

شناسایی موجودیت‌های نام‌دار در متون فارسی

پونه سادات مرتضوی

آزمایشگاه تحقیقاتی پردازش زبان طبیعی، دانشکده

مهندسی برق و کامپیوتر، دانشگاه شهید بهشتی، تهران، ایران

pooneh_mortazavi@yahoo.com

مهرنوش شمس فرد

آزمایشگاه تحقیقاتی پردازش زبان طبیعی، دانشکده

مهندسی برق و کامپیوتر، دانشگاه شهید بهشتی، تهران، ایران

m-shams@sbu.ac.ir

چکیده: شناسایی موجودیت‌های نام‌دار^۱ در پردازش زبان طبیعی به عملیاتی گفته می‌شود که در طی آن کلیه‌ی اسامی خاص موجود در متن و متعلق به مقوله‌های معنایی مختلف، شناسایی و استخراج می‌گردند.

در این مقاله به معرفی سیستمی توسعه یافته به منظور تشخیص اسامی نام‌دار و دسته‌بندی آنها در زبان فارسی پرداخته می‌شود. این سیستم با بکارگیری ساختار واژوی اسامی خاص و نیز الگوهای متنی ممکن برای اسم‌های خاص متعلق به یک دسته، سعی در شناسایی موجودیت‌های نام‌دار می‌کند. علاوه بر آن، با بکارگیری برچسب نحوی و معنایی برای هر کلمه و توجه به سایر رخداد‌های آن در متن، عملیات رفع ابهام برای بهبود شناسایی را انجام می‌دهد.

واژه‌های کلیدی: شناسایی موجودیت‌های نام‌دار، پردازش زبان طبیعی، مقوله بندی موجودیت‌های نام‌دار، برچسب معنایی، الگو شناسایی، ریخت شناسی موجودیت‌های نام‌دار.

۱- مقدمه

موجودیت نام‌دار به کلمه و یا عبارتی گفته می‌شود که برای ارجاع به نمونه‌های یک مقوله‌ی مشخص مانند شخص، شرکت یا موسسه، تاریخ، بیماری، گونه‌ای باکتری و سایر بکار می‌رود. نیاز به شناسایی موجودیت‌های نام‌دار، در دنیای امروز که عصر ارتباطات و اطلاعات است رو به رشد می‌باشد. شناسایی موجودیت‌های نام‌دار برای جستجوهای معنادار، ترجمه‌ی خودکار، استخراج خودکار مفاهیم متن، کشف ارجاعات در متن و بسیاری دیگر از زمینه‌های مربوط به پردازش زبان‌های طبیعی کاربرد دارد.

اینکه سیستم چه نوع موجودیتی را تشخیص دهد و یا به بیان دیگر دسته‌های معنایی مورد نظرش چه باشند، وابسته به زمینه‌ی کاربردی سیستم می‌باشد. شناسایی موجودیت نام‌دار در علم زیست‌شناسی می‌تواند تشخیص اسامی وابسته به انواع «پروتئین‌ها»، «DNAها»، «نوع سلول» و ...، در حوزه پزشکی

می‌تواند تشخیص انواع بیماری‌ها، داروها، مراکز درمانی و مانند این‌ها و در حوزه تجارت نام شرکت‌ها و موسسات، تراکنش‌های مالی، بورس و غیره باشد. همچنین این امر می‌تواند به صورت خیلی خاص مثلاً فقط برای کشف اسامی شرکت‌های تولید کننده فولاد از روی متون مربوطه بکار رود. یک دسته بندی عام که در بسیاری تحقیقات NER مورد استفاده قرار می‌گیرد و در MUC-6 [۲] نیز ملاک مقایسه بوده است، دسته بندی معنایی بر اساس «شخص»، «سازمان» و «ناحیه» است. در این مقاله پس از مروری بر فعالیت‌های انجام شده در این حوزه به معرفی سیستم NER پیاده سازی شده برای زبان فارسی می‌پردازیم. در این معرفی ابتدا به دسته‌های معنایی مورد نظر در این سیستم اشاره نموده سپس با ارائه معماری کلی سیستم، رهیافت مورد استفاده در آن را شرح خواهیم داد.

۲- کارهای مرتبط

در زبان‌های مختلف بر روی مبحث «تشخیص و دسته بندی اسامی نام‌دار (NER)» به عنوان زیر شاخه‌ای از پردازش زبان طبیعی کارهایی انجام شده است. اگرچه نسبت بالایی از تحقیقات NER متعلق به زبان انگلیسی بوده است، فعالیت‌هایی برای زبان‌هایی چون آلمانی، یونانی، ژاپنی، فرانسوی، ایتالیایی، هلندی، چینی، روسی، کره‌ای، رومانیایی و ترکی نیز انجام شده است. اکثر کارهای انجام شده متعلق به زبان، دامنه و یا گونه نوشتاری خاص است و بزرگ‌ترین مشکل این گونه سیستم‌ها نیز مربوط به انتقالشان به دامنه‌ی جدید می‌باشد. در این راستا برای زبان فارسی، کار قابل توجهی انجام نگرفته است؛ البته در مواردی در حواشی موضوع کارهایی شده است که از جمله می‌توان به تعیین عبارات اسمی هم-مرجع در فارسی [۱] اشاره نمود.

