

دانشگاه فردوسی مشهد
دانشکده مهندسی

پروژه کارشناسی
مهندسی کامپیوتر

عنوان:

ایجاد گراف دانش با Neo4j

نگارش:

سارا اسدی

استاد راهنما:

دکتر محسن کاهانی

استاد مشاور:

مریم کرابی

زمستان ۱۴۰۱

سپاس

از استاد بزرگوارم، دکتر کاهانی، که با کمک ها و راهنمایی های بی دریغشان، بنده را در انجام این پروژه یاری داده‌اند، تشکر و قدردانی میکنم.

از خانم مریم کرابی که در طول مدت انجام این پروژه، زمان و علم و دانش زیادی را در اختیار من قرار داد نیز بسیار سپاسگزارم.

از جناب دکتر اله بخش که زحمت داوری این پروژه را به عهده داشته اند نیز کمال تشکر را دارم. همچنین تشکر میکنم از دوستان و همکلاسیان گرامی که در طی دوران تحصیل در کنار من بوده اند. و در آخر از پدر، مادر و خواهرانم که در همواره در طی این سالیان پشتیبان من بوده اند.

چکیده

با رشد تصاعدی داده ها، نیاز به مدیریت کارآمد داده ها و سیستم های بازیابی اهمیت فزاینده ای پیدا کرده است. گراف های دانش به عنوان یک راه حل امیدوارکننده برای مقابله با این چالش با نمایش داده ها به شیوه ای ساختاریافته و به هم پیوسته پدیدار شده اند. در این پایان نامه، فرآیند ایجاد یک گراف دانش از منابع داده بدون ساختار را بررسی می کنیم و بر جنبه های عملی توسعه چنین سیستمی تمرکز می کنیم. ما با بحث در مورد معرفی گراف دانش و هستی شناسی و اهمیت آنها در زمینه مدیریت داده ها شروع می کنیم. هدف اصلی این پایان نامه توسعه یک گراف دانش برای یک حوزه خاص و نمایش اثربخشی آن در سازماندهی و استخراج اطلاعات سودمند از مجموعه داده های بزرگ است. ما از یک مجموعه داده دنیای واقعی برای نشان دادن اثربخشی گراف دانش استفاده خواهیم کرد که برای بنیاد امام رضا (ع) می باشد. در ایجاد این گراف دانش مراحل از جمله جمع آوری، پیش پردازش داده و شناسایی موجودیت، استخراج رابطه و تولید گراف نهایی ذکر شده است. امیدواریم که این پژوهش به مجموعه دانش رو به رشد در این زمینه کمک کند و الهام بخش کار آینده برای محققان و متخصصان علاقه مند به استفاده از نمودارهای دانش برای مدیریت و تجزیه و تحلیل داده ها باشد.

کلیدواژه ها: Neo4j، گراف دانش، پایگاه داده گرافی، هستی شناسی .

فهرست مطالب

۸	۱ مقدمه
۹	۱-۱ تعریف پروژه
۱۰	۲-۱ اهداف پروژه
۱۰	۳-۱ ساختار پروژه
۱۱	۲ پیش زمینه پروژه و تعاریف کلی
۱۱	۱-۲ گراف دانش چیست؟
۱۲	۲-۲ آشنایی مختصر با زبان OWL
۱۳	۱-۲-۲ تفاوت آنتولوژی با گراف دانش
۱۵	۳ روش پیشنهادی برای انجام پروژه
۱۵	۱-۳ جمع آوری داده
۱۷	۲-۳ ایجاد گراف دانش
۱۷	۳-۳ معرفی Neo4j
۱۸	۱-۳-۳ نصب و راه اندازی Neo4j
۲۲	۲-۳-۳ تعامل با رابط کاربری Neo4j
۲۴	۴ نتایج بدست آمده

۱-۴ پیاده‌سازی طرح کلی گراف دانش ۲۴

۲-۴ گراف دانش پروژه بنیاد ۳۰

۳-۴ رابط کاربری ۳۳

۵ جمع بندی و کارهای پیشرو ۳۵

فهرست شکل‌ها

۱۶	۱-۳
۱۶	۲-۳
۱۹	۳-۳
۱۹	۴-۳
۱۹	۵-۳
۲۰	۶-۳
۲۰	۷-۳
۲۰	۸-۳
۲۱	۹-۳
۲۱	۱۰-۳
۲۱	۱۱-۳
۲۵	۱-۴
۲۵	۲-۴
۲۵	۳-۴
۲۶	۴-۴
۲۷	۵-۴

۲۸	۶-۴
۲۹	۷-۴
۳۰	۸-۴
۳۱	۹-۴
۳۲	۱۰-۴
۳۲	۱۱-۴
۳۳	۱۲-۴
۳۴	۱۳-۴
۳۴	۱۴-۴

فصل ۱

مقدمه

در سال‌های اخیر، حجم داده‌های تولید شده توسط مشاغل، افراد و ماشین‌ها افزایش بسیار یافته و به پدیده‌ای منجر شده که معمولاً به عنوان پدیده کلان داده شناخته می‌شود. این داده‌ها می‌توانند از طیف گسترده‌ای از منابع، از جمله پلتفرم‌های رسانه‌های اجتماعی، دستگاه‌های IoT، تراکنش‌های مالی و غیره به دست آیند، و اغلب بدون ساختار یا نیمه‌ساختارمند هستند، که مدیریت و تجزیه و تحلیل با استفاده از ابزارهای مدیریت داده سنتی را دشوار می‌کند. حجم زیاد و پیچیدگی کلان داده چالش‌هایی را ایجاد می‌کند، از جمله:

- ذخیره سازی و بازیابی داده ها: سیستم های ذخیره سازی سنتی داده ها ممکن است توانایی مدیریت حجم داده های تولید شده توسط منابع بزرگ داده را نداشته باشند، و ذخیره و بازیابی به موقع داده ها را دشوار می‌کند.
- پردازش و تجزیه و تحلیل داده‌ها: با وجود داده‌های زیادی برای پردازش، ابزارهای سنتی تجزیه و تحلیل داده‌ها ممکن است برای ارائه بینش و هوش عملی در مدت زمان معقول با مشکل مواجه شوند.
- کیفیت داده: داده‌های بزرگ می‌توانند کثیف باشند، با داده‌های ناسازگار یا ناقص از منابع متعدد. این می‌تواند دریافت یک تصویر کامل و دقیق از آنچه داده ها به ما می‌گویند دشوار باشد.
- حریم خصوصی و امنیت داده ها: با جمع آوری و ذخیره داده های بیشتری، نگرانی هایی در مورد حفظ حریم خصوصی و امنیت داده ها وجود دارد، از جمله خطر دسترسی غیرمجاز یا نقض

داده‌ها.

عدم استفاده از تجزیه و تحلیل داده‌ها برای مشاغل و سازمان‌ها می‌تواند منجر به مشکلات متعددی مانند عدم شناسایی ناکارآمدی‌ها، ناتوانی در رقابت، از دست دادن مشتریان، تصمیمات ضعیف، از دست دادن فرصت‌های بهبود و ... شود که موجب هدر رفتن منابع، زمان، پول و عواقب منفی بسیاری گردد.

برای مقابله با این چالش‌ها، کسب و کارها و سازمان‌ها به فناوری‌های جدیدی مانند پایگاه‌های داده NoSQL روی آورده‌اند که برای مدیریت حجم زیادی از داده‌های بدون ساختار یا نیمه ساختار یافته طراحی شده‌اند. Neo4j یکی از این فناوری‌هاست، و چندین مزیت کلیدی را نسبت به پایگاه‌های داده سنتی و دیگر پایگاه‌های داده NoSQL در مدیریت و تجزیه و تحلیل داده‌های بزرگ ارائه می‌دهد. اول و مهمتر از همه، مدل پایگاه داده گراف Neo4j به ویژه برای برنامه‌های کاربردی کلان داده مناسب است، زیرا امکان مدیریت کارآمد و پرس و جو از مجموعه داده‌های بسیار به هم پیوسته را فراهم می‌کند. این امکان استخراج بینش از روابط پیچیده بین نقاط داده را فراهم می‌کند که ممکن است به راحتی با استفاده از ابزارهای مدیریت داده سنتی قابل شناسایی نباشند. علاوه بر این، زبان پرس و جو آن یعنی Cypher، به گونه‌ای طراحی شده است که بصری و آسان برای استفاده باشد و به توسعه دهندگان و تحلیلگران داده اجازه می‌دهد تا به سرعت و به راحتی پرس و جوهای پیچیده را در مجموعه داده‌های بزرگ انجام دهند. و از آنجایی که Neo4j بسیار مقیاس پذیر است، به راحتی می‌تواند رشد کند و با نیازهای کسب و کارها و سازمان‌ها سازگار شود زیرا نیازهای داده آنها در طول زمان تکامل می‌یابد. به طور کلی، نقاط قوت Neo4j در مدل‌سازی داده‌ها، پرس و جو و مقیاس‌پذیری، آن را به ابزاری قدرتمند برای مدیریت و تجزیه و تحلیل داده‌های بزرگ تبدیل می‌کند و آن را به گزینه‌ای ایده‌آل برای کسب و کارها و سازمان‌هایی تبدیل می‌کند که به دنبال کسب بینش از مجموعه داده‌های پیچیده و بسیار به هم مرتبط هستند.

در ادامه به بیان شرح کل صورت پروژه، دلایل انجام و بررسی ساختار کل این گزارش می‌پردازیم.

