مطالبی راجع به مباحث پایه‌ای که در ادامه تحقیق از آن‌ها استفاده شده است را مطرح خواهیم کرد.

۲-۱- شبکه‌های همگون و ناهمگون

در سال‌های اخیر شبکه‌های اجتماعی پی‌شرفت زیادی داشته است و در زمینه‌های مختلفی شبکه سازی شده است. عمدتاً در این شبکه‌ها روابط خاصی مد نظر و قابل استخراج است، مثل روابط دوستی، همکاری و ... که اگر گراف این گونه شبکه‌ها را رسم کنیم در این گراف‌ها یال‌ها از یک نوع است و نشانگر یک رابطه‌ی خاص است. مثلاً گرافی هست که همه‌ی گره‌های آن انسان‌ها هستند و یال‌های بین گره‌ها نشان‌دهنده‌ی وجود یا عدم وجود دوستی بین اشخاص است. این گونه شبکه‌ها که در آن‌ها یال و گره‌ها از یک نوع است را شبکه‌های همگون می‌نامیم [X]. در شکل + یک نمونه شبکه‌ی همگون که در آن اشخاص و رابطه‌ی دوستی بین آن‌ها به تصویر کشیده شده است را مشاهده می‌کنیم.

عکس شبکه‌ی معمولی

اما همانطور که در مقدمه نیز اشاره شد، اخیراً شبکه‌های بزرگتر و پیچیده‌تری معرفی شده است که فقط یک نوع رابطه را پوشش نمی‌دهد و روابط مختلفی را بین موجودیت‌های مختلف در بر می‌گیرد، به این شبکه‌ها به علت یکسان نبودن نوع روابط و موجودیت‌ها شبکه‌های ناهمگون می‌گوییم، همچنین

به خاطر وجود دانشی که در این شبکه‌ها نهفته و قابل استخراج است، آن را پایگاه‌دانش^۱ نیز می‌نامیم. در شکل + یک قسمت کوچک از یک شبکه‌ی ناهمگون را مشاهده می‌کنیم که موجودیت‌های آن از دو نوع انسان و شهر هستند و روابط موجود در این شبکه از دو نوع «والد بودن» و «متولد شهر ... بودن» است.

عکس شبکه‌ی ناهمگون

۲-۲- پیش‌بینی پیوند

یکی از روش‌های یادگیری ماشین که در زمینه‌های مختلف به کمک انسان آمده است و کارهای انسانی را تسهیل کرده است پیش‌بینی پیوند است. در پیش‌بینی پیوند یک گراف از روابط بین موجودیت‌ها به عنوان ورودی مساله دریافت می‌کنیم و وجود یا عدم وجود یک یال بین دو موجودیت را پیش‌بینی می‌کنیم. در شکل + یک گراف از یک شبکه‌ی ناهمگون را مشاهده می‌کنیم. مساله‌ی پیش‌بینی لینک تلاش می‌کند که بررسی کند که رابطه‌ی + بین گره‌های + و + قرار دارد یا نه. این پیش‌بینی عمدتاً از روی دیگر روابط مرتبط بین موجودیت‌ها انجام می‌شود و با بررسی شباهت‌ها و معیارهایی که در ادامه‌ی این پژوهش توضیح خواهیم داد تصمیم بگیرد که پیوند رابطه‌ی + بین موجودیت‌های + و + برقرار هست یا خیر.

شکل گراف با سوال

راه‌کارهای مختلفی برای حل این مساله مطرح شده است که به طور کلی می‌توان این راه‌کارها را به سه دسته‌ی ۱- روش‌های مبتنی بر ویژگی‌های گراف ۲- روش‌های مدل تصادفی مارکوف ۳- روش‌های مبتنی بر ویژگی‌های پنهان تقسیم کرد که در این پژوهش تمرکز ما روی دسته‌ی سوم یعنی

¹ Knowledgebase

روش‌های مبتنی بر ویژگی‌های پنهان خواهد بود و در بخش + این روش‌ها را به تفصیل توضیح خواهیم داد.

۳-۲- فرض جهان باز^۱ و فرض جهان بسته^۲

فرض‌های جهان باز و جهان بسته در سیستم رسمی منطق^۳ برای بازنمایی دانش^۴ مورد استفاده قرار می‌گیرد. در فرض جهان بسته در نظر می‌گیریم که اگر داده‌ای در پایگاه دانش نبود، آن داده را غلط فرض می‌کنیم [x]. برای مثال اگر در پایگاه دانش حقیقت «x, r, y» که به معنی این است که x با y رابطه‌ی r را دارد وجود نداشته باشد، می‌توانیم در نظر بگیریم که این حقیقت اشتباه است و مطمئنیم که x با y رابطه‌ی r ندارد.

اما در فرض جهان باز اینگونه نیست و اگر حقیقتی در پایگاه دانش وجود نداشته باشد نمی‌توانیم با اطمینان بگوییم که آن حقیقت اشتباه است، و ممکن است صحیح باشد [x]. وجود این فرض از کامل نبودن پایگاه‌های دانش ناشی می‌شود که وقتی نمی‌توانیم همه‌ی اطلاعات موجود در مورد موضوع مربوط به پایگاه دانش را جمع‌آوری و در آن قرار دهیم پس نمی‌توانیم در مورد حقیقت‌هایی که در پایگاه دانش نیست نظری بدهیم و این حقیقت را ناشناس^۵ در نظر می‌گیریم.

دسته‌ی دیگر فرض جهان بسته‌ی محلی است، که این فرض میانه‌ای بین دو فرض قبل است...

¹ Open world assumption (OWA)

² Closed world assumption (CWA)

³ Formal system of logic

⁴ Knowledge representation

⁵ unknown

۴-۲- گرادینان نزولی تصادفی^۱ (SGD)

¹ Stochastic gradient descent

فصل ۳: پیشینه پژوهش

در این بخش به معرفی روش های موجود که به حل مساله ی پیش بینی پیوند در شبکه های ناهمگون می پردازند، خواهیم پرداخت. روش هایی با راه کارهای مختلفی در حوزه های مختلف تلاش به حل این مساله کرده اند. این روش ها می توانند به سه دسته ی مختلف تقسیم می شوند: ۱- روش های مبتنی بر ویژگی های گراف ۲- روش های مدل تصادفی مارکوف ۳- روش های مبتنی بر ویژگی های پنهان. در ادامه توضیح مختصری در مورد هر دسته از روش ها می دهیم.

- در روش های مبتنی بر ویژگی های گراف از روی ویژگی های ساختاری گراف داده ها استفاده می شود مانند دسته بندی گره ها، دسته بندی نوع یال ها، تعداد گره های مشابه و ...
- روش های مدل تصادفی مارکف که در آن دید بالایی از گراف داده نداشته و سعی می کنیم مساله را به صورت محلی حل کنیم به این صورت که روابط هر موجودیت را با موجودیت های اطراف آن بررسی می کنیم و جواب های محلی را بدست می آوریم.
- روش های مبتنی بر ویژگی های پنهان که در این روش ها هر موجودیت و نوع رابطه بین آنها به صورت برداری از ویژگی های پنهان تعریف می شود که ویژگی های پنهان نام دارد. برای مثال ویژگی هایی که یک موجودیت می تواند داشته باشد، محل به وجود آمدن آن، سن آن، جاندار یا بی جان بودن آن و ... است.

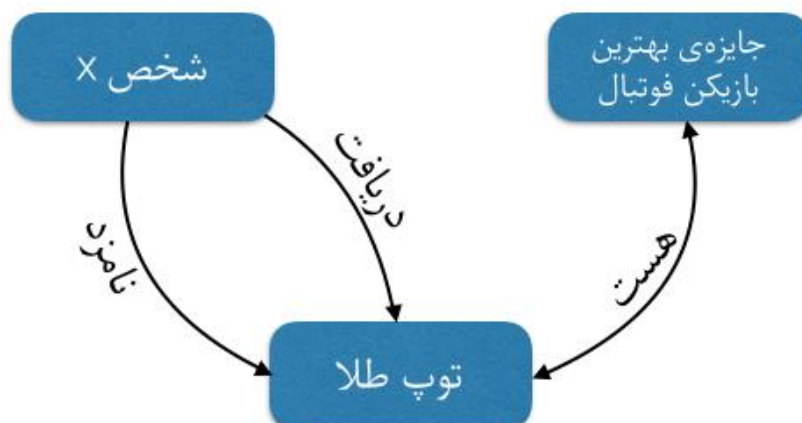
همان طور که در بخش XX گفتیم تمرکز ما در این تحقیق روی روش های مبتنی بر ویژگی های پنهان است، در ادامه این روش ها را به صورت کامل توضیح داده و تعدادی از آنها را به طور مختصر معرفی کرده و نحوه ی کار آنها را توضیح می دهیم و نتایج بدست آمده از آنها را توضیح می دهیم.

۱-۳- روش های مبتنی بر ویژگی های پنهان

روش های مبتنی بر ویژگی های پنهان از جدید ترین راه کارهایی است که در حوزه ی پیش بینی پیوند روی پایگاه های دانش موجود معرفی شده است. همانطور که در بخش قبل گفتیم این روش از ویژگی هایی که در نگاه اول از پایگاه دانش برداشت نمی شود استفاده می کند که ویژگی های پنهان نام دارند، در همه ی روش های مبتنی بر ویژگی پنهان این ویژگی ها را به صورت برداری تعریف می کنیم که هر مولفه از این بردار نشان دهنده ی یک ویژگی می باشد.

برای مثال در رابطه ی دوستی فاکتورهایی تاثیر گذار هستند و اگر در پایگاه دانش همه ی رابطه های دوستی موجود را بررسی کنیم به یک بردار از فاکتورها می رسیم که بردار رابطه ی دوستی را تشکیل می دهد. مثلاً فاکتورهایی مانند شهر محل زندگی، سن، دانشگاه، رشته ی دانشگاهی، جنسیت، مذهب و ... در شکل گیری رابطه ی دوستی می تواند موثر باشد، اما اینکه هر کدام از این روابط چقدر در ایجاد رابطه ی دوستی تاثیر دارند و اهمیت هر کدام چقدر است و این میزان اهمیت را چگونه در تشخیص این رابطه تاثیر دهیم به مدل یادگیری بستگی دارد که در ادامه در معرفی هر یک از روش ها به صورت کامل توضیح داده خواهد شد.

برای مثال برای اینکه بررسی کنیم شخص X بازیکن فوتبال خوبی هست یا خیر از دیگر روابط موجود استفاده می کنیم و میزان ارتباط این شخص را با معیارهای بازیکن خوب فوتبال بودن بررسی می کنیم. در مثال شکل - شخص x هم برای توپ طلا نامزد شده و هم توپ طلا را دریافت کرده و از اطلاعات دیگر پایگاه دانش می دانیم که توپ طلا به بهترین بازیکن فوتبال هر سال داده می شود. پس شخصی که این جایزه را دریافت کرده بازیکن فوتبال خوبی است.



روش های زیادی برای مدل های مبتنی بر ویژگی های پنهان معرفی شده است، روش هایی که در این تحقیق مورد مطالعه و بررسی قرار گرفته این به صورت زیر دسته بندی می شوند.

- روش های خطی
- روش های دو خطی
 - RESCAL
- روش های ادراک چندلایه ای
 - Neural Tensor Network (NTN)
- روش های فاصله ی پنهان
 - Structured Embedding
 - Translating Embedding (TransE)
 - Translating on Hyperplane (TransH)

در ادامه ابتدا پایگاه دانش Freebase که مطالعات روی آن انجام می گیرد و نحوه ی تقسیم بندی آن را توضیح خواهیم داد. سپس روش کلی آموزش مدل های مبتنی بر ویژگی های پنهان را شرح می دهیم و در ادامه روش هایی که در بالا نام برده شدند را توضیح داده و نقاط قوت و ضعف آن ها را بررسی می کنیم و نتایج بدست آمده از هر یک از این روش ها را نیز بررسی خواهیم کرد و در فصل بعد تلاش بر بهبود این روش ها خواهیم کرد.

۲-۳- تقسیم بندی داده های پایگاه دانش

برای آموزش دادن مدل های نام برده شده در بخش قبل از پایگاه دانش Freebase که قسمتی از گراف دانش تولید شده توسط گوگل است استفاده می کنیم. در این پایگاه دانش ۸۰ میلیون موجودیت، ۲۰ هزار نوع رابطه مختلف و ۱,۲ میلیارد حقیقت وجود دارد. حقیقت هایی که در این پایگاه دانش وجود دارد به صورت سه تایی مرتب

(subject, predicate, object)

مشخص شده اند به این صورت که موجودیت subject رابطه ی predicate دارد با موجودیت object. به این نوع ذخیره ی داده اصطلاحاً SPO گفته می شود. برای مثال یک نمونه حقیقت موجود در این پایگاه دانش به صورت:

(Barack Obama, place_of_birth, Hawaii)

است که بیانگر حقیقت «باراک اوباما متولد هاوایی است» می باشد.

این پایگاه دانش شامل تعداد زیادی حقیقت است که عملیات آموزش روی آن هزینه ی زیادی از نظر زمان و منابع خواهد داشت. برای سادگی و تسریع کار از یک نمونه ی نرمال کوچک این پایگاه دانش به نام FB15k استفاده می کنیم که در آن ۱۴۹۵۱ موجودیت، ۱۳۴۵ رابطه ی مختلف و ۵۹۲۲۱۳ حقیقت وجود دارد.

	Entities (n_e)	Rel. (n_r)	Train. Ex.	Valid. Ex.	Test Ex.
FB13	75,043	13	316,232	5,908	23,733
FB15k	14,951	1,345	483,142	50,000	59,071
FB1M	1×10^6	23,382	17.5×10^6	50,000	177,404

روش های مبتنی بر ویژگی های پنهان به صورت تکراری انجام می شوند و نیاز است در هر تکرار بررسی کنیم که به آستانه ی مناسب برای قطع تکرار الگوریتم رسیده ایم یا خیر. همچنین پس از انجام عملیات آموزش نیاز است که مدل آموزش دیده را آزمایش کنیم و میزان دقت آن را بدست آوریم.

برای انجام عملیات آموزش، بررسی کیفیت آموزش در هر مرحله و بررسی کیفیت کلی آموزش به سه دسته مختلف از داده نیاز داریم.

- داده های آموزش: قسمت عمده ی داده ها برای عملیات آموزش استفاده می شود، که الگوریتم اصلی هر روش روی آن اعمال می شود و سعی می کنیم پارامترهایی که همان ویژگی های پنهان هستند را تنظیم کنیم.
- داده های ارزیابی: قسمتی از داده که داده های ارزیابی نام دارند برای بررسی میزان بهبود یا تخریب مدل در هر تکرار استفاده می شوند، این داده ها کاملاً از داده های آموزش جدا هستند و در مرحله ی آموزش اصلاً به مدل نشان داده نمی شود و مدل تحت تاثیر این داده ها قرار نمی گیرد. در انتهای هر مرحله به کمک این داده ها بررسی می شود که تغییراتی که روی پارامترهای این مدل انجام شده باعث بهتر یا بدتر شدن این مدل شده است، در صورتی که بهبودی مشاهده شود تغییرات انجام شده در این مرحله نگه داشته می شود و به سراغ مراحل بعدی می رویم اما اگر نتایج بدتر شده باشد پارامترها را به مقادیر قبلی برگردانده و مرحله ی بعدی را شروع می کنیم.
- داده های آزمون: قسمت دیگری از داده که داده های آزمون نام دارند برای بررسی کیفیت عملکرد کلی مدل به کار می روند. این داده ها نیز کاملاً از داده های آموزش تفکیک شده اند و در زمان آموزش روی مدل تاثیری نمی گذارند و کلاً در هیچ یک از مراحل آموزش استفاده نمی شوند و تنها پس از آموزش مدل استفاده می شوند تا کیفیت مدل آموزش دیده شده را بررسی کنند.

در این تحقیق پایگاه دانش Freebase15k به صورت شکل - تقسیم و استفاده شده است، به این صورت که برای قسمت آموزش 483,142 حقیقت استفاده کرده ایم، برای قسمت ارزیابی 50k حقیقت استفاده کرده ایم و برای قسمت آزمون 59k حقیقت استفاده کرده ایم.



۳-۳- مراحل آموزش روش های مبتنی بر ویژگی های پنهان

همانطور که در بخش قبل گفته شد پایگاه دانش را به سه قسمت آموزش، ارزیابی و آزمون تقسیم می کنیم. نحوه ی آموزش کلی همه ی روش های مبتنی بر ویژگی های پنهان در سه مرحله انجام می شود:

۱- مرحله ی آموزش ۱- مرحله ی ارزیابی ۳- مرحله ی آزمون. مراحل آموزش و ارزیابی به صورت تکراری و معمولاً با تکرار بالا انجام می شوند و در هر تکرار بررسی می شود که بهبودی اتفاق افتاده است یا خیر، اگر بهبودی داشتیم نتایج این مرحله تکرار را نگه داشته و مرحله ی تکرار بعدی را شروع می کنیم و اگر بهبودی اتفاق نیافتاده بود نتایج این مرحله را تاثیر نمی دهیم و مرحله ی تکرار بعدی را شروع می کنیم.

در ادامه این سه مرحله را توضیح می دهیم.

۱-۳-۳- مرحله ی آموزش:

همانطور که قبل تر تو ضیح داده شد آموزش به صورت تکراری انجام می شود و در هر تکرار الگوریتم آموزش روش مورد نظر روی قسمتی یا همه ی داده های مجموعه ی آموزش انجام می شود. در این قسمت سعی می شود که پارامترهای همه ی ویژگی ها جهت دهی شده و آموزش ببینند تا کمترین خطا در پاسخ به سوالاتی که از مدل پرسیده می شود را داشته باشند.

۲-۳-۳- مرحله ی ارزیابی:

پس از هر تکرار مدل آموزش دیده شده را روی داده های ارزیابی اجرا می کنیم و معیارهایی که برای بررسی کیفیت مدل در نظر گرفته ایم را بدست می آوریم و از روی آن میزان بهبود مدل در این تکرار را بررسی می کنیم.

برای مثال فرض کنید که حقیقت زیر در مجموعه داده های ارزیابی وجود دارد و در مرحله ی آموزش مدل این حقیقت مشاهده نشده است:

(WALL-E, has_genre, Fantasy)

این حقیقت به این معناست که «ژانر فیلم WALL-E فانتزی است». در مرحله ی ارزیابی موجودیت اول یا آخر این حقیقت را حذف می کنیم و قسمت حذف شده را از مدل سوال می پرسیم و انتظار داریم که قسمت حذف شده را حدس بزند. سوالی که از این مدل پرسیده می شود به این صورت است:

(WALL-E, has_genre, ?)

به این معنی که «ژانر فیلم وال-ای چیست؟».

در مدل های مبتنی بر ویژگی های پنهان پاسخ به اینگونه سوالات به صورت مجموعه ای مرتب شده ی همه موجودیت هاست. یعنی این مدل میزان نزدیکی همه ی موجودیت ها با پاسخ این سوال را بررسی می کند و به ترتیب نزدیک به دور مرتب می کند و به عنوان پاسخ به ما می دهد. برای مثال پاسخ به سوال بالا به صورت زیر می باشد:

(WALL-E, has_genre, Fantasy)

WALL-E has the genre?!	1- Animations
(WALL-E, has_genre, ?)	2- Computer Animation
	3- Comedy film
	4- Adventure film
	5- Science Fiction
	6- Fantasy
	7- Stop motion
	8- Satire
	...

همانطور که می بینیم مدلی که این سوال از آن پرسیده شده پاسخ درست را در ششمین حدس به ما داده است.

در مرحله ی ارزیابی نیاز به معیارهایی داریم که بررسی کنیم که مدل با توجه به این معیارها بهبود داشته یا خیر. در این تحقیق از دو معیار زیر استفاده شده است:

- رتبه ی میانگین (Mean rank): میانگین رتبه ی جواب های درستی که مدل داده است.
- حدس زیر ۱۰: درصد سوالاتی که پاسخ درست مدل به آن زیر رتبه ی ۱۰ بوده است که در ادامه به آن hit@10 می گوئیم.

همانطور که قبل تر توضیح داده شد ۵۰ هزار حقیقت در دسته ی ارزیابی وجود دارد، ما پس از هر تکرار الگوریتم دو معیار بالا را به دست آورده و میزان بهبود مدل را اندازه می گیریم. پس از پرسیدن این ۵۰ هزار سوال میانگین رتبه ای که جواب های درست داشته معیار اول را به ما می دهد و درصد سوالاتی که جواب درست آن زیر رتبه ی ۱۰ بوده است معیار دوم را به ما می دهد.

۳-۳-۳- مرحله ی آزمون:

پس از انجام کامل مراحل آموزش و ارزیابی و متوقف شدند الگوریتم از داده ای آزمون که در بخش – توضیح دادیم و معیارهای کیفیت که در بخش – توضیح دادیم استفاده می کنیم و کیفیت مدل را بررسی می کنیم. در این بخش هم مانند بخش ارزیابی داده ها به مدل در حال آموزش نشان داده نشده و برای مدل جدید هستند. در این مرحله هم یک قسمت از هر حقیقت موجود در داده های آزمون را حذف کرده و آن را از مدل سوال می پرسیم، دقیقاً به مانند مرحله ی ارزیابی. پس از پرسیدن سوالات دو معیار رتبه ی میانگین و $hit@10$ را بدست می آوریم که این دو معیار نشان دهنده ی میزان کیفیت و دقت روش است.

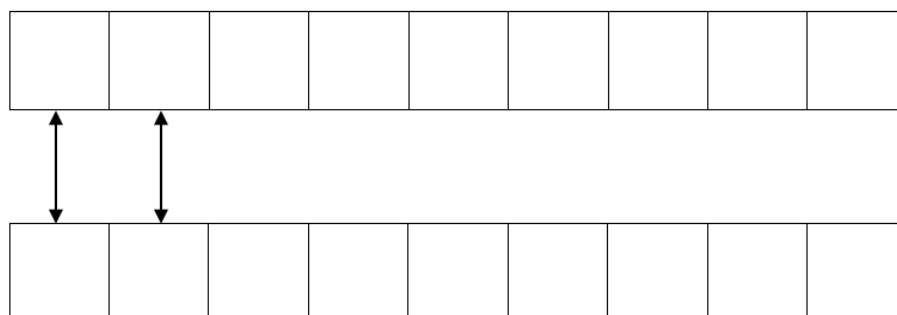
۳-۴- بررسی روش های موجود

در ادامه نحوه ی عملکرد روش های موجود را توضیح داده و بررسی می کنیم.