روش‌های اولیه‌ای که برای کشف موجودیت‌های نام‌دار پیشنهاد می‌شد بیشتر مبتنی بر قانون بود، در حالی که امروزه اکثر کارها به سمت استفاده از روش‌های یادگیری پیش می‌رود [۳] و [۴].

تولید نتیجه تشکیل شده است. شمای کلی از معماری سیستم در شکل (۱) نمایش داده شده است.

در این سیستم پس از دریافت متن، با ورودی مورد پذیرش سیستم تطبیق داده می‌شود، واحدهای خاص شناسایی شده و سپس بر روی آن برچسب POS زده می‌شود. بعد از این برچسب گذاری نحوی، با توجه به ترکیب‌های اعداد، کلمات و علامات به شناسایی اولیه اعداد و برچسب گذاری‌های معنایی پرداخته می‌شود و برچسب گذاری ثانویه‌ای از نام‌های خاص مکمل این مرحله است. در انتهای مرحله پیش پردازش، با استفاده از ساختار واژگانی اسامی خاص، سعی در شناخت این واژگان مستقل از متن می‌شود.

در مرحله شناسایی، با توجه به نحوه چیدمان واژگان در متن و الگوهای موجود برای هر مقوله، موجودیت‌های نامدار شناخته شده و مرز آنها با توجه به تکرارها و برچسب نحوی و معنایی کلمات مشخص می‌شود.

در انتها با حذف خطاها و تکرارهای شناخته شده، هر موجودیت در مقوله‌ی مناسب خود دسته بندی شده و خروجی با برچسب گذاری معنایی موجودیت‌های نامدار با عنوان مقوله‌ای که در آن قرار می‌گیرند تولید می‌شود.

در ادامه پس از معرفی مجموعه برچسب‌های معنایی مورد استفاده در این سیستم به ارائه نحوه عملکرد سیستم در بخش‌های مجزا خواهیم پرداخت.

۴- مقوله بندی برچسب‌های معنایی

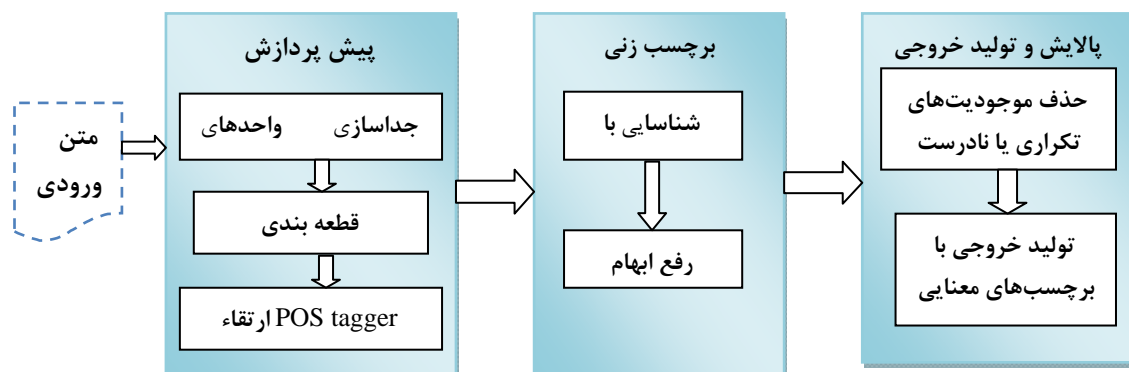
موجودیت‌های نامدار در مقوله‌های مختلفی دسته بندی می‌شوند. این مقوله‌ها می‌توانند بسته به هدف سیستم تشخیص اسامی نامدار، در طیفی از زمینه‌های تخصصی چون تشخیص اسامی بیماری‌ها و «DNA» و ... تا زمینه‌های عامی چون نام افراد یا سازمان‌ها رده بندی شوند. همچنین در این طبقه بندی هر مقوله‌ی کلی هم می‌تواند به زیر مقوله‌های جزئی‌تر تقسیم شود.

سیستم‌های توسعه یافته بر پایه یادگیری برای NERC سه دسته تقسیم می‌شوند [۵]: یادگیری با نظارت^۲ (SL)، یادگیری نیمه نظارت^۳ (SSL) و یادگیری بی نظارت^۴ (UL). در یادگیری با نظارت، نیاز به حجم زیادی از متون حاشیه نویسی شده است و از تکنیک‌هایی نظیر درخت تصمیم‌گیری و مدل مخفی مارکوف (HMM) استفاده می‌کند. این روش مناسب برای سیستم‌های وابسته به دامنه است (۳). در یادگیری نیمه نظارت از مجموعه‌ای از کلمات نمونه برای یادگیری استفاده می‌شود (مثل مجموعه‌ای از نام کتاب‌ها). این روش‌ها ابتدا سعی در یافتن موجودیت‌های نمونه در متون کرده و سپس با توجه به ساختار متن در بر دارنده کلمات، بدنبال سایر موجودیت‌هایی از آن نوع می‌گردد. یادگیری بدون نظارت، با استفاده از منابع لغوی چون WordNet^۵ و الگوهای لغوی و محاسبات آماری به شناسایی می‌پردازد. به عبارتی تکنیک معمول برای یادگیری بدون نظارت بر پایه دسته بندی موجودیت‌ها است [۶].

