





دانشکده مهندسی فناوری اطلاعات و کامپیوتر

رشته مهندسی کامپیوتر، گرایش هوش مصنوعی و رباتیکز

پایان نامه کارشناسی ارشد

تخمین شباهت معنایی متن با استفاده از روش‌های آماری

نگارنده: فرشته ریاحی

استاد راهنما

دکتر مرتضی زاهدی

استاد مشاور

دکتر علیرضا تجری

بهمن ۱۳۹۹



فرم شماره (۳) صورتجلسه نهایی دفاع از پایان نامه دوره کارشناسی ارشد

با نام و یاد خداوند متعال، ارزیابی جلسه دفاع از پایان نامه کارشناسی ارشد خانم - فرشته ریاحی با شماره دانشجویی ۹۶۰۷۷۲۴ رشته مهندسی کامپیوتر گرایش هوش مصنوعی و رباتیکز تحت عنوان تخمین شباهت معنایی متن با استفاده از روش های آماری که در تاریخ ۱۳۹۹/۱۱/۲۹ با حضور هیأت محترم داوران در دانشگاه صنعتی شاهرود برگزار گردید به شرح ذیل اعلام می گردد:

الف) درجه عالی: نمره ۲۰-۱۹ ب) درجه خیلی خوب: نمره ۱۸/۹۹-۱۸

ج) درجه خوب: نمره ۱۷/۹۹-۱۶ د) درجه متوسط: نمره ۱۵/۹۹-۱۴

ه) کمتر از ۱۴ غیر قابل قبول و نیاز به دفاع مجدد دارد

نوع تحقیق: نظری عملی

اعضای	مرتبه علمی	نام و نام خانوادگی	عضو هیأت داوران
	استادیار	دکتر مرتضی زاهدی	۱- استاد راهنمای اول
	-----	-----	۲- استاد راهنمای دوم
	استادیار	دکتر علیرضا تجری	۳- استاد مشاور
	عربی	مهندس محسن فرهادی	۴- نماینده تحصیلات تکمیلی
	استادیار	دکتر منصور فاتح	۵- استاد منتحن اول
	استادیار	دکتر محسن رضوانی	۶- استاد منتحن دوم

تایید

تاریخ و امضاء و مهر دانشکده:

دکتر علیرضا الفی

نام و نام خانوادگی رئیس دانشکده

تقدیم

به پدرم که زندگی را از دستانش،
مادرم که محبت را از نگاهش آموختم
و خانواده‌ام به پاس ایثار و فداکاری‌هایشان.

تشکر و قدردانی

سپاس خدایی را که به من این توفیق را عطا فرمود تا بتوانم قدمی در راه علم و دانش بردارم.

بدین وسیله از زحمات و تلاش بی دریغ استاد محترم جناب آقای دکتر مرتضی زاهدی و جناب آقای دکتر علیرضا تجری و خانواده عزیزم صمیمانه سپاسگزاری می‌نمایم و همچنین از دوست عزیزم خانم الهام احمدی و دیگر دوستانی که هر کدام به نحوی در تهیه این مجموعه با اینجانب همکاری داشته‌اند تشکر نموده و موفقیت همه آن‌ها را از خداوند متعال خواهانم.

فرشته ریاحی

بهمن ۱۳۹۹

تعهد نامه

اینجانب فرشته ریاحی دانشجوی کارشناسی ارشد رشته مهندسی کامپیوتر مهندسی فناوری اطلاعات و کامپیوتر دانشگاه شاهرود، نویسنده پایان نامه با عنوان تخمین شباهت معنایی متن با استفاده از روش های آماری، تحت راهنمایی مرتضی زاهدی متعهد می شوم:

- تحقیقات در این پایان نامه توسط اینجانب انجام شده است و از صحت و اصالت برخوردار است.
- در استفاده از نتایج پژوهش های دیگر پژوهش گران، به مرجع مورد استفاده استناد شده است.
- مطالب این پایان نامه، تا کنون توسط خود، یا فرد دیگری برای دریافت هیچ نوع مدرک یا امتیازی در هیچ جا ارایه نشده است.
- حقوق معنوی این اثر، به دانشگاه صنعتی شاهرود تعلق دارد، و مقالات مستخرج با نام “دانشگاه صنعتی شاهرود” یا “Shahrood University of Technology” به چاپ خواهد رسید.
- حقوق معنوی تمام افرادی که در به دست آوردن نتایج اصلی پایان نامه تاثیرگذار بوده اند، در مقالات مستخرج از پایان نامه رعایت می گردد.
- در تمام مراحل انجام این پایان نامه، در مواردی که از موجود زنده (یا بافت های آنها) استفاده شده است، ضوابط و اصول اخلاقی رعایت شده است.
- در تمام مراحل انجام این پایان نامه، در مواردی که به حوزه اطلاعات شخصی افراد دسترسی یافته (یا استفاده شده است)، اصل رازداری و اصول اخلاق انسانی رعایت شده است.