۳-۴-۱- روش های خطی

همانطور که در بخش – گفتیم در روش های مبتنی بر ویژگی های پنهان موجودیت ها و رابطه ها به صورت بردارهایی در یک فضای n بعدی تبدیل می شوند که به کمک معیارهای مختلف شباهت،

رابطه ی بین دو موجودیت را بدست می آوریم. در روش های خطی در زمان آموزش و بررسی میزان شباهت بردارها را به صورت خطی با یکدیگر مقایسه می کنیم، به این صورت که هر اندیس از بردار موجودیت اول را با اندیس متناظر آن در بردار رابطه یا موجودیت دیگر بررسی می کنیم، شکل +.

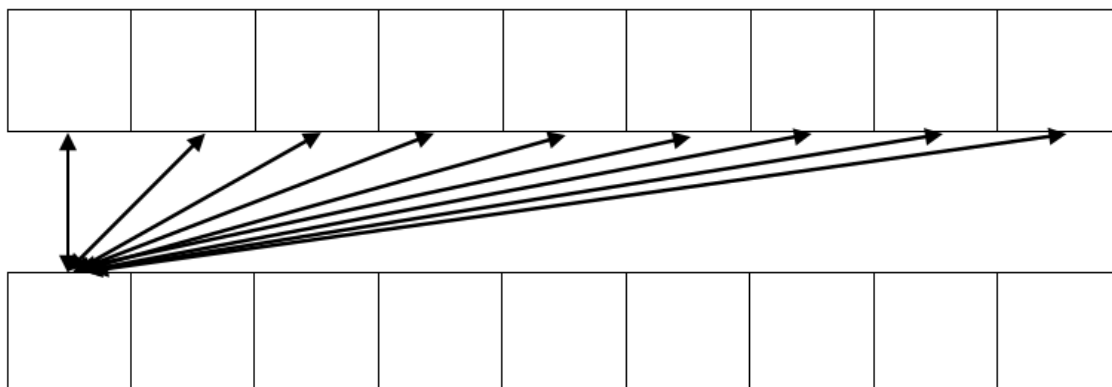


با توجه به نتایج ضعیفی که روش های این دسته در آموزش مدل و پیش بینی پیوند بدست آورده اند به این روش ها نمی پردازیم و به همین معرفی کفایت می کنیم.

۲-۴-۳- روش های دوخطی - رسکال

روش RESCAL [۱] یک روش ویژگی های پنهان رابطه ای^۱ است که حقیقت های پایگاه دانش را به صورت تراکنش های بین جفت ویژگی های پنهان در نظر می گیرد. یعنی بر خلاف روش های خطی هر ویژگی پنهان از هر بردار را با همه ی ویژگی های پنهان دیگر بردار بررسی می کند، به صورت شکل +. به همین دلیل این روش را روش دوخطی نیز می نامیم.

¹ Relational latent factor



در این روش امتیاز هر سه تایی را از رابطه ی + بدست می آوریم که در آن H_e تعداد ابعاد بردار در نظر گرفته شده برای موجودیت ها است (تعداد ویژگی های پنهان هر موجودیت). W_k یک بردار وزن با ابعاد $H_e * H_e$ است که هر اندیس w_{abk} نشانگر این است که ویژگی پنهان a و b در رابطه ی k چقد با هم کنش می کنند.

$$f_{ijk}^{\text{RESCAL}} := \mathbf{e}_i^T \mathbf{W}_k \mathbf{e}_j = \sum_{a=1}^{H_e} \sum_{b=1}^{H_e} w_{abk} e_{ia} e_{jb}$$

همانطور که در رابطه ی + مشاهده می شود در این روش هر رابطه به یک ماتریس تبدیل شده و هر موجودیت به یک بردار، اگر حقیقت (a,k,b) به این معنی که موجودیت a رابطه ی k را با b دارد را در نظر بگیریم و بخواهیم بررسی کنیم که این رابطه برقرار هست یا خیر، احتمال وجود این رابطه را از روی امتیازی که تابع امتیاز + به ما می دهد بدست می آوریم. این امتیاز به این صورت محاسبه می شود که بردار موجودیت a در ماتریس مربوط به k ضرب شده و پس از آن در بردار b ضرب می شود که نتیجه ی آن یک مقدار حقیقی است که امتیاز این حقیقت را به ما می دهد.

در ادامه به برخی از نکات مورد توجه این مدل به صورت موردی اشاره می کنیم.

آموزش رابطه ای^۱ از طریق بازنمایی مشترک^۲: در رابطه ی + هر موجودیت به صورت یک بردار بازنمایی شده است بدون توجه به اینکه در قسمت اول حقیقت می آید یا قسمت دوم آن. همچنین این موجودیت ها برای همه ی رابطه ها یک بازنمایی مشترک دارند و در هر رابطه نیاز به تعریف جدید ندارند. برای مثال موجودیت i در قسمت اول حقیقت x_{ijk} با رابطه ی k آمده است و همین موجودیت در رابطه ی x_{piq} به عنوان موجودیت دوم در رابطه ی q ظاهر شده است. هر دو تابع امتیاز $f_{ijk} = e_i W_k e_j$ و $f_{piq} = e_p W_q e_i$ از یک بازنمایی برای موجودیت i (e_i) استفاده می کنند. بنابراین همه ی پارامترها به صورت مشترک^۳ آموزش دیده می شوند و این بازنمایی مشترک باعث می شود که اطلاعات روی همه ی حقیقت ها به وسیله ی بازنمایی موجودیت ها و ماتریس وزن دار رابطه ها پخش شوند و بتوانیم وابستگی های جهانی^۴ در داده ها را تشخیص دهیم [1].

ارتباط معنایی بردارها^۵: خاصیت بازنمایی مشترک در این روش کمک می کند که میزان شباهت موجودیت ها در فضای رابطه ای^۶ نیز بدست بیاید. برای مثال موجودیت هایی که با رابطه های مشابه به موجودیت های مشابه متصل هستند به یکدیگر شبیه هستند. به عنوان نمونه اگر بازنمایی e_i و e_p شبیه به هم باشد، تابع امتیاز f_{ijk} و f_{piq} باید مقادیر نزدیک به هم داشته باشند پس موجودیت ها با تعداد زیادی

¹ Relational Learning

² Shared representations

³ jointly

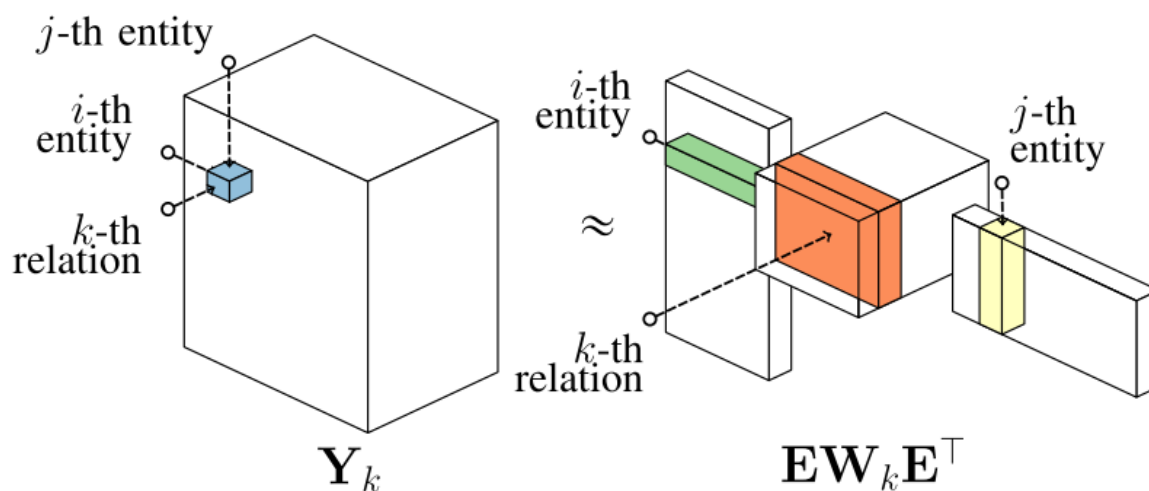
⁴ global dependency

⁵ Semantic embeddings

⁶ relational domain

رابطه ای مشترک بازنمایی یکسانی خواهند داشت. این خصوصیت می تواند در بازنمایی^۱ و خوشه بندی های با مقیاس بالای^۲ موجودیت ها روی داده های رابطه ای مورد استفاده قرار گیرد. [63,64,1]

ارتباط با عامل بندی تنسور^۳: این روش شباهت زیادی به روش های استفاده شده در سیستم های توصیه گر [66] و عامل بندی تنسور سنتی دارد [۶۷]. ضرب ماتریس که در معادله ی + معرفی کردیم می تواند به صورت $F_k = EW_kE^T$ نوشته شود که در آن $F \in R^{Ne*Ne}$ ماتریسی است که همه ی امتیازات مربوط به رابطه ی k را در خود جا داده است و سطر i ام از ماتریس $E \in R^{Ne*He}$ بازنمایی موجودیت e_i است که برداری از ویژگی های پنهان این موجودیت است. در شکل + این تبدیل ماتریس به تنسور نمایش داده شده است. [1]



¹ resolution

² large-scale hierarchical clustering

³ Tensor factorization

برازش مدل: اگر بخواهیم از یک مدل آماری استفاده کنیم، پارامترهای این روش می توانند به صورت یک مدل کمینه سازی بر مبنای گرادیان¹ تخمین زده شوند مانند گرادیان نزولی تصادفی (SGD)² [68]. این روش می توانیم به صورت یک روش بر مبنای امتیاز حل کرد که می تواند از پارامترهای مدل را به صورت بسیار بهینه ای تخمین بزند: با توجه به ساختار تنسور که در بالا توضیح دادیم و همچنین با توجه به تنک بودن داده های موجود، نشان می دهیم که روش رسکال می تواند به کمک توالی³ ای از بروزرسانی های بسته ی کارآمد⁴ محاسبه شود. [63,64] به صورت تحلیلی می توان نشان داد که در این راهکار با هر بروزرسانی در E و W_k به صورت خطی با تعداد موجودیت ها N_e ، تعداد رابطه ها N_r و تعداد حقیقت های مشاهده شده توسط مدل رشد می کند. [64]

پیش بینی مجزا: در رابطه ی + احتمال وجود یک رابطه از روی یک ضرب ماتریسی ساده از مرتبه ی $O(H_e^2)$ بدست می آید. بنابراین، زمانی که پارامترهای مدل تخمین زده شدند، پیچیدگی محاسبات برای پیش بینی امتیاز یک حقیقت فقط به تعداد ویژگی های پنهان وابسته است و مستقل از اندازه ی کل گراف است. با این حال به لطف بازنمایی مشترک که قبل تر توضیح داده شده، این مدل می تواند در زمان تخمین پارامترها، وابستگی های جهانی بین موجودیت ها و رابطه ها را فهمیده و در فرایند آموزش تاثیر دهد. [1]

نتایج یادگیری رابطه ای: رسکال در زمینه های مختلف مدل های یادگیری مدرن⁵ مورد استفاده قرار گرفته است. مثلاً در [63] نشان داده شده است که رسکال موفق شده در پیش بینی رابطه نتایج نزدیک

¹ gradient-based minimization

² stochastic gradient descent

³ sequence

⁴ efficient close-form update

⁵ state-of-the-art

یا بهتر روی چندین مجموعه داده ی معیار نسبت به روش های [70] Markov Logic Networks و Bayesian Clustered Tensor و the Infinite (Hidden) Relational model [71, 72] Factorization [73] بدست آورد. همچنین رسکال برای پیش بینی پیوند روی کل پایگاه دانش مانند YAGO و DBpedia مورد استفاده قرار گرفته است. [64, 74] فارق از پیش بینی پیوند، رسکال در روش های یادگیری رابطه ای تصادفی (SRL)^۱ مانند بازنمایی موجودیت ها^۲ و خوشه بندی بر مبنای پیوند^۳ عملکرد موفق داشته است. برای مثال این روش عملکرد موفق در دسته بندی نویسندگان، ناشران و سالن های انتشار روی مجموعه داده های ناشران داشته است [63, 65]. علاوه بر این، ارتباط معنایی موجودیت ها که در این روش محاسبه شد در ایجاد طبقه بندی^۴ به کمک خوشه بندی سلسله مراتبی^۵ روی داده های دسته بندی نشده^۶ به کار برده شده است [75].

¹ stochastic relational learning

² entity resolution

³ link-based clustering

⁴ taxonomies

⁵ hierarchical clusterings

⁶ uncategorized data via

۳-۴-۳- روش های ادراک چندلایه ای^۱

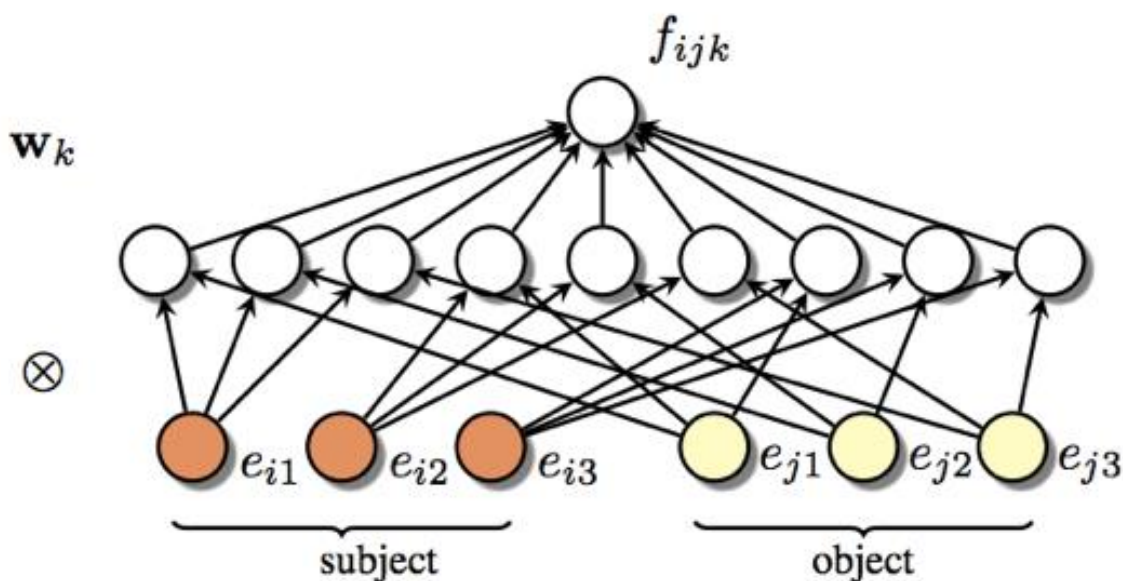
می توانیم رسکال را به صورت مدلی که برای هر حقیقت یک بازنمایی تولید می کند و از روی این بازنمایی ها وجود یا عدم وجود این حقیقت ها را پیش بینی می کند تفسیر کنیم. به طور خاص می توانیم رسکال را به صورت رابطه های + و + بازنویسی کنیم.

$$f_{ijk}^{\text{RESCAL}} := \mathbf{w}_k^T \phi_{ij}^{\text{RESCAL}}$$

$$\phi_{ij}^{\text{RESCAL}} := \mathbf{e}_j \otimes \mathbf{e}_i,$$

در روابط + ضرب داخلی \mathbf{e}_j و \mathbf{e}_i که در محاسبه ی مجموع امتیازها شرکت می کردند را از مجموع بیرون کشیده و به صورت ضرب خارجی نوشتیم. بنابراین رسکال بازنمایی جفت موجودیت i و j را به صورت ضرب تنسور ویژگی های پنهان این دو موجودیت بدست آورد رابطه ی + و وجود حقیقت x_{ijk} را از روی بازنمایی جفت موجودیت ها (ϕ_{ij}) و ماتریس وزن دار رابطه (w_k) بدست آورد. این تقسیم بندی رابطه ی رسکال در شکل + به تصویر کشیده شده است. برای توضیح بیشتر در مورد ایجاد بازنمایی پنهان به وسیله ی ضرب تنسور به [88, 89, 90] مراجعه شود.

¹ Multi-layer perceptrons



از آنجایی که ضرب تنسور تراکنش بین همهی جفت موجودیت ها را مدل می کند، وقتی تعداد ویژگی های پنهان زیاد باشد، رسکال در این راهکار به تعداد زیادی پارامتر نیاز پیدا خواهد کرد. این موضوع می تواند باعث مشکل در مقیاس پذیری روش روی پایگاه های دانش بزرگ با تعداد زیاد رابطه شود.

در ادامه ی این بخش به روش های ادراک چندلایه ای که به شبکه های عصبی پیش خور¹ نیز معروف هستند می پردازیم. این راهکار این امکان را به ما می دهد که مدل های جایگزینی برای ساختن بازنمایی حقیقت ها در نظر بگیریم و همچنین بتوانیم از توابع غیر خطی در پیش بینی وجود پیوندها استفاده کنیم.

مدل $E\text{-MLP}^2$:

¹ feedforward neural networks

² Entity-MLP

در ابتدا مدل Entity-MLP را معرفی می کنیم. تابع امتیاز این مدل به صورت معادلات + و + و + محاسبه می گردد.

$$\begin{aligned} f_{ijk}^{\text{E-MLP}} &:= \mathbf{w}_k^\top \mathbf{g}(\mathbf{h}_{ijk}^a) \\ \mathbf{h}_{ijk}^a &:= \mathbf{A}_k^\top \phi_{ij}^{\text{E-MLP}} \\ \phi_{ij}^{\text{E-MLP}} &:= [\mathbf{e}_i; \mathbf{e}_j] \end{aligned}$$

در روابط فوق $\mathbf{g}(u) = [\mathbf{g}(u_1), \mathbf{g}(u_2), \dots]$ یک تابع \mathbf{g} است که بر روی تک تک المان های بردار u اعمال می شود. که می تواند یک تابع غیر خطی باشد، مانند $\mathbf{g}(u) = \tanh(u)$.

در این روابط h_a یک لایه ی مخفی^۱ اضافه شده است که ماتریس وزن دار دیگر روی بازنمایی موجودیت ها اعمال می کند. در اصل ما در رابطه ی + دو موجودیت e_i و e_j را ترکیب کردیم و هیچ کنشی بین این دو موجودیت محاسبه و تاثیر داده نشده است، بنابراین نیاز به یک ماتریس وزن دار برای محاسبه ی تاثیر این دو موجودیت نیاز بود و h_a وارد معادله شد.

بزرگترین تفاوتی که این روش با روش های ضرب تنسور مانند رسکال دارد این است که در اینجا بجای محاسبه ی همه ی تراکنش های ممکن بین دو موجودیت، فقط تراکنش های موجود در A_k مورد بررسی قرار می گیرد. این راهکار به طور قابل ملاحظه ای تعداد پارامتری که باید آموزش دیده شوند را کاهش می دهد.

روش ER-MLP^2 :

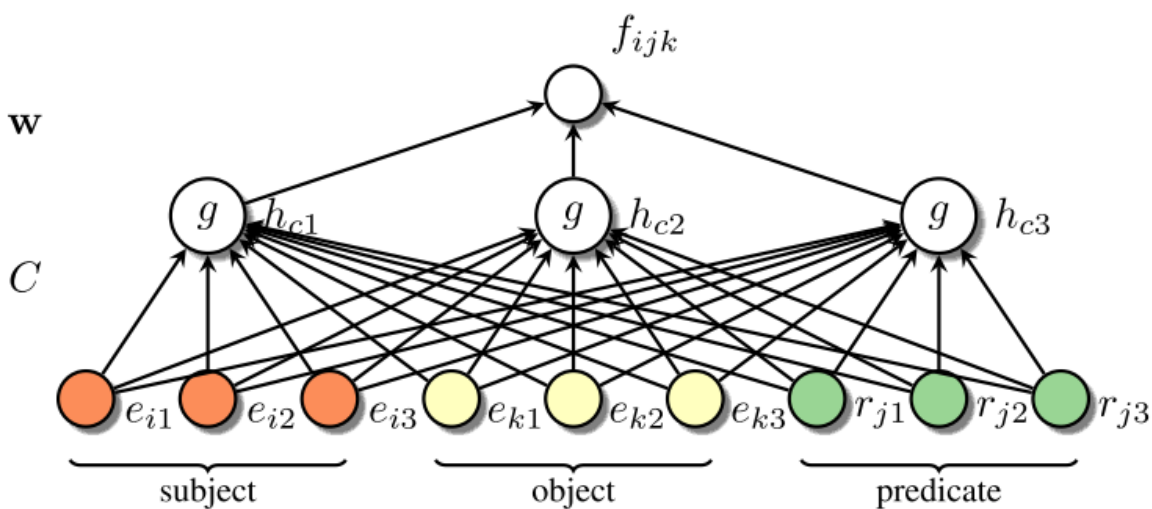
¹ hidden layer

² Entity-Relation-MLP

یکی از اشکالات روش E-MLP این است که باید برای هر رابطه ی ممکن، یک بردار w_k و یک ماتری A_k تعریف شود، که تعداد پارامترها را زیاد می کند. برای حل این مشکل روش ER-MLP معرفی شده است که در این روش رابطه را نیز به صورت برداری در کنار موجودیت ها در نظر می گیریم، و بجای ماتریس A_k می توانیم از یک بردار ثابت C استفاده کنیم. بنابراین روابط ER-MLP به صورت $+ + +$ خواهند بود.

$$\begin{aligned} f_{ijk}^{\text{ER-MLP}} &:= \mathbf{w}^\top \mathbf{g}(\mathbf{h}_{ijk}^c) \\ \mathbf{h}_{ijk}^c &:= \mathbf{C}^\top \phi_{ijk}^{\text{ER-MLP}} \\ \phi_{ijk}^{\text{ER-MLP}} &:= [\mathbf{e}_i; \mathbf{e}_j; \mathbf{r}_k]. \end{aligned}$$

توجه کنید که در این روش از یک بردار وزن دار جهانی برای همه روابط استفاده شده است و بردار C تعریف شده مستقل از رابطه ی r_k است که باعث کاهش تعداد پارامترهای مساله می گردد. نحوه ی کار روش ER-MLP در شکل + نمایش داده شده است.