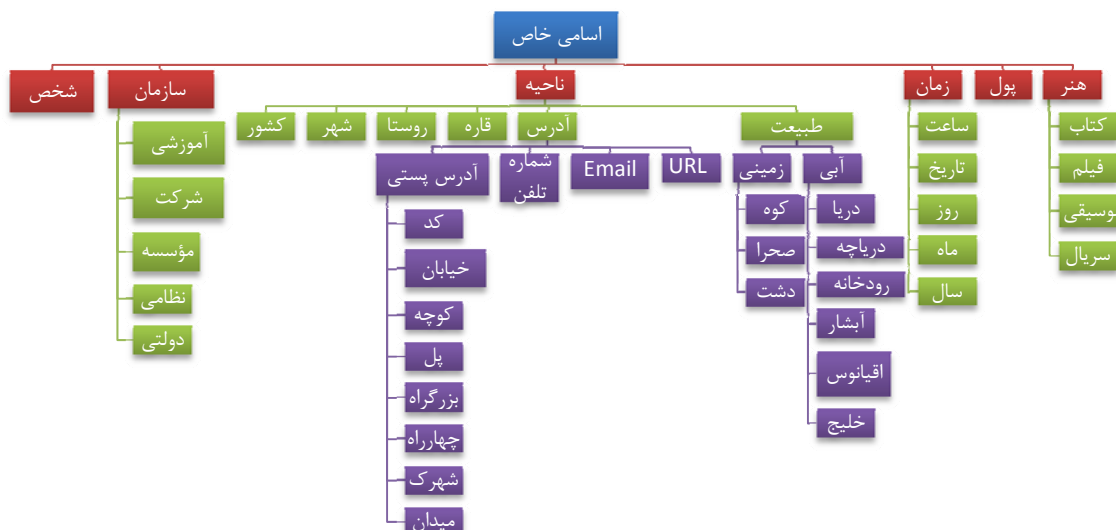
در زبان فارسی با توجه به عدم وجود پیکره‌های برچسب خورده با برچسب‌های معنایی یا طبقات نامدار و همچنین دردسترس نبودن واژگان‌های کاملی چون وردنت به ناچار باید تمرکز بیشتر بر روش‌های مبتنی بر قاعده باشد. سیستم معرفی شده در این مقاله از لحاظ استفاده از قواعد زبان شناسی، مشابه سیستم توسعه یافته [۷] عمل می‌کند با این تفاوت که در زبان فارسی همچون زبان انگلیسی امکان درک موجودیت نامدار با توجه به بزرگی حرف اول کلمات وجود ندارد. پس بجای آن برای بهبود سیستم از برچسب‌گذاری‌های نحوی و مجموعه‌ای از اسامی خاص مانند مجموعه اسامی مناطق، افراد، روزهای هفته، واحدهای پولی، عناوین و ... استفاده می‌شود. استفاده از ریخت-شناسی کلمات از دیگر نقاط قوت این سیستم می‌باشد که در مقاله [۸] هم به نوعی به آن اشاره شده است.

۳- معماری سیستم پیشنهادی

سیستم NER معرفی شده در این مقاله از سه بخش اصلی (۱) پیش پردازش، (۲) تشخیص و برچسب زنی و (۳) پالایش و



شکل ۱: شمای کلی از معماری سیستم



شکل ۲: بخشی از سلسله مراتب برچسب‌های معنایی

در این پیمان‌نامه پس از برچسب‌زنی نحوی معمولی، برای اصلاح برچسب‌های نحوی ایجاد شده برای هر کلمه، از نوع دیگری از برچسب‌های اسم خاص ذخیره شده در پایگاه داده استفاده می‌شود. همچنین برچسب‌های معنایی اولیه، برای اعداد و ترکیب آنها با علامات اطراف و نیز ساخت شناسی کلمات خاص، در این مرحله زده می‌شود تا در مرحله تشخیص، اصلاح و تکمیل شوند.

۵-۱-۱- استخراج واحدهای خاص

برخی از واحدهای متنی لازم است به همان شکل وارد شده در متن و بدون هیچ پردازش یا تغییر شکلی استخراج شوند. از جمله این واحدها می‌توان به آدرس صفحات وب (URL) و آدرس نامه‌های الکترونیکی اشاره نمود. این آدرس‌ها لازم است قبل از انجام قطعه‌بندی متن جدا شوند چراکه بدلیل پیش پردازش‌هایی که بر روی متن ورودی صورت می‌گیرد و ایجاد فواصل میان علائم نگارشی با سایر کلمات، یک آدرس URL و یا Email که دنباله‌ی چسبیده‌ی حروف در آنها حائز اهمیت است، از هم جدا شده و درصد درستی تشخیص را پایین می‌آورد. به همین منظور کشف این دو موجودیت نامدار که جزو دسته آدرس در ساختار سلسله‌مراتبی برچسب‌های معنایی قرار می‌گیرند، باید پیش از این جداسازی‌ها انجام گیرد. پس از جمله اولین اقداماتی که انجام می‌گیرد شناسایی این دو مقوله در متن و جداسازی آنها پس از شناسایی از متن می‌باشد.

