

رفع ابهام معنایی واژگان مبهم فارسی

با مدل موضوعی LDA

بابک مسعودی^۱ و سعید راحتی قوچانی^۲

^۱گروه مهندسی کامپیوتر و فناوری اطلاعات، دانشگاه پیام نور زابل، سیستان و بلوچستان، ایران

^۲گروه برق، دانشگاه آزاد اسلامی واحد مشهد، مشهد، ایران

چکیده

ابهامزدایی از واژگان مبهم و دارای معنای متعدد، موضوع مهمی در حوزه پردازش زبان‌های طبیعی است. در این مقاله، مدلی برای رفع ابهام از واژگان مبهم فارسی با استخراج ویژگی‌های جدید پیشنهاد شده است. برای ایجاد این مدل دو دسته ویژگی واژگان و نشانه‌های همراه واژه مبهم و ویژگی‌هایی که با به‌کاربردن روش‌های مدل‌سازی موضوع به‌دست می‌آید، استفاده شده است. یک مدل موضوعی، مدلی آماری برای استخراج چکیده موضوع‌های موجود در اسناد یک پیکره است. در مقاله حاضر ما از روش بدون سربرستی تخصیص پنهان دریکله (LDA) برای این منظور استفاده کردیم. نتایج آزمایش‌ها برای پانزده واژه مبهم پردازشگر در زبان فارسی که از پیکره پژوهشکده پردازش شومند علائم استخراج شد، دقّت حدود ۹۷٪ را نشان می‌دهد که بیان‌گر تأثیر این روش در یافتن معنی مناسب واژگان مبهم است.

واژگان کلیدی: تخصیص پنهان دریکله، چندمعنایی، رفع ابهام معنایی، مدل‌سازی موضوع.

موضوع پژوهشی می‌تواند راه‌گشای این مباحث در زبان فارسی باشد.

۲- مروری بر پژوهش‌های پیشین

روش‌های بسیاری برای رفع ابهام از واژگان مبهم در زبان‌های مختلف پیشنهاد شده است؛ همچنین به‌منظور ایجاد انگیزه در بین پژوهش‌گران، جهت پژوهش در این زمینه از سال ۱۹۸۸، هر سه سال رقابت‌هایی تحت عنوان SemEval^۳ برگزار می‌شود.

در (یاروسکی، ۲۰۰۰)، الگوریتمی با سربرستی بر پایه فهرست‌های تصمیم‌گیری^۴ پیشنهاد شده است. در این الگوریتم مجموعه کاملی از ویژگی‌های همسایگی، ساخت واژی و نحوی استفاده و دقّت سامانه ۷۸/۱٪ گزارش شده است. (ونسترا و همکاران، ۲۰۰۰) روشی مبتنی بر حافظه برای رفع ابهام معنایی لغات پیشنهاد می‌کند که دقّت بدست آمده از این سامانه ۷۵/۱٪ گزارش شده است. در

۱- مقدمه

در اکثر زبان‌های طبیعی، وجود برخی کلمات باعث ایجاد ابهام در متون آن زبان می‌شود. به این کلمات مبهم که دارای ساختار نوشتاری یکسان و معانی متفاوت هستند، همنگاره (همنویسه)^۱ گفته می‌شود. ابهامزدایی از کلمات مبهم^۲ موضوع مهمی در حوزه پردازش زبان‌های طبیعی است و پایه بسیاری از دیگر مباحث مطرح در این حوزه نیز هست. تاکنون پژوهش‌های گسترده‌ای با رویکردهای متعدد در زبان‌های دیگر انجام و نتایج چشم‌گیری در این زبان‌ها حاصل شده است. این موضوع سنگبنای بسیاری از مباحث در حوزه پردازش زبان طبیعی، همچون ترجمة ماشینی، خطایاب، تبدیل متن به گفتار، تشخیص گفتار و درک متن است. بهعنوان مثال مترجم زبان فارسی به انگلیسی برای ترجمه واژه «شیر» با سه معنی "Lion"، "Milk" و "Valve" مواجه است که باید با توجه به بافت متن، مناسب‌ترین واژه را انتخاب کند. فعالیت بر روی این

³ <http://aclweb.org/>

⁴ Decision lists

^۱ Homograph

^۲ Word-sense disambiguation

موضوع، مدل‌های آماری از متون هستند و فرض می‌کنند موضوع‌های پنهان درون پیکره نهفته شده‌اند، که می‌توانند برای استخراج، دسته‌بندی متون و بخش‌بندی گفتگوها مورد استفاده قرار گیرند. از آن جا که رفع ابهام از واژگان بهم مبنوعی یک مسئله دسته‌بندی مجزا است، می‌توان با ایجاد مدلی برای هریک از واژه‌های مبهم، به کمک روش دسته‌بندی بیشینه‌بی‌نظمی^۴ از آن‌ها رفع ابهام کرد.

ادامه مقاله به این صورت بخش‌بندی شده است: بخش سوم اصول روش پیشنهادی را با توصیف مدل‌های موضوعی، مدل‌های زایا، تخصیص پنهان دریکله و دسته‌بندی بیشینه‌بی‌نظمی بیان می‌کند. مدل پیشنهادی در بخش چهارم ذکر شده و آزمایش‌ها، نتیجه‌گیری و کارهای آینده در بخش‌های پنجم و ششم بیان شده‌اند.

۳- اصول روش پیشنهادی

۱-۳- مدل‌سازی موضوع

یک مدل موضوعی، مدلی است که موضوع‌های درون پیکره‌ای از اسناد را استخراج می‌کند و به هر سند برخی از این موضوع‌ها را نسبت می‌دهد. در واقع می‌توان گفت یک مدل موضوعی، جعبه‌ای سیاه با دو خروجی است: انتساب واژگان به موضوع‌ها و انتساب چندین موضوع به اسناد. نخستین خروجی، توزیع موضوع بر روی واژگان است. ما به‌طور عمومی موضوع‌ها را به صورت فهرست‌هایی از واژه‌ها نشان می‌دهیم. شکل(۱) سه موضوع استخراج شده از پیکره New York Times را نشان می‌دهد. هر واژه با درصدی احتمال، در یک موضوع حضور دارد. در اینجا تعدادی واژه با احتمال بیشتر نشان داده شده است. خروجی دیگر یک مدل موضوعی، امکان انتساب یک سند به چند موضوع مرتبط است.