در [91] نشان داده شده است که روش های MLP کلماتی که قرابت معنایی¹ دارند را به درستی نزدیک به یکدیگر تشخیص می دهند در حالی که برای چنین کاربردی آموزش دیده نشده اند. همچنین در [28] این ویژگی نشان داده شده است، برای مثال به کمک MLP مساله ی نزدیک ترین همسایه² برای بازنمایی پنهان چندین رابطه ی انتخاب شده از پایگاه دانش Freebase را حل کرده اند. در نتایج می توان مشاهده کرد که روابطی که ارتباط معنایی دارند نزدیک به یکدیگر قرار گرفته اند.

¹ semantically similar

² nearest neighbors

۴-۳-۴ شبکه های عصبی تنسور ۱

در [92] با ترکیب روش های ادراک چند لایه ای (MLP) و روش های دوخطی (Bilinear) روش جدیدی به نام شبکه های عصبی تنسور (NTN) معرفی شده است. روابط محاسبه ی تابع امتیاز این روش به صورت + و + و + است.

$$\begin{aligned} f_{ijk}^{NTN} &:= \mathbf{w}_k^\top \mathbf{g}([\mathbf{h}_{ijk}^a; \mathbf{h}_{ijk}^b]) \\ \mathbf{h}_{ijk}^a &:= \mathbf{A}_k^\top [\mathbf{e}_i; \mathbf{e}_j] \\ \mathbf{h}_{ijk}^b &:= [\mathbf{e}_i^\top \mathbf{B}_k^1 \mathbf{e}_j, \dots, \mathbf{e}_i^\top \mathbf{B}_k^{H_b} \mathbf{e}_j] \end{aligned}$$

در اینجا \mathbf{B}_k یک تنسور است، که در آن میزان رابطه ی بین جفت موجودیت ها نگاشت شده است. H_{ijk}^b را یک لایه ی مخفی دوخطی^۲ می نامیم، زیرا هم به صورت یک لایه ی مخفی از مدل های MLP هست و هم به کمک یک تابع وزن دار میزان کنش جفت موجودیت ها را مشخص می کند مانند آنچه در روش رسکال توضیح دادیم.

با توجه به معادله های معرفی شده در +و+و+ مشخص است که این روش مشکلات مقیاس پذیری که در هر دو روش E-MLP و RESCAL وجود داشت را دارد. همچنین در [95] و [28] نشان داده شده است که این روش به بیش برآزش^۳ میل می کند (حداقل روی مجموعه داده هایی که در این مقالات استفاده شده است) [1].

¹ Neural tensor networks

² bilinear hidden layer

³ overfit

۵-۴-۳- روش های فاصله ی پنهان^۱

کلاس دیگری از مدل ها، مدل های فاصله ی پنهان هستند (همچنین در تحلیل شبکه های اجتماعی به مدل های فضای پنهان معروفند) که احتمال وجود رابطه بین بین موجودیت ها را از فاصله ی بین بازنمایی پنهان آن ها در فضا می سنجد: موجودیت ها با یکدیگر رابطه دارند اگر بازنمایی پنهان آن ها با یک معیار فاصله نزدیک به هم باشد [۱].

در [96] مدلی برای داده های تک-رابطه ای^۲ برای اولین بار راهکاری در این زمینه معرفی شده است. این راهکار احتمال وجود پیوند در شبکه های اجتماعی را با تابع امتیاز $f(e_i, e_j) = -d(e_i, e_j)$ محاسبه می کند، که در آن $d(.,.)$ یک تابع اندازه گیری فاصله ی دلخواه مانند فاصله ی اقلیدسی^۳ است.

مدل Structured embedding

در [93] مدلی به نام structured embedding (SE) معرفی شده که در آن ایده ای که در بند قبل توضیح داده شد برای داده های چند-رابطه ای^۴ گسترش داده شده است. در این روش تابع امتیاز برای حقیقت های x_{ijk} به صورت + مدل شده است.

$$f_{ijk}^{SE} := -\|\mathbf{A}_k^s \mathbf{e}_i - \mathbf{A}_k^o \mathbf{e}_j\|_1 = -\|\mathbf{h}_{ijk}^a\|_1$$

¹ Latent distance models

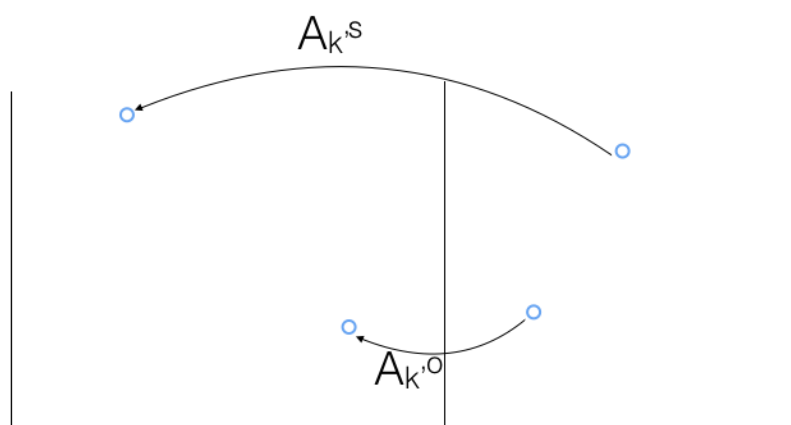
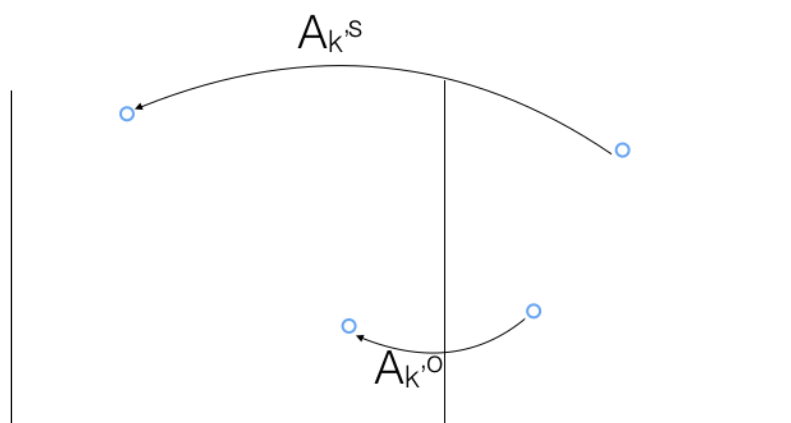
² uni-relational data

³ Euclidean distance

⁴ multi-relational

در رابطه ی $A_k +$ از ماتریس های مربوط به موجودیت های اول و دوم $[A_{ks}, A_{ko}]$ تشکیل شده است. ماتریس های A_{ks} و A_{ko} بازنمایی ویژگی های پنهان موجودیت ها را به فضای مخصوص به رابطه ی k منتقل می کند. این انتقال به صورتی آموزش دیده می شود که جفت رابطه هایی که برقرار هستند، نسبت به جفت رابطه هایی که برقرار نیستند به یکدیگر نزدیکتر باشند.

در شکل + مشاهده می شود که در صورتی که بین دو موجودیت i و j رابطه ی k برقرار باشد، پس از نگاشت این دو موجودیت توسط ماتریس های نگاشت A_{ks} و A_{ko} این دو موجودیت به یکدیگر نزدیکتر شده اند. همینطور در شکل + مشاهده می شود که در صورتی که رابطه ی k' بین دو موجودیت i و j برقرار نباشد، پس از نگاشت این دو موجودیت توسط ماتریس های نگاشت $A_{k'o}$ و $A_{k's}$ دو موجودیت در فاصله ی بیشتری از یکدیگر قرار می گیرند.



یکی از نکات منفی که در این روش به چشم می خورد یادگیری ماتریس های وزن دار جدا برای موجودیت های اول و دوم است، که باعث افزایش تعداد پارامتر مورد نیاز برای آموزش است.

۶-۴-۳- مدل TransE^۱

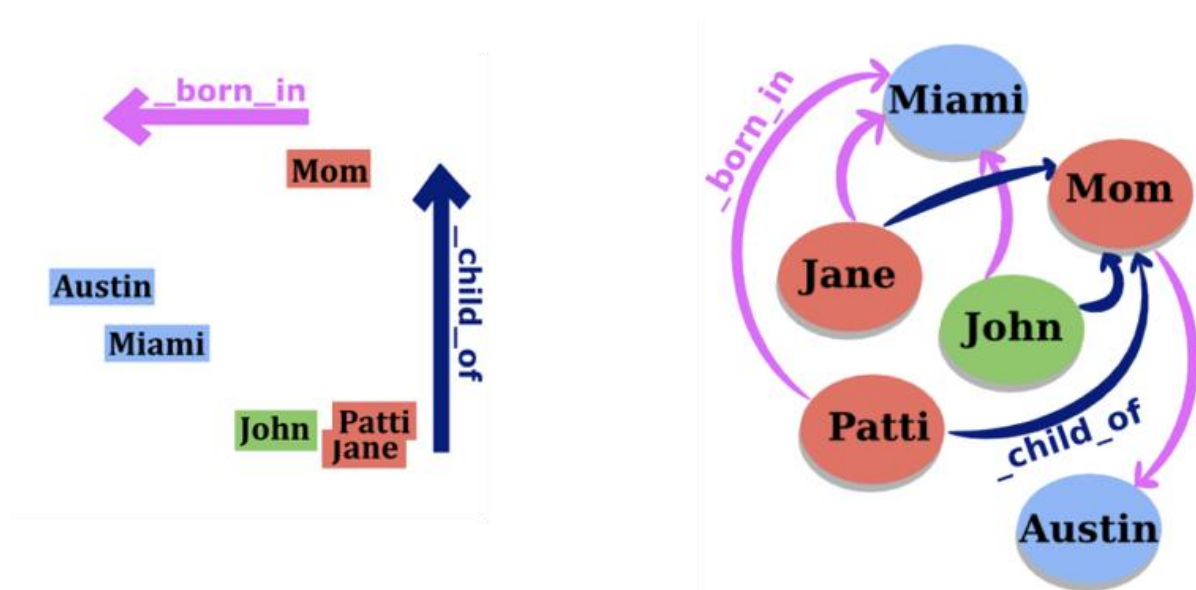
در [94] مدلی برای حل مشکل تعداد پارامتر زیاد در روش SE مطرح شده است که بجای اینکه از ماتریس های Aks و Ako برای تاثیر دادن رابطه ها استفاده شود، رابطه را یک بردار همانند بردار موجودیت ها در نظر گرفته شده و به عنوان یک آفست در کنار موجودیت اول در محاسبات از آن استفاده می کند. [1] به طور خاص امتیاز حقیقت xijk از رابطه ی + بدست می آید.

$$f_{ijk}^{\text{TransE}} := -d(\mathbf{e}_i + \mathbf{r}_k, \mathbf{e}_j).$$

در این روش رابطه ها به صورت یک بردار انتقال استفاده می شوند، به این صورت که فقط روی موجودیت اول اعمال می شوند و در صورت وجود رابطه بین این دو موجودیت، موجودیت اول را به موجودیت دوم نزدیک می کنند. در شکل + یک گراف را مشاهده می کنیم که در آن روابط فرزند و مکان تولد بین ۶ موجودیت نمایش داده شده است. در شکل + یک مثال از اعمال روش TransE روی این گراف را مشاهده می کنیم. مشاهده می شود که موجودیت ها بردارهایی ثابت در نظر گرفته شده اند و بردارهای رابطه به صورت یک بردار انتقال به تصویر کشیده شده است که موجودیت ها را به هدف های مورد نظر نزدیک می کنند.

¹ Translation Embedding

مثلا اگر john را توسط بردار Child_of منتقل کنیم به موجودیت Mom نزدیک می شود که حقیقت (john, Child_of, Mom) را تشکیل می دهد. اما اگر همین موجودیت را توسط بردار رابطه ی born_in منتقل کنیم، به موجودیت Miami نزدیک می شود که حقیقت (john, born_in, Miami) را تشکیل می دهد.



روش TransE این مساله را به صورت یک مساله ی کمینه سازی حل می کند به این صورت که سعی می کند فاصله ی مجموع بردارهای موجودیت اول و رابطه (e1+r) را با موجودیت دوم (e2) کمینه کند. همانطور که قبل تر گفتیم تابع امتیاز در این روش ها یک تابع فاصله است، پس در اینجا هدف کاهش فاصله بین (e1+r) و e2 است که بصورت + نوشته می شود.

$$d(e1, r, e2) = ||e1 + r - e2||$$

برای رسیدن به نتیجه ی بهینه نوآوری دیگری که در این روش معرفی شده است این است که تابع کمینه سازی را به گونه ای تغییر داده است که فاصله ی نمونه های اشتباه را نیز زیاد کرده است.

در زمانی آموزش پارامترها علاوه بر در نظر گرفتن این نکته که باید تابع فاصله ی نمونه های مثبت موجود در پایگاه دانش کمینه شود، سعی شده تا فاصله ی نمونه های منفی را نیز افزایش دهد. از آنجایی که

پایگاه های دانش مورد استفاده از فرض جهان باز پیروی می کنند وقتی حقیقت (e1, r, e2) در پایگاه دانش موجود نیست نمی توانیم نتیجه بگیریم که این سه گانه غلط است و آن را نمونه ای منفی تلقی کنیم.

در این روش برای ساختن نمونه های منفی، نمونه های مثبت مورد استفاده قرار گرفته اند به این صورت که یک بار موجودیت اول حذف شده و یک موجودیت تصادفی جای آن قرار داده شده است و یک بار موجودیت دوم حذف شده و یک موجودیت تصادفی جای آن قرار گرفته است، به این صورت به ازای هر نمونه ای مثبت دو نمونه ای منفی نیز تولید شده است (رابطه ی +). در رابطه ی + مشاهده می شود که علاوه بر کاهش فاصله ی موجودیت های مثبت، یک جریمه هم برای نمونه های منفی در نظر گرفته شده است و همچنین یک حاشیه γ بین نمونه های مثبت و منفی اضافه شده است.

$$S' = \{(\text{sub}', \text{rel}, \text{obj}) | \text{sub}' \in \mathcal{E}\} \cup \{(\text{sub}, \text{rel}, \text{obj}') | \text{obj}' \in \mathcal{E}\}$$

$$\sum_{pos} \sum_{neg \in S'} [\gamma + ||\mathbf{s} + \mathbf{r} - \mathbf{o}||_2^2 - ||\mathbf{s}' + \mathbf{r} - \mathbf{o}'||_2^2]_+$$

در این روش نیز مانند روش رسکال از SGD استفاده شده است که هم امکان آموزش به کمک نمونه برداری دسته ای^۱ را فراهم می کند و هم از مشکل گیر کردن در نقاط بهینه ای^۲ محلی جلوگیری می کند. الگوریتم کامل این روش در + آمده است.

¹ Batch Sampling

² local optimom

- 1: **input:** Training set $S = \{(\text{sub}, \text{rel}, \text{obj})\}$, margin γ , learning rate λ
- 2: **initialize** $\mathbf{r} \leftarrow \text{uniform}(-\frac{6}{\sqrt{k}}, \frac{6}{\sqrt{k}})$ for each rel
- 3: $\mathbf{r} \leftarrow \ell / \|\ell\|$ for each ℓ
- 4: $\mathbf{e} \leftarrow \text{uniform}(-\frac{6}{\sqrt{j}}, \frac{6}{\sqrt{k}})$ for each entity ent(sub or obj)
- 5: **loop**
- 6: $\mathbf{e} \leftarrow \mathbf{e} / \|\mathbf{e}\|$ for each entity ent
- 7: $S_{\text{batch}} \leftarrow \text{sample}(S, b)$ //sample minibatch of size b
- 8: $T_{\text{batch}} \leftarrow \emptyset$ //initialize set of pairs
- 9: **for** $(\text{sub}, \text{rel}, \text{obj}) \in S_{\text{batch}}$ **do**
- 10: $(\text{sub}', \text{rel}, \text{obj}') \leftarrow \text{sample}(S'(\text{sub}, \text{rel}, \text{obj}))$ //sample negative triplet
- 11: $T_{\text{batch}} \leftarrow T_{\text{batch}} \cup \{((\text{sub}, \text{rel}, \text{obj}), (\text{sub}', \text{rel}, \text{obj}'))\}$
- 12: **end for**
- 13: Update embeddings w.r.t. $\sum_{T_{\text{batch}}} \nabla [\gamma + \|\mathbf{s} + \mathbf{r} - \mathbf{o}\|_2^2 - \|\mathbf{s}' + \mathbf{r} - \mathbf{o}'\|_2^2]_+$
- 14: **end loop**

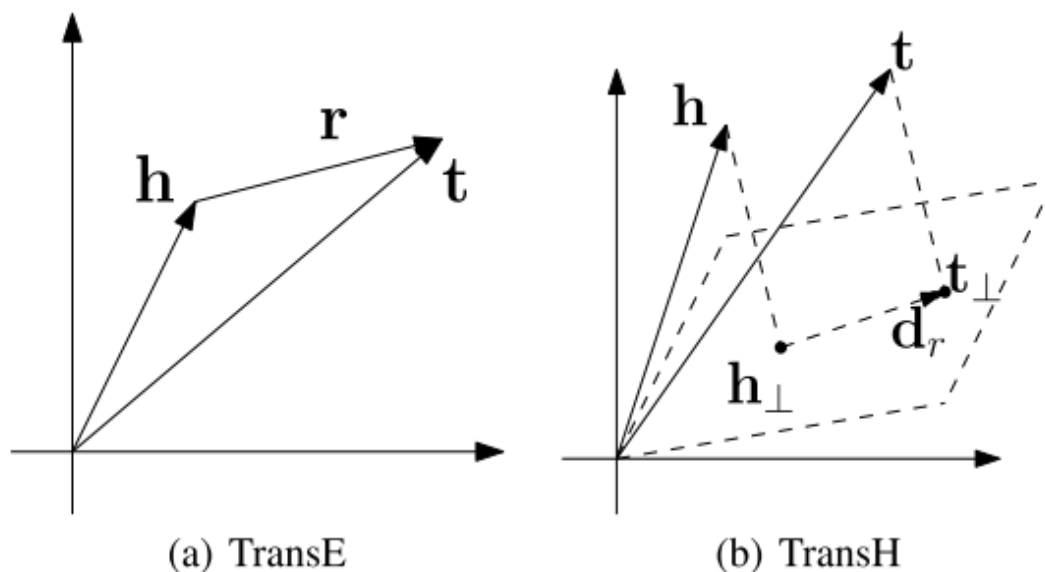
TransH^۱ مدل ۳-۴-۷

در بخش قبل مشاهده کردیم که TransE روش مناسبی را ارائه داد با تعداد پارامتر کم و مقیاس پذیر که قادر به آموزش مدل است. در [۱] روشی معرفی شده است که به بررسی برخی از نگاشته ها مانند یک به چند، چند به یک، چند به چند و انعکاسی پرداخته است. در این روش رابطه ها به صورت یک ابرصفحه و یک بردار انتقال روی آن در نظر گرفته می شود، از همین رو به این روش روش انتقال روی ابر صفحه یا به اختصار TransH گفته می شود.

در این روش موجودیت ها همچون روش قبل به صورت یک بردار از ویژگی های پنهان تعریف می شوند ولی رابطه ها به صورت دو بردار تعریف می شوند، یک بردار برای انتقال فضای مساله به ابرصفحه آن

¹ Translating on Hyperplanes

رابطه و دیگری بردار انتقالی است که در روش TransE نیز داشتیم و موجودیت های اول و دوم را به یکدیگر وصل می کرد. در روش قبل هدف این بود که دو موجودیت مشارکت کننده در یک حقیقت درست به وسیله ی بردار انتقال رابط به یکدیگر وصل شوند و اگر حقیقت صحیح نیست انتظار داشتیم که این اتصال برقرار نباشد. در این روش نیز هدف به همین صورت است با این تفاوت که این انتقال برای هر رابطه روی یک ابرصفحه انجام می شود که نگاشتی از موجودیت های اول و دوم در آن قرار دارد. به کمک این روش ساده می توانیم نگاشت های یک به چند، چند به یک، چند به چند و انعکاسی را نیز در عملیات آموزش تاثیر دهیم در حالی که پیچیدگی و هزینه ی این روش TransE نداریم. با هزینه و پیچیدگی کمی مانند روش TransE، را نیز در نظر گرفت.



همانطور که در شکل - مشاهده می شود روش TransE بردارهای h و t را به کمک بردار r به یکدیگر متصل می کرد اما در روش TransH همانطور که در شکل - نشان داده شده است ابتدا بردارهای h و t به صفحه ی مربوط به رابطه ی مورد نظر منتقل شده اند و توسط بردار d_r که بردار رابطه است به یکدیگر متصل می شوند.

مقادیر h_{\perp} و t_{\perp} به صورت معادله ی - بدست می آیند که در آن بردار انتقال به ابرصفحه ی مربوط به رابطه ی r است و dr معادل بردار r در ابرصفحه ی رابطه است.

$$\mathbf{h}_{\perp} = \mathbf{h} - \mathbf{w}_r^{\top} \mathbf{h} \mathbf{w}_r, \quad \mathbf{t}_{\perp} = \mathbf{t} - \mathbf{w}_r^{\top} \mathbf{t} \mathbf{w}_r.$$

پس تابع امتیاز این روش به صورت زیر خواهد بود:

$$f_r(\mathbf{h}, \mathbf{t}) = \|(\mathbf{h} - \mathbf{w}_r^{\top} \mathbf{h} \mathbf{w}_r) + \mathbf{d}_r - (\mathbf{t} - \mathbf{w}_r^{\top} \mathbf{t} \mathbf{w}_r)\|_2^2.$$

در این روش نیز مانند روش TransE برای کمینه کردن خطا از حقیقت های صحیح و غلط و یک فاصله بین آن ها استفاده می شود که قصد دارد فاصله ی بین بردار $h + r$ حقیقت های صحیح را با t کم و همین فاصله را در حقیقت های غلط زیاد کند. انتخاب حقیقت های غلط در این روش بمانند روش قبل است اما در این روش بجای استفاده از تابع امتیاز fr که در معادله ی - مطرح کردیم از تابع امتیاز fr که در - معرفی کردیم استفاده می کنیم و تابع کمینه سازی به صورت - خواهد شد.

$$\mathcal{L} = \sum_{(h,r,t) \in \Delta} \sum_{(h',r',t') \in \Delta'_{(h,r,t)}} [f_r(\mathbf{h}, \mathbf{t}) + \gamma - f_{r'}(\mathbf{h}', \mathbf{t}')]_+$$

فصل ۴: قوانین انجمنی

در این قسمت سعی بر این داریم که نقطه ضعفی که در بخش + از آن صحبت کردیم را به کمک قوانین انجمنی حل کنیم. در ابتدار ...