از جمله فواید جداسازی‌ها افزایش دقت در شناسایی و نیز افزایش سرعت در tagger و مراحل پس از آن می‌باشد. اما از

در این سیستم، ساختار سلسله‌مراتبی برای مقوله‌های موجودیت‌های خاص در دسته‌ی عام طراحی شده است که در شکل (۲) به تصویر کشیده شده است.

یک موجودیت نامدار می‌تواند در هر رده از سلسله‌مراتب بالا قرار گیرد. هر چه در عمق پایین‌تر، برچسب معنایی آن دقیق‌تر است.

۵- شناسایی موجودیت‌های نامدار در جمله فارسی

تشخیص موجودیت‌های نامدار به طور کلی به دو صورت انجام می‌شود، یکی به صورت درونی که با استفاده از برچسب‌های معنایی و برچسب‌های مقوله نحوی نسبت داده شده به هر کلمه انجام می‌گیرد و شیوه‌ی دوم به صورت برونی است و تشخیص موجودیت نامدار با توجه به متنی که کلمه در آن قرار گرفته است و کلمات اطراف آن صورت می‌گیرد.

این سیستم، از هر دو روش استفاده می‌کند و در هر کدام برای افزایش صحت جواب حاصله، از چندین راه حل ممکن بهره می‌گیرد. در ادامه به بیان این راه‌حل‌ها و نحوه‌ی بکارگیری آنها در سیستم پرداخته می‌شود.

۵-۱- پیش‌پردازش

پس از دریافت متن ورودی و پیش از انجام هر گونه عملیاتی بر روی متن، ابتدا واحدهای خاص (چون URL) شناسایی می‌شوند و سپس پردازش‌های اولیه مورد نیاز بر روی متن، اجرا می‌شوند. این پردازش‌ها شامل قطعه‌بندی و تحلیل ساختاروی احتمالی جهت آماده‌سازی متن برای ورود به پیمان‌نامه برچسب‌زن مقوله نحوی ارتقاء یافته است.

طرف دیگر، ورود آن به کد از ویژگی سادگی در ارتقاء آن می-کاهد.

پیش از انجام شناسایی باید فرمت نگارشی و علائم این دو آدرس بررسی شده و سپس با در نظر گرفتن آنها رشته‌ای را که دارای آن ویژگی‌ها است در یکی از مقوله‌های ذکر شده قرار داد و آن را از متن برای جلوگیری از پردازش‌های اضافی حذف کرد. از آنجایی که بهترین کاراکتری که در شناسایی Email مؤثر است، قابل استفاده در URL نیز می‌باشد، مهم است که پیش از شناسایی Email‌ها به شناسایی URL پرداخته شود.

۵-۱-۱-۱- فرمت آدرس Email

برای تشخیص آدرس Email در متن باید دانست که از چه اجزایی تشکیل شده است و در هر جزء آن چه کاراکترهایی مجاز به استفاده می‌باشند. هر آدرس Email به طور کلی از دو جزء «بخش محلی»^۷ و «دامنه اینترنتی»^۸ تشکیل شده است و به این شکل می‌باشد:

«local-part "@" domain». «بخش محلی» که آغازگر آدرس است با علامت @ از «دامنه» جدا می‌شود. هر یک از این دو جزء هم فرمت و محدودیت‌های کاراکتری خود را دارد و گرامر صحیح آن در RFC آمده است. برای کشف صحیح Email باید از این گرامر پیروی کرد.^[۹]

۵-۱-۱-۲- فرمت URL [۱۰]

URL به رشته‌ای از کاراکترها گفته می‌شود که بدون هیچ فاصله‌ای از یکدیگر آمده‌اند و نشان دهنده‌ی منبعی موجود در اینترنت می‌باشد. ویژگی‌های آن از مفاهیم بیان شده توسط درک اطلاعات سراسری شبکه گسترده جهانی مشتق شده است.

به طور کلی طبق استاندارد بیان شده در rfc1738 و rfc4395، URL‌ها به این صورت نوشته می‌شوند: <بخش خاص شما>: <شما>^۹

تفسیر «بخش خاص شما» به «شما» بستگی دارد. اسامی شما، از دنباله‌ای از کاراکترها تشکیل می‌شود. کاراکترهای لاتین می‌توانند به صورت بزرگ هم نوشته شوند که مفسر آنها را به حروف کوچک تبدیل می‌کند. (مثلاً «HTTP» هم همچون «http» ممکن است). شناسایی URL تنها به کاراکترهایی که برای آن استفاده شده است بستگی دارد. پس برای شناسایی URL تنها لازم است کاراکترهای مجاز برای استفاده در آن را بیابیم.

هر شیمای مورد استفاده در سطح اینترنت، پیش از معرفی در مخزنی در IANA^۱ ثبت می‌شود. به طور کلی تمامی URL‌ها در IANA وجود دارند و هر شیمایی که جدید تعریف می‌شود باید علاوه بر نام، شیوه و علائم دسترسی به منابع درون آن شیمای معرفی کند.