۱-۱ تعریف پروژه

پروژه شامل چهار بخش کلی می‌باشد که عبارت‌اند از:

۱. استخراج موجودیت ها و روابط از داده های جمع آوری شده از فایل OWL
۲. ساخت گراف دانش از تمام کلاس های موجود و تهیه یک طرح جامع
۳. ایجاد پایگاه داده گرافی بر اساس داده های پروژه بنیاد امام رضا و مقادیر گردآوری شده از سایت
۴. طراحی یک رابط کاربری مناسب برای نمایش خروجی پروژه

۲-۱ اهداف پروژه

هدف از انجام این پروژه استفاده از پایگاه داده گرافی و ایجاد گراف دانش برای داده های غیر ساختار یافته حوزه خاص و نمایش اثربخشی آن در سازماندهی و استخراج اطلاعات سودمند از مجموعه داده های بزرگ است. ما از یک مجموعه داده دنیای واقعی برای نشان دادن اثربخشی نمودار دانش استفاده خواهیم کرد که برای بنیاد امام رضا (ع) می باشد. در ایجاد این گراف دانش مراحل از جمله جمع آوری، پیش پردازش داده و شناسایی و استخراج موجودیت و روابط و تولید گراف نهایی شرح داده شده است.

۳-۱ ساختار پروژه

این گزارش شامل پنج فصل می باشد که شامل محتوی زیر می باشد. فصل اول (این فصل) شامل تعریف کلی صورت پروژه و بیان کلیات می باشد. فصل دوم دربرگیرنده معرفی گراف دانش - آشنایی با OWL و تفاوت آن با گراف دانش می باشد. در فصل سوم روش پیشنهادی برای انجام پروژه و پیاده سازی به صورت مفصل شرح داده می شود. سپس در فصل چهارم به نتایج به دست آمده از این عملیات اشاره میکنیم. فصل پنجم شامل خلاصه و جمع بندی پروژه به همراه کارهای پیش رو می باشد که به بیان چشم انداز ما برای ادامه پروژه و اقدامات قابل انجام توسط دیگران می پردازد.

فصل ۲

پیش زمینه پروژه و تعاریف کلی

در این بخش به توضیح برخی موارد و اصطلاحاتی که در طول گزارش استفاده شده است و ممکن است برای همگان ناآشنا باشد می‌پردازیم.

۱-۲ گراف دانش چیست؟

گراف دانش نوعی پایگاه داده است که دانش را به عنوان مجموعه‌ای از موجودیت‌های به هم پیوسته و روابط آنها با یکدیگر نشان می‌دهد. گراف‌های دانش معمولاً حاوی اطلاعاتی در مورد حوزه‌ها یا حوزه‌های موضوعی خاص، مانند مفاهیم علمی، فرآیندهای تجاری یا شبکه‌های اجتماعی هستند و یک ابزار قدرتمند برای نمایش و تجزیه و تحلیل داده‌های پیچیده ارائه می‌کنند و به طور فزاینده‌ای در طیف گسترده‌ای از کاربردها در صنایع مورد استفاده قرار می‌گیرند. در گراف دانش، موجودیت‌ها به صورت گره و روابط بین موجودیت‌ها به صورت لبه‌ها یا پیوندها نمایش داده می‌شوند. به عنوان مثال، در یک نمودار دانش از افراد و روابط آنها، هر فرد به عنوان یک گره نشان داده می‌شود، و روابط آنها (مانند عضو خانواده، دوستان و همکاران) به عنوان یال‌هایی نشان داده می‌شود که گره‌ها را به هم متصل می‌کنند. استفاده از نمودار دانش امکان نمایش انعطاف‌پذیرتر و ظریف‌تر داده‌ها را نسبت به پایگاه‌های داده سنتی فراهم می‌کند. از نمودارهای دانش می‌توان برای ثبت پیچیدگی و به هم پیوستگی سیستم‌های دنیای واقعی استفاده کرد و امکان پرس و جو و تحلیل کارآمدتر حجم زیادی از داده‌ها را فراهم کرد و به ویژه برای برنامه‌هایی که نیاز به استدلال و استنباط دارند، مانند پردازش زبان طبیعی، سیستم‌های توصیه

و کشف تقلب مناسب هستند. گراف‌های دانش معمولاً با استفاده از فناوری های وب معنایی، مانند مدل داده RDF (چارچوب توصیف منابع) و OWL (زبان هستی شناسی وب) برای نمایش هستی شناسی ها و ساختارهای دانش ساخته می شوند. علاوه بر گره ها و روابط، نمودارهای دانش ممکن است دارای ویژگی ها یا ویژگی هایی باشند که موجودیت ها را با جزئیات بیشتری توصیف می کنند. به عنوان مثال، یک گره شخصی در یک نمودار دانش ممکن است ویژگی هایی مانند سن، شغل و علایق داشته باشد. گراف دانش را می توان برای نمایش طیف وسیعی از انواع داده ها، از جمله داده های ساخت یافته (مانند داده های پایگاه داده های رابطه ای) و داده های بدون ساختار (مانند متن و محتوای چند رسانه ای) استفاده کرد و می توان با استفاده از تکنیک های مختلف، از جمله الگوریتم های نمودار، پردازش زبان طبیعی، و یادگیری ماشین، پرس و جو و تحلیل کرد. یکی از مزیت های کلیدی نمودارهای دانش، توانایی آن ها در یکپارچه سازی داده ها از منابع متعدد، مانند مجموعه های داده عمومی، پایگاه های داده داخلی و پلتفرم های رسانه های اجتماعی است. برخی از شناخته شده ترین نمودارهای دانش عبارتند از نمودار دانش گوگل که برای ارائه پرس و جو و پاسخ در نتایج جستجو استفاده می شود و نمودار دانش DBpedia که از ویکی پدیا مشتق شده و حاوی اطلاعات ساختاریافته درباره افراد، مکان ها و چیزها است. استفاده از نمودارهای دانش به طور فزاینده ای در زمینه هایی مانند علم داده، هوش مصنوعی و مدیریت اطلاعات رایج می شود، زیرا سازمان ها به دنبال به دست آوردن بینش عمیق تر در مورد داده های خود و استفاده بهتر از آن هستند.

۲-۲ آشنایی مختصر با زبان OWL

همانطور که می دانید ما تا به حال ۳ نسل از وب را پشت سر گذاشته ایم:

۱. وب ۱ یا وب سنتی همان وبی که آقای «تیم برنرزی» اختراع کرد و حاوی یک سری صفحات ساده که در آن ها فقط متن و عکس بود.

۲. وب ۲ یا اشتراکی: وبی که ظهور شبکه های اجتماعی (Social-Networks) را شاهد بودیم و کاربران امکان اشتراک هر نوع داده را داشتند.

۳. وب ۳ یا وب معنایی (Semantic-Web): در وب ۳ صفحات به طور کامل برای موتورهای جستجو و هر نوع موتور تحلیل گر دیگر باید شناخته شده باشد و متون برای آن فقط یک سری

کلمه‌ی ساده در نظر گرفته نشود.

روش‌های معنا بخشیدن به صفحات وب :

۱. XML

۲. URI

۳. OWL ، RDF

زبان هستی‌شناسی وب یا همان OWL یک زبان وب معنایی می‌باشد که برای بیان مفاهیم و دانش پیچیده در رابطه با اشیا موجودیت‌ها گروه‌های اشیا و روابط بین آنها ایجاد شده است OWL یک زبان محاسباتی منطقی می‌باشد به نحوی که دانشی که توسط OWL بیان می‌گردد می‌تواند توسط برنامه‌های کامپیوتری مورد استفاده قرار گیرد. برای مثال برای بررسی کردن سازگاری و پایداری یک ساختار دانش و یا آشکار کردن دانش ضمنی کاربرد دارد. اسناد OWL که آن را آنتولوژی یا هستان‌شناسی می‌نامیم می‌تواند در فضای وب منتشر شود می‌تواند به سایر هستیان‌شناسی‌های OWL ارجاع دهد و یا توسط آنها مورد ارجاع قرار گیرد. OWL بخشی از مجموعه تکنولوژی وب معنایی می‌باشد که شامل RDF، RDFS، و SPARQL می‌باشد. یکی از ابزارهای مرتبط با هستی‌شناسی که در این پروژه از آن استفاده شده است protégé می‌باشد که از مزایای آن موارد زیر می‌باشد:

- طراحی و ساخت آنتولوژی‌های owl بدون نیاز به کدنویسی
- دارای موتور استنتاج جهت ارزیابی سازگاری هستان‌شناسی (آنتولوژی)
- خروجی گرافیکی از ساختار هستان‌شناسی
- قابلیت تعریف انواع ویژگی‌ها برای آنتولوژی

۲-۲-۱ تفاوت آنتولوژی با گراف دانش

از آنجایی که برنامه‌های معنایی به سرعت در حال تبدیل به موضوعات مهم در صنعت می‌شوند، سوالاتی غالباً در رابطه با آنتولوژی‌ها و گراف‌های معنایی پیش می‌آید به خصوص در رابطه با تفاوت این دو.