شکل(۱) نشان می‌دهد چگونه یک سند می‌تواند با یک یا چند موضوع شناخته شود.

روش‌های متنوعی برای یافتن موضوع‌ها و انتساب آن‌ها به اسناد پیشنهاد شده‌اند. در این مقاله ما بر روی تخصیص پنهان دریکله (LDA) تمرکز کرده‌ایم. LDA یک مدل احتمالاتی و زایا است. یک مدل احتمالاتی است، زیرا با زبان احتمال بیان شده و زایا است؛ زیرا بیان می‌کند که داده‌ها چگونه به وجود آمده‌اند. در این مدل همه جزئیات شناخته شده نیستند؛ این قطعات گمشده را متغیرهای

⁴Maximum Entropy

(فلورین و همکاران، ۲۰۰۲) روشی با ترکیب چند دسته‌بندی کننده نامتجانس (بردار شباهت کسینوسی، مدل‌های بیزی^۱ و فهرست‌های تصمیم‌گیری) پیشنهاد شده و مجموعه‌ای از ویژگی‌های استخراج شده از بافت متن نظیر روابط دستوری و برجسب‌های POS را برای مشخص کردن معنای مناسب واژه، استفاده کرده است. سامانه GAMBEL (دکادت و همکاران، ۲۰۰۴) بر پایه بادگیری با سرپرستی مبتنی بر حافظه پیشنهاد شد که از ویژگی‌های بافت متن و منابع دانشی نظیر WordNet استفاده می‌کند. در NUS-MI SemEval-2007 سامانه NUS-MI معرفی شد که از تخصیص پنهان دریکله^۲ (LDA) - یک مدل احتمالاتی بیزی سلسه‌مراتبی سه‌لایه و ویژگی‌های نحوی و موضوعی استفاده می‌کرد. (استیونسون و همکاران، ۲۰۱۲) نشان دادند کارایی یک الگوریتم بدون سرپرستی WSD با افزودن اطلاعات بافت متن بهبود می‌یابد. در (پریسیس و همکاران، ۲۰۱۳) به ویژگی‌های استخراج شده توسط سه سامانه مرسوم رفع ابهام معنایی بدون سرپرستی و یک سامانه با سرپرستی، ویژگی موضوع بافت استخراج شده توسط افزوده شد، که نتایج به دست آمده نشان داد کارایی سامانه‌ها به طور قابل ملاحظه‌ای افزایش می‌یابد.

از فعالیت‌هایی که در زمینه رفع ابهام واژگان فارسی صورت گرفته، به (مکی، ۲۰۰۸) می‌توان اشاره کرد. در این مقاله روشی مبتنی بر پیکره و یک واژگان برای امتیازدهی به دسته تعلاق مفهومی هر معنی واژه مبهم، پیشنهاد شده است. دقّت میانگین این روش برای پانزده واژه مبهم ۹۱/۴۶٪ گزارش شده است. همچنین یک روش مبتنی بر قاعده در (سعیدی، ۲۰۰۹) پیشنهاد شده که با درنظر گرفتن این قواعد دقّت آن برای یکی از پیکره‌های فارسی ۹۱٪ است.

در مقاله حاضر روشی مبتنی بر پیکره و با سرپرستی جهت ابهام‌زدایی از واژگان مبهم فارسی پیشنهاد شده است. بر این اساس از یک پیکره فارسی برای هر واژه مبهم، دو دسته ویژگی استخراج می‌شود. دسته نخست ویژگی‌ها، واژگان و نشانه‌های همنشین واژه هدف و دسته دوم موضوع بافت متن مورد نظر است. برای استخراج موضوع بافت در این مقاله از روش تخصیص پنهان دریکله (LDA) (بلی، ۲۰۰۳) که یکی از روش‌های بدون سرپرستی مدل‌سازی موضوع^۳ است، استفاده می‌شود. روش‌های مدل‌سازی

¹ Bayes model

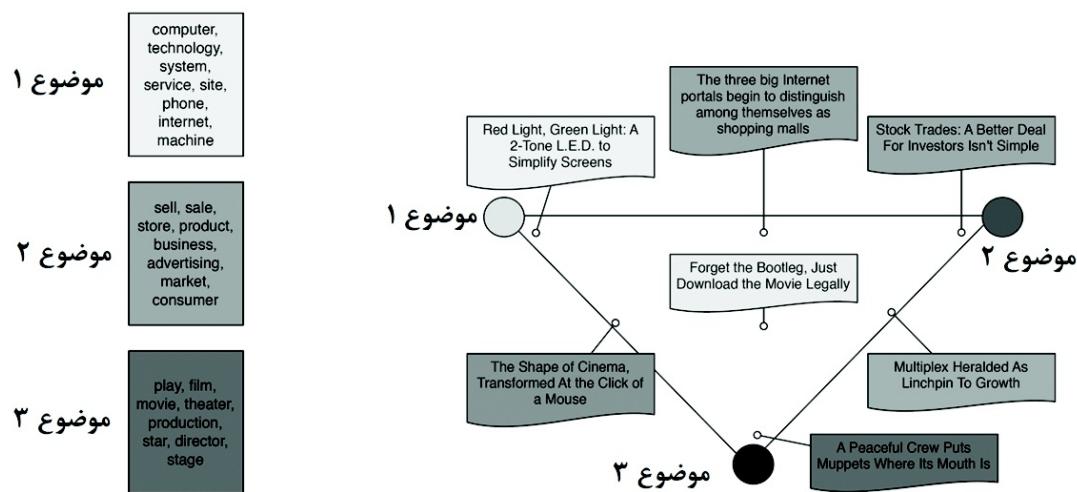
² Latent Dirichlet Allocation

³ Topic Model



فراهام آورند. این مدل از مجموعه فن‌های بدون سرپرستی برای یافتن موضوع‌ها از داده‌های خام محسوب می‌شود.