فرشته ریاحی

بهمن ۱۳۹۹

مالکیت نتایج و حق نشر

- تمام حقوق معنوی این اثر و محصولات آن (مقالات مستخرج، کتاب، برنامه های رایانه ای، نرم افزارها و تجهیزات ساخته شده) متعلق به دانشگاه صنعتی شاهرود می باشد. این مطلب باید به نحو مقتضی، در تولیدات علمی مربوطه ذکر شود.
- استفاده از اطلاعات و نتایج موجود در این پایان نامه بدون ذکر منبع مجاز نمی باشد.

چکیده

در دنیای امروز به دلیل افزایش حجم اطلاعات اینترنتی و فراوانی فرهنگ‌های مختلف دیجیتالی، وجود تشابه در داده‌های متنی به دلایل مختلفی افزایش پیدا کرده است، به همین دلیل تخمین میزان شباهت میان آن‌ها امری ضروری به حساب می‌آید. همچنین یافتن شباهت میان داده‌های متنی در زمینه‌هایی همچون سیستم‌های بازیابی اطلاعات، سرقت ادبی، داده‌کاوی، طبقه‌بندی اسناد و موارد بسیار دیگری مورد استفاده قرار می‌گیرد.

با ورود جمله یا سند یا متنی به سیستم، میزان شباهت آن با اسناد موجود مورد بررسی قرار می‌گیرد و در زمینه مورد نظر استفاده می‌شود. تاکنون روش‌های مختلف پردازش زبان طبیعی و یادگیری ماشین برای محاسبه میزان شباهت داده‌های متنی معرفی شده‌اند. این روش‌ها به میزان دقت متفاوتی دست یافته‌اند، برای بهبود این معیار نیاز است تا پژوهش بهتری صورت گیرد. از جمله روش‌های آماری که جز روش‌های مبتنی بر پیکره می‌باشد و از ترکیب این روش با روش‌های دیگر می‌توان به نتایج جالب دست یافت.

در این پژوهش، قبل از تبدیل متن به فرمت قابل خواندن برای ماشین، پیش‌پردازش روی داده‌ها اعمال می‌شود تا ویژگی‌ها کاهش یابد و دقت اندازه‌گیری را بالا ببرد. پس از آن با استفاده از تخصیص پنهان دیریکلت هر کلمه به یک نمره همبستگی اختصاص می‌یابد و این روند براساس احتمالات تکرار می‌شود تا اختصاص نمره همبستگی به کلمات بهبود یابد و جمله در دسته مورد نظر قرار گیرد. برای بازنمایی جملات نیز از doc2bow استفاده شده است. سپس برای محاسبه شبیه‌ترین جملات از فاصله جنسن شانون که با مقایسه واگرایی توزیع برچسب مورد نظر بدست می‌آید استفاده شده است. در انتها ماشین بردار پشتیبان خطی و تابع پایه شعاعی را برای دسته‌بندی جملات مشابه به کار برده است. روش پیشنهادی به دقت دسته‌بندی ۸۹ درصد و نمره همبستگی ۰,۹۲ درصد روی پایگاه داده SICK دست پیدا کرده است و توانسته است ۳,۸ درصد دقت بهتری نسبت به پژوهش MaLSTM داشته باشد.