آیا آنتولوژی ها و گراف های دانش یکسان هستند؟ اگر نیستند تفاوتشان چیست؟ ارتباط بین این دو چیست؟ آنتولوژی ها مدل های داده معنایی هستند که نوع اشیایی که در حوزه مورد نظر ما وجود دارند و ویژگی هایی که می توانند برای توصیف آنها مورد استفاده قرار بگیرند را تعریف می کنند. در واقع می توان گفت آنتولوژی ها مدل های داده عمومی یا کلی هستند و فقط به توصیف کلی یا عمومی چیزهایی که ویژگیهای معینی دارند می پردازد، اما به اطلاعات مشخص در رابطه با اعضا و موجودیت های حاضر در آن حوزه نمی پردازد. از سه بخش اصلی تشکیل شده است که به صورت زیر توصیف می شود:

- کلاس ها: یعنی انواعی از چیزها که دارای مشخصات معین و منحصر به فرد هستند که در داده های ما وجود دارد
- ارتباطات: ویژگی هایی که دو کلاس را به یکدیگر مرتبط می کنند
- خصیصه ها: ویژگی هایی که کلاس یک عضو را مشخص می کنند

پس هستی شناسی به مشخصات رسمی یک مفهوم سازی اشاره دارد، یعنی نمایشی ساختاریافته و استاندارد شده از مفاهیم و روابط در یک حوزه خاص که معمولاً شامل سلسله مراتبی از مفاهیم، مجموعه ای از روابط بین آن مفاهیم، و مجموعه ای از بدیهیات یا قوانینی هستند که بر رفتار آن مفاهیم و روابط حاکم هستند ولی گراف دانش یک پایگاه داده یا ساختار داده ای است که دانش را به عنوان مجموعه ای از موجودیت های به هم پیوسته و روابط آنها با یکدیگر نشان می دهد و یک پیاده سازی مشخص از مشخصات هستی شناسی در یک پایگاه داده یا ساختار داده است. هستی شناسی ها برای تعریف ساختار و معناشناسی یک گراف دانش استفاده می شوند، در حالی که نمودارهای دانش برای ذخیره و پرس و جو از داده های ارائه شده توسط هستی شناسی استفاده می شوند. یک گراف دانش هنگامی ایجاد می شود که شما یک آنتولوژی (همان دیتا مدل) را به یک دیتاست از داده های معین اعمال کنید. به بیان دیگر می توان گفت:

$$\text{Ontology} + \text{Data} = \text{Knowledge Graph}$$

اما به طور کل هستی شناسی ها و نمودارهای دانش مفاهیم مکملی هستند که اغلب با هم در برنامه هایی مانند یکپارچه سازی داده ها، پردازش زبان طبیعی و هوش مصنوعی استفاده می شوند.

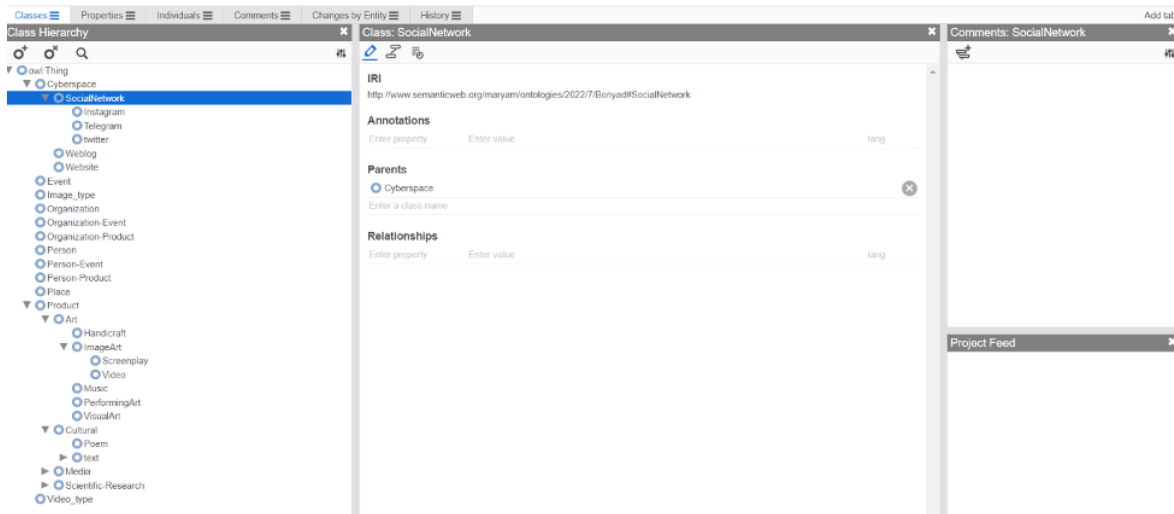
فصل ۳

روش پیشنهادی برای انجام پروژه

۳-۱ جمع آوری داده

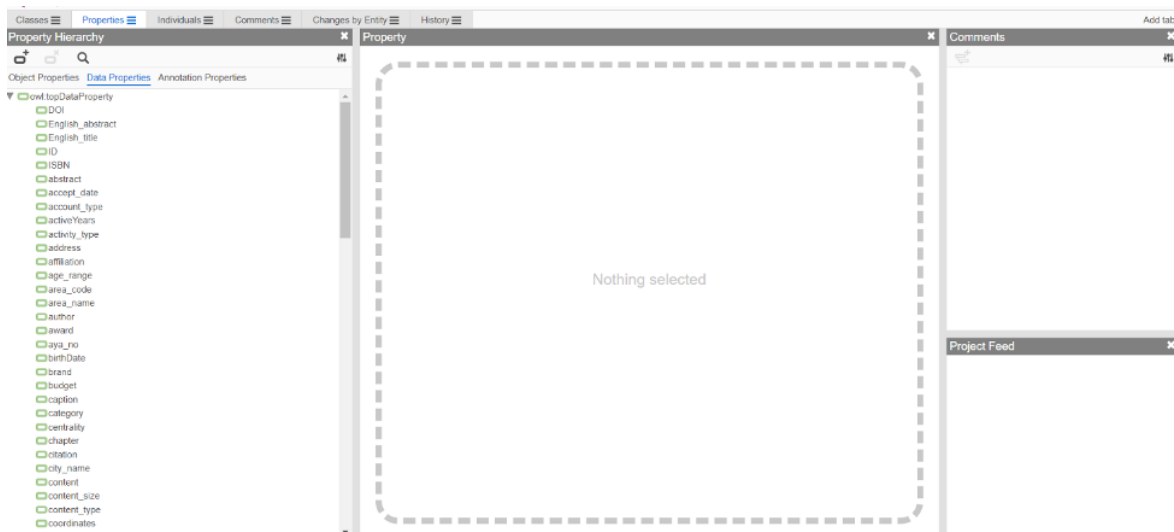
در انجام این پژوهش داده های جمع آوری شده در قالب فایل OWL ارائه شده است که همانطور که گفته شد OWL نوعی فرمت فایل است که در زمینه فناوری های وب معنایی برای نمایش هستی شناسی ها، که توصیف های رسمی حوزه های دانش هستند، استفاده می شود که معمولاً به عنوان سلسله مراتبی از مفاهیم و روابط آنها نشان داده می شوند که می توانند برای توصیف ساختار یک حوزه خاص و روابط بین عناصر مختلف آن استفاده شوند. در زمینه Neo4j، یک فایل OWL می تواند منبع اطلاعاتی ارزشمندی برای ایجاد یک پایگاه داده گراف باشد. با استخراج داده ها از یک فایل OWL و استفاده از آن برای ایجاد گره ها و روابط در یک پایگاه داده گراف، می توان یک مدل داده قدرتمند و منعطف ایجاد کرد که به طور دقیق ساختار یک حوزه دانش خاص را نشان می دهد. این رویکرد امکان پرس و جو و تجزیه و تحلیل کارآمد داده های پیچیده و بهم پیوسته را فراهم می کند و آن را برای برنامه هایی که با حجم زیادی از داده ها از منابع متعدد سروکار دارند ایده آل می کند. به طور کلی، استفاده از یک فایل OWL برای ایجاد یک پایگاه داده گراف در Neo4j می تواند مزایای متعددی از جمله مدل سازی و تحلیل داده ها کارآمدتر، کیفیت بهتر داده ها و انعطاف پذیری بیشتر در انطباق با نیازهای متغیر داده در طول زمان را ارائه دهد. Neo4j با توانایی مدیریت مجموعه داده های پیچیده و به هم پیوسته، ابزاری قدرتمند برای مدیریت و تجزیه و تحلیل داده ها از طیف گسترده ای از منابع است و استفاده از فایل های OWL می تواند به باز کردن بینش های بیشتر از این مجموعه داده ها کمک کند. در تصویر زیر به کمک

نرم افزار protégé محتویات فایل را مشاهده میکنیم تا از آن مفاهیم مورد نیاز را استخراج نماییم:



شکل ۳-۱

همانطور که مشاهده میکنید سلسله مراتبی از کلاس ها را داریم که هر کدام دارای زیرشاخه و روابط با یکدیگر هستند.



شکل ۳-۲

در این قسمت سلسله مراتب ویژگی ها را به صورت مجزا داریم که با بررسی هرکدام، ویژگی هر موجودیت یا کلاس را استخراج کرده و در گراف دانش مان به عنوان خصیصه های گره اضافه کردیم. سپس روابط بین کلاس ها را بر اساس Domain و Range که گره شروع و پایان برای آن یال را مشخص میکند را جدا کردیم تا در گراف دانش مان پیاده سازی کنیم.