پنهان می‌نامیم و به روش‌های استنتاج آماری جستجو می‌شوند تا بهترین توصیف را برای داده‌های قابل مشاهده



(شکل-۱): مثالی از توزیع کلمات بر روی موضوع‌ها و توزیع موضوع‌ها بر روی اسناد در پیکره (بودید-گرابر ۲۰۱۰ New York Times)

اطلاعات وابسته به مدل، دفعات تولید واژه است که به نام سبد واژگان^۳ شناخته می‌شود و در بسیاری از مدل‌های آماری زبان، رایج است. البته اطلاعات ترتیب واژه ممکن است، حاوی اشاراتی به محتوای یک سند باشند که این اطلاعات توسط مدل استفاده نمی‌شوند گریفیتس و همکاران، ۲۰۰۵).

سمت راست شکل (۲) مسئله استنتاج آماری را نشان می‌دهد؛ به طوری که با واژگان مشاهده شده در یک مجموعه از اسناد، تمایل داریم بدانیم کدام مدل موضوعی به توزیع موضوع‌های هر سند، شبیه‌تر است.

۳-۳- تخصیص پنهان دریکله (LDA)

مدل LDA فرایند زایایی زیر را برای ایجاد پیکرهای از M سند با N_d واژه در سند d با k موضوع $\{\beta_1 \dots \beta_k\}$ فرض می‌کند:

- For each document $d \in \{1 \dots\}$
- Choose the document's topic weights
 $\theta_d \sim Dir(\alpha)$
 - For each word $n \in \{1 \dots n_d\}$
 - Choose topic assignment $z_{d,n} \sim Mult(\theta_d)$

²Bag-of-Words

۲-۳- مدل‌های زایا^۱

یک مدل زایا برای اسناد، براساس قوانین نمونه‌برداری ساده احتمالاتی توصیف می‌شود. در مدل زایا، هدف یافتن بهترین مجموعه از متغیرهای پنهانی است که می‌توانند داده‌های مشاهده شده (واژگان موجود در اسناد) را بهترین وجه توصیف کنند. در شکل (۲) مدل موضوعی با رویکرد مجزای مدل زایا و استنتاج آماری نشان داده می‌شود.

در سمت چپ، مدل زایا با دو موضوع مختلف نشان داده شده است. موضوع‌های ۱ و ۲ وابسته به «شیر معادل Lion» و «شیر معادل Milk» هستند و به عنوان سبد‌های حاوی توزیع‌های مختلف روی واژگان در نظر گرفته می‌شوند. اسناد مختلف را می‌توان با برداشتن واژگان از یک موضوع، بسته به وزن داده شده، انتخاب کرد. به عنوان مثال اسناد یک و سه به ترتیب با نمونه‌برداری از موضوع‌های یک و دو ایجاد می‌شوند؛ در حالی که سند دو با ترکیبی از هر دو موضوع ساخته می‌شود. توجه کنید، اعداد بالاترین نشان می‌دهد از کدام موضوع برای واژه، نمونه‌برداری شده است.

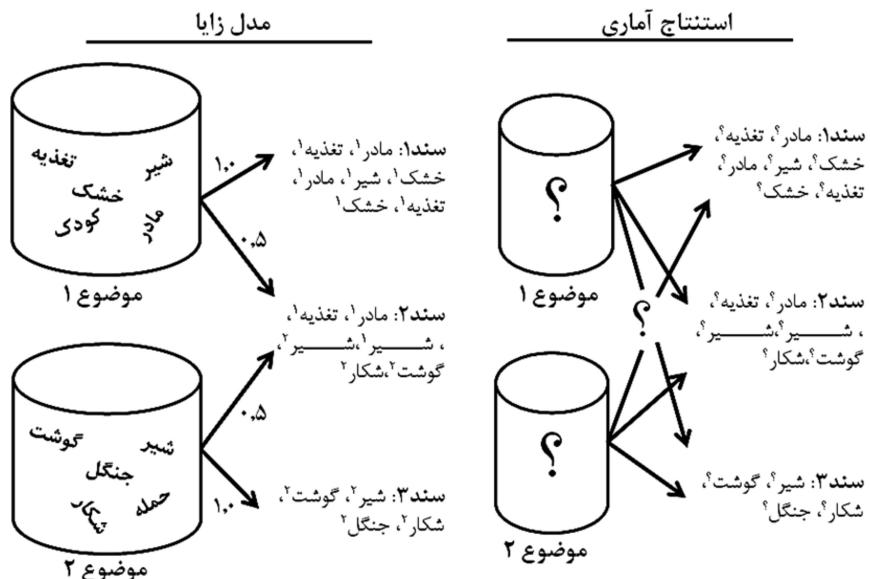
مدل زایا که در اینجا توصیف شد، هیچ فرضی در مورد ترتیب واژگانی که در سند ظاهر شده‌اند، ندارد. تنها

¹Generative Models

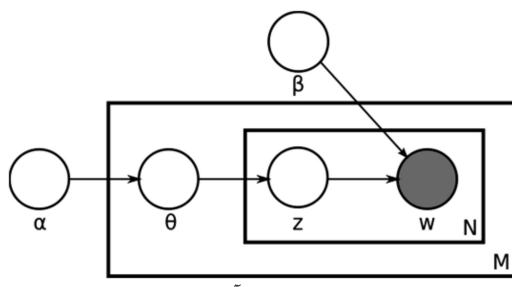
ii) Choose word $w_{d,n} \sim \text{Mult}(\beta_{z_{d,n}})$