کلمات کلیدی: شباهت معنایی متن، تخصیص پنهان دیریکلت، سند به کیسه کلمات، فاصله جنسن شانون، ماشین بردار پشتیبان

فهرست مطالب

ق	فهرست تصاویر
ش	فهرست جداول
۱	۱ مقدمه
۲	۱.۱ مقدمه
۲	۲.۱ تعریف مسئله
۵	۳.۱ کاربردهای تخمین شباهت معنایی متن
۷	۴.۱ لزوم استفاده از تخمین شباهت معنایی متون
۸	۵.۱ اهداف تحقیق
۸	۶.۱ ساختار تحقیق
۹	۲ پیشینه تحقیق
۱۰	۱.۲ مقدمه
۱۰	۲.۲ تاریخچه شباهت معنایی متن
۱۱	۳.۲ پردازش زبان طبیعی
۱۲	۴.۲ ویژگی‌های مهم و مبانی نظری در شباهت معنایی
۱۲	۱.۴.۲ شبکه واژگان
۱۳	۲.۴.۲ مفهوم اطلاعات
۱۳	۳.۴.۲ ویکی‌پدیا
۱۴	۱.۳.۴.۲ تحلیل معنایی صریح مبتنی بر ویکی‌پدیا
۱۵	۴.۴.۲ مدل N-گرم
۱۵	۵.۴.۲ طولانی‌ترین پیامد مشترک
۱۵	۶.۴.۲ فرکانس مدت - فرکانس اسناد معکوس
۱۶	۷.۴.۲ چارچوب کلمه به بردار
۱۷	۸.۴.۲ مدل کیسه ای از کلمات

۱۸	تخصیص پنهان دیریکلت	۹.۴.۲
۱۸	فاصله جنسن شانون	۱۰.۴.۲
۱۹	ماشین بردار پشتیبان	۱۱.۴.۲
۲۰	تجزیه و تحلیل آماری	۱۲.۴.۲
۲۰	حافظه کوتاه مدت طولانی	۱۳.۴.۲
۲۱	پیشینه تحقیق	۵.۲
۲۱	روش‌های مبتنی بر پیکره	۱.۵.۲
۲۱	روش‌های آماری	۱.۱.۵.۲
۲۳	روش‌های مبتنی بر یادگیری عمیق	۲.۱.۵.۲
۲۶	روش‌های مبتنی بر معنا	۲.۵.۲
۲۹	روش‌های مبتنی بر رشته	۳.۵.۲
۳۰	روش‌های ترکیبی	۴.۵.۲
۳۴	سایر روش‌ها	۵.۵.۲
۳۵	جمع‌بندی	۶.۲
۳۹	روش تحقیق	۳
۴۰	مقدمه	۱.۳
۴۰	طرح مسئله	۲.۳
۴۰	روش پیشنهادی	۳.۳
۴۱	ارائه‌ی مدل کلی	۴.۳
۴۱	پیش پردازش متن	۱.۴.۳
۴۵	استخراج ویژگی از متن	۲.۴.۳
۴۷	جزئیات روش جنسن شانون و دسته‌بندی	۳.۴.۳
۴۹	جمع‌بندی	۵.۳
۵۱	نتایج و تحلیل آن‌ها	۴
۵۲	مقدمه	۱.۴
۵۲	مجموعه دادگان مورد استفاده	۲.۴
۵۳	پیاده‌سازی	۳.۴
۵۵	طبقه‌بندی نمرات شباهت به‌دست‌آمده	۱.۳.۴
۵۵	آزمون و پارامترهای مورد نظر	۴.۴
۵۷	معیار ارزیابی	۵.۴
۵۹	نتایج	۶.۴
۶۲	جمع‌بندی	۷.۴

۶۵	نتیجه‌گیری	۵
۶۶	نتیجه‌گیری	۱.۵
۶۷	کارهای آینده	۲.۵
۶۹	مراجع	

فهرست تصاویر

۴	روش‌های مختلف برای محاسبه تشابه کلمه به کلمه [۱]	۱.۱
۱۷	برآورد کارآمد از بازنمایی کلمات در فضای بردار	۱.۲
۲۴	نمای کلی مدل	۲.۲
۴۲	ساختار کلی روش پیشنهادی	۱.۳
۴۵	مراحل پیش‌پردازش متن	۲.۳
۴۷	مدل تخصیص پنهان دیریکلت [۷۴]	۳.۳
۵۰	مرز تصمیم برای بردار پشتیبان خطی [۷۵]	۴.۳
۵۰	مرز تصمیم برای داده‌های غیرخطی [۷۵]	۵.۳
۵۴	نمودار داده‌های پیش‌پردازش شده	۱.۴
۶۰	نمودار دسته‌بندی داده‌های آموزش	۲.۴
۶۳	نمودار دقت روش‌های بیان شده	۳.۴