شیمای ثبت شده در IANA به سه دسته تقسیم می‌شوند: «شماهای دائمی»، «شماهای موقتی» و «شماهای قدیمی». اعضای هر یک از این سه دسته، طبق آخرین به روز رسانی مخزن IANA در پایگاه داده‌های سیستم ذخیره شده است و برای شناسایی URL‌ها بکار می‌رود. [۱۱]

۵-۱-۲- قطعه بندی و برچسب گذاری مقوله نحوی ارتقاء یافته

پس از استخراج واحدهای خاص متن تحت عمل قطعه بندی قرار می‌گیرد تا مرز کلمات و جملات آن مشخص شود. برای عمل قطعه بندی، کلیه عملیات لازم اعم از جدا کردن کلمات چسبیده، برداشتن فاصله‌ها و نیم‌فاصله‌های اضافه، چسباندن کلمات جدا را سیستم انجام می‌دهد. بعد از این پردازش‌های اولیه، متن حاصل به برچسب زن مقوله نحوی^{۱۱} داده می‌شود. این برچسب زن [۱۲] در آزمایشگاه NLP دانشگاه شهید بهشتی توسعه یافته و متن را بر اساس مجموعه‌ای از ۱۹ نوع برچسب^{۱۲} برچسب گذاری می‌کند.

پس از برچسب گذاری نحوی، محل هر واژه در جمله هم به عنوان یک نمایه قابل بازیابی به کلمه الصاق و نگهداری می‌شود. این امر موجب می‌شود در صورت نیاز (در مراحل بعدی) بتوان به کلمات پس از آن در متن دسترسی داشت و از خاصیت تکرار واژه‌های پشت هم در یک متن بهره برد. همچنین در مرحله‌ی نهایی برای تشخیص و اصلاح خطاهایی نظیر شناسایی اسامی خاص موجود در اسم خاص دیگر، نیاز به نگهداری جای هر واژه در متن است.

علاوه بر آن حروفی را که تحلیلگر ساختوازی مورد استفاده از کلمات حذف کرده است (مثل «ها» یا ضمیرهای متصلی چون «م») برای هر کلمه نگهداری می‌شود تا در مراحل آتی از آنها استفاده شود. این حروف خود می‌توانند عامل تشخیص دهنده‌ی عدم تعلق یک واژه به یک نوع موجودیت باشند.

۵-۱-۲-۱- برچسب زنی اعداد

در سیستم شناسایی موجودیت های نامدار، علاوه بر برچسب زنی معمول نیاز به برچسب های پیشرفته تری نیز هست

که بر اساس ترکیب برچسب های قبلی بدست می آید. یک دسته از این برچسب ها مربوط به اعداد صحیح و اعشاری نوشته شده به صورت حرفی و عددی و همچنین نمایش های مختلف زمان و تاریخ است. کار دیگری که در این مرحله انجام می شود ترکیب و تشخیص اعداد است. در این مرحله چند نوع برچسب جدید تعریف می شوند که با ترکیب علامات (!، '، '، ' و '-') با اعداد ترکیب جدید با برچسب جدید ایجاد می کند. همچنین اعداد فارسی را که بینشان 'و' می باشد را با هم ترکیب کرده و یک عدد مرکب ایجاد می کند.

در جدول (۱) برچسب های بکار رفته در این مرحله آورده شده است.

جدول ۱: برچسب های مقوله های نحوی اضافه

کد مقوله نحوی	توضیح	کد مقوله نحوی	توضیح
FNo	عدد اعشاری	timeNo	ترکیب زمان (۸:۵۶)
moneyNo	عدد پول (۱۲۳.۲۳۴.۰۰۰)	dateNo	ترکیب تاریخ (۱۹۹۹/Jun/۱۲)

بر اینکه نوع هر وند را نشان می دهد محدودیت وند را هم در پذیرش نوع ریشه ی کلمه (مثلاً اسم خاص فقط می تواند باشد و یا ...) و نیز نقش سازنده وند را هم نشان می دهد.

پس از مرحله ثانویه برچسب گذاری، به بررسی ساختار کلمات به دو شکل مجزا و داخل متن پرداخته می شود. در مرحله ی اول هر کلمه از لحاظ اجزای تشکیل دهنده اش و ترکیب حرفی خاص که در ابتدا و انتهای آن آمده است بررسی می شود. مثلاً کلمه «گلشیفته» که در مجموع کلمات ذخیره شده وجود ندارد و عبارتست از «گل + شیفته»، گل هم که به عنوان پیشوند چسبیده اسم ساز شناخته شده است و برچسب شیفته هم جز برچسب های مجاز ترکیب با گل است پس اسم خاص تشخیص داده می شود.

در مرحله دوم، کلمه ی قبل و بعد کلمه مذکور - که اینک ممکن است به دلیل مرحله قبل، نوع آن تغییر کرده باشد - را نگاه می کند. در فارسی گاهی اعضای یک کلمه به صورت کلمات مجزا با فاصله نوشته می شوند. بنابراین اسم خاصی چون «معین زاده» به صورت مجزا خوانده و چه بسا نادرست برچسب زده می شود. در این مرحله اینگونه خطاها تصحیح می شوند.

پس از انجام این پیش پردازش ها وارد مرحله پردازش اصلی می شویم که در ادامه به آن پرداخته می شود.