در این فرآیند توزیع دریکله و توزیع چندجمله‌ای به صورت

زیر تعریف می‌شود:



(شکل - ۲): نمایشی از مدل زایا و مسئله استنتاج آماری در تفسیر مدل موضوعی



(شکل - ۳): نمودار فرآیند زایایی LDA

متغیرهای پنهان، داده‌های قابل مشاهده و پارامترها هر کدام با یک گره در این نمودار نشان داده می‌شوند. هر وابستگی آماری توسط خط بین گره‌ها مشخص می‌شود، داده‌های قابل مشاهده با گره سایه‌دار و تکرار داده با زمینه مستطیل نشان داده می‌شود، پایین‌نویس سمت راست پایین نشان‌دهنده تعداد تکرار متغیرهای درون مستطیل است. از آن‌جا که فقط واژگان جمع‌آوری شده از اسناد قابل مشاهده است، یافتن موضوعات $\beta_{1:k}$ ، انتساب موضوع $z_{1:d}$ و توزیع موضوع $\theta_{1:D}$ برای هر سند، یک مسئله استنتاج آماری است.

$$\text{Dir}(\theta|\alpha_1, \dots, \alpha_k) = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_k \theta_k^{\alpha_k - 1} \quad (1)$$

$$\text{Mult}(n|\Theta_1, \dots, \Theta_k) = \frac{(\sum_k n_k)!}{\prod_k n_k!} \prod_k \Theta_k^{n_k} \quad (2)$$

α و β نیز پارامتر هستند.

توزیع (θ_d) برداری به طول k (تعداد موضوعات در

مدل) است و برای توصیف موضوع منتبشده به هر Z_n

یک از n واژه از سند d استفاده می‌شود.

امین موضوع β_k برداری به طول V است، که هر

بخش آن متناظر با احتمال حضور واژه در آن موضوع است.

در شکل ۱ (ب)، β_{zn} یکی از لغات ذخیره‌شده در فهرست‌ها

و θ_d موقعیت قرارگرفتن اسناد را نشان می‌دهد.

ابزار دیگری که برای توصیف این مدل به کار برده

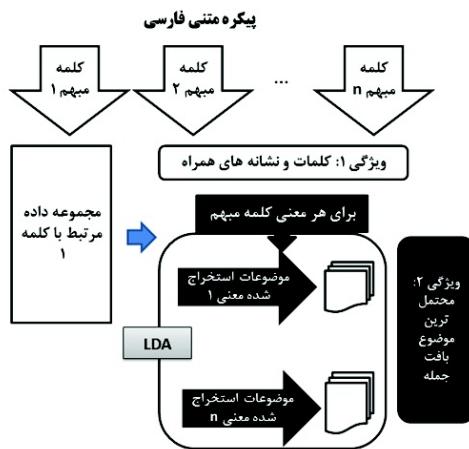
می‌شود، نمودار شکل (۳) است. این نمودار فرایند LDA را

نشان می‌دهد.

فصل نهم



شکل (۴)، نمودار استخراج ویژگی برای ایجاد مدل را نشان می‌دهد.



شکل - (۳): نمودار استخراج ویژگی ها برای ایجاد مدل پیشنهادی

استخراج دو دسته ویژگی به شرح زیر است:

دسته نخست - واژگان و نشانه های همراه: این ویژگی ها با درنظر گرفتن خصوصیات زبان فارسی انتخاب شده است. ابهام بسیاری از واژگان زبان فارسی، ناشی از نبود اعراب در متون نوشته ای است. از این جهت واژگان و نشانه گذاری هایی که بلافاصله قبل و بعد از واژه هدف قرار می گیرند، نقش مؤثری در رفع ابهام از این واژگان دارند. در مثال زیر با دقت در نشانه گذاری قبل و بعد از واژه «شیر» معنی صحیح آن را می توان مشخص کرد:

شیر آب
شیر، آب

برای این کار با پنجره ای به مرکز واژه مورد نظر و عرض ± 2 کلیه واژگان و نشانه گذاری های قبل و بعد از آن را به عنوان دسته ویژگی نخست استخراج می کنیم. بدیهی است در صورت قرار گرفتن واژه در ابتدا یا انتهای جملات، مقادیر این ویژگی فضای خالی خواهد بود، که در نظر گرفته نمی شوند.

دسته دوم - مدل سازی موضوع: جهت یافتن ویژگی مدل سازی موضوع در مرحله آموزش مدل، ابتدا سبد واژگان همراه هر واژه م بهم استخراج می شود. برای این منظور یک پنجره به مرکز واژه مورد نظر و عرض ± 10 در نظر می گیریم. کلیه واژگان به استثنای حروف ریط، اضافه و اعداد را به عنوان واژگانی که بافت جمله مورد نظر را نشان می دهند، استخراج می کنیم؛ سپس به کمک روش LDA استخراج شده

۴-۳- دسته بندی بیشینه آنتروپی

روش دسته بندی آماری استفاده شده در این مقاله دسته بندی بیشینه آنتروپی است. بیشینه آنتروپی فنی عمومی برای تخمین توزیع احتمال، از یک مجموعه داده است. اگر هیچ دانش نخستینی از داده ها در دسترس نباشد، آنتروپی بیشینه، یک نواخت ترین توزیع را انتخاب می کند؛ به طوری که همه وقایع دارای احتمال برابر باشند.