۴-۱- معیارهای اطمینان

برای اینکه بتوانیم به قانون هایی که استخراج کرده ایم اعتماد کنیم و از آن ها در آموزش مدل ها استفاده کنیم، نیاز به یک معیار اطمینان داریم. معیارهای رایج برای اندازه گیری دقت یک قانون میزان پشتیبانی قانون^۱، پوشش سر^۲ و معیار اطمینان استاندارد^۳ است.

¹ support rule

² head coverage

³ Standard confidence

۱-۴-۱- معیار پشتیبانی قانون:

این معیار برای میزان اهمیت یک قانون استفاده می شود و به این صورت تعریف می شود: تعداد قوانین یافت شده در پایگاه دانش که یک حقیقت را نتیجه می دهند. برای مثال در + قانون R معرفی شده است که مشخص می کند هر شخص x که در شهر y زندگی می کند، در آن شهر متولد شده است. به تعداد دفعاتی که این دو حقیقت در پایگاه دانش اتفاق بیافتد و قانون R را نقض نکند، پشتیبان این قانون گفته می شود.

$$R: \text{livesIn}(x, y) \Rightarrow \text{wasBornIn}(x, y)$$

این معیار به صورت رابطه ی + تعریف می شود.

$$\text{supp}(\vec{B} \Rightarrow r(x, y)) := \#(x, y) : \exists z_1, \dots, z_m : \vec{B} \wedge r(x, y)$$

در رابطه ی + B مجموعه ای از حقیقت ها است که نتیجه میدهد حقیقت $r(x, y)$ را، به تعداد باری که این اتفاق می افتد معیار پشتیبانی قانون گفته می شود. در جدول + یک پایگاه دانش متشکل از دو رابطه و ۵ حقیقت وجود دارد. قانون + را در نظر بگیرید، میزان پشتیبان قانون R در این جدول برابر ۱ است، بخاطر وجود حقیقت (Adam, LivesIn, Paris) و (Adam, wasBornIn, Paris) که از قانون R پیروی می کنند [AMIE].

<i>livesIn</i>	<i>wasBornIn</i>
(Adam, Paris)	(Adam, Paris)
(Adam, Rome)	(Carl, Rome)
(Bob, Zurich)	

۲-۱-۴- معیار پوشش سر:

معیار پشتیبان قانون یک مقدار مطلق است و برای تعریف کردن یک آستانه برای این معیار نیاز است تا اندازه ی پایگاه دانش را نیز بدانیم. مثلاً اگر پشتیبان یک قانون در یک پایگاه دانش با ۲۰۰۰ حقیقت ۵۰۰ با شد، عدد بسیار بزرگی است اما اگر در یک پایگاه دانش با چندین میلیون حقیقت ۵۰۰ با شد عدد قابل توجهی نیست. برای حذف این وابستگی معیار پوشش سر را به صورت زیر معرفی می کنیم.

$$hc(\vec{B} \Rightarrow r(x, y)) := \frac{supp(\vec{B} \Rightarrow r(x, y))}{size(r)}$$

$$size(r) := \#(x', y') : r(x', y')$$

در این معیار بررسی می شود که چند درصد از $r(x, y)$ هایی که اتفاق افتاده است بخاطر وجود زنجیره قانون B بوده است. در جدول + مقدار معیار پوشش سر بخاطر وجود (Carl, Rome) برابر با ۰,۵ است [AMIE].

۳-۱-۴- معیار اطمینان استاندارد:

معیارهای قبل اهمیت قانون را بررسی می کردند و فقط پیش بینی های درست از قانون را در نظر می گرفتند، و پیش بینی های غلطی که این قانون تولید می کند را در نظر نمی گیرند. پس ما به معیاری نیاز داریم که کیفیت قانون را نیز بررسی کند.

در این معیار میزان پشتیبان هر قانون را بر تعداد باری که قسمت بدنه ی این قانون در پایگاه دانش دیده شده است تقسیم می کنیم. این کار باعث می شود که معیاری داشته باشیم برای اینکه بدانیم در چند درصد مواقع که قسمت بدنه ی این قانون برقرار بوده است منجر به تولید قسمت نتیجه شده است.

$$conf(\vec{B} \Rightarrow r(x, y)) := \frac{supp(\vec{B} \Rightarrow r(x, y))}{\#(x, y) : \exists z_1, \dots, z_m : \vec{B}}$$

رابطه ی بالا به زبان ساده برابر است با تعداد رخداد قانون بخش بر تعداد دفعاتی که می توانست رخ دهد. برای مثال در پایگاه دانش معرفی شده در جدول + معیار اطمینان استاندارد برابر $\frac{1}{3}$ است. زیرا $wasBornIn(Adam, Paris)$ به عنوان نمونه ی مثبت^۱ در نظر گرفته شده و $wasBornIn(Adam, Rome)$ و $wasBorn(Bob, Zurich)$ که در پایگاه دانش وجود ندارند نمونه ی منفی^۲ در نظر گرفته شده است. که معیار اطمینان استاندارد برابر می شود با تعداد نمونه ی مثبت یعنی ۱ بخش بر تعداد کل نمونه ها یعنی ۳.

این معیار در پایگاه های دانش فرض جهان بسته معیار خوبی است و به خوبی دقت قانون استخراج شده را مشخص می کند زیرا همان طور که در مثال بالاتر هم مشاهده کردیم، عدم وجود یک نمونه ی مثبت از حقیقت در پایگاه دانش دلیلی بر غلط بودن آن حقیقت بوده است. اما همانطور که در بخش + اشاره کردیم فضای این مساله فرض جهان باز است و حقیقت هایی که در این پایگاه های دانش وجود ندارند لزوماً غلط نیستند [AMIE]. برای درک بیشتر این موضوع به مثال زیر دقت کنید:

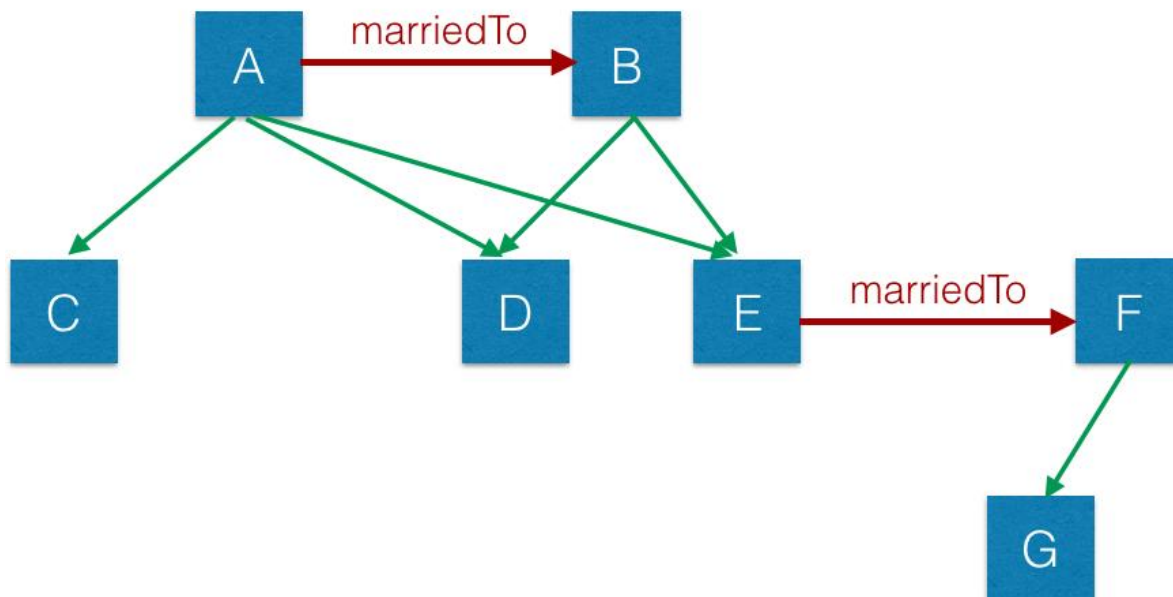
به گراف ارائه شده در شکل + دقت کنید. در این گراف رابطه های افقی رابطه ی ازدواج ($marriedTo$) و روابط عمودی رابطه ی داشتن فرزند ($hasChild$) را مشخص می کنند. در این گراف قصد داریم که قانون R با تعریف زیر را بررسی کنیم:

$$hasChild(y, x), marriedTo(y, z) \Rightarrow hasChild(z, x)$$

¹ positive example

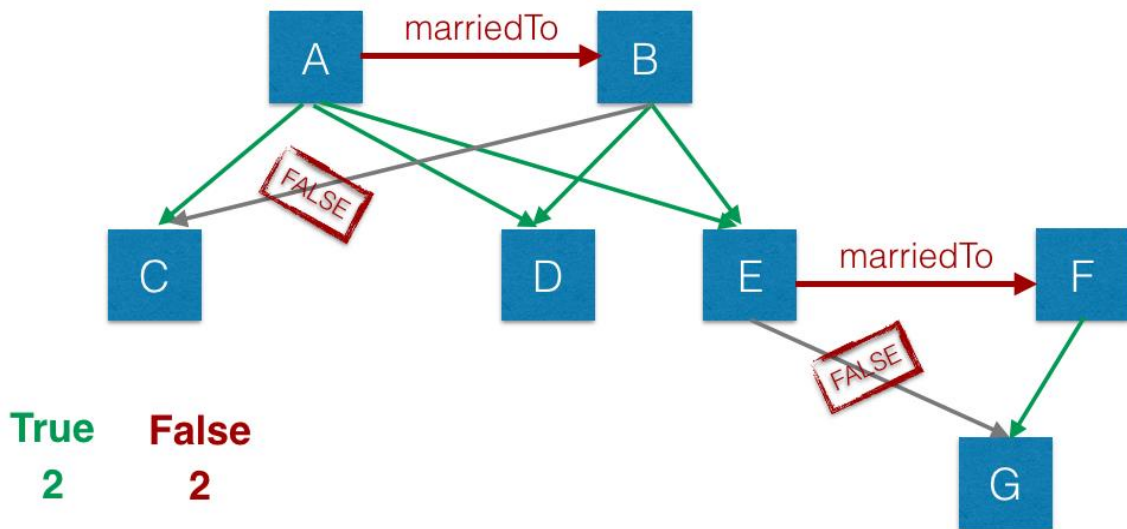
² negative example

این قانون به این معناست که اگر شخصی y فرزندی به نام x داشته باشد و همین شخص با شخص دیگری به نام z ازدواج کرده باشد، میتوان نتیجه گرفت که شخص z هم فرزندی به نام x دارد.

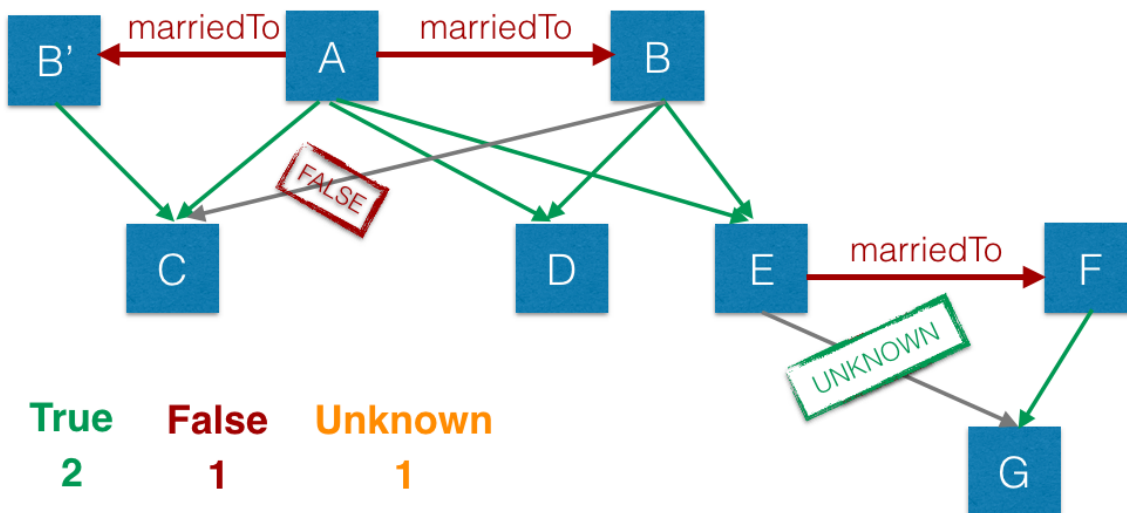


اگر بخواهیم معیار اطمینان استاندارد در این گراف را بررسی کنیم، مشاهده می شود که این قانون ۲ بار در این گراف صدق کرده است در حالی که می توانست ۴ بار اتفاق بیافتد. در شکل + دفعاتی که این قانون باید اتفاق می افتاده است و اتفاق نیافتاده است نمایش داده شده و به عنوان نمونه ی منفی در نظر گرفته شده است. با این اوصاف معیار اطمینان استاندارد قانون R در گراف معرفی شده برابر ۰,۵ می شود.

$$\text{Standard Confidence} = \frac{\sup(B \Rightarrow r(x,y))}{\#(x,y):B} = 2 / 4 = 0.5$$



اما همانطور که توضیح دادیم عدم وجود حقایق در پایگاه های دانش فرض جهان باز دلیل بر غلط بودن این حقیقت ها نیست. برای مثال فرض کنید که در گراف معرفی شده یک گرهی B' داشته باشیم که رابطه ای ازدواج داشته باشد با گرهی A و رابطه ای فرزندنی داشته باشد با گرهی C پس می دانیم که رابطه ای فرزندنی بین گرهی B و C برقرار نیست اما در مورد وجود و عدم وجود این رابطه بین E و G اطلاعی نداشته باشیم (شکل +).



پس معیار اطمینان استاندارد برای پایگاه های دانش جهان باز معیار مناسبی نیست و نیاز به تعریف معیار بهتری داریم. در [AMIE] معیار اطمینان با فرض نیمه کامل^۱ (PCA-Confidence) را معرفی کرده است.

۴-۱-۴- معیار اطمینان با فرض نیمه کامل

در این معیار سعی شده نمونه هایی که در پایگاه دانش وجود ندارند و در معیار اطمینان استاندارد آنها را نمونه ی غلط در نظر می گرفتیم را در اینجا بیشتر بررسی کنیم و با احتمال بهتری غلط بودن یا نبودن آن را مشخص کنیم.

در این معیار اگر حقیقت $r(x,y)$ در پایگاه دانش موجود نبود، بررسی می شود که y وجود دارد که برای آن حقیقت $r(x,y')$ برقرار باشد یا خیر، اگر وجود داشته باشد در نظر می گیرد که $r(x,y)$ غلط بوده و در غیر این صورت این حقیقت را ناشناخته فرض می کند.

$$conf_{pca}(\vec{B} \Rightarrow r(x,y)) := \frac{supp(\vec{B} \Rightarrow r(x,y))}{\#(x,y) : \exists z_1, \dots, z_m, y' : \vec{B} \wedge r(x,y')}$$

(1)

به مثال جدول + برگردیم، در این جدول wasBornIn(Adam,Paris) را یک نمونه ی مثبت در نظر می گیریم و wasBornIn(Adam,Rome) را یک نمونه ی منفی در نظر می گیریم اما اینکه livesIn (Bob, Zurich) داریم و در مورد محل تولد آن اطلاعی نداریم دلیل نمی شود که این نمونه را یک نمونه ی منفی در نظر بگیریم زیرا ممکن است bob اصلاً متولد نشده باشد. پس در این مثال مقدار $PCA-Confidence = \frac{1}{2}$ می شود.

یا در گرافی که در شکل + معرفی شده است، در معیار اطمینان استاندارد هر دو رابطه ی بین (B,C) و (E,G) را غلط در نظر گرفتیم، اما در اینجا داریم که C مادری به نام B' دارد و می توانیم مطمئن

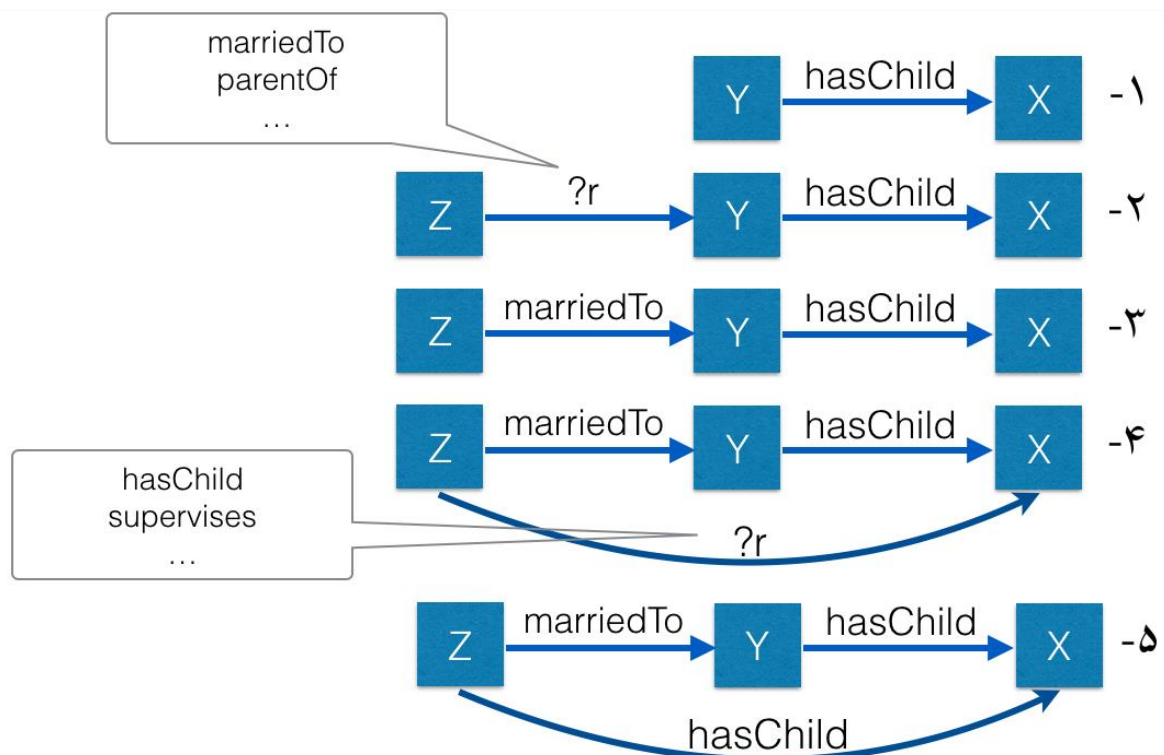
¹ Partial Completeness Assumption

شویم که رابطه ی فرزندی بین B و C برقرار نیست و این رابطه را غلط در نظر بگیریم. اما برای رد رابطه ی E و G هیچ مدرکی نداریم و این رابطه را غلط در نظر نمی گیریم بلکه ناشناخته در نظر می گیریم و از آن در فرومول PCA-confidence استفاده نمی کنیم. پس PCA-Confidence در این مثال برابر با X است.

۲-۴- الگوریتم ها

برای استخراج قوانین انجمنی از الگوریتم های معرفی شده در [۱] استفاده شده است که در ادامه توضیح مختصری در مورد بخشی از این الگوریتم ها خواهیم داد.

نحوه ی استخراج قوانین در شکل + نشان داده شده است که هر مرحله را در ادامه توضیح خواهیم داد.



به ازای هر حقیقت این روند یک بار تکرار می شود:

- 1 - حقیقت انتخاب شده را در نظر می گیریم. در این مثال حقیقت $(Y, \text{hasChild}, X)$ در نظر گرفته شده است.
 - 2 - تمامی روابطی که بین یکی از موجودیت ها و موجودیت دیگری وجود دارد را بررسی می کنیم. در این مثال روابط موجود برای موجودیت اول یعنی Y نمایش داده شده است، برای مثال روابط ازدواج و والد بودن با این موجودیت مورد بررسی قرار گرفته است.
 - 3 - همه ی روابطی که در مرحله ی ۲ کاندید شده بودن را در کنار Y قرار می دهیم و یک زنجیر از قوانین را ایجاد می کنیم. زنجیر ایجاد شده در این قسمت به صورت رابطه ی ازدواج بین Z و Y و رابطه ی داشتن فرزند بین Y و X است.
 - 4 - در این قسمت همه ی روابطی که بین موجودیت اضافه شده به زنجیر و موجودیت اول را بررسی می کنیم یعنی روابط بین Z و X . روابط کاندید در این مثال روابط “داشتن فرزند” و “ناظر بودن” است.
 - 5 - در این مرحله روابطی که کاندید شده اند را بررسی می کنیم و معیارهای اطمینان و اهمیت که در بخش قبل معرفی کردیم را برای آن ها بدست می آوریم، در صورتی که به معیارها مقدار قابل قبولی داشته باشند قانون تولید شده را ذخیره می کنیم و در غیر این صورت از آن رد می شویم.
 - 6 - پس از بررسی همه ی روابط کاندید تولید شده در بخش ۴ کل زنجیر تولید شده را به مرحله ی ۱ ارسال می کنیم و با اضافه کردن یک حقیقت دیگر به ابتدا یا انتهای آن روند رشد زنجیر را تا آستانه ی تعریف شده ادامه می دهیم.
- به کمک الگوریتم معرفی شده در بالا همه ی قانون های ممکن در پایگاه دانش با توجه به معیارهای کیفیت و اهمیت استخراج خواهند شد. اما با توجه به اندازه ی بزرگ پایگاه های دانش مورد استفاده این روش برای بررسی همه ی حالت های ممکن بسیار وقت گیر خواهد بود، پس بهتر است در

مرحله ی ۲ و ۴ که به انتخاب کاندید می پردازیم بجای بررسی همه ی حقیقت های موجود در پایگاه داده یک عملیات هرس^۱ روی کاندیدها با توجه به همرخدای با حقیقت موجود انجام شود و همه ی حقیقت ها مورد بررسی قرار نگیرند.