۵-۲-۱- شناسایی موجودیت و برچسب معنایی آن

پس از مرحله پیش پردازش نوبت به شناسایی موجودیت های نامدار با استفاده از متن برچسب گذاری شده می رسد. در این مرحله ابتدا NEها با استفاده از روش های بکار برده شده که در ادامه مطرح شده اند حدس زده می شوند، سپس صحت حدس های زده شده بررسی می شود و نهایتاً NEها و نوعشان مشخص می شود.

۵-۲-۱- تشخیص مبنتی بر الگو

برای تشخیص یک کلمه در متنش، می توان از الگو استفاده کرد. در این سیستم نیز توجه به الگویی که هر کلمه در آن قرار گرفته به عنوان مبنای تشخیص بکار می رود. الگوها به صورت مجموعه ای از کلمات و برچسب هایی هستند که در متن بدنبال هم می آیند. برای هر نوع از موجودیت های نامدار که در ساختار سلسله مراتبی بخش قبل آورده شده است می توان یک یا چندین الگو در نظر گرفت.

الگوها به طور عام به صورت «دسته x» و یا «دسته xx» ظاهر می شوند. x در الگوی اول یک موجودیت یک کلمه ای و

۵-۲-۱- برچسب زنی اسامی خاص

برای ارتقاء و افزایش کارایی برچسب گذاری ها از پایگاه داده دیگری استفاده می شود که در بر دارنده ماه ها، روزها، عناوین اشخاص، اماکن، واحدهای پولی و اسامی خاص اند. این پایگاه داده کلمات ترکیبی (مثل «آلیل پیروپیل دی سولفید») را نیز در بر دارد که در صورت وجود چنین کلمه ای در متن، آن جایگزین چندین کلمه برچسب گذاری شده قبلی با توجه به مکانش می شود. به این ترتیب به با اضافه شدن این برچسب های جدید، برچسب گذاری ارتقاء می یابد.

۵-۲-۱- ریخت شناسی^{۱۳} اسامی خاص

در زبان فارسی بعضی از اسامی خاص از ساختار خاص خود پیروی می کنند مثلاً نام های افراد و یا نواحی، ممکن است از وندهای خاصی نظیر آنچه که در «یعقوب + زاده»، «پور + حسن»، «پونه سادات مرتضوی» یا «علی آباد» مشاهده می شود پیروی کنند. علاوه بر موارد ذکر شده، یک اسم خاص هم می تواند خود، ترکیبی از چند اسم خاص دیگر باشد.

این سیستم برای تشخیص اسامی خاص از روی ساختار آنها، از پایگاه داده ای از وندها بهره می برد. این پایگاه داده علاوه

×× در الگوی دوم جایگاه یک موجودیت چند کلمه ایست. مثلاً «دولت ×» برای نام کشور و «نهاد ××» برای نام یک نهاد بکار می‌رود. در اینصورت عباراتی مثل «...دولت ایران..» و یا «...نهاد نمایندگی مقام معظم رهبری در دانشگاه ها...» با این الگوها مطابقت خواهند کرد.

مجموعاً در این سیستم، جستجو برای تشخیص موجودیت نامدار با کمک الگوها به این صورت انجام می‌گیرد که برای تک- تک عناصر برچسب گذاری شده در متن (با در نظر داشتن ضمایر چسبیده به کلمات و نیز ادات جمع)، هم برای خود کلمه و هم برای برچسب آن بررسی می‌شود که آیا در هیچ یک از الگوها صدق می‌کنند یا نه و به این ترتیب مجموعه‌ای از اسامی خاص و نوع آنها (که عنوان مقوله هر الگو است) مشخص می‌شود.

۵-۲-۲- تأثیر کلمات جمع در اسامی خاص

گاهی مجموعه‌ای از چندین موجودیت نامدار متعلق به یک مقوله بدنبال یکدیگر ذکر می‌شوند. مثلاً داریم «شرکت‌های ایزایران و سافت نت» با همکاری هم نرم افزار جدید ایجاد نمودند. در اینجا هر یک از مجموعه کلماتی که زیرشان خط کشیده شده است به عنوان یک اسم خاص برای شرکت بکار می‌رود.

در دو حالت مجموعه‌ای از موجودیت‌های خاص بدنبال یکدیگر را سیستم در یک دسته قرار می‌دهد. حالت اول زمانی است که کلمه‌ای که به عنوان شناساینده موجودیت، در همسایگی آن آمده است به صورت جمع آورده شده باشد و حالت دیگر وجود صفت شمارشی پیش از آن، مثلاً «دو شرکت ایزایران و سافت نت». در حالت دوم سیستم بدنبال اعداد شمارشی استفاده نشده در الگوها پیش از شناساینده‌ها می‌گردد تا آنها را به صورت جمع دیده و بجای یک اسم نامدار همچون شناساینده‌های جمع، بدنبال مجموعه‌ای از اسامی نامدار بگردد.

۵-۲-۳- نا کافی بودن الگو

این گونه استخراج موجودیت نامدار سه مشکل دارد، اول اینکه چه تعداد از دنباله کلمات، موجودیت نامدار درست را تشکیل می‌دهند یا به عبارتی مرز موجودیت نامدار کجاست. مثلاً برای نام شرکت که در الگو به صورت * مشخص شده است چه دنباله‌ای نام شرکت است.