۲-۳ ایجاد گراف دانش

در پایگاه داده های گراف، داده ها در گره ها و روابط ذخیره می شوند. گره ها موجودیت ها یا اشیاء را نشان می دهند، در حالی که روابط نشان دهنده ارتباطات یا تعاملات بین آن موجودیت ها هستند. پس از شناسایی و استخراج موجودیت ها و روابط، موجودیت ها، هر موجودیت یا کلاس در فایل OWL باید با یک گره در گراف دانش مطابقت داشته باشد و هر رابطه باید با یک یال بین گره ها ترسیم شود. پس ویژگی ها را به گره ها و یال ها در نمودار دانش اختصاص دادیم. نمودار دانش به دست آمده همان اطلاعات فایل OWL را نشان خواهد داد اما در قالبی ساختاریافته و قابل خواندن توسط ماشین. هر گره می تواند یک یا چند برچسب برای دسته بندی نوع موجودیتی که نشان می دهد داشته باشد. برچسب ها در اصل برچسب هایی هستند که به ما کمک می کنند داده های خود را سازماندهی کرده و پرس و جو را آسان تر کنند. برای مثال، برچسب هایی مانند «شخص»، «محصول» یا «مکان» داریم. علاوه بر برچسب ها، گره ها همچنین می توانند یک یا چند ویژگی داشته باشند که جفت های کلید-مقدار هستند که اطلاعات مربوط به گره را ذخیره می کنند. به عنوان مثال، یک گره "Person" دارای ویژگی هایی مانند name، age، email است. روابط گره ها را به هم متصل می کنند و همچنین می توانند برچسب هایی برای توصیف نوع تعامل بین آنها داشته باشند. برای مثال، برچسب های رابطه ای مانند has_Participant، has_location، knows داریم. هنگامی که پایگاه داده گراف را پرس و جو می کنید، می توانید از برچسب ها، ویژگی ها و روابط برای یافتن و تجزیه و تحلیل زیرمجموعه های خاصی از داده ها استفاده کنید.

۳-۳ معرفی Neo4j

Neo4j یک سیستم مدیریت پایگاه داده گراف قدرتمند و همه کاره است که در سال های اخیر به دلیل توانایی آن در مدیریت ساختارهای داده پیچیده و به هم پیوسته محبوبیت پیدا کرده است. Neo4j با تاکید بر روابط بین نقاط داده و توانایی آن برای انجام پرس و جوهای بسیار کارآمد در مجموعه داده های بزرگ، به گزینه ای برای توسعه دهندگان و تحلیلگران داده در صنایع مختلف تبدیل شده است. برخلاف پایگاه های داده رابطه ای سنتی، که داده ها را در جداول دارای ردیف و ستون ذخیره می کنند، Neo4j داده ها را به عنوان مجموعه ای از گره ها و یال ها نشان می دهد، جایی که گره ها موجودیت ها و یال ها

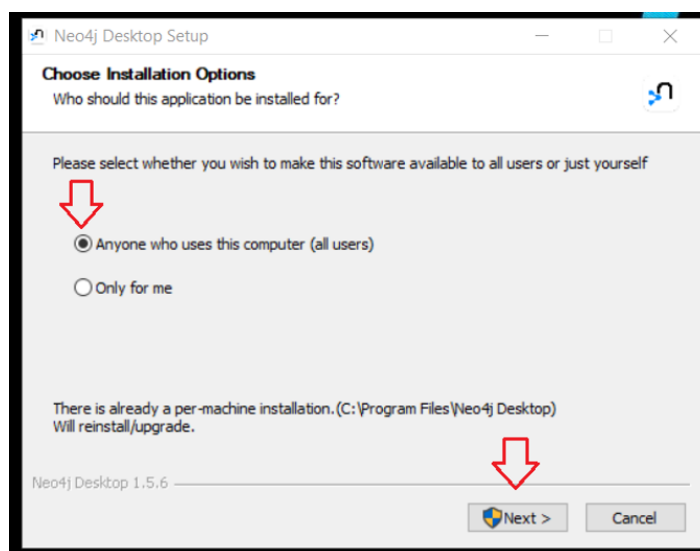
روابط بین آن موجودیت‌ها را نشان می‌دهند. این رویکرد امکان ایجاد مدل‌های داده بسیار انعطاف‌پذیر را فراهم می‌کند که می‌توانند به راحتی در صورت نیاز اصلاح و به روز شوند، و آن را برای برنامه‌هایی که با داده‌های بسیار بهم پیوسته سروکار دارند ایده‌آل می‌کند.

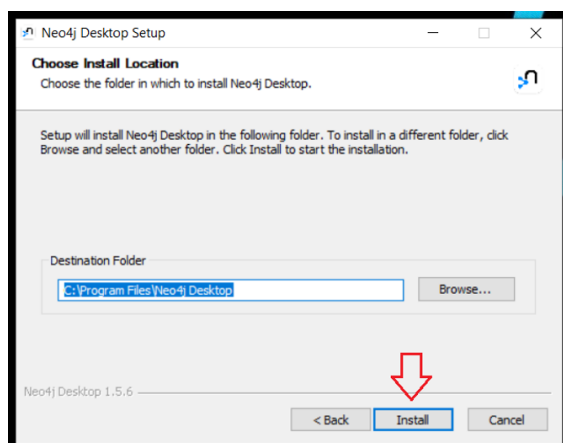
۱-۳-۳ نصب و راه‌اندازی Neo4j

ابتدا به نصب و راه‌اندازی نسخه ۶.۵.۱ آن می‌پردازیم: به سایت www.neo4j.com مراجعه کرده و از قسمت neo4j.com/download وارد بخش Started Get می‌شویم که باید ابتدا فرم مشخصات را پر کنیم و موافقت نامه مجوز Neo4j را برای نرم افزار دسکتاپ Neo4j اعلام کنیم تا فایل مورد نظر شروع به دانلود شود. سیستم مورد نیاز توصیه شده برای نصب

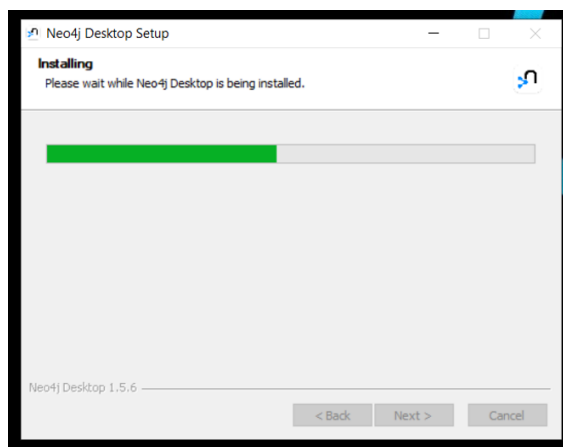
MacOS 10.10 (Yosemite)+, Windows 8.1+ with Powershell 5.0+, Ubuntu 12.04+,
Fedora 21, Debian 8

تعیین شده است. در صفحه بعد Key Activation تهیه شده است که باید از این کلید برای فعال کردن کپی DesktopNeo4j استفاده کنید. بعد از دانلود شدن و اجرای installer طبق تصاویر زیر پیش می‌رویم:

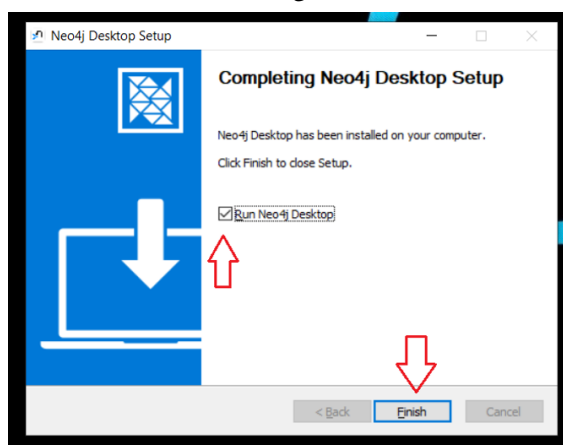




شکل ۳-۳



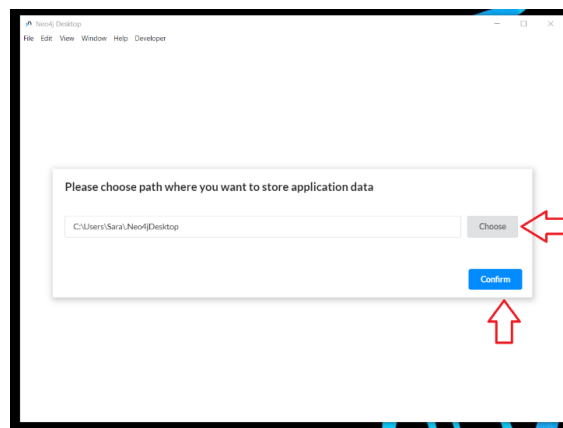
شکل ۴-۳



شکل ۵-۳

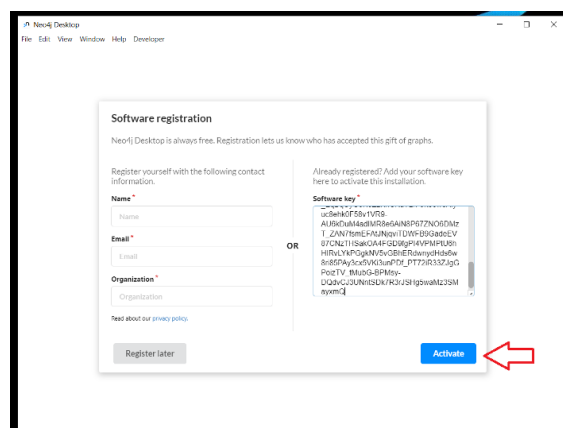


شکل ۳-۶

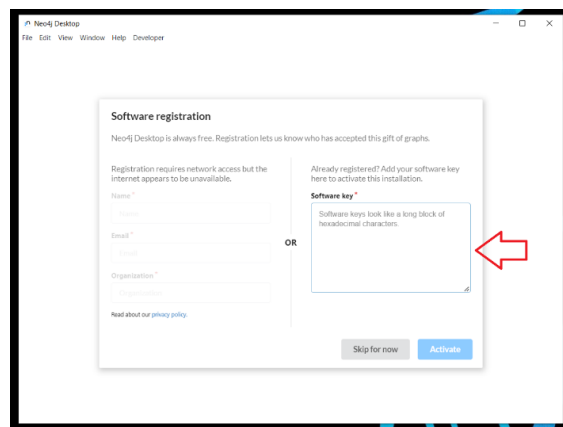


شکل ۳-۷

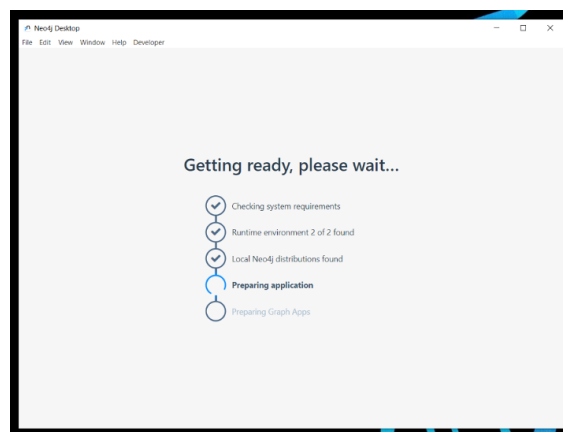
در این قسمت Key Activation داده شده در قبل را وارد کنید. و یا می‌توانید با پر کردن فرم سمت راست صفحه برنامه، کلیدی را از داخل برنامه ایجاد کنید.



شکل ۳-۸

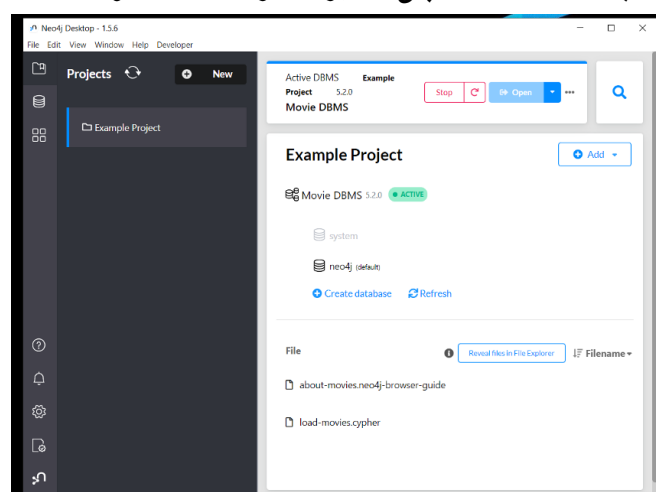


شکل ۳-۹



شکل ۳-۱۰

برنامه با موفقیت نصب شده و صفحه پنل کاربری زیر نمایان میشود:



شکل ۳-۱۱

۲-۳-۳ تعامل با رابط کاربری Neo4j

رابط کاربری Neo4j دارای سه بخش زیر است:

۱. Sidebar

۲. Editor

۳. Stream

برای ایجاد پایگاه داده، در سمت راست بر روی دکمه Add و Graph New کلیک کرده گزینه Create graph local a را انتخاب کنید سپس نام دیتابیس مورد نظر و پسورد را وارد کرده و روی Create کلیک نمایید. پس از ایجاد پایگاه داده، روی دکمه Start کلیک کنید. Open را که بزنید صفحه Neo4j Browser باز میشود. مرورگر Neo4j یک رابط بصری غنی برای تعامل با پایگاه داده گراف فراهم می کند. این شامل یک پنل در سمت چپ صفحه است که به شما امکان می دهد طرحواره پایگاه داده گراف را بررسی کنید، پرس و جوها را اجرا کنید و تنظیمات پایگاه داده را مدیریت کنید. این ویرایشگر به شما اجازه می دهد تا کوئری های Cypher را بنویسید و اجرا کنید. راه دیگر برای داشتن مرورگر Neo4j این است که آدرس <http://127.0.0.1:7474> را در مرورگر خود وارد کنید و سپس وارد ویرایشگر کوئری شوید. برای بار اول، از شما نام کاربری و رمز عبوری را درخواست می کند که همان نام و رمز عبور، هنگام نصب است.

مرورگر یک Cypher shell command تعاملی است که به شما امکان می دهد با گراف خود تعامل داشته باشید و اطلاعات موجود در آن را visualize کنید. حال کمی با زبان cypher آشنا شویم:

برای ساخت یک آبجکت از CREATE استفاده میکنیم که به فرمت زیر است:

```
CREATE (VariableName:objectType (property1 :value1 ,property2 : value2 ,... ))
```

این معادل ساخت جدول در بانکهای اطلاعاتی رابطه ای است.

برای بازیابی یک آبجکت، از کوئری زیر استفاده می شود:

```
MATCH (VariableName:objectType ) RETURN VariableName
```

که میتوان با where برای آن شرط تعیین شود.

```
MATCH (n:Person)-[:KNOWS]->(m:Person)
WHERE n.name = 'Alice'
RETURN m AS person
```

برای یافتن و بازنمایی همه گره ها از عبارت های زیر استفاده میکنیم :

```
MATCH (n) RETURN n
```

برای ایجاد رابطه بین شی a و b که شروع از a و به سمت گره b است به شکل زیر کوئری مینویسیم که نوع رابطه را داخل براکت مشخص میکنیم :

```
CREATE (a)-[:relationship_type]->(b)
```

برای update یا ایجاد یک property از set استفاده میشود:

```
SET e.property1 = $value1
```

برای حذف رابطه با مشخصات داده شده:

```
MATCH (n:Label)-[r]->(m:Label)
WHERE r.id = 123
DELETE r
```

طبیعی است گستردگی این زبان به این موارد محدود نمیشود برای اطلاعات کامل تر به داکيومنت خود سایت Neo4j میتوان مراجعه کرد.

فصل ۴

نتایج بدست آمده

۴-۱ پیاده‌سازی طرح کلی گراف دانش

پس از راه اندازی Neo4j و جمع آوری و استخراج داده ها آغاز به ایجاد گراف دانش مان میکنیم. پایگاه داده گراف مان از گره ها و یال ها تشکیل شده است که در آن گره ها موجودیت هایی مانند افراد، سازمان ها، محصولات و رویدادها را نشان می دهند و یال ها روابط بین این موجودیت ها را نشان می دهند. ابتدا از نوع کلاس thing یک کلاس Cyberspace میسازیم. سپس سه گره با برچسب های «وبلاگ» و «وب سایت» و یک گره دیگر با برچسب «شبکه اجتماعی» ایجاد می کنیم که زیرکلاس های فضای مجازی هستند. این گره ها با استفاده از رابطه "IS_A" به گره Cyberspace متصل می کنیم.

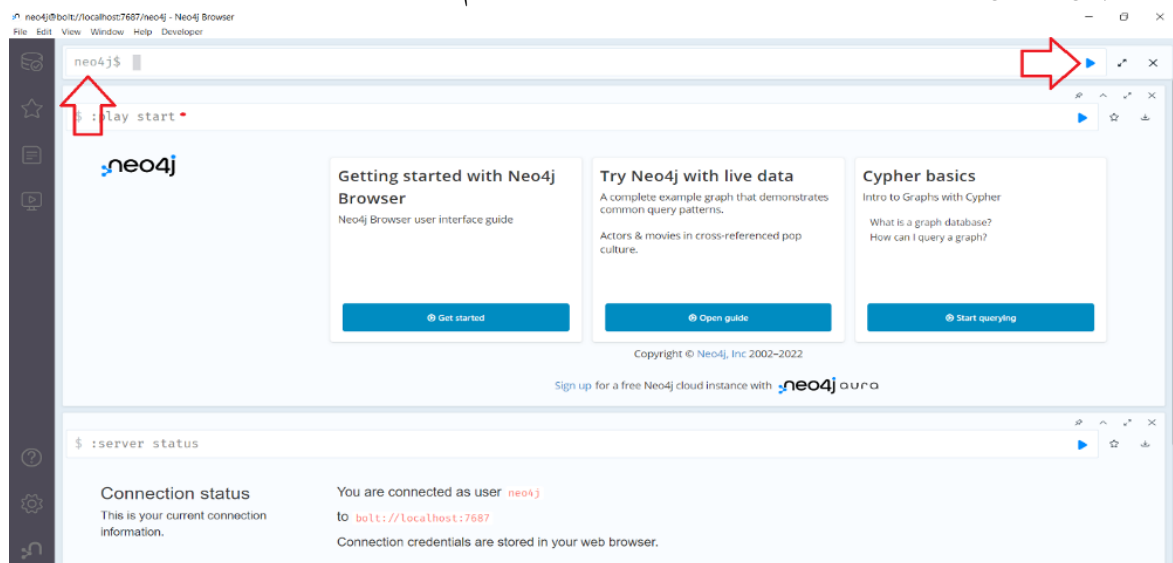
سپس به ساخت زیرکلاس های Social-Network میپردازیم که Instagram و Telegram و Twitter هستند. از آنجایی که برای property های تعریف شده مقادیری در فایل owl داده نشده است و مقدار خالی در Neo4j به منزله عدم وجود آن ویژگی است، از مقدار ۰ برای مقادیر عددی و از NULL برای رشته استفاده کردیم که بعداً طبق داده های بدست آمده کامل، آن ها را آپدیت کنیم.


```
CREATE (Cyberspace:thing {title: "Cyberspace"}),
  (Weblog:Cyberspace {title: "Weblog", url:"NULL", description:"NULL", person_name:"NULL"}),
  (Website:Cyberspace {title: "Website", url:"NULL", description:"NULL", affiliation:0}),
  (SocialNetwork:Cyberspace {title: "SocialNetwork", account_type:"NULL", description:"NULL"})
CREATE (Weblog)-[:is_a]->(Cyberspace),
  (Website)-[:is_a]->(Cyberspace),
  (SocialNetwork)-[:is_a]->(Cyberspace)

CREATE (Instagram:SocialNetwork {title:"Instagram", InstagramID:"NULL"})-[:is_a]->(SocialNetwork),
  (Telegram:SocialNetwork {title:"Telegram", TelegramID:"NULL"})-[:is_a]->(SocialNetwork),
  (Twitter:SocialNetwork {title:"Twitter", TwitterID:"NULL"})-[:is_a]->(SocialNetwork)
```