در صورتی که مجموعه داده در دسترس باشد، داده های دارای برچسب مجموعه آموزش، به عنوان تعدادی ویژگی که محدودیت های مدل از آن ها استخراج می شود، در نظر گرفته می شوند. بنابراین در حالی که توزیع باید به گونه ای باشد تا آنتروپی را بیشینه کند، مدل نیز باید محدودیت های اعمال شده توسط برچسب های داده آموزشی را برآورده کند. بنابراین یک مدل بیشینه آنتروپی، مدلی برای اراضی تمامی محدودیت های اعمال شده روی مجموعه داده ها است. این مدل براساس رابطه زیر بنا می شود:

$$p(c|x) = \frac{1}{Z} \exp\left(\sum_i \lambda_i f_i(x, c)\right) \quad (3)$$

که در آن تابع ویژگی $f_i(x, c)$ تعداد دفعاتی است که ویژگی i استفاده شده تا برچسب c را برای رویداد x پیدا کند و وزن های λ_i برای بیشینه کردن درست نمایی در فرایند آموزش انتخاب می شوند، تا مقدار آنتروپی p بیشینه شود. مهم ترین برتری مدل بیشینه آنتروپی آن است که این مدل چارچوبی برای یک پارچه کردن داده های ناهمگون به دست آمده از منابع مختلف اطلاعاتی، فراهم می کند (ویکام، ۲۰۰۲).

مدل احتمالاتی بیشینه آنتروپی در بسیاری از کاربردهای حوزه پردازش زبان های طبیعی نظیر برچسب گذاری اجزا کلام^۱ (POS) و تعیین مرز جملات با موفقیت استفاده شده است (رتناپارکی، ۱۹۹۸).

۴- مدل پیشنهادی

برای ایجاد مدلی مؤثر جهت رفع ابهام واژگان همنگاره فارسی از دو دسته ویژگی استفاده می شود. دسته نخست، واژگان و نشانه های همراه واژه هدف و دسته دوم، موضوع بافت جمله است که به کمک روش LDA استخراج شده است.

^۱Part of speech

(جدول-۱): جزئیات واژگان مبهم مورد استفاده در آزمایش‌ها

واژه	معانی	تعداد جملات	تعداد
شیر	Lion	۱۵۷	۶۸۰
	Milk	۴۴۸	
	Valve	۱۵۷	
سیر	Garlic	۱۰۷۲	۱۴۹۸
	Journey	۴۲۶	
مهر	Punch	۲۰۹	۷۲۹
	Love	۱۲۰	
	A Solar Month	۴۰۰	
جو	Atmosphere	۳۰۷	۴۱۷
	Barley	۱۱۰	
شکر	Sugar	۲۶۰	۳۷۱
	Thankfulness	۱۱۱	
حلال	Lawful	۸۲	۱۰۹
	Resolvent	۲۷	
خود	Wisdom	۲۰۹	۴۰۷
	Tiny	۱۹۸	
نفس	Breath	۲۰۴	۷۶۵
	Oneself	۵۶۱	
محرم	ALunar Month	۲۳۱	۲۵۵
	Intimate	۲۴	
تن	Ton	۲۸۹۲	۴۳۴۹
	Body	۱۴۵۷	
اشکال	Figures	۲۶۴	۵۳۹
	Defect	۲۷۵	
اشراف	Dminance	۵۱	۱۰۷
	Knight,Gentleman	۵۶	
قسم	Oath	۷۱	۱۳۰
	Sort	۵۹	
سرم	My head	۹۳	۱۳۱
	Serum	۳۸	
شرف	Honor	۸۹	۱۳۶
	Soon Expected	۴۷	

۲-۵- نتایج آزمایش‌های مدل پیشنهادی

به منظور ارزیابی مدل پیشنهادی، برای هر کلمه مبهم سبد لغات همنشین آن از مجموعه آموزشی استخراج می‌شود. آزمایش‌ها برای رفع ابهام هر کلمه چندین مرتبه و هر بار با درنظر گرفتن مقداری متفاوت برای تعداد موضوع‌های موجود در هر معنی آن کلمه صورت گرفته و احتمال حضور کلمات همنشین موجود در سبد لغات توسط LDA ، محاسبه می‌شود. جداول ۴-۲ ده واژه (انتخاب شده از کل سبد

پنهان دریکله) که در جعبه ابزار 'Matlab Toolbox' مبادله‌سازی شده است برای هر یک از معانی واژه‌های مبهم تعدادی موضوع در نظر گرفته و احتمال وقوع هر یک از واژگان استخراج شده و موجود در سبد واژگان را در هر موضوع محاسبه می‌کنیم. جداول ایجادشده برای هر کلمه مبهم جهت رفع ابهام در مراحل بعدی نگهداری خواهد شد. در مرحله آزمایش با جمع جدایانه احتمال حضور واژگان همنشین کلمه مبهم در هر موضوع، هر یک که احتمال وقوع در آن بیشتر باشد، مقدار ویژگی آن موضوع را برابر یک و برای سایر موضوعات این مقدار برابر صفر در نظر گرفته می‌شود. به این ترتیب برای واژه‌ای مانند شیر با سه معنی، تعداد $3 \times$ ویژگی استخراج می‌شود؛ که T تعداد موضوعات درنظر گرفته شده برای هر معنی است. به عنوان مثال در یکی از آزمایش‌ها تعداد موضوعات را سه در نظر گرفته‌ایم؛ بنابراین برای واژه «شیر» تعداد نه ویژگی دوستایی درنظر گرفته می‌شود که تنها یکی از آن‌ها یک و بقیه صفر خواهد بود.

پس از استخراج ویژگی‌ها، جداساز بیشینه آنتروپی برای ایجاد سامانه رفع ابهام آموزش می‌بیند. در اینجا برای هر کلمه مدلی با استفاده از روش دسته‌بندی بیشینه آنتروپی و با استفاده از ابزار Mallet (مک‌کالم، ۲۰۰۲) که بر پایه جاوا تهیه شده است، ایجاد می‌شود.