مشکل دوم این است که آیا کلمه تشخیص داده شده همان موجودیت نامدار است یا موجودیت ناخواسته^{۱۴} می‌باشد. مثلاً

برای «تاریخ هجری قمری» با استفاده از الگو «تاریخ ×»، هجری را به عنوان تاریخ در نظر می‌گیرد در حالی که تاریخ فرمت خاص خود را دارد که با برچسب محاسبه شده dateNo یا ترکیبی از برچسب‌های مشخص حاصل می‌شود.

مشکل سوم، ناشناس ماندن برخی از موجودیت‌هایی است که در هیچ الگویی صدق نکرده‌اند. در ادامه روش‌هایی برای ابهام زدایی و رفع این مشکلات احتمالی آورده شده است.

۵-۲-۴- ابهام زدایی

۵-۲-۴-۱- ابهام زدایی با توجه به تکرار رخداد

برای اولین مشکل مطرح شده یعنی تعیین مرز اسم نامدار، کار ابتدایی که انجام شد بررسی بلندترین دنباله اسمی تکرار شده در متن است. چرا که اسامی خاص به احتمال زیاد جای دیگری از متن باز تکرار می‌شوند. در این صورت، اگر هر دو متعلق به یک مقوله بودند، احتمال صحت افزایش می‌یابد. از آنجایی که یک موجودیت ممکن است در جای دیگر به صورت خلاصه شده‌ای تکرار شود پس در صورت ترکیبی بودن، ترکیب- های مختلف آن با حفظ ترتیب، برای تکرار آزمایش می‌شود.

۵-۲-۴-۲- ابهام زدایی با ترکیب کلمات

برای رفع دومین مشکل مطرح شده که در بهبود مشکل اول هم دخیل است، با بهره‌وری از پایگاه داده‌ای از ترکیب‌های برچسبی و کلمه‌ای ممکن برای هر مقوله به عنوان مکملی برای الگوها، مشکل تقلیل می‌یابد. به عبارتی هر کلمه پیش از اینکه جزوی از یک موجودیت نامدار محسوب شود امکان موجودیت بودن آن بررسی می‌شود. پیش از این در مورد پایگاه داده طراحی شده برای ترکیب کلمات توضیح داده شد. در جایی که سیستم، بدنبال موجودیت نامدار مناسب برای جایگزینی در الگو می‌گردد، آن موجودیتی را می‌پذیرد که خود کلمه و یا برچسب نحوی که الگو برای آن مطرح است، برای موجودیت صدق کند. مثلاً برای «مؤسسه آموزشی دانشجو و ایران» با توجه به الگوی «+مؤسسه +آموزشی ××» برای مقوله سازمان-آموزشی تک کلمات پس از آموزشی را تا جایی که خود کلمه و یا برچسب آن برای مقوله آموزشی صدق می‌کنند می‌پذیرد.

۵-۳- پالایش و تولید خروجی نهایی

پس از کشف موجودیت‌ها و اعمال برچسب معنایی به آنها، سیستم وارد مرحله پردازش نهایی می‌شود. در این مرحله با توجه به ساختار سلسله مراتبی سیستم، موجودیت‌های نامدار تکرار شده در یک مقوله و یا در مقوله‌های سطح پایین‌تر (گره

شده است و قابل استفاده به عنوان مکملی برای کاربردهایی از قبیل برچسب‌گذاری‌های نحوی، استخراج اطلاعات و درک مطلب و خلاصه سازی ماشین می‌باشد.

نتایج حاصله از آزمایش سیستم بجز شرایط بروز ابهام در متن خوب گذارش شده است. برای رفع مشکل مرز موجودیت نامدار که از روش‌های رفع ابهام مطرح شده در متن رفع نمی‌شود، سیستم به عنوان مکمل، نیاز به قطعه‌ساز^{۱۷} دارد که اجزای جمله را از یکدیگر تفکیک نماید. علاوه بر آن استفاده از الگوریتم‌های یادگیری ماشین که از متون حاشیه نویسی استفاده می‌کنند می‌تواند در افزایش دقت سیستم تأثیر چشم‌گیری داشته باشد.