شکل ۴-۱

سپس در این قسمت از browser کوثری هارا وارد میکنیم:



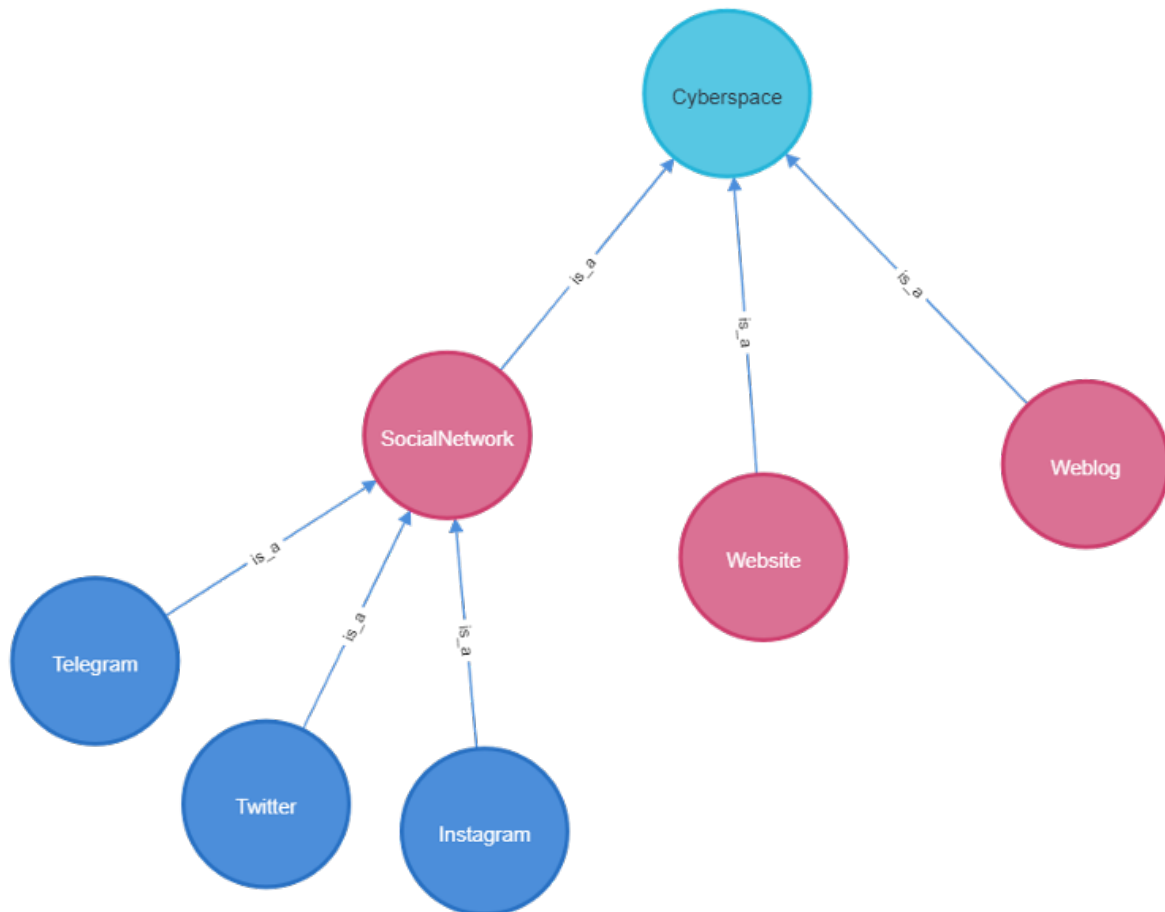
شکل ۴-۲

که برای کوثری بالا خروجی زیر را میدهد و گره ها و روابط را ایجاد میکند:



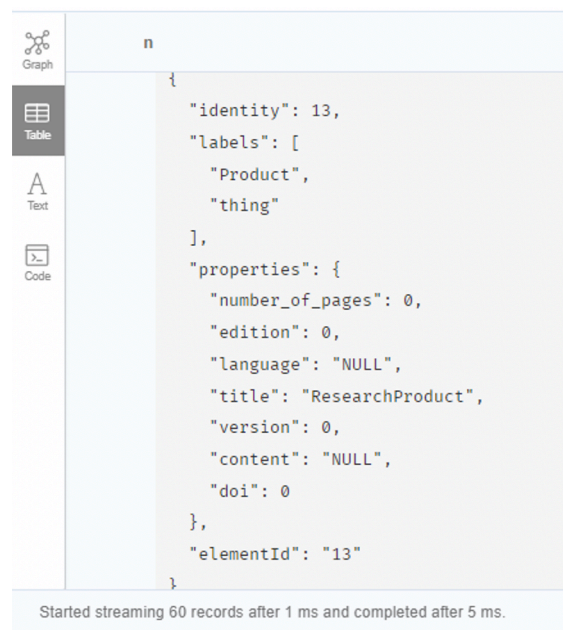
شکل ۴-۳

با نوشتن کوئری $MATCH(n)RETURN n$ کل آنچه تا کنون ساخته‌ایم را می‌توانیم به صورت گراف مشاهده کنیم که خروجی کوئری بالا به شکل زیر است:



شکل ۴-۴

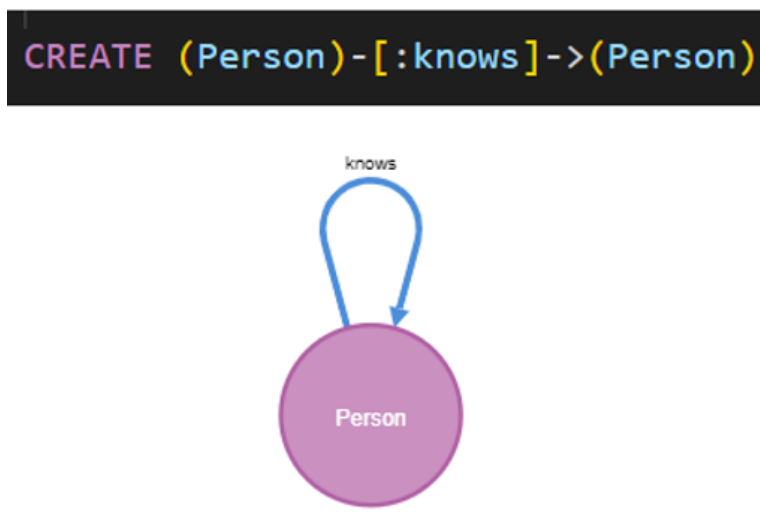
در قسمت view table می‌توانیم property های هر گره را مشاهده کنیم:



شکل ۴-۵

سپس، گره های دارای برچسب های Person، Event، Organization، Place Product و Person_Product را ایجاد کردیم. سپس چندین زیرگروه از این گره ها را ایجاد و با استفاده از رابطه "IS_A" آنها را به هم متصل نمودیم.