فهرست جداول

۳۷	مقایسه برخی کارهای صورت گرفته در زمینه شباهت معنایی متن	۱.۲
۵۳	طبقه‌بندی مجموعه دادگان مورد استفاده براساس امتیازات	۱.۴
۵۳	طبقه‌بندی مجموعه دادگان مورد استفاده براساس برچسب گذاری موجود	۲.۴
۵۳	نمونه‌ای از مجموعه داده‌ها همراه با ضریب همبستگی	۳.۴
۵۶	جملات مورد مقایسه ۴۱۴ و ۷۳۰ مجموعه دادگان	۴.۴
۵۷	روش‌های مقایسه شده با روش پژوهش	۵.۴
۵۹	نتایج طبقه‌بندی مجموعه دادگان آموزش براساس امتیازات	۶.۴
۶۰	ماشین بردار پشتیبان خطی	۷.۴
۶۱	ماشین بردار پشتیبان کرنل چند جمله‌ای	۸.۴
۶۱	ماشین بردار پشتیبان کرنل تابع پایه شعاعی	۹.۴
۶۲	مقایسه دقت بدست آمده روش‌های دسته‌بندی	۱۰.۴
۶۲	مقایسه روش ارائه شده با سایر روش‌ها	۱۱.۴

فصل ١

مقدمه

۱.۱ مقدمه

با رشد روزافزون اینترنت و افزایش قابلیت استفاده از آن، توانایی تعیین تشابه معنایی بین متون برای بسیاری از برنامه‌های کاربردی موردنیاز است. روزانه حجم عظیمی از داده‌های زبان طبیعی در اینترنت اضافه می‌شود. علاوه بر این، ادبیات انسانی در فرهنگ‌های مختلف دیجیتالی شده و در کتابخانه‌های دیجیتالی در دسترس قرار می‌گیرد. در میان این حجم انبوه از داده‌ها، داده‌های مشابه دیگری نیز به دلایل مختلفی مانند متفاوت بودن منابع خبری، زمان و مکان ایجاد گردیده است. مقدار بسیار زیادی از این داده‌ها به زبان طبیعی قالب‌بندی شده‌اند. این امر باعث می‌شود فن‌های پردازش زبان طبیعی^۱ برای استفاده از این مقدار داده بسیار مهم باشند [۱]. بنابراین تخمین شباهت بین داده‌ها یک مسئله کاربردی و مهم می‌باشد. علاوه بر این، صحت سنجش تشابه داده‌ها نیز مهم است. ارزیابی دو بخش از متن توسط معیار ارزیابی مرتبی انجام می‌شود که اختلاف معنایی کامل تا همبستگی کامل معنایی را تشخیص می‌دهد [۲]. به‌عنوان حامل اصلی اطلاعات، متون نقش مهمی در عصر داده‌های بزرگ دارند. تعداد زیادی از داده‌های مبهم مانند واژه‌هایی که چندمعنا دارند، در متون وب چالش‌های بیشتری را برای درک خودکار متن به وجود می‌آورند. برای پرداختن به این مشکل، اندازه‌گیری تشابه معنایی پیشنهاد شده است تا ارتباط معنایی بین کلمات را از نظر شباهت معنایی و واژگانی به‌جای تشابه سطحی کلمات پیدا کند [۳].

۲.۱ تعریف مسئله

در پردازش زبان طبیعی، شباهت معنایی^۲ نقش مهمی دارد و یکی از وظایف اساسی برای بسیاری از برنامه‌های پردازش زبان طبیعی و مناطق مرتبط با آن است [۴]. شباهت معنایی متن^۳ درصد تخمین میزان هم‌ارزی معنایی بین دو قطعه متن است. هدف از یافتن شباهت بین متون، محاسبه میزان نزدیک بودن دو جمله است؛ هم از نظر معنایی و هم واژگانی. شباهت معنایی متن می‌خواهد درجه همگرایی بین دو جمله را ارزیابی کند [۵]. با محاسبه شباهت یک جفت جمله، نمره سطح همبستگی بین آن دو جمله به دست می‌آید [۶]. هدف این است که شباهت معنایی بین دو جمله داده‌شده در یک مقیاس از ۰ تا ۵ را اندازه‌گیری کنیم؛ ۰ نشان می‌دهد که معنایی جملات کاملاً مستقل و ۵ به این معنی است که یک معادل معنایی کامل وجود دارد [۷]. اندازه‌گیری شباهت جملات به سه دسته تقسیم شده‌اند: روش‌های مبتنی بر آمار، روش‌های مبتنی بر معنا و روش‌های ترکیبی. روش‌های آماری با استفاده از محاسبه کلمات مشترک در دنباله رشته‌ها شباهت جملات را

¹Natural language processing(NLP)

²semantic similarity(SS)

³semantic textual similarity(STS)

محاسبه می‌کنند [۸]. روش‌های مبتنی بر معنا، معنای واقعی پشت کلمات یا کل عبارت را در نظر می‌گیرد و روش‌های ترکیبی از ترکیب این دو روش استفاده می‌کند.