۵- آزمایش‌ها

۱-۱- دادگان مورد استفاده

مجموعه داده مورد استفاده در این مقاله از پیکره متنی پژوهشکده پردازش هوشمند علائم^۱ جمع‌آوری شده است. این پیکره متنی، مجموعه‌ای از متنون نوشتاری و گفتاری زبان فارسی رسمی است که از منابع واقعی همچون روزنامه‌ها، پایگاه‌ها و ... جمع‌آوری شده، تصحیح شد و برچسب خورده است. حجم این دادگان در حدود یکصد میلیون واژه است و تنوع بسیار زیادی دارد.

برای ایجاد مجموعه‌های داده مورد نیاز این پیکره متنی جملاتی را که حاوی پانزده کلمه پر تکرار مبهم فارسی هستند، استخراج می‌کنیم.

جدول (۱) جزئیات کلمات مبهم مورد استفاده در آزمایش‌ها را نشان می‌دهد:

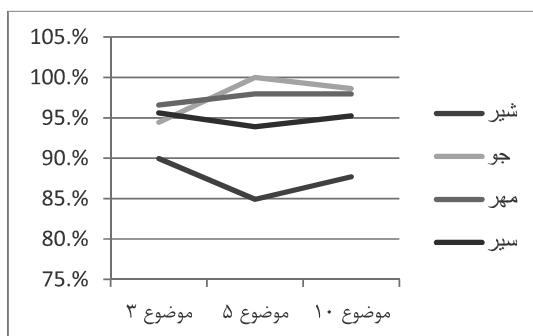
¹http://psiexp.ss.uci.edu/research/programs_data/toolbox.htm
²www.rcisp.com



(جدول - ۴): ده واژه با بالاترین احتمال استخراج شده توسط LDA برای معنی ۳ واژه «شیر»

موضوع ۷	احتمال	موضوع ۸	احتمال	موضوع ۹	احتمال
نتفیط	۰/۰۴۹۱	آب	۰/۰۴۹۱	فشار	۰/۰۸۳۸۳
فشار	۰/۰۲۳۳۰	شیلنج	۰/۰۲۳۳۰	باز	۰/۰۲۷۹۶
منلا	۰/۰۲۳۳۰	برقی	۰/۰۲۷۹۶	یک	۰/۰۴۲۸۹
تخلیه	۰/۰۱۸۶۵	گاز	۰/۰۲۷۹۶	استفاده	۰/۰۲۵۷۴
اصلی	۰/۰۱۸۶۵	کار	۰/۰۲۲۳۸	کترول	۰/۰۲۲۸۹
صفحه	۰/۰۱۴۰۰	هنگام	۰/۰۱۹۵۸	آب	۰/۰۲۰۰۳
رانش	۰/۰۱۴۰۰	سیلندر	۰/۰۱۹۵۸	طرف	۰/۰۱۷۱۷
دیگ	۰/۰۱۴۰۰	دو	۰/۰۱۹۵۸	دهانه	۰/۰۱۶۲۶
عمل	۰/۰۱۴۰۰	کردن	۰/۰۱۶۷۹	قرار	۰/۰۱۴۲۲
دما	۰/۰۱۴۰۰	سه	۰/۰۱۴۰۰	هوا	۰/۰۱۴۲۲
دکمه	۰/۰۱۴۰۰	تحریک	۰/۰۱۴۰۰	قطع	۰/۰۱۴۲۲
تنها	۰/۰۱۴۰۰	سمت	۰/۰۱۴۰۰	عنوان	۰/۰۱۴۲۲

در آزمایشی دیگر اثر تغییر در انتخاب تعداد موضوع برای هر معنی واژه بررسی شد. شکل زیر دقت سامانه برای ۳، ۵ و ۱۰ موضوع را نشان می‌دهد.



(شکل - ۴): مقایسه اثر تغییر تعداد موضوع در دقت دسته‌بندی

با توجه به نمودار بالا می‌توان دریافت تغییر در انتخاب موضوع برای معنی مختلف یک واژه، می‌تواند اثرات مختلفی داشته باشد؛ اما به طور کلی می‌توان گفت با زیاد شدن تعداد موضوع، واژه در موضوع‌هایی با جزئیات بیشتر قرار می‌گیرد و با کم شدن تعداد آن‌ها معنی واژه در دسته‌های کلی‌تر قرار داده می‌شوند. افزایش تعداد موضوع‌ها در بعضی واژگان ممکن است باعث افزایش دقیقت و در بعضی باعث کاهش دقیقت شود. علت این مسئله را می‌توان در نوع مجموعه داده مورد استفاده و یا به دسته مفهومی که واژه به آن تعلق دارد، یافت؛ که نیاز به پژوهش بیشتر دارد.

واژگان) با بالاترین احتمال وقوع را در تعداد سه موضوع برای هر معنی واژه «شیر» نشان می‌دهد.

(جدول - ۲): ده واژه با بالاترین احتمال استخراج شده توسط LDA برای معنی ۱ واژه «شیر»

موضوع ۱	احتمال	موضوع ۲	احتمال	موضوع ۳	احتمال
سر	۰/۰۴۱۶۰	رویاه	۰/۰۲۷۰۲	یک	۰/۰۴۷۰۰
آن	۰/۰۳۳۳۰	حق	۰/۰۲۱۶۳	شکار	۰/۰۴۱۱۳
نر	۰/۰۳۳۳۰	نه	۰/۰۲۱۶۳	گرگ	۰/۰۳۵۲۶
هوا	۰/۰۲۴۹۹	فیل	۰/۰۲۱۶۳	خدود	۰/۰۲۹۴۰
جسممه	۰/۰۲۴۹۹	مثل	۰/۰۲۱۶۳	روی	۰/۰۲۲۵۳
البته	۰/۰۱۶۶۹	حیوانات	۰/۰۲۱۶۳	فیلم	۰/۰۱۷۶۶
زبانی	۰/۰۱۶۶۹	حمله	۰/۰۱۶۲۳	مأمور	۰/۰۱۷۶۶
فرار	۰/۰۱۶۶۹	سنگی	۰/۰۱۶۲۳	فعل	۰/۰۱۷۶۶
پارک	۰/۰۱۶۶۹	خورشید	۰/۰۱۶۲۳	جایزه	۰/۰۱۷۶۶
پیدا	۰/۰۱۶۶۹	دین	۰/۰۱۶۲۳	ستگ	۰/۰۱۱۷۹