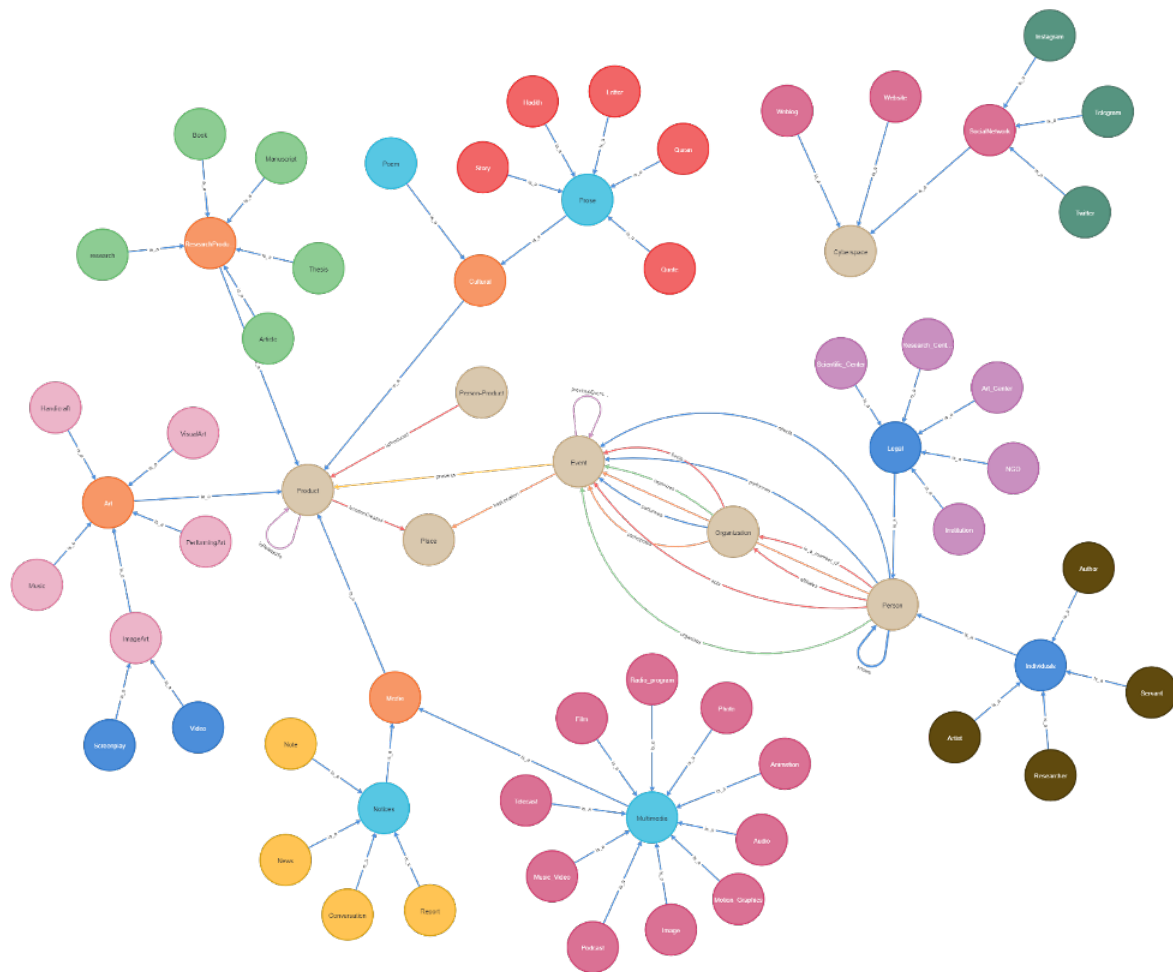
پس از ایجاد گره ها و زیرگروه ها، روابط و لبه هایی را بین آنها ایجاد کردیم تا روابط غیر زیرکلاس را نشان دهیم. به عنوان مثال، لبه ای ایجاد کردیم تا نشان دهد که یک فرد شخص دیگری را می شناسد یعنی یک self-relation به خود کلاس Person، که این امکان در Neo4j ساپورت میشود و به صورت زیر است:



شکل ۴-۶

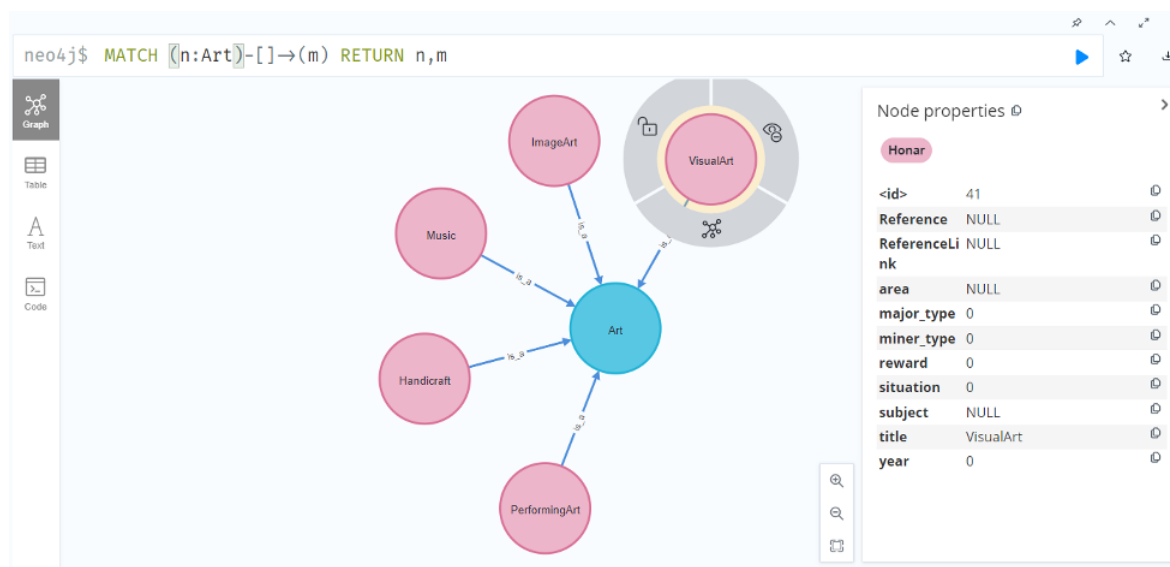
از روابط دیگری که استخراج شد یک نمونه این هم بود که انسان یک رویداد را سازماندهی می کند، در رویداد اقدام می کند، در رویداد شرکت می کند، رویداد را اجرا می کند، رویداد را هدایت می کند، به یک سازمان وابسته است، و عضو یک سازمان است. پس از آن همچنین لبه هایی ایجاد شد تا نشان دهد که سازمان بودجه رویداد را تامین می کند، رویداد را سازماندهی می کند، رویداد را اجرا می کند و در رویداد شرکت می کند.

برای رویداد اضافه کردیم که رویداد دارای مکان است، محصول را ارائه می دهد و به رویداد دیگری مرتبط است. لبه های دیگری نیز ایجاد شد تا نشان دهد که یک محصول فرهنگی یک نوع فرعی از محصول است و برای مثال یک شعر و یک نثر زیرگونه های محصول فرهنگی هستند. در نهایت بعد از ساخت همه گره ها و روابط شان که حدود ۹۱ گره و ۳۴۱ ویژگی و ۷۴ رابطه شد، نمای کلی گرافی که بدست آمد به شکل زیر است:



شکل ۴-۷

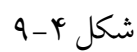
در گراف دانش ساخته شده، با کمک پرس و جوهای متعدد می توان اطلاعات متفاوتی استخراج کرد. مثلا اگر بخواهیم لیستی از تمام هنر هایی که داریم مجزا مشاهده کنیم کوئری زیر را نوشته و گره های زیرمجموعه بدست می آید که با کلیک روی هر گره در سمت راست قسمت properties Node ویژگی های آن گره را به اختصاص نشان میدهد.



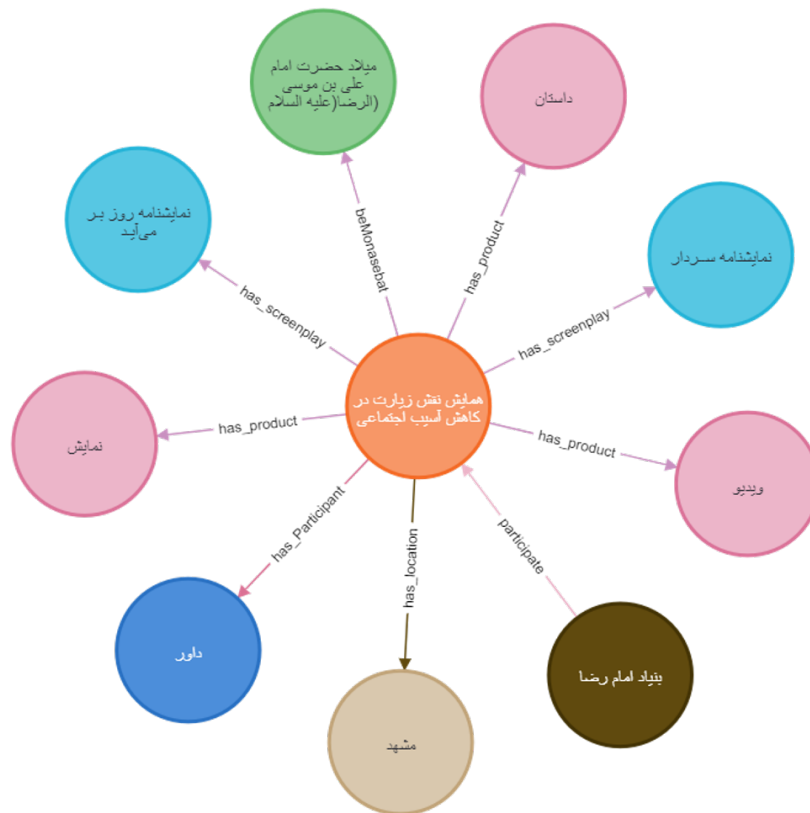
شکل ۴-۸

۴-۲ گراف دانش پروژه بنیاد

بعد از ایجاد این طرح کلی، از سایت بنیاد امام رضا (ع) جزییات بیشتری راجع به جشنواره ها استخراج کردیم و این گراف را با داده های این بنیاد بازنمایی نمودیم. ۱۲ جشنواره داریم اعم از یادمان شهری رضوی، همایش ایتار و شهادت، همایش نقش زیارت در کاهش آسیب اجتماعی و... که در شهر های مشهد و سبزوار به مناسبت میلاد حضرت امام علی بن موسی الرضا (علیه السلام) توسط این بنیاد برگزار شده است. این داده ها را وارد پایگاه داده گرافی مان کرده و گراف دانش زیر را بدست آوردیم:



```
MATCH (e:Event {Ev:1})-[r]-(n)
RETURN e, r, n
```