با اعمال این الگوریتم روی پایگاه دانش Freebase15k که شامل حدود ۵۰۰ هزار حقیقت است، تعداد ۴۱۱۹۶ قانون استخراج شد. در ادامه چند مثال از این قانون ها را بررسی می کنیم.

در قانون زیر داریم که اگر b یک تیم فوتبال باشد که بازیکنی به نام a داشته باشد، می توانیم نتیجه بگیریم که تیم فعلی بازیکن a تیم b است. معیار اطمینان PCA در این مثال ۰,۸۹ است که مقدار قابل اعتمادی است.

?b /sports/soccer/team/player ?a =>

?a /soccer/current_team/team ?b

در مثال دیگر داریم که اگر بازیگر سینمای a جایزه ی b را برنده شده باشد، نتیجه می گیریم که این بازیگر برای جایزه ی b نامزد شده است.

?b /award/awards_won ?a =>

?a /award/award_nomination/nominated_for ?b

قوانین پیچیده تری نیز استخراج شده است که بدنه ی قانون^۲ از چند حقیقت تشکیل شده است که یک نمونه از این قوانین به صورت زیر است. در این قانون داریم که اگر a یک موسسه ی آموزشی باشد که در مکان f قرار دارد و داشته باشیم که f در مکان b واقع شده است، می توانیم نتیجه

¹ Pruning

² rule body

بگیریم که موسسه ی آموزشی a نیز در مکان b قرار دارد. علی رقم پیچیدگی ای که این قانون نسبت به دو قانون قبلی دارد، اما معیار اطمینان PCA این قانون ۰,۹۳ بوده و بسیار قابل اعتماد است و می توانیم از آن در آموزش مدل ها استفاده کنیم.

?a /educational_institution/located_in ?f

?f /location/location/containedby ?b

=>

?a /location/location/containedby ?b

در بخش بعدی انواع این قوانین را بررسی کرده و توضیح خواهیم داد که چگونه از این قوانین در بهبود مدل های مبتنی بر ویژگی های پنهان استفاده خواهیم کرد.

۱-۲-۴- قانون R-equivalence

این قانون رابطه شبیه به قانون قبلی است با این تفاوت که رابطه ی برگشت نیز بین دو طرف قانون برقرار است.

$$r(x, y) \Leftrightarrow r'(x, y)$$

این قانون در آموزش مدل ها کمک بسیاری می تواند بکند زیرا ماهیت این قانون به این صورت است که روابط یکسان (روابطی که به دو صورت در پایگاه دانش استفاده شده اند اما یک معنی می دهند) را شناسایی می کند، برای نمونه مثالی که در بخش قبل زدیم، دو رابطه ی located_in و contained_by در توالی یکدیگر آمده بودند که یک معنی را می دهند، پس می توانیم نتیجه بگیریم که اگر حقیقت (x, located_in, y) داشته باشیم، حقیقت (x, contained_by, y) را نیز داریم و در عملیات آموزش روش

از یکی از این دو مفهوم استفاده کنیم که باعث کاهش تعداد روابط موجود در پایگاه دانش و در نتیجه کاهش پیچیدگی مساله می شود.

۳-۴- قوانین هدف

در بخش قبل مشاهده شد که چگونه قوانین از پایگاه های دانش استخراج می شوند و چند نمونه از این قوانین را مرور کردیم. در این بخش دسته بندی از قوانینی که استخراج می شوند را معرفی می کنیم و در ادامه بررسی های انجام شده روی روش های موجود را روی این دسته بندی ها اعمال می کنیم و نقاط قوت و ضعف هر روش را نسبت به این دسته بندی ها می سنجیم.

۱-۳-۴- قانون R-subsumption

این قانون به صورت زیر تعریف می شود، به این معنی که اگر x و y رابطه ی r را باهم داشته باشند، رابطه ی r' نیز بین آن ها برقرار است.

$$r(x, y) \Rightarrow r'(x, y)$$

برای مثال اگر در پایگاه دانش رابطه ی پدر بودن را داشته باشیم، می توانیم از آن رابطه ی والد بودن را نتیجه بگیریم، مثلاً اگر حقیقت «علی، پدر، حسن» به معنی «علی پدر حسن است» را داشته باشیم، می توانیم نتیجه بگیریم که حقیقت «علی، والد، حسن» نیز حقیقت درستی است.

۲-hope translation قانون ۴-۳-۲

این قانون، همان قانون تعدی است، به این صورت که اگر x و y رابطه ی $r1$ را باهم داشته باشند و همچنین y و z نیز رابطه ی $r2$ را داشته باشند، می توانیم نتیجه بگیریم که x و z نیز با یکدیگر در ارتباطند.

$$r1(x, y), r2(y, z) \Rightarrow r'(x, z)$$

نمونه ای از این قانون را در مورد مناطق جغرافیایی در بخش قبل مشاهده کردیم. مثلاً می دانیم که شهر تهران در کشور ایران قرار دارد، و کشور ایران در منطقه ی خاورمیانه قرار دارد، پس می توانیم نتیجه بگیریم که شهر تهران در منطقه ی خاورمیانه قرار دارد.

Triangle alignment قانون ۴-۳-۳

این قانون نیز مانند قانون قبلی در سمت بدنه ی قانون دو حقیقت را بررسی می کند و از روی آن ها حقیقت جدید را نتیجه می گیرد.

$$r(x, z), r(y, z) \Rightarrow r'(x, y)$$

مشاهده می شود که در این قانون اشتراک موجودیت های دوم مد نظر قرار گرفته است و در صورتی که دو حقیقت در موجودیت دوم مشترک باشند بررسی می شود که بین موجودیت های اول نیز رابطه ای برقرار است یا خیر. برای مثال اگر داشته باشیم که $hasChild(x, z)$ و $hasChild(y, z)$ به این معنی که هم x و هم y فرزندی به نام z داشته باشند، می توان نتیجه گرفت که $married(x, y)$ یعنی x و y زن و شوهر هستند.

۴-۳-۴ قانون Specific R-subsumption

این قانون حالت کامل تری از قانون اول یعنی R-sub است به این صورت که بجز بررسی رابطه ها، ویژگی های موجودیت اول را نیز بررسی می کند. برای مثال در قانون اول داشتیم اگر شخصی پدر x باشد می توان نتیجه گرفت که آن شخص والد x نیز هست، اما عکس این قضیه را نمی توان نتیجه گرفت، یعنی نمی توان نتیجه گرفت که اگر شخصی والد x است پس پدر x است زیرا ممکن است مادر x باشد. در این قانون یک ویژگی از موجودیت اول نیز بررسی می شود.

$$r_1(x, y), r_2(x, v) \Rightarrow r'(x, y)$$

مثلا اگر داشته باشیم که شخصی والد x است و جنسیت آن شخص مذکر است، می توان نتیجه گرفت که آن شخص پدر x است و اگر جنسیت مونث داشت مادر x است.

در جدول + قوانین معرفی شده به اختصار آمده اند.

name	Target rule		Body rule
R-subsumption	$r'(x, y)$	\Rightarrow	$r(x, y)$
R-equivalence	$r'(x, y)$	\Leftrightarrow	$r(x, y)$
2-hope translation	$r'(x, z)$	\Rightarrow	$r_1(x, y), r_2(y, z)$
Triangle alignment	$r'(x, y)$	\Rightarrow	$r(x, z), r(y, z)$

$r_1(x, y), r_2(x, v)$	\Rightarrow	$r'(x, y)$	Specific R-sub
------------------------	---------------	------------	----------------

۴-۴- جمع بندی

متن

فصل ۵: روش پیشنهادی

متن

۱-۵- چالش ها

متن

۲-۵- عملکرد روش های موجود

در این بخش نتایج بدست آمده از روش های معرفی شده در بخش قبل را مورد بحث و بررسی قرار می دهیم. روش هایی که در این پایان نامه مورد بررسی و بهبود قرار گرفته اند ۱- روش RESCAL از روش های دوخطی ۲- روش NTN از روش های ادراک چندلایه ای ۳- روش SE ۴- روش TransE و ۵- روش TransH از روش های فاصله ی پنهان هستند.

در ادامه آزمایشاتی روی این ۵ روش انجام شده است و این روش ها از نظر مقیاس پذیری و میزان کیفیت خروجی که به کمک معیارهای رتبه ی میانگین و $hit@10$ محاسبه می شود مورد بررسی قرار گرفته اند. همه ی این آزمایشات روی پایگاه دانش Freebase15k انجام شده است که شامل حدود ۱۵ هزار موجودیت، ۱۳۰۰ رابطه و در مجموع ۶۰۰ هزار حقیقت است. تقسیم بندی داده های آموزش، ارزیابی و آزمون به شکلی که در + توضیح داده شد انجام شده است.

۱-۲-۵- روش رسکال:

تعداد پارامترها: در این روش برای هر موجودیت یک بردار d بعدی و برای هر رابطه یک ماتریس وزن دار $d \times d$ بعدی در نظر گرفتیم. پس تعداد پارامترهایی که در این روش نیاز است که آموزش دهیم از مرتبه $O(n_e d + n_r d^2)$ است. در آزمایشات تعداد ابعادی که بهترین نتیجه را بدست می دهد $d=250$ است. اگر ابعاد مساله را ۲۵۰ در نظر بگیر باید به تعداد ۸۸ میلیون پارامتر آموزش داده شود $(n_e d + n_r d^2 = 15k \times 250 + 1.3k \times 250 \times 250 = 88m)$.

کیفیت خروجی: پس از اعمال مدل آموزش دیده رسکال روی مجموعه داده ی آزمون ۵۰ هزار حقیقتی، میزان ۴۲٫۱٪ پاسخ های داده شده توسط این مدل زیر رتبه ی ۱۰ بوده $(hit@10)$ و میانگین رتبه ی کل پاسخ هایی که داده شده $(mean rank)$ ۶۸۳ است.

۲-۲-۵- روش NIN:

تعداد پارامترها: در این روش نیز برای هر موجودیت یک بردار d بعدی در نظر گرفته شده و برای هر رابطه یک ماتریس B_k و دو بردار A_k و w_k در نظر گرفته می شود که در مجموع برای هر رابطه d^3 بعد در نظر گرفته می شود و تعداد کل پارامترهایی که باید آموزش دیده شوند از مرتبه $O(n_e d + n_r d^3)$ است. تعداد پارامتری که در آزمایشات برای این روش نتیجه ی مناسبی بدست می دهد $d=50$ است. برای این تعداد ویژگی پنهان تعداد کل پارامتری که باید آموزش داده شوند تقریباً ۱۶۵ میلیون پارامتر می شود که نسبت به روش رسکال با اینکه تعداد ابعاد بردارها بسیار کمتر است، تقریباً دو برابر است $(n_e d + n_r d^3 = 15k \times 50 + 1.3k \times 50 \times 50 \times 50 = 163,250k)$.

کیفیت خروجی: مقدار $hit@10$ در این روش ۲۷٪ و میانگین رتبه ی همه ی پاسخ های این روش ۱۶۴ است. همانطور که در معرفی این روش گفته شد این روش به سرعت به بیش برآزش میل می کند و تاثیر این امر در میزان $hit@10$ مشخص است اما میانگین رتبه این روش به میزان خوبی از روش رسکال بهتر است که نشان می دهد روش رسکال برای سوالاتی که به میزان کافی نمونه ی آموزش ندیده است و نتایج خیلی پرتی بدست می دهد ولی برای ۴۲٫۱٪ حقیقت که روابط بین آن ها به خوبی آموزش دیده شده است و پاسخ های درست در ۱۰ پاسخ اول بوده است.

۳-۲-۵- روش Structured Embedding

تعداد پارامتر: در این روش برای هر موجودیت یک بردار d بعدی و برای هر رابطه دو ماتریس $d \times d$ بعدی آموزش دیده می شود، یکی برای انتقال موجودیت اول و دیگری برای انتقال موجودیت دوم. بنابراین تعداد پارامتر این روش از مرتبه ی $O(n_e d + 2n_r d^2)$ است. برای آموزش این مدل نیز بردارهایی با اندازه $d=50$ کفایت می کند و برای آموزش این مدل باید حدود ۸ میلیون پارامتر آموزش دیده شوند. $(n_e d + 2n_r d^2 = 15k \times 50 + 2 \times 1.3k \times 50 \times 50 = 7.250k)$

کیفیت خروجی: در این روش $hit@10$ برابر با ۳۹٫۸٪ است و مقدار میانگین رتبه ۱۶۲ بوده است. مشاهده می شود که علی رقم کاهش بسیار زیادی که در تعداد پارامترهای مساله نسبت به دو روش قبل داشته ایم، نتیجه ی خوبی حاصل شده است و در معیار $hit@10$ با روش رسکال و در معیار میانگین رتبه با روش NTN رقابت می کند.

۴-۲-۵- روش TransE:

تعداد پارامتر: در این روش تعداد پارامترها به میزان بسیار خوبی کاهش پیدا کرده است، همانطور که در معرفی این روش مشاهده شد، برای هر موجودیت یک بردار d بعدی در نظر گرفته شد و برای روابط از ماتریس استفاده نشده و به عنوان یک بردار انتقال در کنار موجودیت اول در نظر گرفته شده است که باعث می شود تعداد پارامتری که برای رابطه ها نیز داریم نیز d بعد باشد. پس تعداد پارامتری که در این روش باید آموزش دیده شود از مرتبه ی $O(n_e d + n_r d)$ است که در مجموع تعداد پارامترهای مساله را بسیار کاهش می دهد و همین مساله روی همین پایگاه دانش را با حدود ۸۰۰ هزار پارامتر حل می کند.

$$(n_e d + n_r d = 15k * 50 + 1.3k * 50 = 815k)$$

کیفیت خروجی: این روش علی رقم کاهش چشم گیری که در تعداد پارامترها داشت و مقیاس پذیری بسیار بالایی که ایجاد کرده است، در نتایج بدست آمده نیز بهتر از روش های معرفی شده ی قبلی کار کرده است. در این روش $hit@10$ برابر ۴۵,۱٪ و میانگین رتبه برابر ۱۲۵ است که مشاهده می شود در این روش هم جواب های درست بسیار بیشتر بوده و هم داده های پرت بسیار کمتر شده است و جواب سوال ها به جواب های منطقی نزدیک تر شده است.

۵-۲-۵- روش TransH:

تعداد پارامتر: همانطور که در بخش - دیدیم این روش از نظر عملیات آموزش بسیار شبیه به روش TransE است و در تعداد پارامترهایی که باید آموزش داده شود فقط یک بردار انتقال به ابرصفحه ی

مربوط به رابطه را بیشتر دارد که به ازای هر رابطه k پارامتر به عملیات آموزش اضافه می کند، پس پارامترهایی که باید آموزش دیده شوند از مرتبه $O(n_k + 2n_r k)$ هستند. تعداد پارامتر این روش روی پایگاه دانش freebase15k حدود 880 خواهد بود $(n_d + 2n_r d = 15k * 50 + 2 * 1.3k * 50 = 880k)$. مشاهده می شود که این روش افزایش چشمگیری در تعداد پارامترهای مساله نداشته اما نتایج بسیار دقیق تری بدست آورده است.

کیفیت خروجی: این روش با در نظر گرفتن روابط پیچیده تری مانند روابط یک به چند و چند به چند و چند به یک و بازگشتی توانسته است که دقت بهتری نسبت به روش های نام برده شده ی قبلی به دست آورد و بدون افزایش غیر منطقی تعداد پارامترها نسبت به روش TransE به مقدار $\text{Hit}@10$ برابر با 64.4% و میانگین رتبه ی ۸۷ برسد.

۳-۵- نقاط قوت و ضعف روش های موجود

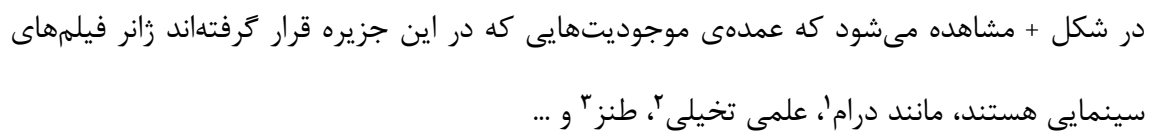
در این بخش به بررسی نقاط قوت و نقاط ضعف روش های موجود می پردازیم. در بخش + به تفصیل در مورد مقیاس پذیری و معیارهای کیفیت روش ها صحبت شد، در این بخش میزان همبستگی و ارتباط موجودیت ها و رابطه ها را بررسی می کنیم که چه مقدار این همبستگی و ارتباط ها در مدل های آموزش دیده شده درک شده است.

روش هایی که معرفی شدند در بهترین حالت تا 64.4% دقت در یافتن پاسخ در ۱۰ جواب اول ($\text{hit}@10$) را داشتند. نمونه هایی از این سوال و ۱۰ جواب اول مدل TransE برای آن را در جدول + مشاهده می کنیم.

Lil Wayne born in?!	New Orleans, Atlanta, Austin, St. Louis, Toronto, New York City, Wellington, Dallas, Puerto Rico
(Lil_Wayne, born_in, ?)	
WALL-E has the genre?!	Animations, Computer Animation, Comedy film, Adventure film, Science Fiction, Fantasy, Stop motion, Satire, Drama
(WALL-E, has_genre, ?)	

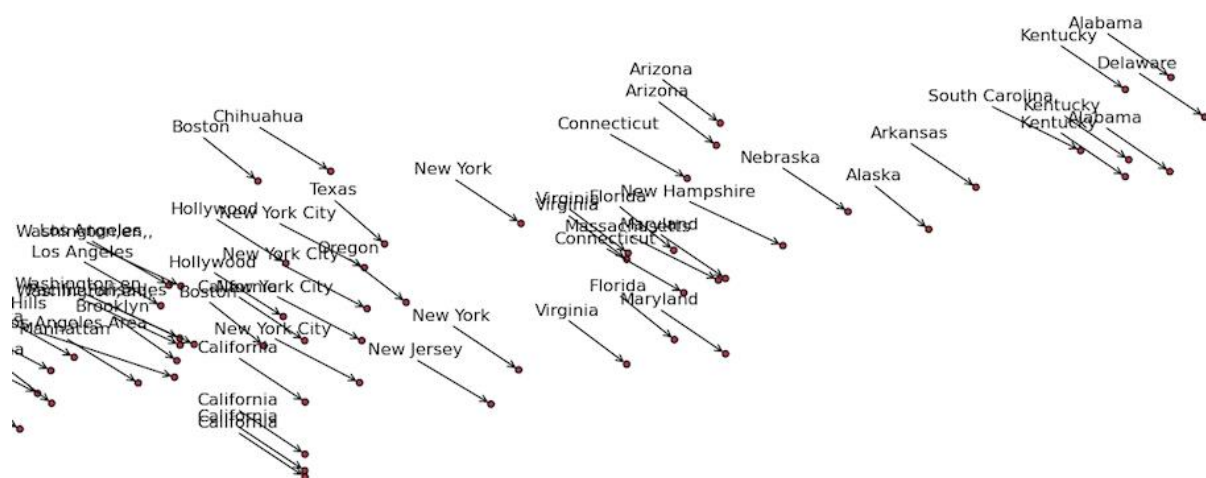
همانطور که مشاهده می شود جواب هایی که داده شده است همبستگی خوبی دارند، برای مثال در نمونه ی اول که از مدل پرسیده شده است Lil Wayne که یک خواننده ی آمریکایی است متولد کجاست، تمام ۱۰ جواب اول همانطور که انتظار می رود ایالت های مختلف آمریکا هستند و پاسخ اول یعنی New Orleans جواب درست می باشد و جواب ها منطقی هستن و مدل TransE این رابطه که پاسخ سوال “متولد کجا است؟” باید یک شهر باشد را خوب فهمیده است. همینطور در مثال دوم وقتی سوال اینکه “ژانر یک فیلم چیست؟” از مدل پرسیده شده است، جواب ها همبستگی خوبی دارند و جواب پرت داخل آن ها نیست.

برای مشاهده ی بهتر این موضوع، پس از آموزش دیدن مدل TransE، بردارهای ۵۰ بعدی ۵ هزار موجودیت را به وسیله ی ابزار tnse در شکل + و در دو بعد نمایش دادیم.



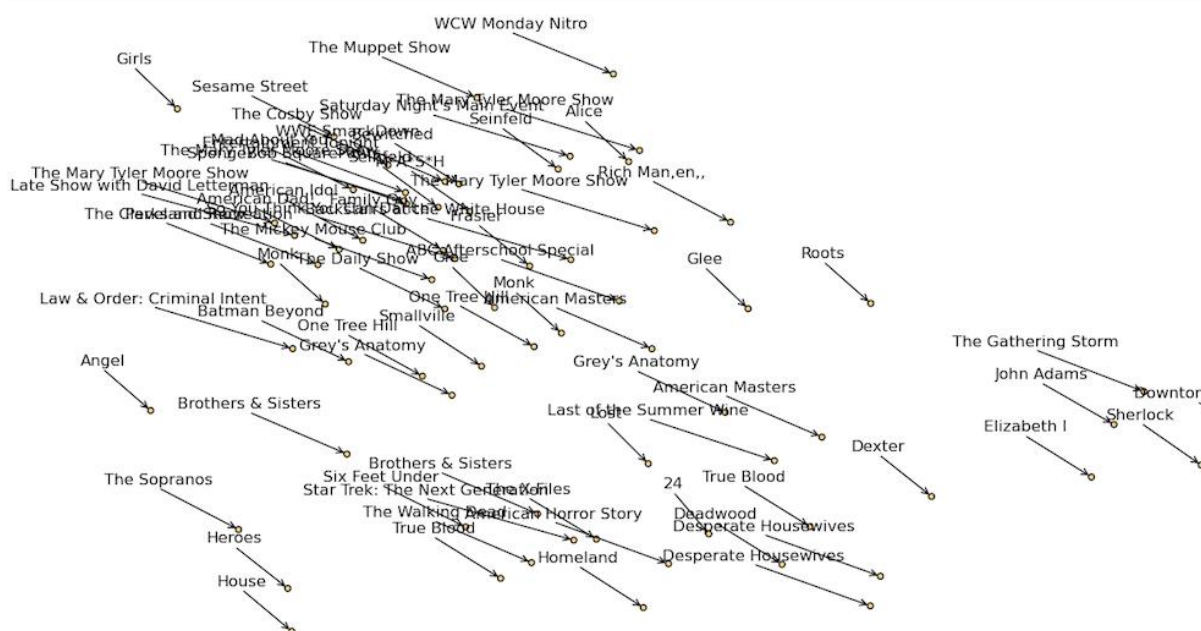
³ Comedy

Table 1



در شکل + مشاهده می شود که ایالت های آمریکا یک جزیره را تشکیل داده اند، مانند آریزونا، کالیفرنیا،

تگزاس و ...