مراجع

- [1] موسوی، نفیسه سادات؛ ثانی، غلامرضا قاسم؛ "بکارگیری دسته بندی کننده و رتبه بندی کننده آنروپی بیشینه در فرایند تعیین"، کنفرانس ملی انجمن کامپیوتر ایران، دوره ۱۴، تهران، ۱۳۸۷.
- [2] Grishman, R., Sundheim, B. *Message Understanding Conference - 6: A brief history*, COLINGS, 1996.
- [3] Zhang, Li, Pan, Yue, Zhang, Tong. *Focused Named Entity Recognition using Machine Learning*. ACM Press, pp. 281-288, 2004.
- [4] Cohen, William W., Sarawagi, Sunita. *Exploiting Dictionaries in Named Entity Extraction: Combining Semi-Markov Extraction Processes and Data Integration Methods*. Conference on Knowledge Discovery in Data, pp. 89-98, 2004.
- [5] Nadeau, David, Sekine, Satoshi. "A survey of named entity recognition and classification.", National Research Council Canada/New York University, 2006.
- [6] Negri, Matteo, Magnini, Bernardo. *Using WordNet Predicates for Multilingual Named Entity Recognition*, The Second Global Wordnet Conference, pp. 169-174, 2004.
- [7] mikheev, Andrei, Moens, Marc, Grover, Claire. *Named Entity Recognition without Gazetteers*, Proceedings of EACL, Vol. 19, Bergen, Norway, pp. 1-8, 1999.
- [8] Cucerzan, Silviu, Yarowsky, David. *Language Independent Named Entity Recognition Combining Morphological and Contextual Evidence*, Proceedings of EMNLP, pp. 90-99, 1999.
- [9] Resnick, Ed., P. *Internet Message Format (rfc5322)*, Network Working Group, 2008.
- [10] Berners-Lee, T., Masinter, L., McCahill, M. *RFC1738 - Uniform Resource Locators (URL)*. Network Working Group, December 1994, <http://www.faqs.org/rfcs/rfc1738.html>
- [11] Uniform Resource Identifier (URI) Schemes, IANA, ICANN, ptc/10-11-2009, <http://www.iana.org/assignments/uri-schemes.html>.
- [12] Shamsfard, Mehrnoush, Fadaee, Hakimeh. *A Hybrid Morphology-Based POS Tagger for Persian*, European Language Resources Association (ELRA), Marrakech, Morocco, 2008.

فرزند) کشف شده و حذف می‌شوند. همچنین موجودیت‌هایی که درون موجودیتی دیگر قرار گرفته‌اند با توجه به مکان آنها در متن شناسایی و حذف می‌گردند. مثلاً «انجمن شرکت‌های بازرگانان نمونه»، ممکن است «بازرگانان نمونه» را از مقوله شرکت ببیند در حالی که «شرکت‌های بازرگانان نمونه» از مقوله انجمن، موجودیت خاص مورد نظر است. پس سیستم با مشاهده چنین شرایطی موجودیت نامدار تکرار شده در موجودیت نامدار دیگر را حذف می‌کند. نهایتاً مجموعه موجودیت‌های شناسایی شده در مقوله مناسب‌شان دسته بندی شده و خروجی با برچسب‌های معنایی تولید می‌شود.

۶- آزمون و ارزیابی

سیستم به ازای متون مختلف و در بر دارنده تقریباً تمامی برچسب‌های معنایی ممکن آزمایش شد. پس از مقایسه‌ی نتایج حاصله از برچسب‌گذاری سیستم بر داده‌های تستی با برچسب‌های معنایی تولید شده به صورت دستی، دقت^{۱۵} سیستم ۷۲٪ و یادآوری^{۱۶} آن ۷۶٪ گزارش شد.

علت عمده عدم تشخیص، به تشخیص ندادن مرز درست برخی از موجودیت‌های نامدار در متن باز می‌گشت. به عنوان نمونه، رفتار سیستم در مقابل عبارت «شرکت فرارسانه اندیشه ساز» تشخیص «فرارسانه» به عنوان موجودیت نامداری از «شرکت» بود. علت رخداد این خطا، عدم تکرار عبارت در جای دیگری از متن و نیز نامرتب بودن برچسب‌های نحوی واژگان تشکیل دهنده عبارت بود. سایر رخدادهای عدم تشخیص موجودیت نامدار، بدلیل صدق نکردن آن در هیچ گونه از الگوها می‌باشد.

اشتباه در تشخیص موجودیت‌های نامدار هم اکثراً یا به میزان تشابه‌شان با ظاهر موجودیت‌های نامدار بازمی‌گشت که ریخت شناسی منجر به این خطا می‌شد و یا به استفاده از موجودیت در جایی که انتظار مشاهده موجودیت نامدار می‌رفت. به عبارتی سیستم با ابهام در تشخیص مواجه می‌شد. گاهی هم خطا در برچسب زن نحوی، منجر به انتشار خطا در نتیجه نهایی سیستم می‌شد.

۷- نتیجه گیری

در اینجا سیستمی برای تشخیص و دسته بندی موجودیت‌های نامدار در زبان فارسی معرفی شد که برای رسیدن به مقصود خود از برچسب‌گذاری نحوی و ویژگی‌های درونی و متنی واژگان بهره می‌برد. این سیستم به طور کامل پیاده سازی و آزمایش

Named Entity Recognition(NER) ^۱

Supervised Learning ^۲

Semi-supervised Learning ^۳

Unsupervised Learning ^۴

^۵ پایگاه داده لغوی انگلیسی است که در آن اسم‌ها، فعل‌ها، صفات و قیود به مجموعه‌ای از مترادف‌های ادراکی (synsets) گروه‌بندی می‌شوند که هر کدام بیانگر مفهوم مجزایی‌اند.

Tokenization ^۶

Local-Part ^۷

Internet Domain ^۸

Scheme ^۹

Internet Assigned Numbers Authority ^{۱۰}

POS Tagger ^{۱۱}

Tagset ^{۱۲}

Morphology ^{۱۳}

noise ^{۱۴}

Precision ^{۱۵}

Recall ^{۱۶}

Chunker ^{۱۷}