(جدول - ۳): ده واژه با بالاترین احتمال استخراج شده توسط LDA برای معنی ۲ واژه «شیر»

موضوع ۴	احتمال	موضوع ۵	احتمال	موضوع ۶	احتمال
خشک	۰/۰۱۴۳۲	مادر	۰/۰۴۲۶۸	تولید	۰/۰۲۹۵۸
گرم	۰/۰۱۳۱۳	یک	۰/۰۳۰۸۸	هزار	۰/۰۲۰۷۰
گوشتش	۰/۰۱۳۱۳	آن	۰/۰۱۸۱۷	تن	۰/۰۱۹۷۲
دیگر	۰/۰۱۱۹۴	قیمت	۰/۰۱۷۲۶	افزایش	۰/۰۱۷۷۵
تولیدی	۰/۰۱۱۹۴	تجذیه	۰/۰۱۷۲۶	آب	۰/۰۱۷۷۵
آنها	۰/۰۱۰۷۵	استفاده	۰/۰۱۶۳۵	خود	۰/۰۱۶۷۷
قرار	۰/۰۱۰۷۵	کارخانه	۰/۰۱۵۴۴	درصد	۰/۰۱۶۷۷
مثال	۰/۰۰۹۵۵	دادن	۰/۰۱۴۵۴	گاو	۰/۰۱۵۷۸
سیریز	۰/۰۰۸۳۶	صنایع	۰/۰۱۴۵۴	شرکت	۰/۰۱۵۷۸
کیلو	۰/۰۰۸۳۶	دامداران	۰/۰۱۷۲۷	پنیر	۰/۰۱۴۸۰

جدول (۵) دقت دسته‌بندی روش پیشنهادی با معیارهای دقت و فراخوانی را نشان می‌دهد. لازم به ذکر است در این آزمایش‌ها تعداد پنج موضوع برای هر معنی واژه در نظر گرفته شده است. علت وجود اختلاف در میزان دقت هر واژه را می‌توان میزان تکرار واژه در پیکره و فراوانی هر کدام از معنی آن و همچنین میزان همپوشانی موضوع‌های مختلف در برگزینده معنی واژه عنوان کرد.

LDA برای استخراج موضوع‌های مختلف هر معنی واژه استفاده شده است و به این ترتیب موضوع بافت جمله نیز به عنوان یک عامل مؤثر در یافتن معنی واقعی کلمه مورد استفاده قرار می‌گیرد. نتایج آزمایش‌ها نشان می‌دهد استخراج ویژگی‌های ذکر شده و استفاده از روش دسته‌بندی پیشینه آنتروپی با دقّت بالایی قادر به مشخص کردن معنی مناسب واژه است، علاوه‌بر این استفاده از مدل موضوعی LDA برای مجموعه آموزشی عملی وقت‌گیر نبوده و پس از آماده‌سازی مجموعه‌های مورد نیاز کار رفع ابهام، به سادگی انجام خواهد گرفت. در مقایسه‌می‌توان گفت روش پیشنهادی با کاهش فضای ویژگی‌ها و استفاده از ویژگی‌های مؤثر، قابلیت ابهام‌زدایی از واژگان مبهم فارسی را با دقّت بالایی دارد.

در کارهای آینده از وردنت فارسی (FarsNet)^۱ در کنار مدل‌های موضوعی (همان‌طور که در برخی زبان‌ها صورت گرفته است) جهت رفع ابهام معنایی از واژگان فارسی استفاده خواهد شد. علاوه‌بر این مدل‌های موضوعی دیگر نظیر... plsa, hdp... را نیز باید آزمود تا مناسب‌ترین مدل برای زبان فارسی به دست آید.

۷- مراجع

Blei D. M. and et al., "Latent Dirichlet allocation", Journal of Machine Learning Research., (2003) , vol. 3, pp. 993-1022.

Boyd-Graber, Jordan. L., "Linguistic Extensions of Topic Models", Princeton University, (2010), pp 2-4.

Decadt B. and et al., "GAMBL, Genetic Algorithm Optimization of Memory-Based WSD", Proceedings of the Third International Workshop on the Evaluation of Systems for the Semantic Analysis of Text (Sens-eval-3), (2004).

Florian R and, et al., "Combining Classifiers for word sense disambiguation", Natural Language Engineering, (2002), vol. 8, pp. 327-341.

Griffiths, T. L., D. M. Blei, et al. "Integrating topics and syntax.", Advances in Neural Information Processing, (2005), 17: 537-544.

Makki R. and Homayoun pour M. M., "Word Sense Disambiguation of Farsi Homographs Using Thesaurus and Corpus", presented at the Proceedings of the 9th international conference on Advances in Natural Language Processing, (2008), Gothenburg, Sweden.