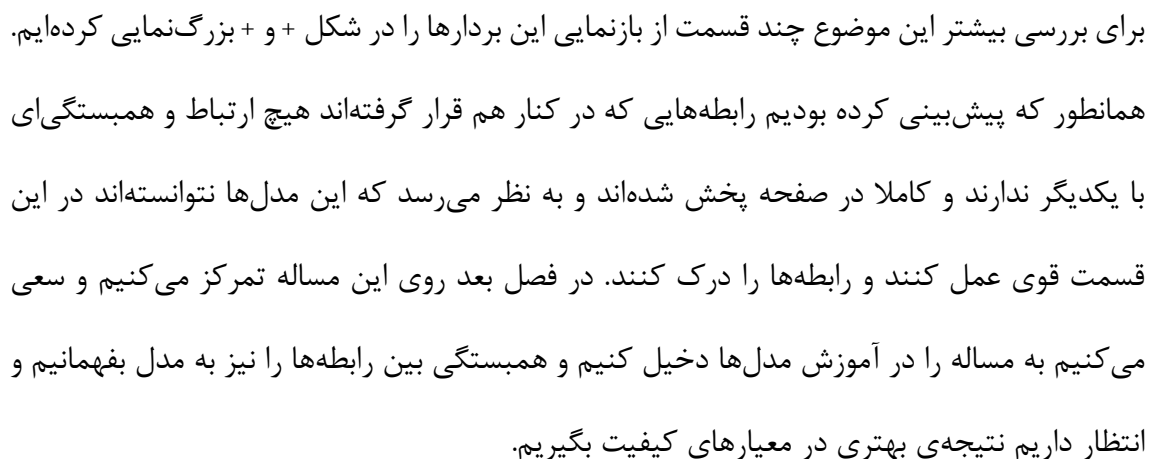
در شکل + مشاهده می شود که نام سریال های تلویزیونی یک جزیره را تشکیل داده اند.

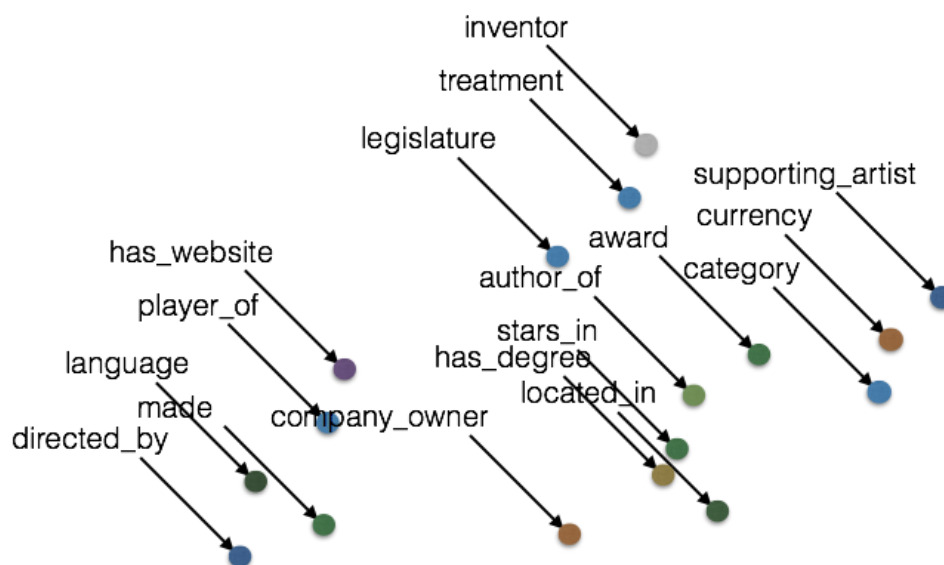
در مثال هایی که گفته شد، مشخص است که رابطه ی بین موجودیت ها به خوبی در این روش ها درک شده است و یک همبستگی بین بردارهای آموزش دیده شده وجود دارد.

می دانیم که رابطه ها هم مانند موجودیت ها می توانند ارتباطات و همبستگی های خود را داشته باشند. برای مثال رابطه ی “پدر بودن” نزدیک تر است به رابطه ی “فرزندی” نسبت به رابطه ی “نویسنده ی کتاب” بودن و انتظار داریم بردارها و ماتریس های آموزش دیده شده برای رابطه ها این همبستگی ها را درک کرده باشد و در فاصله ی کمی از یکدیگر قرار گرفته باشند و همان دسته بندی و جزیره شدنی که برای موجودیت ها اتفاق افتاده بود را اینجا نیز شاهد باشیم.

حال همین مساله را در بردارهای آموزش دیده شده بررسی می کنیم. برای بررسی این امر در شکل + بردارهای همه رابطه های استفاده شده در Freebase15k را در دو بعد به نمایش در آورده ایم. همانطور که در شکل مشخص است، برای رابطه ها اتفاق جزیره ای شدن که در بازنمایی موجودیت ها اتفاق افتاده بود، به وجود نیامده است و بردارها روی صفحه بخش¹ است. به نظر می رسد که در بردارهای آموزش دیده برای رابطه ها این همبستگی و ارتباط وجود ندارد.

¹ Scatterd





۴-۵- بررسی عمل کرد روش ها در یک نگاه

در این قسمت روش های موجود را ...

در جدول + تعداد پارامترهایی که هر یک از روش های بالا برای آموزش دیدن نیاز دارند، آورده شده است، همانطور که مشاهده می شود روش های Rescal و NTN تعداد پارامتر خیلی زیادی باید آموزش دهند و به مشکل مقیاس پذیری بر خواهند خورد و روی پایگاه دانش های بزرگ خیلی کند خواهند بود. در مقابل روش های SE و TransE و TransH با ایده هایی که برای کاهش تعداد پارامترهای مساله پیاده سازی کردن این تعداد را بسیار کاهش داده و مقیاس پذیری خوبی دارند، به صورتی که مدل TransE روی پایگاه دانش Freebase 1M در مدت تقریباً یک روز با $hit@10$ برابر با ۳۴٪ به جواب رسیده است در صورتی که روش های Rescal و NTN روی این پایگاه دانش به جواب نمی رسند، که این موضوع بر عدم مقیاس پذیری روش های Rescal و NTN تاکید می کند.

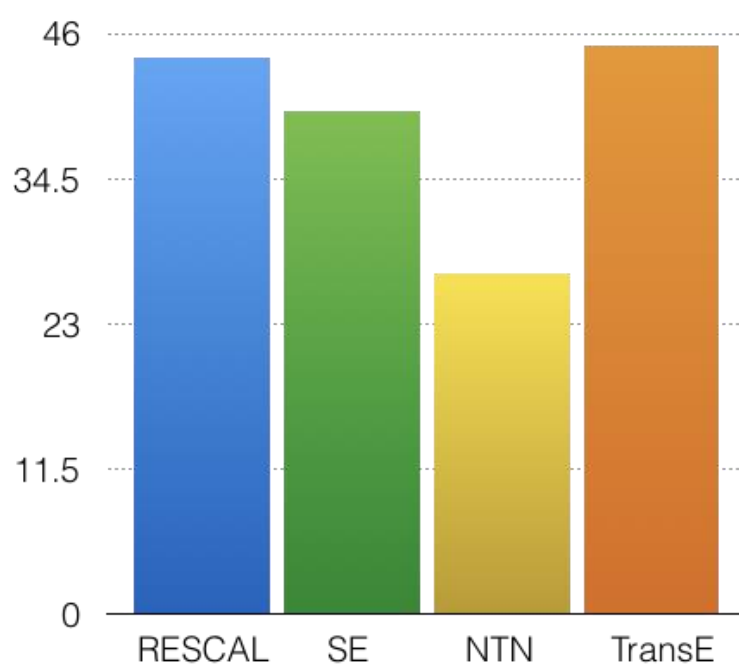
Method	#Params	On FB15K
RESCAL	$O(n_e d + n_r d^2)$	88M (d=250)
MLP (NTN)	$O(n_e d + n_r d^3)$	165M (d=50)
SE	$O(n_e d + 2n_r d^2)$	8M (d=50)
TransE	$O(n_e d + n_r d)$	0.8M (d=50)

در جدول + میزان Hit@10 و میانگین رتبه ی هر یک از روش های معرفی شده نمایش داده شده است که در این جدول مشاهده می شود که روش NTN در hit@10 نتیجه ی مطلوبی نگرفته است که در معرفی این روش اشاره کردیم که این روش به سرعت به بیش برآزش میل می کند و همین امر باعث شده است که پاسخ های درستی در مورد مساله نداشته باشد.

نکته ی مهم دیگری که در این جدول دیده می شود کارایی خوب روش های SE و TransE هست که علاوه بر کاهش پارامتری که داشته اند، هم از نظر hit@10 و هم از نظر میانگین رتبه نتایج خوبی داشته اند.

در روش RESCAL هم hit@10 خوب بوده است اما رتبه ی میانگین مقدار بسیار بالایی نسبت به روش های دیگر داشته که نشان از این دارد که علی رقم بسطی از روابط را درست فهمیده است و نتایج خوبی در آن ها گرفته است، اما بسطی از روابط را نیز اصلا درک نکرده و پاسخ هایی که به سوالات در مورد این روابط داده است جواب های تقریباً تصادفی ای داشته است و رتبه های پرتی گرفته اند که باعث شده میانگین رتبه ی این روش چنین افزایشی داشته باشد.

	Hit@10	Mean
RESCAL	42.1%	683
SE	39.8%	162
NTN	27%	164
TransE	45.1%	125



۵-۵- استفاده از قوانین انجمنی برای بهبود نتایج

در این فصل به طور کامل قوانین انجمنی و چگونگی استخراج آن ها را شرح دادیم و پنج قانون هدف نیز مشخص کردیم که به کمک آن ها تصمیم داریم مدل های موجود مبتنی بر ویژگی پنهان را بهبود دهیم. در این بخش به چگونگی انجام این بهبود می پردازیم.

در برخی از روش های معرفی شده در بخش + ابتدا همه ی این قانون ها را روی حقیقت های مجموعه ی آموزش اعمال کردیم و حقایق جدید بدست آمده را دوباره به پایگاه دانش اضافه کردیم و سپس از پایگاه دانش جدید که بزرگ تر، کامل تر و دقیق تر شده است استفاده کردیم.

اما در روش های TransE و TransH این قوانین را مستقیماً وارد فرایند آموزش کردیم. به این صورت که پایگاه دانش را گسترش ندادیم و با همان حدود ۵۰۰ هزار حقیقت اولیه عملیات آموزش را شروع کردیم اما در حین آموزش از این قوانین استفاده کردیم تا نتایج دقیق تری بگیریم که در ادامه به این موضوع خواهیم پرداخت که این قوانین چگونه استفاده شده اند.

عدم اعمال این قوانین روی پایگاه دانش باعث شد که هم مساله ی پیش پرداز برای اعمال قوانین روی پایگاه دانش را نداشته باشیم و هم مساله را با پایگاه دانش کوچک تری شروع کنیم که در دو مرحله باعث کاهش زمان محاسبات می شود.

در ادامه توضیحات استفاده از این قوانین را روی روش TransE می دهیم. در روش TransE هدف کاهش فاصله ی بین بردار موجودیت اول + بردار رابطه و بردار موجودیت دوم بود که برای حقیقت $r(e1, e2)$ این رابطه را به صورت زیر می نوشتیم.

$$\| e_1 + r - e_2 \|$$

۱-۵-۵- قانون R-subsumtion

برای قانون R-subsumtion که به صورت $r(x, y) \Rightarrow r'(x, y)$ است در عملیات آموزش هر جا عملیات کمینه سازی روی حقیقت $r(x, y)$ انجام گرفت آن را روی روابط هم ارز آن که از این قانون بدست می آیند نیز انجام می دهیم، یعنی روابط $r'(x, y)$. پس در کمینه سازی ها علاوه بر کمینه کردن رابطه ی $x + ||$ رابطه ی $r - y ||$ را نیز کمینه می کنیم.

۲-۵-۵- قانون R-equi val ence

در قانون R-equivalence که همانند رابطه ی قبل است با این تفاوت که قانون برگشت پذیر نیز هست، یعنی $r(x, y) \Leftrightarrow r'(x, y)$ می توان بصورت بالا عمل کرد و هر قسمت (سمت چپ و راست قانون) از این قانون در حقیقت ها دیده شد قسمت دیگر را نیز در معادله ی کمینه سازی قرار دهیم. اگر معیار اطمینان این قانون بالا باشد، عملاً اتفاقی که می افتد این است که رابطه های r و r' به یک شکل آموزش دیده می شوند و به یک صورت عمل خواهند کرد، که با توجه به هم معنا بودن رابطه های مثل `located_in` و `contained_by` این اتفاق، اتفاق منطقی ای خواهد بود و از آن امید بهبود در مدل ها را داریم.

۳-۵-۵- قانون تعدی (2-hope translation)

قانون تعدی یا 2-hope translation که به صورت

$$r_1(e_1, e_2), r_2(e_2, e_3) \Rightarrow r'(e_1, e_3)$$

تعریف می شود، نسبت به قانون های قبلی پیچیدگی محاسباتی بیشتری خواهد داشت، زیرا پس از پیدا کردن دو حقیقت در سمت بدنه ی قانون قادر به اعمال قانون خواهیم بود. پس در هر مرحله از آموزش

که به یکی از حقیقت های سمت چپ قانون رسیدیم، عملیات جستجو برای حقیقت دیگر را شروع می کنیم و در صورت پیدا شدن حقیقت دوم این قانون را اعمال می کنیم. (عملیات جستجو برای حقیقت دوم به صورت موازی انجام می پذیرد و باعث کاهش سرعت آموزش نخواهد شد).

پس از یافتن هر دو حقیقت سمت بدنه ی قانون یعنی $r_1(e_1, e_2)$ و $r_2(e_2, e_3)$ برای اعمال قانون در آموزش باید عملیات کمینه سازی را برای حقیقت $r'(e_1, e_3)$ نیز انجام دهیم یعنی عبارت $\|e_1 + r' - e_3\|$ را نیز کمینه کنیم.

از طرفی روابط را در روش TransE به صورت یک بردار انتقال در نظر گرفتیم پس انتظار داریم که مجموع انتقالی که بردارهای r_1 و r_2 ایجاد می کنند برابر باشد با انتقال بردار r' یعنی:

$$r_1 + r_2 = r'$$

پس از عکس این رابطه نیز می توانی استفاده کنیم و اگر جایی به قانون r' را مشاهده کردیم در کنار کمینه کردن مقدار

$$\|e_1 + r' - e_3\|$$

مقدار مجموع دو بردار دیگر را نیز کمینه کنیم:

$$\|e_1 + r_1 + r_2 - e_3\|$$

۴-۵-۵- Triangle alignment قانون

در این قانون نیز مانند قانون قبل می توان عمل کرد و با پیدا کردن هر یک از حقیقت های موجود در بدنه ی قانون، دومین حقیقت را جستجو کرد و در صورت پیدا کردن آن، علاوه بر کمینه سازی که روی خود حقیقت های اصلی انجام می دهیم کمینه سازی را روی حقیقت بدست آمده از قانون نیز اعمال کنیم.

$$r_1(e_1, e_3), r_2(e_2, e_3) \Rightarrow r'(e_1, e_2)$$

$$\Rightarrow ||e_1 + r' - e_2||, ||e_1 + r_1 - e_3||, ||e_2 + r_2 - e_3||$$

۵-۵-۵- Specefic R-subsumption قانون

در این قانون همچون دو قانون قبل باید دو حقیقت که در سمت بدنه ی قانون آمده اند را یافته و پس از آن نتیجه ی قانون را در عملیات آموزش تاثیر دهیم. این قانون را به این صورت تاثیر می دهیم که اگر دو قانون $r_1(x, y)$ و $r_2(x, V)$ در حقایق وجود داشت، علاوه بر کمینه کردن فاصله ی بین این دو حقیقت عبارت نتیجه ی این قانون را نیز در عملیات کمینه سازی شرکت می دهیم به صورت

$$||x + r' - y||$$

فصل ۶: ارزیابی

۱-۶- مقدمه

در فصل قبل قوانین انجمنی و نحوه ی استخراج، و چگونگی استفاده از آن ها را توضیح دادیم. در این فصل این قوانین انجمنی را وارد عملیات آموزش کرده و نتایج بدست آمده را بررسی و مقایسه می کنیم. در بخش + به تشریح آزمایش ها و نتایج بدست آمده می پردازیم، در بخش + به بررسی نتایج و مقایسه آن ها اختصاص می دهیم و در بخش + به کارهایی که در آینده می توان برای بهبود بیشتر نتایج انجام داد را معرفی می کنیم.

۲-۶- آزمایش ها

در این تحقیق هدف فقط بهبود نتایج روش های مبتنی بر ویژگی های پنهان نبوده و به دنبال ارائه ی چارچوبی هستیم که به کمک آن بتوانیم قدرت و نقاط قوت و ضعف روش ها را نیز مشخص کنیم. برای ارائه ی این چارچوب از بخش بندی ای که در فصل + روی قوانین انجمنی داشتیم استفاده می کنیم. و بررسی می کنیم که هر یک از این دسته قوانین چه مقدار توسط مدل درک شده اند و خوب جواب می دهند.

برای انجام آزمایش ها باید پایگاه دانش و قوانین انجمنی که می خواهیم عملیات آموزش را روی آن ها انجام دهیم را مشخص کنیم. پایگاه دانش استفاده شده همان پایگاه دانش Freebase15k است که در

بخش + معرفی شد (حدود ۶۰۰ هزار حقیقت، ۱۵ هزار موجودیت و ۱۳۰۰ رابطه) است. در ادامه قوانین انجمنی را روی این پایگاه دانش اجرا می کنیم و قوانین بدست آمده را دسته بندی می کنیم.

همانطور که در بخش + اشاره شد، پس از اجرای قوانین انجمنی روی پایگاه دانش ۴۱۱۹۶ Freebase15k قانون استخراج شد اما همه ی این قانون ها از نظر معیار کیفیت و اطمینان، شرایط استفاده در عملیات آموزش را ندارند، نیاز است تا قوانین با کیفیت را مشخص کنیم و فقط از آن ها استفاده کنیم. برای مشخص کردن قوانین با کیفیت با آزمون و خطا به معیار اطمینان PCA برابر ۰,۸ و پوشش سر ۲۰۰ رسیدیم و فقط قوانینی که معیار اطمینان بالای ۰,۸ و پوشش سر بیش از ۲۰۰ داشتند را انتخاب و در عملیات آموزش شرکت دادیم. پس از اعمال این محدودیت ها تعداد قوانین کاندید بدست آمده برابر با + قانون شد. این قوانین را بر اساس تقسیم بندی های بخش + دسته بندی کردیم که این دسته بندی در جدول + مشخص شده است.

Rules		
name	rule	#rule
R-Subsumption	$r(x, y) \Rightarrow r'(x, y)$	1127
R-Equivalence	$r(x, y) \Leftrightarrow r'(x, y)$	782
2-hope	$r_1(x, y), r_2(y, z) \Rightarrow r'(x, z)$	810
Triangle alignment	$r(x, z), r(y, z) \Rightarrow r'(x, y)$	695
SR-Subsumption	$r_1(x, y), r_2(x, v) \Rightarrow r'(x, y)$	779
all rules		4193

همانطور که قبل تر هم اشاره شد تصمیم داریم که علاوه بر بهبود نتایج یک چارچوب برای سنجش کیفیت روش ها نیز ارائه دهیم، برای این کار طبق آنچه در بخش + گفته شد این قوانین را به صورت دسته ای به روش های آموزش اضافه می کنیم و نتایج را مشاهده می کنیم و در آخر نیز همه ی قوانین را باهم به روش ها اضافه می کنیم و میزان بهبود انجام شده توسط این قانون ها را در روش های مختلف بررسی می کنیم.

در نمودار + معیار $hit@10$ و در نمودار + معیار میانگین رتبه برای روش های معرفی شده در +، قبل و بعد از اضافه کردن قوانین استخراج شده در دسته ی R-Subsumption را مشاهده می کنیم.

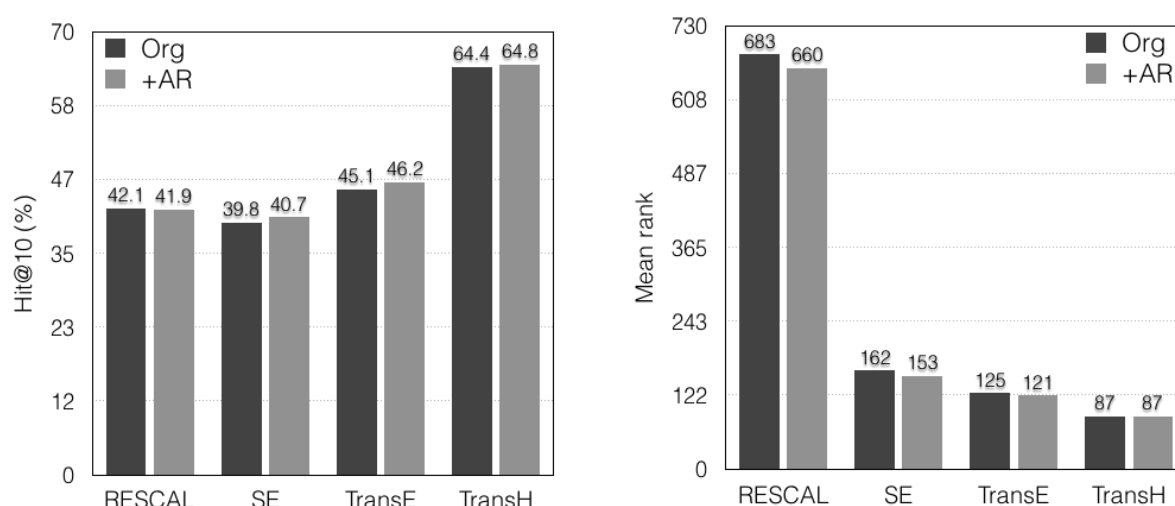


Figure 1 - R-sub

در نمودارها مشاهده می شود که اضافه کردن قوانین R-subsumption در هیچ یک از معیارها تاثیر چندان چشم گیری نداشته و بهبودهای جزئی روی نتایج روش ها داشته است. این مساله نشانگر این است که هر ۴ روش معرفی شده در فهمیدن قوانین R-subsumption مشکلی نداشته اند و در زمان آموزش این روابط توسط مدل شناسایی و درک شده است.