شباهت معنایی متن و استنباط زبان طبیعی دو وظیفه محبوب زبان طبیعی است که برای محک زدن مدل‌های بیان جمله درجایی که دو جمله باهم جفت شده‌اند، استفاده می‌شود [۹]. جملات می‌تواند از نظر لغوی یا معنایی مشابه باشد. تشابه واژگانی به معنای شباهت مبتنی بر رشته است. با این حال، شباهت معنایی نشانگر معنای مشابه در جمله است. شباهت یابی مبتنی بر معنا که به معنی جملات بستگی دارد، رویکردهای مختلفی دارد. این روش‌ها فن‌های مختلفی را برای مقایسه دو جمله از نظر معنایی به کار می‌گیرند. رویکرد اول، مبتنی بر مجموعه نوشته‌ها (پیکره)^۱ شباهت کلمات را بر اساس تجزیه و تحلیل آماری متون بزرگ می‌یابد. علاوه بر این، برای تجزیه و تحلیل یک گروه بزرگ برای نشان دادن معانی کلمات از یادگیری عمیق می‌تواند استفاده کند. رویکرد دوم، مبتنی بر دانش، به شبکه معنایی برای کلمات بستگی دارد. معنی کلمات و روابط بین کلمات در این شبکه معنایی گنجانده شده است. بیشترین کاربرد معنای خالص شبکه واژگان^۲ است. رویکرد سوم، مبتنی بر ساختار، از اطلاعات ساختار یک جمله برای به دست آوردن معنی این جمله استفاده می‌کند. جمله مشابه باید دارای ساختار اساسی مشابه باشد [۱].

در پردازش زبان طبیعی، شباهت معنایی نقش مهمی دارد و یکی از وظایف اساسی برای بسیاری از برنامه‌های پردازش زبان طبیعی و مناطق مرتبط با آن است. شباهت متنی معنایی را می‌توان با استفاده از متریک بیش از مجموعه‌ای از اسناد تعریف کرد و ایده این است که شباهت معنایی بین آن‌ها را پیدا کنید. شباهت بین اسناد مبتنی بر روابط مستقیم و غیرمستقیم بین آن‌هاست. این روابط را می‌توان با وجود روابط معنایی در بین آن‌ها سنجید و شناخت [۴].

شباهت معنایی یک مورد خاص از ارتباط است که به شباهت بین مفاهیم مربوط می‌شود. علاوه بر این، روش‌های اندازه‌گیری شباهت معنایی دو مفهوم را به عنوان ورودی در نظر می‌گیرند و یک نمره عددی را برمی‌گردانند که میزان آن‌ها مشابه است. در حقیقت، چندین اثر در مورد اقدامات شباهت معنایی در سال‌های گذشته ارائه شده است. چهار خانواده مختلف به شرح زیر ذکر شده است:

۱. اقدامات مبتنی بر مسیر [۱۰]؛ [۱۱]؛ [۱۲] شباهت‌ها را با فاصله معنایی مفاهیم در یک هستی‌شناسی معین ارزیابی می‌کنند. این نوع اقدامات در محاسبات ساده در بسیاری از مناطق مورد استفاده گسترده قرار گرفته است.

۲. اقدامات مبتنی بر محتوای اطلاعات [۱۳]؛ [۱۴]؛ [۱۵]؛ [۱۶]؛ [۱۷] شباهت‌های مفهوم اطلاعات^۳ را در یک هستی‌شناسی مشخص می‌کنند. مفهوم اطلاعات بر پایه این فرض استوار است که موجودات خاص و ویژه، مفهوم اطلاعات بیشتری نسبت به موارد کلی و انتزاعی ارائه می‌دهند.

