

بهبود مدل های مبتنی بر ویژگی های پنهان پایگاه های دانش با استفاده از قوانین انجمنی

## تعریف مساله و هدف و ضرورت

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. پیش بینی وقوع لینک ها، یک مساله اساسی و بنیادین در شبکه های اجتماعی می باشد. در موضوع پیش بینی لینک، نمایی از یک شبکه به ما داده می شود و ما مایل هستیم که بدانیم در آینده نزدیک، احتمالا چه تراکنش هایی میان اعضای فعلی شبکه روی می دهد و یا اینکه کدام یک از تراکنش های موجود را از دست می دهیم.

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 یا ل های گراف شبکه به صورت «الف، ب» خواهد بود به این معنی که شخص الف با شخص ب رابطه دوستی دارد. در سال های اخیر پایگاه های دانشی ایجاد شده اند که داده ها در آن به صورت چند رابطه ای ذخیره شده اند و اطلاعات بیشتری از یک رابطه دوتایی بلی یا خیر به ما می دهند. ما در این تحقیق از داده های چند رابطه ای استفاده می کنیم. منظور از داده های چند رابطه ای گراف جهت داری است متشکل از موجودیت ها و روابط بین آن ها که بصورت «مبدأ 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 WordNet، Google Knowledge Graph، FreeBase وجود دارند که شامل تعداد زیادی نمونه چند رابطه ای می باشند و تعداد زیادی موجودیت و روابط بین آن ها را می توان در آن ها یافت و از آن برای یادگیری مدل استفاده کرد. شناخت پایگاه دانش و دریافت اطلاعات از آن نیز یکی از مسایل پیش بینی لینک روی داده های چند رابطه ای است. پیش بینی لینک مقیاس پذیر روی داده های چند رابطه ای عمدتا با استفاده از ویژگی های پنهان صورت می گیرد. مشکل این روش ها تعداد بسیار زیاد پارامترها در فاز یادگیری مدل است به صورتی که انجام عملیات یادگیری مدل را غیرممکن می کند. کاهش تعداد پارامترهای مساله نیز که در برخی روش ها پیشنهاد شده نیز برخی از الگوهای موجود در پایگاه دانش را در نظر نمی گیرند و باعث کاهش دقت مدل خواهند شد. در این تمرین قصد داریم که بدون افزایش تعداد پارامترهای روش های الگوی پنهان شده و با استفاده از قوانین انجمنی موجود در پایگاه دانش عمل کرد این روش ها را بهبود بخشیم.

در این تحقیق قصد داریم که با کار روی داده‌های چند رابطه‌ای به سؤالاتی که از پایگاه‌دانش می‌شود پاسخ دهیم، به این صورت که یک مدل از روی داده‌های پایگاه دانش ساخته و آموزش داده می‌شود و پس از آن سؤالات به گونه‌ای که مقصد یا مبدأ آن مجهول باشد از مدل ایجاد شده پرسیده می‌شود. برای نمونه در مثال بالا اگر قسمت کشور مجهول باشد از پایگاه دانش پرسیده می‌شود «تهران، واقع در، ؟» و انتظار می‌رود که مدل آموزش دیده پاسخ سؤال را با دقت نسبتاً خوبی بیابد.

نتایج این مساله کاربردهای زیادی خواهد داشت به عنوان مثال در ادامه تعدادی از این کاربردها را مطرح خواهیم کرد:

- پیشبینی لینک‌های احتمالی در شبکه‌های اجتماعی، برای مثال در شبکه‌های 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 به عنوان سیستم‌های توصیه‌گر، برای مثال کاربران و کالاها یا اشیاء موجودیت‌ها هستند و خریدن، امتیاز دادن، بررسی کردن و ... رابطه‌ها هستند که می‌توان از بررسی این روابط و موجودیت‌ها اطلاعاتی نظیر کاربر  $x$  کالای  $y$  را خواهد خرید یا خیر یا اینکه کالای  $x$  به تعداد بالا فروش خواهد رفت یا خیر بدست آورد.
- کامل‌تر کردن پایگاه‌های داده، از نتایج این تحقیق می‌توان استفاده کرد و با اضافه کردن رابطه‌های انجمنی بدست آمده به پایگاه‌داده آن را کامل‌تر کرد.