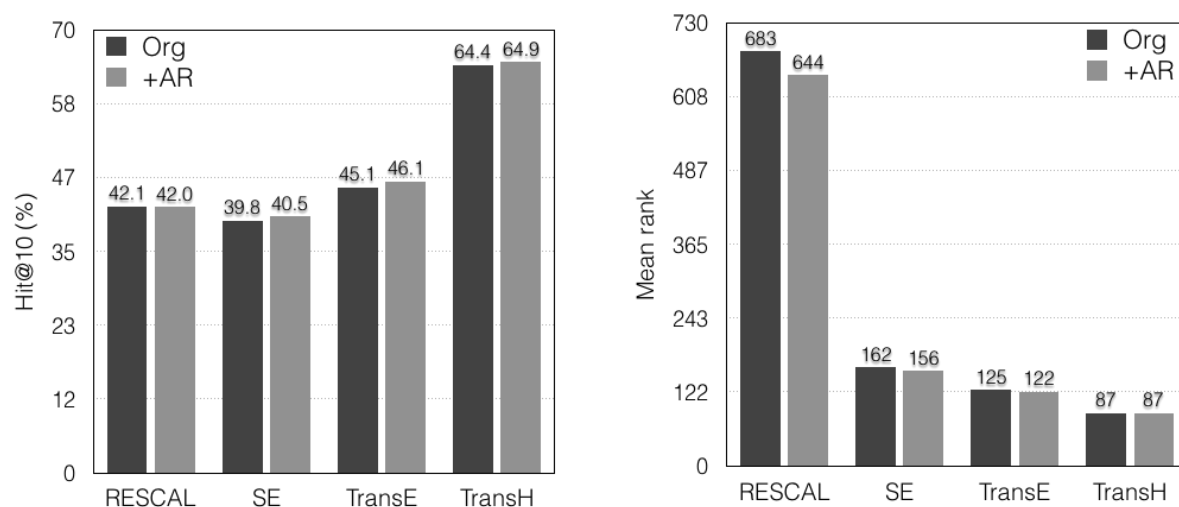
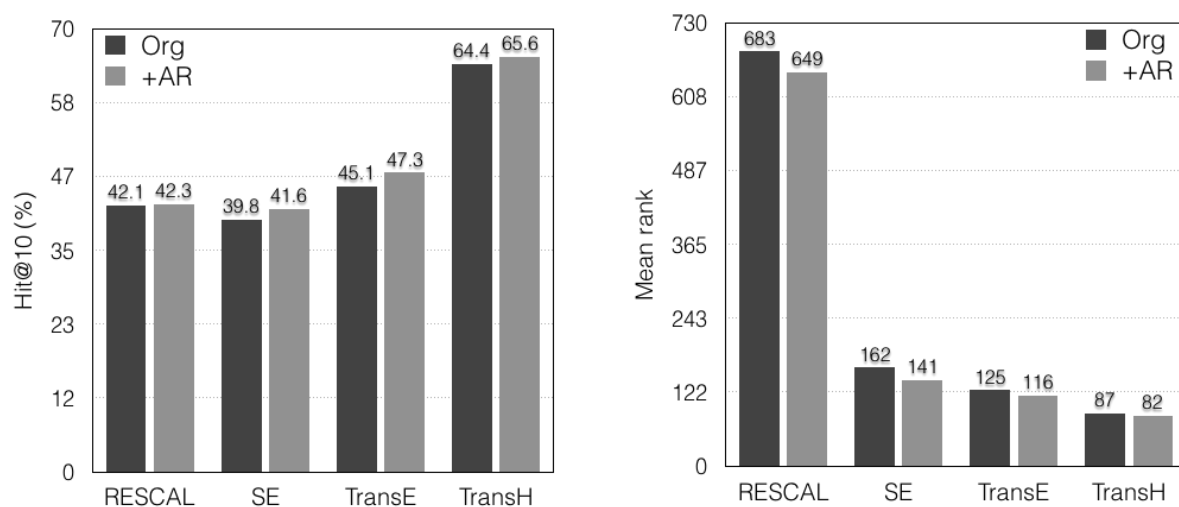


Figure 2- SR-sub

در نمودارهای + و + نتایج معیارهای کیفیت روش های آموزش را قبل و بعد از اضافه کردن قوانین SR-subsumption مشاهده می کنیم. این قانون ها نیز مانند قانون های R-subsumption تاثیر چندانی در معیارهای hit@10 و میانگین رتبه نداشته است.

می توان نتیجه گرفت که مدل های مبتنی بر ویژگی های پنهان در کل قوانینی که به طور مستقیم از روی یکدیگر نتیجه گرفته می شوند را به خوبی درک می کنن و مدل آموزش دیده این روابط رو به خوبی تشخیصی می دهد و در پاسخ به سوالات به اشتباه نمی افتد.

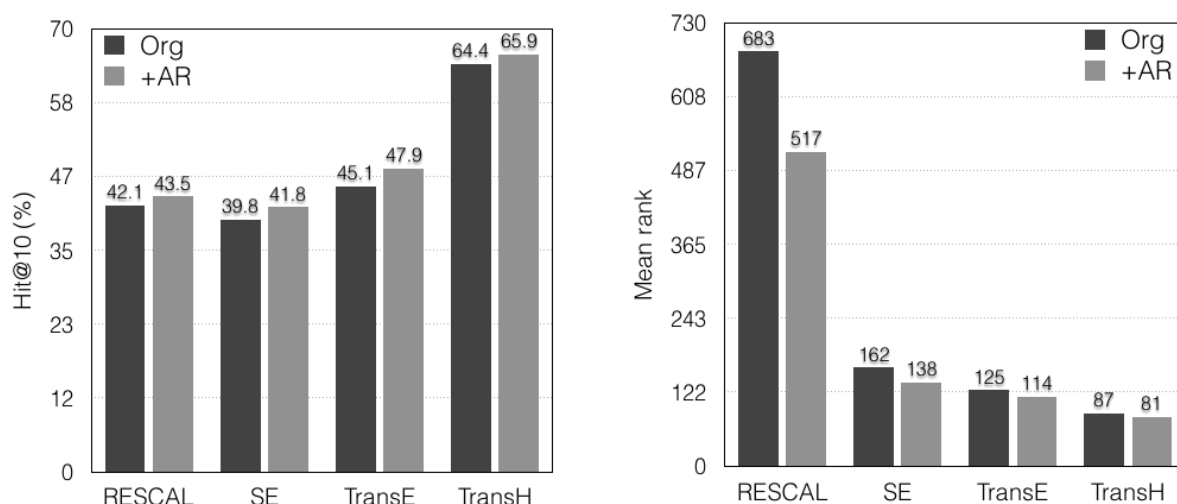


R-Equ-Figure 3

در نمودارهای + و + مشاهده می شود اضافه کردن قوانین R-equivalence تاثیر چندانی روی نتایج RESCAL نداشته است و بهبود جزئی داشته است، اما در روش های دیگر تاثیر نسبتا خوبی داشته است و باعث بهبود نتایج شده است.

همانطور که در بخش + عنوان کردیم این قوانین به کشف روابطی می پردازند که به دو صورت بیان شده اند ولی معنای یکسانی دارند، مانند روابط `located_in` و `contained_by` که هر دو معنی قرار داشتن یک مکان در مکان دیگر را دارد اما به دو صورت بیان شده است. اینجا مشاهده می شود که تعداد زیاد این چنین رابطه ها در پایگاه دانش و عدم شناسایی آن ها توسط مدل های آموزش تاثیر منفی ای در نتایج داشته است که با اضافه کردن قوانین R-equivalence این مشکل رفع شده و بهبود نسبتا خوبی حاصل شده است.

عدم بهبود در روش RESCAL نیز به این دلیل است که این روش همانطور که در توضیح آن در بخش + توضیح دادیم همه ی جفت ویژگی های پنهان از دو موجودیت را با یکدیگر مقایسه و بررسی می کند و می تواند به مقدار خوبی این روابط یکسان را ببیند و در ماتریس های رابطه قرار دهد.



در نمودارهای + و + با اضافه کردن قوانین 2-hope translation نیز در هر دو معیار بهبود محسوسی دیده می شود. مشاهده می شود که در روش RESCAL نیز که با اضافه کردن قانون های قبلی تغییر چندانی نکرده بود، با اضافه کردن این قانون بهبود بسیار زیادی داشته ایم و مقدار hit@10 به ۴۳,۵ و مقدار میانگین رتبه با ۱۶۶ رتبه کاهش به ۵۱۷ رسیده است که بهبود چشم گیری است. همچنین در دیگر مدل ها بهبودهای زیادی را داشته ایم، در روش های SE و TransE و TransH هر کدام به ترتیب ۲ و ۲,۸ و ۱,۵ درصد رشد را داشته ایم. همچنین در این آزمایش میانگین روش TransH از ۸۷ به ۸۱ رسیده است که با توجه به پایین بودن میانگین اولیه و سخت بودن تغییر در آن، بهبود بزرگی محسوب می شود.

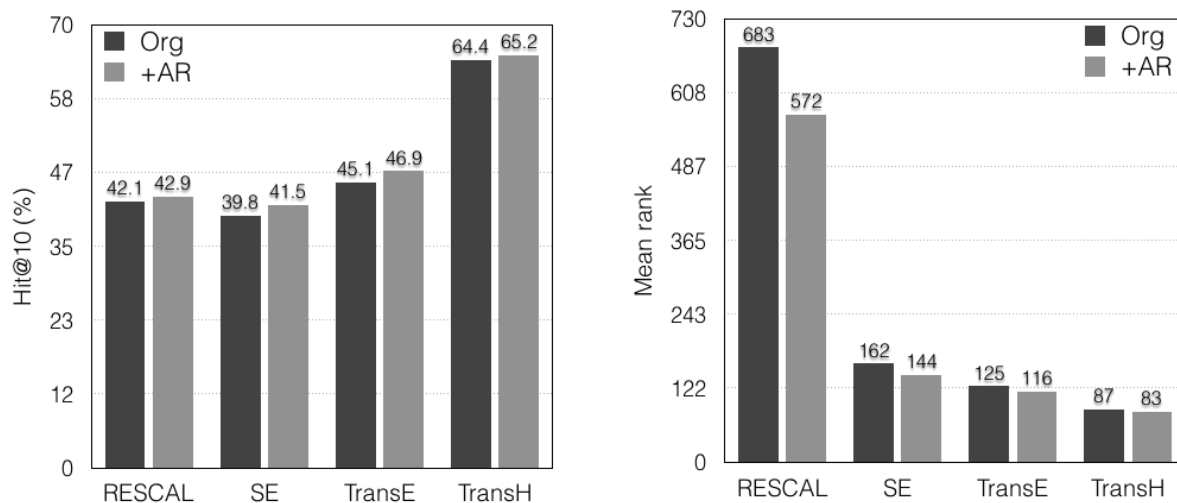
می توانیم نتیجه بگیریم که هیچ یک از این مدل ها قادر به شناسایی قوانین غیر مستقیم پیچیده نبوده است. در آزمایش های قبلی دیدیم که قوانینی غیر مستقیم که با استنتاج از روی یک رابطه ساخته شده بودند (در قسمت بدنه ی قانون فقط یک حقیقت وجود داشت) را مدل ها بهتر درک کرده بودن و اضافه کردن قوانین مربوط به این آزمایش ها تاثیر چندانی در نتیجه نداشت، اما اضافه کردن قوانین غیر مستقیم پیچیده تری مانند قانون 2-hope translation (که در سمت بدنه ی قانون از دو حقیقت استفاده شده است) بهبودهای زیادی را در همه ی مدل ها اعمال کرده است. پس به کمک این قانون ها توانستیم چند

قدم جلوتر از آنچه مدل ها قادر به دیدن آن بودند را به آن ها نشان دهیم و قوانین پیچیده ی موجود در پایگاه های دانش را نیز در امر آموزش دخیل کنیم.

بهبود در معیار $hit@10$ نشان می دهد که اضافه کردن این قوانین باعث شده که در جواب های دقیق که به سوالات داده شده است بهبود داشته باشیم و تعداد جواب هایی صحیح که زیر ۱۰ کاندید اول بودند را افزایش دهیم و از طرفی بهبود در معیار میانگین رتبه نشان می دهد که در رابطه هایی که مدل برای آن ها درست آموزش ندیده است و جواب های پرتی برای آن ها در نظر گرفته است، اوضاع پاسخ ها بهتر شود.

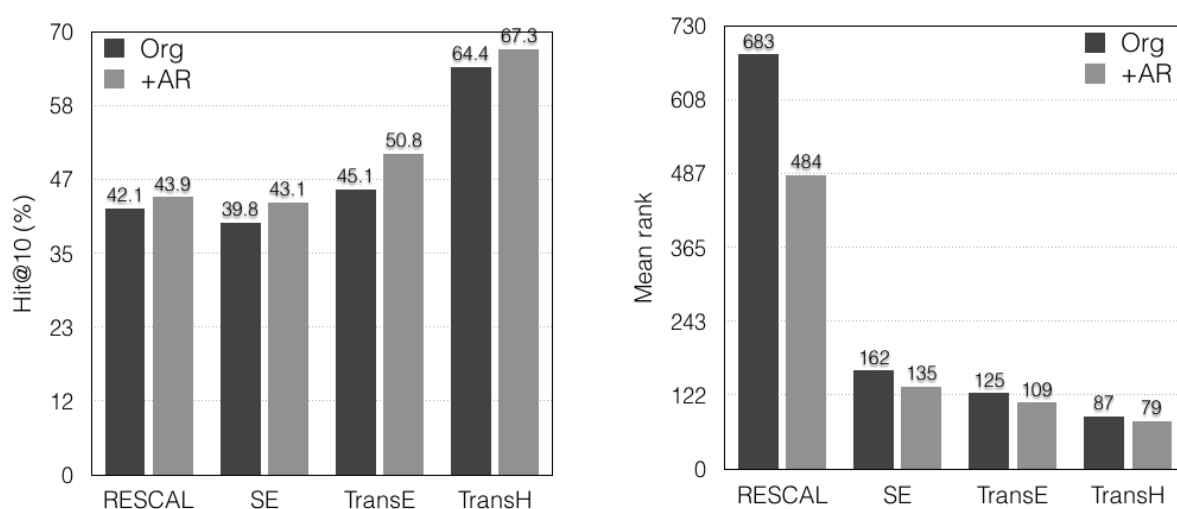
به طور خاص یکی از دلایل بالا بودن معیار میانگین رتبه این است که برای برخی روابط حقیقت های کمی در پایگاه دانش موجود است و مدل ها از روی این تعداد کم رابطه قادر به شناسایی الگو برای پاسخگویی به سوالات در مورد این حقیقت ها و روابط نیستند. همین امر باعث می شود که پاسخ هایی که به سوالات در مورد این حقیقت ها می دهند به صورت تصادفی باشد و رتبه های پرتی بگیرد. این رتبه های پرت باعث افزایش زیادی در معیار میانگین رتبه می شود، در روش RESCAL این مساله را به خوبی مشاهده می کنیم که با وجود اینکه قریب به ۴۲٪ جواب ها زیر رتبه ی ۱۰ قرار می گیرند اما میانگین رتبه ای که برای این روش اعلام شده است ۶۸۳ است.

این مشکل عدم وجود تعداد نمونه ی کافی برای آموزش مدل را تا حدی می توان با اضافه کردن نمونه های غیر مستقیم موجود در پایگاه دانش حل کرد، به این معنی که برای خیلی از روابط نمونه هایی وجود دارد که از روی بقیه ی حقایق درون پایگاه دانش می توان به آن ها پی برد و از آن ها استفاده کرد. در اینجا هم ما با اضافه کردن قانون های 2-hope translation چنین کاری صورت داده ایم و حقیقت هایی که از به وسیله ی قانون تعدی قابل شناسایی بوده اند را به پایگاه دانش اضافه کرده ایم و در عملیات آموزش تاثیر داده ایم و نهایتا نتایج خوبی نیز در پاسخ به سوالاتی که از مدل پرسیده می شود بدست آمده است.



در نمودارهای + و + نیز تاثیر استفاده از قوانین Triangle alignment را مشاهده می کنیم که بهبودهای این نتایج نیز به مانند آزمایش قبل قابل قبول بوده است و همان دلیل پیچیدگی قانون که مدل قادر به درک آن نبوده است در این مورد هم صدق می کند و باعث بهبود نتایج شده است.

در ادامه نتایج اعمال کل قانون هایی که در قسمت + معرفی کردیم را روی مدل های معرفی شده را بررسی می کنیم. نتایج بدست آمده در نمودارهای + و + نشان داده شده است.



مشاهده می شود که اعمال همه ی قانون ها در کنار هم نیز باعث بهبود در نتایج همه ی روش های معرفی شده، شده است و علاوه بر چارچوبی که برای سنجش نقاط قوت و ضعف روش ها معرفی کردیم توانستیم در همه ی روش ها بهبودهای قابل قبولی نیز اعمال کنیم.

۳-۶- جمع بندی

در جدول + و + جمع بندی ای داریم روی نتایج اعمال قوانین روی روش های معرفی شده که در بخش قبل به تفصیل توضیح داده شد و این نتایج را به تفکیک معیار $hit@10$ و میانگین رتبه را در کنار هم قرار داده ایم.

Hit@10 on Freebase 15K (%)							
Method	original	R-Subsumption	R-Equivalence	2-hope	Triangle alignment	SR-Subsumption	all rules (~improve)
RESCAL	42.1	41.9	42.3	43.5	42.9	42.0	43.9 (1.8%)
NTN	27	-	-	-	-	-	-
SE	39.8	40.7	41.6	41.8	41.5	40.5	43.1 (3.3%)
TransE	45.1	46.2	47.3	47.9	46.9	46.1	50.8 (5.7%)
TransH	64.4	64.8	65.6	65.9	65.2	64.9	67.3 (2.9%)

Mean rank on Freebase 15K							
Method	original	R-Subsumption	R-Equivalence	2-hope	Triangle alignment	SR-Subsumption	all rules (~improve)

RESCAL	683	660	649	517	572	644	484 (30%)
NTN	164	-	-	-	-	-	-
SE	162	153	141	138	144	156	135 (17%)
TransE	125	121	116	114	116	122	109 (13%)
TransH	87	87	82	81	83	87	79 (10%)

نکاتی که در این جداول قابل توجه هستند:

- اعمال همه ی قوانین روی هر چهار روش مورد آزمایش باعث بهبود قابل قبول نتایج گشته است.
- قوانین پیچیده تر مانند قانون 2-hope و قانون Triangle alignment که به طور مستقیم از روی داده های موجود در پایگاه دانش قابل برداشت نیستند توسط مدل ها به خوبی درک نشده اند و از نقاط ضعف روش ها به شمار می روند که در جدول مشاهده می کنیم اضافه کردن این قوانین بهبود خوبی در نتایج داشته است.
- قانون R-equivalence مانند دو قانون بند قبل قانون پیچیده ای نیست، اما به حل یکی از مشکلات اساسی پایگاه های دانش پرداخته است که وجود روابط مختلف با تعاریف یکسان است. این خاصیت این قانون که روابط یکسان را تشخیص و در عملیات آموزش تاثیر می دهد نیز باعث تاثیر خوبی در نتایج همه ی روش ها بجز روش RESCAL شده است که در بخش + دلیل عدم بهبود برای روش RESCAL را توضیح دادیم.

- قوانین ساده تر مانند R-subsumption و SR-subsumption تقریباً توسط همه ی مدل ها به خوبی درک شده است و اضافه کردن این قانون ها بهبود چندانی را در نتایج حاصل نکرده است.
 - به طور خاص در روش RESCAL بهبود در معیار $hit@10$ نسبت به دیگر روش ها کمتر بوده است (حدود ۱,۸٪) اما بهبود در معیار میانگین رتبه مقدار زیادی بوده است و این معیار را از ۶۸۳ به ۴۸۴ (حدود ۳۰٪) کاهش داده است. این مساله نشان می دهد که روش RESCAL برای روابطی که به میزان کافی نمونه برای آموزش دیدن داشته است خوب عمل کرده و جواب حدود ۴۰٪ از سوال ها را در رتبه ی زیر ۱۰ پاسخ داده است، اما برای مواردی که به میزان کافی نمونه برای آموزش نداشته است خیلی بد عمل کرده و باعث شده است که میانگین رتبه به مقدار زیادی بالا برود. اضافه کردن قوانین انجمنی به این روش نیز تاثیر چندانی در $hit@10$ نداشته است و برای روابطی که نمونه ی کافی از آن ها موجود بوده است مفید واقع نشده است، اما کاهش بسیار زیاد میانگین رتبه ی پاسخ ها نشان می دهد که کمک بسیاری به روابطی که نمونه ی کافی نداشته اند شده است و جواب های بیشتری به سمت منطقی شدن پیش رفته است.
 - متأسفانه پیاده سازی روش NTN کامل نشده و نتایج اعمال این قوانین روی این روش را برای مقایسه در اختیار نداریم.
- در این تحقیق ابتدا روش های موجود در زمینه ی پیش بینی پیوند که از ویژگی های پنهان استفاده می کردند معرفی شد و سپس با معرفی و استفاده از قوانین انجمنی سعی شد تا در این روش ها بهبودهایی ایجاد شود. با توجه به نتایجی که مشاهده کردیم و نکاتی که گفته شد در کل دست آوردهای این تحقیق را می توان به صورت زیر خلاصه کرد:
- استفاده از قوانین انجمنی در بهبود روش های مبتنی ویژگی های پنهان

- ارایه ی چارچوبی برای مقایسه و بررسی عملکرد و نقاط ضعف و قوت روش های موجود مبتنی بر ویژگی های پنهان که قابل گسترش به روش های دیگر در زمینه ی پیش بینی پیوند نیز هست.