۳. اقدامات مبتنی بر ویژگی [۱۸]؛ [۱۹]؛ [۲۰]؛ [۲۱]؛ [۲۲] شباهت‌ها را با استفاده از مجموع وزنی مشترک و ویژگی‌های غیرمعمول مفاهیم ارزیابی می‌کنند.

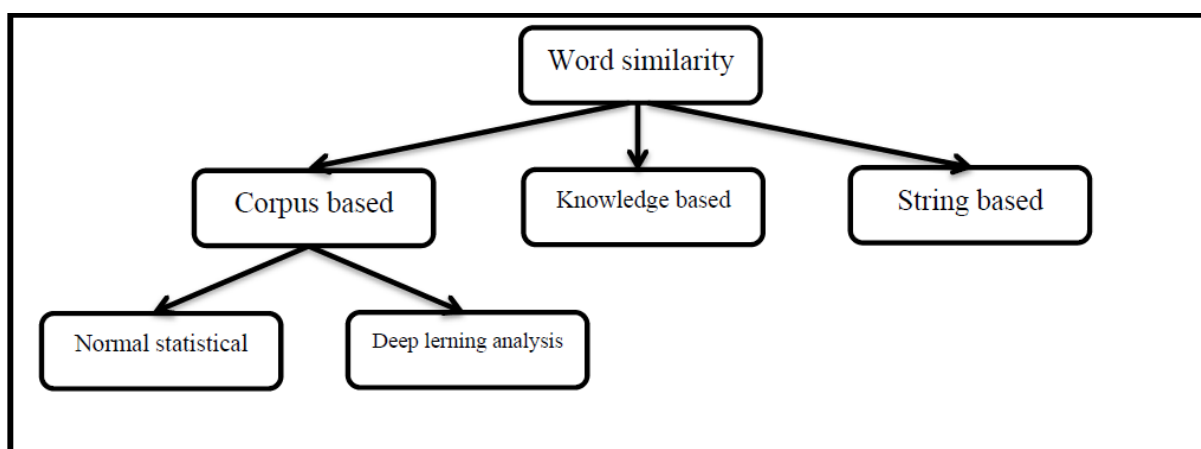
¹corpuse

²wordnet

³Information Content(IC)

۴. اقدامات ترکیبی حداقل دو تا از اقدامات فوق را برای ادغام مزایا ترکیب می‌کنند که منجر به وابستگی به پارامترها می‌شود [۲۳].

سه روش اصلی برای اندازه‌گیری شباهت بین کلمات وجود دارد ۱. تشابه کلمات مبتنی بر مجموعه نوشته‌ها ۲. دانش مبتنی بر کلمات ۳. شباهت کلمات مبتنی بر رشته که تشابه مبتنی بر پیکره خود دو روش مختلف دارد اول استفاده از آنالیز آماری عادی و دوم استفاده از یادگیری عمیق می‌باشد. پیدا کردن شباهت بین کلمات، هسته تشابه جمله است. در ادبیات معیارهای زیادی برای محاسبه شباهت کلمه به کلمه وجود دارد. سلسله‌مراتب برای روش‌های مورد استفاده برای اندازه‌گیری شباهت کلمه در شکل ۱.۱ نشان داده شده است [۱].



شکل ۱.۱: روش‌های مختلف برای محاسبه تشابه کلمه به کلمه [۱]

با توجه به منابع لغوی مورد استفاده، رویکردهای موجود برای محاسبه شباهت معنایی را می‌توان به رویکردهای مبتنی بر پیکره و رویکردهای دانش‌محور طبقه‌بندی کرد. رویکردهای مبتنی بر پیکره معنای واژگانی را با استخراج بردارهای کلمه‌ای از خصوصیات توزیع کلمات متنی آن‌ها در یک کالبد، مانند پیکره ملی انگلیس، ویکی‌پدیا و موتور جستجوی وب یاد می‌گیرند. با توجه به داده‌های متن گسترده، چنین رویکردهایی برای تصفیه اطلاعات معنایی واژگان پنهان در یک فضای بردار، مانند تحلیل معنایی صریح^۱ (گابیلوویچ و مارکوویچ ۲۰۰۷) به فن‌های کاهش بعد تکیه می‌کنند. به‌عنوان یک جایگزین، لو و همکاران (۲۰۱۵، ۲۰۱۶) مدل شبکه عصبی بسته‌های کیسه کلمات^۲ را آموزش می‌دهد تا بردارهای کلمه‌ای کم بعدی برای اندازه‌گیری شباهت معنایی واژگانی یاد بگیرد. با این حال، رویکردهای مبتنی بر پیکره معمولاً از پیچیدگی محاسباتی بالایی برای شرکت‌های پیش‌پردازش و کاهش ابعاد بردارهای کلمه برخوردار هستند. علاوه بر این، رویکردهای مبتنی بر بدنه برای محاسبه ارتباط معنایی نسبت به شباهت معنایی کاربردی‌تر است. در حالی که، رویکردهای دانش‌بنیان که