McCallum, A. K. (2002). MALLET: A Machine Learning for Language Toolkit

(جدول-۵): کارایی روش پیشنهادی به درصد

واژه	معانی	دقّت Precision	فراموشی Recall	دقّت کل
شیر	Lion	۸۹/۴۲	۸۷/۳۹	۹۲/۷۳
	Milk	۹۵/۳	۹۶/۶۴	
	Valve	۹۳/۳۴	۹۰/۵۹	
سیر	Garlic	۹۹/۹۴	۹۸/۲۳	۹۸/۶۴
	Journey	۹۷/۳۵	۹۷/۴۲	
مهر	Punch	۹۹/۳۶	۹۹/۸۵	۹۹/۳۱
	Love	۹۱/۷۱	۹۳/۳۳	
	A Solar Month	۹۹/۶۴	۹۷/۴۱	
جو	Atmosphere	۱۰۰	۱۰۰	۱۰۰
	Barley	۱۰۰	۱۰۰	
شکر	Sugar	۹۸/۶۷	۹۵/۴	۹۷/۳۵
	Thankfulness	۹۸/۷۱	۹۶/۵۶	
حلال	Lawful	۹۹/۸۹	۹۶/۳۲	۹۷/۴۲
	Resolvent	۹۸/۹۷	۹۷/۶۶	
خود	Wisdom	۹۸/۹۴	۹۴/۶۲	۹۷/۲۸
	Tiny	۹۹/۳۱	۹۶/۸۶	
نفس	Breath	۹۹/۳۸	۹۹/۹۲	۹۹/۲۵
	Oneself	۹۷/۵۶	۹۸/۵۶	
محرم	Lunar Month	۸۳/۷۹	۸۹/۹۶	۹۶/۸۴
	Intimate	۹۷/۵۹	۹۶/۵۶	
تن	Ton	۸۹/۵۷	۹۴/۳	۹۳/۵۱
	Body	۹۰/۳۶	۹۳/۷۸	
اشکال	Figures	۹۸/۴۵	۹۸/۹۴	۹۹/۰۱
	Defect	۱۰۰	۹۷/۷۲	
اشراف	Dminance	۹۵/۷۴	۹۷/۳۷	۹۷/۱۳
	Knight , Gentleman	۹۴/۴۵	۹۶/۵۴	
قسم	Oath	۹۹/۶۱	۹۹/۱۶	۹۹/۵۸
	Sort	۹۹/۴۸	۱۰۰	
سرم	My head	۱۰۰	۱۰۰	۱۰۰
	Serum	۱۰۰	۱۰۰	
شرف	Honor	۹۹/۷۲	۹۸/۸۷	۹۹/۶۴
	Soon Expected	۹۸/۲۴	۹۹/۳۱	

۶- نتیجه‌گیری و کارهای آینده

این مقاله به بحث واژگان همنگاره فارسی و ارائه روشی برای ابهام‌زدایی از آن‌ها می‌پردازد. اگرچه در بسیاری از زبان‌های طبیعی به طور گسترده از مدل‌های موضوعی در کاربردهای مختلف استفاده می‌شود، در زبان فارسی کمتر به استفاده از آن‌ها پرداخته شده است. در این مقاله از مدل موضوعی

فصل نهم



^۱ <http://nlp.sbu.ac.ir/site/farsnet/>



سعید راحتی متولد سال ۱۳۴۶ در شهرستان قوچان است. وی در سال ۱۳۶۴ در رشته الکترونیک دانشکده فنی دانشگاه تهران تحصیلات دوره کارشناسی خود را آغاز کرد. پس از فراغت از تحصیل، دوره کارشناسی ارشد خود را در رشته مخابرات دانشگاه آزاد اسلامی واحد تهران جنوب در سال ۱۳۷۲ به اتمام رساند؛ سپس مدرک دکتراخود را در سال ۱۳۷۸ در رشته مخابرات دانشگاه آزاد اسلامی واحد علوم و تحقیقات اخذ کرد. با پایان دوره دکترا، به عنوان استادیار و از سال ۱۳۹۰ در سمت دانشیاری دانشگاه آزاد اسلامی واحد مشهد مشغول به خدمت شد. تاکنون بیش از یکصد مقاله در زمینه سامانه‌های هوشمند از ایشان به چاپ رسیده است. نشانی رایانمۀ ایشان عبارتست از:

Rahati@mshdiau.ac.ir

Preiss J., Stevenson M., "Unsupervised Domain Tuning to Improve Word Sense Disambiguation", in Proceedings of NAACL-HLT, (2013), pp. 680-684.

Ratnaparkhi, A., "MAXIMUM ENTROPY MODELS FOR NATURAL LANGUAGE AMBIGUITY RESOLUTION", University of Pennsylvania. Doctor of Philosophy, (1998).

Saeidi C. and Shamsfard M., "Translating Persian documents into English using knowledge based WSD", in Digital Information Management, ICDIM 2009 .Fourth International Conference, (2009), on, pp. 1-6

Steyvers M., Griffiths T., "Probabilistic Topic Models", In Landauer, T., McNamara, D., Dennis, S., Kintsch, W., Latent Semantic Analysis: A Road to Meaning, ,(2006).

Stevenson M, Agirre, E and Soroa A, "Exploiting domain information for Word Sense Disambiguation of medical documents", Journal of the American Medical Informatics Association, (2012), vol. 19, pp. 235-240.

Veenstra J. and et al., "Memory-Based Word Sense Disambiguation", Computers and the Humanities, (2000), vol. 34, pp. 171-177.

Weikum G., "Foundations of statistical natural language processing", SIGMOD Rec, (2002), vol. 31, pp. 37-38.

Walter D., TiMBL: Tilburg Memory-Based Learner, Tilburg University and CNTS Research Group, University of Antwerp (2006).

Yarowsky D., "Hierarchical Decision Lists for Word Sense Disambiguation", Computers and the Humanities, (2000), vol. 34, pp. 179-186.
(<http://rahati.mshdiau.ac.ir>)



بابک مسعودی مدرک کارشناسی خود را در رشته مهندسی برق - الکترونیک در سال ۱۳۷۹ از دانشگاه سیستان و بلوچستان و مدرک کارشناسی ارشد خود را از دانشگاه آزاد اسلامی مشهد در سال ۱۳۸۹ اخذ کرد. وی هم‌اکنون عضو هیئت علمی دانشگاه پیام نور زابل است. زمینه‌های پژوهشی ایشان پردازش زبان طبیعی و پردازش تصاویر دیجیتال است.

نشانی رایانمۀ ایشان عبارتست از:

B.masoudi@pnu.ac.ir