فصل ۷: نتیجه گیری

۷-۱- نتیجه گیری

متن

۷-۲- کارهای آینده

با روش هایی که معرفی کردیم توانستیم علاوه بر چارچوبی که برای سنجش کیفیت روش های موجود و یافتن نقاط و قوت و ضعف این روش ها معرفی کردیم، بهبودهایی در همه ی روش های معرفی شده داشته باشیم. اما با وجود بهبودهایی که داشتیم بهترین نتیجه ای که بهترین روش به ما داده است مقدار $hit@10$ برابر با ۶۷,۳٪ است که برای استفاده های واقعی از این روش ها مقدار خوب و قابل قبولی نیست. این عدد به این معناست که اگر ما سوالی از این مدل بپرسیم به احتمال ۳۳٪ جواب اشتباه می دهد و اگر جواب درست بدهد نیز باید این جواب درست را از بین ۱۰ جواب کاندید انتخاب کنیم. اینکه این ۱۰ جواب کاندید با یکدیگر در ارتباط هستند این نوید را می دهد که از این روش ها در سیستم های توصیه گر که لیستی از موارد را به کاربر توصیه می کند به خوبی قابل استفاده هستند، اما اگر بخواهیم از این مدل ها برای پاسخ به سوالات استفاده کنیم بجای $hit@10$ باید مقدار $hit@1$ را در نظر بگیریم، زیرا فقط یک جواب صحیح مدنظر است و دیگر جواب ها فاقد اعتبارند و جواب های درست به شدت کاهش پیدا می کند.

با این که این روش ها در حال حاضر در آخرین تکنولوژی های شرکت های بزرگی مانند گوگل و IBM در حال استفاده هستند، اما این روش ها هنوز باید بسیار دقیق تر شوند تا در کاربردهایی مثل پاسخگویی به سوالات نیز قابل استفاده باشند.

بهبودهایی که می توان روی این مساله داد می توانند از طریق راه کارهای زیر باشند:

- استفاده از روشی که در این تحقیق ارائه شده و یافتن نقاط ضعف روش ها و انتخاب راه حل برای حل این ضعف ها
- در این تحقیق فقط قوانین انجمنی با پیچیدگی ۱ و ۲ (قوانینی که در قسمت بدنه ی خود دو حقیقت را بررسی می کردند) در نظر گرفته شدند، مطمئنا درک قوانین با پیچیدگی های بیشتر نیز برای مدل های موجود سخت بوده و ممکن است نقاط ضعف این روش ها باشند، پس اضافه کردن قوانین با پیچیدگی بیشتر از ۲ می تواند بهبودهایی را حاصل کند.
- راهکار دیگری برای پایین آوردن معیار میانگین رتبه، فیلتر کردن نتایج روش ها است. می دانیم پاسخ هایی که روش های موجود به سوالات می دهند به صورت لیستی از جواب های مرتب شده است که ممکن است جواب درست سوال مورد ۵۰ام باشد. با بررسی جواب ها مشخص می شود که خیلی از پاسخ هایی که در ۴۹ پاسخ غلط وجود داشته کلا بی ارتباط با سوال است. مثلا سوال شده که «نویسنده ی کتاب x چه شخصی است؟» و انتظار داریم جواب های دریافتی نام اشخاص باشد، اما پاسخ هایی از جنس های دیگر مانند نام کشورها، وضعیت های آب و هوا و ... در پاسخ های غلط پیش از پاسخ درست وجود دارد. برای حل این مشکل و فیلتر کردن پاسخ های بی ربط می توان از روش های کشف جامعه^۱ استفاده کرد و موجودیت ها را دسته بندی کرد، و به مدل ها آموزش داد که جواب هر سوال باید از جنس چه مجموعه موجودیتی باشد و پس از

¹ community detection

دریافت پاسخ ها، پاسخ هایی که در این مجموعه موجودیت نیستند را حذف کرد و سریع تر به

جواب صحیح رسید.

فصل ۸: مراجع

-
- [1] N. Feamster, J. Rexford, and E. Zegura, "The Road to SDN: An Intellectual History of Programmable Networks," *ACM SIGCOMM Comput. Commun. Rev.*, vol. 44, no. 2, pp. 87–98, 2014.
 - [2] D. Kreutz, E. Verissimo, and S. Azodolmolky, "Software-Defined Networking: A Comprehensive Survey," *Proc. IEEE*, vol. 103, no. 1, pp. 14–76, 2015.
 - [3] R. Masoudi and A. Ghaffari, "Software defined networks: A survey," *J. Netw. Comput. Appl.*, vol. 67, no. C, pp. 1–25, 2016.
 - [4] "Software-Defined Networking : The New Norm for Networks," 2012. [Online]. Available: <https://www.opennetworking.org/images/stories/downloads/sdn-resources/white-papers/wp-sdn-newnorm.pdf>.
 - [5] H. Hata, "A Study of Requirements for SDN Switch Platform," in *Intelligent Signal Processing and Communications Systems (ISPACS), 2013 International Symposium on*, 2013, pp. 79–84.
 - [6] Y. R. Qu, S. Zhou, and V. K. Prasanna, "A Decomposition-based Approach for Scalable Many-field Packet Classification on Multi-core Processors *," *Int. J. Parallel Program.*, vol. 43, no. 6, pp. 965–987, 2015.
 - [7] Y. R. Qu and V. K. Prasanna, "Power-efficient range-match-based packet classification on FPGA," in *Field Programmable Logic and Applications (FPL)*, 2015.
 - [8] S. Banerjee and K. Kannan, "Tag-In-Tag : Efficient Flow Table Management in SDN Switches," in *Network and Service Management (CNSM), 2014 10th International Conference on*, 2014, pp. 109–117.
 - [9] H. Farhadi and A. Nakao, "Rethinking Flow Classification in SDN," in *Cloud Engineering (IC2E), 2014 IEEE International Conference on*, 2014, pp. 598–603.
 - [10] K. G. Pérez, X. Yang, and S. Sezer, "A Configurable Packet Classification Architecture for Software-Defined Networking," in *System-on-Chip Conference (SOCC)*, 2014, pp. 353–358.
 - [11] P. T. Congdon, P. Mohapatra, M. Farrens, and V. Akella, "Simultaneously reducing latency and power consumption in openflow switches," *IEEE/ACM Trans. Netw.*, vol. 22, no. 3, pp. 1007–1020, 2014.
 - [12] Y. R. Qu, H. H. Zhang, S. Zhou, and V. K. Prasanna, "Optimizing many-field packet classification on FPGA, multi-core general purpose processor, and GPU," in *Architectures for Networking and Communications Systems (ANCS)*, 2015, no. 3, pp. 87–98.
 - [13] M. Varvello, R. Laufer, F. Zhang, and T. V Lakshman, "Multi-Layer Packet Classification with Graphics Processing Units Categories and Subject Descriptors," in *Proceedings of the 10th ACM International on Conference on emerging Networking Experiments and Technologies*, 2014, pp. 109–120.
 - [14] C. Hsieh and N. Weng, "Many-Field Packet Classification for Software-Defined Networking Switches," in *Proceedings of the 2016 Symposium on Architectures for Networking and Communications Systems*, 2016, pp. 13–24.
 - [15] H. Lim, N. Lee, G. Jin, J. Lee, Y. Choi, and C. Yim, "Boundary cutting for packet classification," *IEEE/ACM Trans. Netw.*, vol. 22, no. 2, pp. 443–456, 2014.

- [16] S. K. Oumya and C. H. S. E. M, "Selective Boundary Cutting For Packet Classification," *Int. J. Sci. Eng. Technol. Res.*, vol. 4, no. 34, pp. 6786–6790, 2015.
- [17] P. Gupta and N. McKeown, "Packet Classification using Hierarchical Intelligent Cuttings," in *Proceedings of the 2003 conference on Applications, technologies, architectures, and protocols for computer communications*, 2003, pp. 213–224.
- [18] S. Singh, F. Baboescu, G. Varghese, and J. Wang, "Packet Classification Using Multidimensional Cutting," in *Proceedings of the 2003 conference on Applications, technologies, architectures, and protocols for computer communications*, 2003, pp. 213–224.
- [19] D. E. Taylor and J. S. Turner, "ClassBench: A packet classification benchmark," *IEEE/ACM Trans. Netw.*, vol. 15, no. 3, pp. 499–511, 2007.
- [20] W. Xia, Y. Wen, S. Member, C. H. Foh, and S. Member, "A Survey on Software-Defined Networking," *IEEE Commun. Surv. Tutorials*, vol. 17, no. 1, pp. 27–51, 2014.
- [21] "Traditional vs Software Defined Networking." [Online]. Available: www.ipknowledge.net/wp-content/uploads/2014/12/SDN.pdf.
- [22] Y. Gong, W. Huang, W. Wang, and Y. Lei, "A survey on software defined networking and its applications," *Front. Comput. Sci.*, vol. 9, no. 6, pp. 827–845, 2015.
- [23] V. Shamugam, I. Murray, L. J. A, and A. S. Sidhu, "Software Defined Networking challenges and future direction : A case study of implementing SDN features on OpenStack private cloud," in *IOP Conference Series: Materials Science and Engineering*, 2016, vol. 121, no. 1, pp. 1–8.
- [24] M. Karakus and A. Durresi, "A survey: Control plane scalability issues and approaches in Software-Defined Networking (SDN)," *Comput. Networks*, vol. 112, pp. 279–293, 2017.
- [25] W. Li, W. Meng, and F. L. Kwok, "A survey on OpenFlow-based Software Defined Networks: Security challenges and countermeasures," *J. Netw. Comput. Appl.*, vol. 68, no. C, pp. 126–139, 2016.
- [26] N. McKeown, T. Anderson, H. Balakrishnan, G. Parulkar, L. Peterson, J. Rexford, S. Shenker, and J. Turner, "OpenFlow: Enabling Innovation in Campus Networks," *ACM SIGCOMM Comput. Commun. Rev.*, vol. 38, no. 2, p. 69, 2008.
- [27] A. Doria, J. H. Salim, W. Wang, and L. Dong, "Forwarding and Control Element Separation (ForCES) Protocol Specification," *Internet Engineering Task Force*, 2010. [Online]. Available: <https://tools.ietf.org/html/rfc5810>.
- [28] B. Pfaff and B. Davie, "The Open vSwitch Database Management Protocol," *Internet Engineering Task Force*, 2013. [Online]. Available: <https://tools.ietf.org/html/rfc7047>.
- [29] H. Song, "Protocol-oblivious forwarding: unleash the power of SDN through a future-proof forwarding plane," in *Proceedings of the second ACM SIGCOMM workshop on Hot topics in software defined networking*, 2013, pp. 127–132.
- [30] M. Smith, M. Dvorkin, and P. Garg, "OpFlex Control Protocol," *Internet Engineering Task Force*, 2014. [Online]. Available: <https://tools.ietf.org/html/draft-smith-opflex-00>.
- [31] G. Bianchi, M. Bonola, A. Capone, and C. Cascone, "OpenState: Programming Platform-independent Stateful OpenFlow Applications Inside the Switch," *Sigcomm Ccr*, vol. 44, no. 2, pp. 44–51, 2014.
- [32] A. Lara, A. Kolasani, and B. Ramamurthy, "Network Innovation using OpenFlow : A Survey," *IEEE Commun. Surv. Tutorials*, vol. 16, no. 1, pp. 493–512, 2013.
- [33] C. Hao, Chang, and Y.-D. Lin, "OpenFlow Version Roadmap," 2015. [Online]. Available:

- http://speed.cis.nctu.edu.tw/~ydlin/miscpub/indep_frank.pdf.
- [34] P. Gupta and N. McKeown, "Algorithms for Packet Classification," *IEEE Netw.*, vol. 15, no. 2, pp. 24–32, 2002.
- [35] J. Sanders and E. Kandrot, *CUDA by Example: An Introduction to General-Purpose GPU Programming*, 1st ed. Addison-Wesley Professional, 2010.
- [36] M. Arora, "The Architecture and Evolution of CPU-GPU Systems for General Purpose Computing," University of California, San Diego, 2012.
- [37] S. Singh, S. Singh, V. Banga, and C. Durlabh, "CUDA for GPGPU Applications – A Survey," *Natl. Conf. Contemporary Techniques Technol. Electron. Eng.*, pp. 1–4, 2013.
- [38] M. Mukerjee, D. Naylor, and B. Vavala, "Packet Processing on the GPU." [Online]. Available: http://www.cs.cmu.edu/~bvavala/misc/project740/15-740_Project_files/Report.pdf.
- [39] S. Han, K. Jang, K. Park, and S. Moon, "PacketShader : a GPU-Accelerated Software Router," in *Proceedings of the ACM SIGCOMM 2010 conference*, 2010, pp. 195–206.
- [40] G. Vasiliadis, S. Antonatos, M. Polychronakis, and P. Evangelos, "Gnort : High Performance Network Intrusion Detection Using Graphics Processors," in *Proceedings of the 11th international symposium on Recent Advances in Intrusion Detection*, 2008, pp. 116–134.
- [41] Y. Zhu and Y. Chen, "Hermes : An Integrated CPU / GPU Microarchitecture for IP Routing," in *Design Automation Conference (DAC)*, 2011, pp. 1044–1049.
- [42] R. Garg and S. Mittal, "Optimization by genetic algorithm," *Int. J. Adv. Res. Comput. Sci. Softw. Eng.*, vol. 4, no. 4, pp. 587–589, 2014.
- [43] D. E. Taylor and J. S. Turner, "Scalable packet classification using distributed crossproducing of field labels," *Proc. IEEE 24th Annu. Jt. Conf. IEEE Comput. Commun. Soc.*, vol. 1, pp. 1–12, 2005.
- [44] D. Yuan, X. Yang, X. Shi, B. Tang, and Y. Liu, "Multi-protocol query structure for SDN switch based on parallel bloom filter," in *International Conference on ICT Convergence*, 2014, pp. 206–211.
- [45] ب. جعفریان, "دسته بندی بسته ها در شبکه های سرعت بالا," پایان نامه کارشناسی ارشد، دانشگاه تهران، ۱۳۹۱.
- [46] C. Thomas H, C. E. Leiserson, R. Ronald L, and S. Chiffoed, *Introduction to Algorithms*, 3rd ed. London, England: MIT Press, 1988.
- [47] P. Gupta, "Algorithms for routing lookups and packet classification," STANFORD, 2000.
- [48] M. M. Buddhikot, S. Suri, and M. Waldvogel, "Space decomposition techniques for fast layer-4 switching," in *Protocols for High-Speed Networks VI*, 2000, pp. 25–41.
- [49] F. Baboescu, S. Singh, and G. Varghese, "Packet classification for core routers: is there an alternative to CAMs?," in *INFOCOM 2003. Twenty-Second Annual Joint Conference of the IEEE Computer and Communications*, 2003, vol. 1, pp. 53–63.
- [50] B. Pfaff, J. Pettit, T. Koponen, E. Jackson, A. Zhou, J. Rajahalme, J. Gross, A. Wang, J. Stringer, P. Shelar, K. Amidon, A. Networks, and M. Casado, "The Design and Implementation of Open vSwitch," in *Proceedings of the 12th USENIX Conference on Networked Systems Design and Implementation*, 2015, pp. 117–130.

فصل ۹: واژه‌نامه انگلیسی به فارسی

Abstraction	تجريد
Action	عمل
Appliance	دستگاه
Asynchronous	ناهمگام
Autonomous	خودمختار
Batch	دسته
Batch Delay	تاخير دسته‌ای
Benchmarking	محک‌زنی
Biased Distribution	توزیع جهت‌دار
Bit Vector	بردار بیتی
Block	بلوک
Bus	گذرگاه
Centralized	متمرکز شده
Component	جزء

Constant Memory	حافظه ثابت
Control plane	سطح کنترل
Core	هسته
Crossover	تقاطع
Cut	برش
Cycle	چرخه
Data plane	سطح داده
Default	پیش فرض
Deterministic	قطعی
Diversity	تنوع
Driver	راه انداز
Duplicate	تکرار
Effective	موثر
Entry	مدخل
Exact Match	انطباق دقیق
False Positive	مثبت کاذب
Fitness	برازش
Flow	جریان
Full Match	تطابق کامل
Fully Efficient	به طور کامل کارآمد

Granularity	دانه‌بندی
Independence	استقلال
Inertia	لختی
Infrastructure	زیرساختار
Lane	خط
Line rate	نرخ خط
Logical	منطقی
Logically	منطقی، به طور منطقی
Mask	ماسک
Maximum	بیشینه
Merge	ادغام
Metadata	فراداده
Minimum	کمینه
Mutation	جهش
Mutually exclusive	دو به دو ناسازگار
Network Access Point	نقطه دسترسی شبکه
Network state	حالت شبکه
North-bound Interface	واسط شمالی
Offline	برون خط
Online	برخط

Out-of-order	خارج از نوبت
Performance	کارایی
Pipeline	خط لوله
Platform	سکو
Pool	مخزن
Pre-filtering	پیش تصفیه
Prefix Match	انطباق پیشوندی
Preprocessing	پیش پردازش
Program Counter	شمارنده برنامه
Programmable	قابل برنامه ریزی
Range Match	انطباق بازه‌ای
Regular Expression	عبارت منظم
Replication Ratio	نسبت تکرار
Rule Programming	برنامه ریزی قانون
Scalability	مقیاس پذیری
Scheduler	زمان بند
Scheduling	زمان بندی
Scope	گستره
Segmentation	قطعه بندی
Slot	شکاف

Smoothing	هموارسازی
South-bound Interface	واسط جنوبی
Sparsity	خلوتی
Streaming Multiprocessor	چندپردازنده جریان‌ی
String	رشته
Subrange	زیرمحدوده
Subset	زیرمجموعه
Switching	راه‌گزینی
Synthetic	مصنوعی
Target Rule	قانون هدف
Thread	ریسمان
Throughput	گذردهی
Tuple	چندتایی
Uniform	یکنواخت
Unifying	یک‌دست‌سازی
Zone	منطقه

فصل ۱۰: واژه‌نامه فارسی به

انگلیسی

Merge	ادغام
Independence	استقلال
Range Match	انطباق بازه‌ای
Prefix Match	انطباق پیشوندی
Exact Match	انطباق دقیق
Fitness	برازش
Online	برخط
Bit Vector	بردار بیتی
Cut	بریدن
Rule Programming	برنامه‌ریزی قانون
Offline	برون خط
Block	بلوک
Fully Efficient	به طور کامل کارآمد
Maximum	بیشینه
Preprocessing	پیش‌پردازش

Pre-filtering	پیش تصفیه
Default	پیش فرض
Batch Delay	تاخیر دسته‌ای
Abstraction	تجرید
Full Match	تطابق کامل
Crossover	تقاطع
Duplicate	تکرار
Diversity	تنوع
Biased Distribution	توزیع جهت‌دار
Flow	جریان
Component	جزء
Mutation	جهش
Cycle	چرخه
Streaming Multiprocessor	چندپردازنده جریانی
Tuple	چندتایی
Constant Memory	حافظه ثابت
Network state	حالت شبکه
Out-of-order	خارج از نوبت
Lane	خط
Pipeline	خط لوله

Sparsity	خلوتی
Autonomous	خودمختار
Granularity	دانه‌بندی
Appliance	دستگاه
Batch	دسته
Mutually exclusive	دو به دو ناسازگار
Driver	راه‌انداز
Switching	راه‌گزینی
String	رشته
Thread	ریسمان
Scheduler	زمان‌بند
Scheduling	زمان‌بندی
Infrastructure	زیرساختار
Subset	زیرمجموعه
Subrange	زیرمحدوده
Data plane	سطح داده
Control plane	سطح کنترل
Platform	سکو
Slot	شکاف
Program Counter	شمارنده برنامه

Regular Expression	عبارت منظم
Action	عمل
Metadata	فراداده
Programmable	قابل برنامه‌ریزی
Target Rule	قانون هدف
Segmentation	قطعه‌بندی
Deterministic	قطعی
Performance	کارایی
Minimum	کمینه
Throughput	گذردهی
Bus	گذرگاه
Scope	گستره
Inertia	لختی
Mask	ماسک
Centralized	متمرکز شده
False Positive	مثبت کاذب
Benchmarking	محک‌زنی
Pool	مخزن
Entry	مدخل
Synthetic	مصنوعی

Scalability	مقیاس پذیری
Zone	منطقه
Logical	منطقی
Logically	منطقی، به طور منطقی
Effective	موثر
Asynchronous	ناهمگام
Line rate	نرخ خط
Replication Ratio	نسبت تکرار
Network Access Point	نقطه دسترسی شبکه
South-bound Interface	واسط جنوبی
North-bound Interface	واسط شمالی
Core	هسته
Smoothing	هموارسازی
Unifying	یکدست سازی
Uniform	یکنواخت

Abstract:

In Software Defined Networking, a flow can be defined using arbitrary set of header fields of each layer. For example, In OpenFlow 1.3.1, total of 15 fields may be used to define a flow. Allowing to define flows with more fields, on the one hand, causes the length of flow table entries to be enlarged; On the other hand, it enables finer control on network traffic which in turn, causes the number of flow table entries to be increased. Consequently, design of a solution for packet classification in high performance environments has been turned out to be a big challenge. To address this problem in soft switches, researchers try to reduce the number of candidate rules that a tuple must match with them. In the MC-SBC algorithm, a two-level trie based structure is proposed for lookup tables in which a set of effective bits are detected through a statistical method, in order to partition the ruleset and construct the trie structure.

In this dissertation, genetic algorithm has been applied to find the effective bit sets for ruleset partitioning. The results show that the proposed method partitions the rulesets in more balanced subsets, such that, both the replication ratio and the maximum number of rules in a leaf node of the trie are decreased. Moreover, the results of running our prototype on a single NVIDIA GPU, shows that packet classification can be performed up to 10 percent faster for large rule sets.

Keywords: Software Defined Networking, Packet Classification, OpenFlow Switch, Graphics Processing Unit, Genetic Algorithm



University of Tehran

College of Engineering

Faculty of Electrical and Computer Engineering

Packet Classification in Software Defined Networking

A thesis submitted to the Graduate Studies Office
In partial fulfillment of the requirements for
The degree of M.Sc in
Hardware Engineering

By:
Mohammad Reza Piroozi

Supervisor:
Dr. Nasser Yazdani