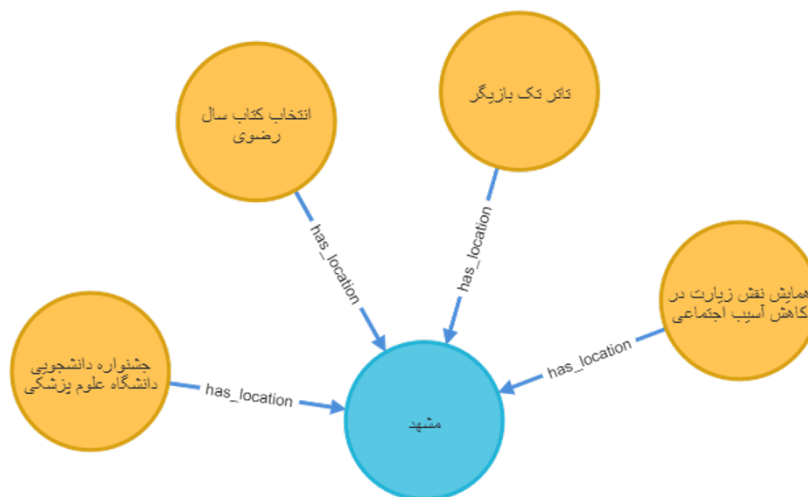


شکل ۴-۱۰

یا مثلا اگر بخواهیم ببینیم فقط در شهر مشهد چه همایش هایی برگزار می شود کوئری زیر را وارد کرده و جشنواره ها را مشاهده میکنیم:

```
MATCH (m)-[]->(n:Place {city_name:"مشهد"}) RETURN n , m
```

شکل ۴-۱۱

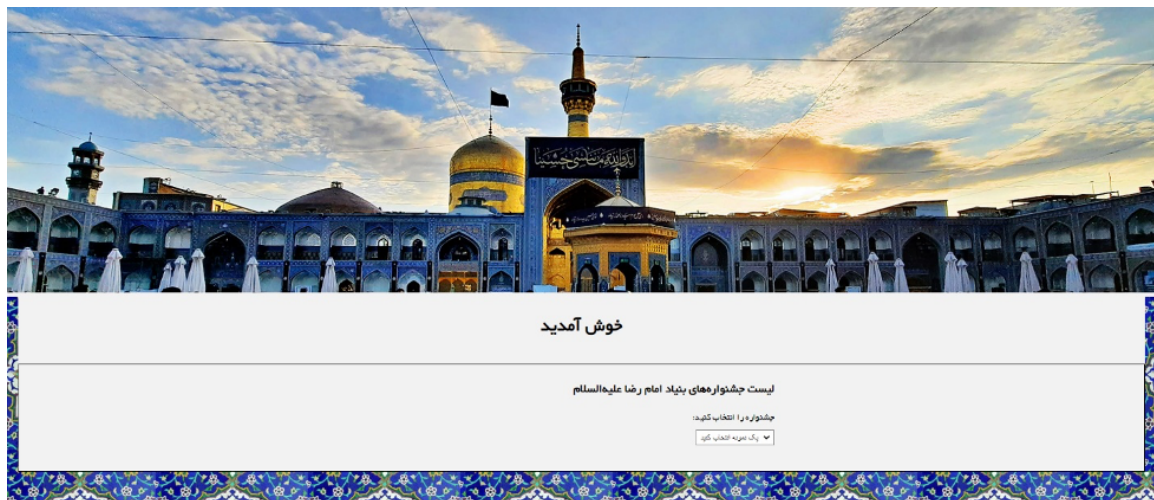


شکل ۴-۱۲

۳-۴ رابط کاربری

جهت داشتن یک رابط کاربری مناسب که کاربر با کوئری های پیچیده درگیر نشود و به آسانی اطلاعات این جشنواره ها را دریافت نماید نیز یک صفحه وب طراحی کرده ایم تا با انتخاب گزینه مراسم ها در بک گراند کوئری ها را نوشته و از پایگاه داده گرافی مان پاسخ پرس و جوی مورد نظر را برای کاربر نشان دهیم که نمونه هایی از آن به شکل زیر است:

صفحه اصلی:



شکل ۴-۱۳

انتخاب جشنواره پنجم توسط کاربر از قسمت گزینه ها و نمایش گراف نتیجه:



شکل ۴-۱۴

فصل ۵

جمع بندی و کارهای پیشرو

این پایان نامه کارشناسی استفاده از Neo4j را برای ایجاد یک گراف دانش نشان داده است. از طریق اجرای این فناوری، مشاهده شد که چگونه می توان از آن برای سازماندهی و ساختار داده ها به گونه ای استفاده کرد که تجزیه و تحلیل و تصمیم گیری مؤثرتر را تسهیل کند.

این پژوهش به بررسی مفاهیم اساسی گراف های دانش و زبان هستی شناسی وب پرداخته است. همچنین در مورد مزایای استفاده از Neo4j توضیح داده شد که بعنوان یک سیستم مدیریت پایگاه داده قدرتمند و همه منظوره برای ذخیره، مدیریت و کوئری گراف های در مقیاس بزرگ استفاده می شود. به طور کلی، این پایان نامه یک نمای کلی از فرآیند ایجاد یک گراف دانش با استفاده از Neo4j، از چارچوب مفهومی تا جزئیات پیاده سازی ارائه کرده است.

در طول فرآیند انجام پروژه، ما با چالش های مختلفی مانند شناسایی منابع داده مربوطه، طراحی یک طرحواره مناسب و بهینه سازی عملکرد کوئری ها و عدم وجود راهنمای کاربری Neo4j مواجه بوده ایم. با این حال، ما استراتژی هایی را برای مقابله با این چالش ها ایجاد کرده ایم، مانند استفاده از تکنیک های پیش پردازش داده ها و تحلیل هستی شناسی های موجود.

در این پژوهش گراف دانش مورد نظر برای داده های بنیاد امام رضا (ع) ساخته شد و برای آن رابط کاربری جهت سهولت استفاده ایجاد گردید. نتایج پیاده سازی های ما بر روی این داده ها نشان داده است که ایجاد یک گراف دانش با استفاده از Neo4j می تواند منجر به بهبود یافته ها، کشف دانش و تصمیم گیری بهتر شود. مطالعات موردی ما نشان داده است که Neo4j یک انتخاب مناسب برای

برنامه‌هایی است که نیاز به تجزیه و تحلیل بلادرنگ و پردازش داده‌های پیچیده و به هم پیوسته دارند، و آن را به یک انتخاب محبوب برای موارد استفاده این‌چنینی تبدیل می‌کند.

ما امیدواریم که این تحقیق به مجموعه دانش رو به رشد در این زمینه کمک کند و الهام بخش کار آینده در مدیریت و تجزیه و تحلیل داده‌های مبتنی بر گراف باشد.

با افزایش داده و به روزرسانی مقادیر گراف دانش، طبیعی است اطلاعات متنوع و سودمندتری به دست خواهیم آورد، که یکی از راه‌های جمع‌آوری دیتا استفاده از خزشگرهای وب است که حاکی از راه طولانی مسیر پیش رو میباشد و نشان می‌دهد کار این پروژه در اینجا به پایان نمی‌رسد.

بعنوان کار پیشرو ادغام با فناوری‌های یادگیری ماشین و هوش مصنوعی می‌توان نام برد. توانایی Neo4j در مدل‌سازی روابط پیچیده، آن را به انتخابی ایده‌آل برای یادگیری ماشین و کاربردهای هوش مصنوعی تبدیل می‌کند. کار آینده می‌تواند بر روی ادغام Neo4j با چارچوب‌های یادگیری ماشین تمرکز کند تا زمینه ساز تجزیه و تحلیل‌ها و پیش‌بینی‌های پیشرفته‌تر باشد.

مراجع

1. neo4j.com
2. <https://neo4j.com/docs/cypher-manual/current/>
3. neo4j.com/developer/cypher/guide-cypher-basics/
4. www.shamstoos.ir/fa/news/bonyad
5. [/www.w3.org/OWL/](http://www.w3.org/OWL/)
6. www.quackit.com/neo4j/tutorial/neo4j_query_language_cypher.cfm
7. documentation.mindsphere.io/MindSphere/howto/howto-create-ontology-owl.html
8. neo4j.com/blog/build-knowledge-graph-from-scratch-even-if-youre-not-full-blown-developer/
9. protege.stanford.edu
10. towardsdatascience.com/how-to-build-a-knowledge-graph-with-neo4j-and-transformers-72b9471d6969
11. farin.academy/google-knowledge-graph/
12. ontology101tutorial.readthedocs.io/en/latest/StartingProtege.html

13. www.javatpoint.com/advantages-of-neo4j
14. https://www.tutorialspoint.com/neo4j/neo4j_data_model.htm



Ferdowsi University of Mashhad
Department of Computer Engineering

B.Sc. Projects

Creating a Knowledge Graph using Neo4J

By:

Sara Asadi

Supervisor:

Dr. Mohsen Kahani

Feb 2023