روش‌ها:

داده‌هایی که قصد کار روی آن‌ها را داریم معمولاً ابعاد بالایی خواهند داشت. برای مثال پایگاه دانش Freebase شامل یک میلیون نوع موجودیت و ۲۵ هزار نوع رابطه و ۱۷ میلیون نمونه برای آموزش و آزمون مدل است. از آنجایی که کار روی داده‌هایی با این تعداد ابعاد بسیار سخت و تقریباً غیر ممکن است باید از روش‌های کاهش ابعاد استفاده کنیم. در این تحقیق سعی می‌کنیم که این داده‌ها را به فضای برداری آورده و مساله را به یک مساله کمینه‌سازی تبدیل کنیم. این داده‌ها را به صورت یک تنسور فرض می‌کنیم که برای هر رابطه یک فضای دوبعدی برداری برای موجودیت‌های مختلف در نظر می‌گیریم، همچنین برای هر رابطه یک بردار بین موجودیت‌های دخیل در آن رابطه خواهیم داشت. سپس تلاش می‌کنیم که از این بردارها در یک مساله کمینه‌سازی استفاده کنیم. برای مثال اگر سه‌تایی « $head, relation, tail$ » را داشته باشیم که هرکدام به یک بردار نگاشت شده‌اند انتظار داریم مجموع بردار  $head$  و  $tail$  به بردار  $relation$  برسد پس باید مقدار زیر را کمینه کنیم:

$$head + relation - tail \Rightarrow 0$$

برای نمایش این قانون از تخمین مرتبه پایین استفاده می‌کنیم به صورت:

$$f_r(h, t) = \|h + r - t\|_2^2$$

در صورتی که تابع  $f$  برای سه‌تایی « $h, r, t$ » مقداری نزدیک به صفر داشته باشد این رابطه برقرار است و در غیر این صورت این رابطه برقرار نیست.

در تحقیقاتی که قبلاً روی این مدل داده‌ها انجام شده است فقط از حقایقی که به صورت مستقیم داده شده استفاده شده است. برای مثال اگر در حقایق موجود در داده‌ها داشته باشیم «الف، مادر، ب» و «الف، همسر، ج» فقط از این دو حقیقت استفاده می‌شود و از این حقیقت که «ج، پدر، ب» استفاده نمی‌شود. در این تحقیق تصمیم داریم از قوانین انجمنی بین داده‌ها نیز استفاده کرده و حقایق ضمنی که بین داده‌ها هست را نیز در آموزش مدل استفاده کنیم.

پیشینه تحقیق:

اولین تحقیقاتی که روی این روش استخراج اطلاع از داده‌های چندرابطه صورت گرفت در سال ۲۰۱۱ شروع شد و معروف ترین مقالاتی که در این موضوع مطرح شده عبارتند از **Unstructured Model (UM)** و **Semantic Matching Energy (SE)** و **Single Layer Model (SLM)** و **Structured Embedding (SE)** و **TransH** و **TransE** و **Latent Factor Model (LFM)** و **Neural Tensor Network (NTN)** و **SME** و **TransR** که در همه موارد فقط از حقایق مستقیم که از پایگاه دانش بدست آمده استفاده شده و از حقایق ضمنی موجود در پایگاه دانش استفاده نشده است. همچنین برای یافتن روابط انجمنی بین داده‌های پایگاه دانش در سال ۲۰۱۳ روش **AMIE** معرفی شده است که قوانین انجمنی را با فرض جهان باز استخراج می‌کند.