^۱ Explicit Semantic Analysis (ESA)

^۲ Continuous Bag-of-Words (CBOW)

عمدتاً به منابع دانش ساختاریافته بستگی دارند، می‌توانند با بهره‌گیری از ساختار ذاتی پایه‌های دانش سلسله مراتبی، مانند هستی‌شناسی، اصطلاحات، فرهنگ لغت، طبقه‌بندی، بر پیچیدگی محاسباتی غلبه کنند. بیشتر رویکردهای هنری مبتنی بر یک پایگاه داده شبکه واژگان برای اندازه‌گیری شباهت معنایی مفهوم ارائه‌شده است (لی و همکاران ۲۰۰۳؛ سانچز و همکاران ۲۰۱۱؛ لیو و همکاران ۲۰۱۲؛ منگ و همکاران ۲۰۱۳؛ گائو و همکاران ۲۰۱۵). در مطالعه لو و همکاران (۲۰۱۶)، شبکه واژگان به‌عنوان کمکی برای تقویت معنایی بردارهای کلمه مورد استفاده قرار می‌گیرد، زیرا ثابت شده است که در مطالعات مبتنی بر دانش قبلی مؤثر است [۲۳].

سنجش شباهت معنایی متن مربوط به تعیین میزان تشابه معنایی بین جفت متون است. چنین اقدامات شباهت معنایی متن، توجه زیادی به تحقیقات در پردازش زبان طبیعی، درک زبان طبیعی و جوامع بازیابی اطلاعات جلب کرده است. اولین روش‌ها برای اندازه‌گیری شباهت بین متون طولانی، به‌عنوان مثال، برای طبقه‌بندی اسناد، یا بین متن کوتاه و متن طولانی، به‌عنوان مثال، برای پرس‌وجو یک سند در یک سیستم جستجو ارائه‌شده است. راحت‌ترین راه برای نشان دادن یک سند توسط آن روش‌های کیف کلمات یک مدل فضای برداری است که در آن هر سند به‌عنوان کیف کلمات در نظر گرفته می‌شود و توسط یک بردار ارائه می‌شود. تشابه کسینوسی، به‌عنوان یک اقدام سنتی بردار که در بازیابی اطلاعات مورد استفاده قرار می‌گیرد، شناخته می‌شود، یک اقدام محبوب برای ارزیابی شباهت بین متون است [۲۴].

۳.۱ کاربردهای تخمین شباهت معنایی متن

شباهت معنایی یک فن گسترده است که در بسیاری از زمینه‌ها از جمله بازیابی اطلاعات، کشف دانش، هوش مصنوعی، علوم شناختی، پاسخ دادن به پرسشنامه سامانه‌ای، سامانه‌های پیشنهاددهنده، ارزیابی عملکرد، داده‌کاوی، مدیریت دانش، پردازش زبان طبیعی، زبان‌شناسی محاسباتی و استخراج اطلاعات کاربرد دارد. علاوه بر این‌ها شباهت معنایی متن قادر به ارائه یک چارچوب یکپارچه برای ارزیابی مؤلفه‌های معنایی بیرونی چندگانه از جمله ابهام در معانی جملات برچسب‌گذاری معنایی متن و کیفیت نام‌گذاری برای مراجع است. بازیابی اطلاعات، اطلاعات و داده‌های از دست‌رفته را بازیابی می‌نماید. داده‌کاوی فرآیند کشف الگوهای موجود در مجموعه‌ی بزرگی از داده‌ها است.

فرآیند کشف دانش از پایگاه داده شامل مراحل زیر می‌باشد:

۱. انتخاب ۲. پیش‌پردازش ۳. انتقال و تبدیل ۴. داده‌کاوی ۵. تفسیر و ارزیابی

می‌توان این فرآیند را به‌صورت زیر نیز بیان نمود:

۱. شناخت کسب‌وکار ۲. شناخت داده‌ها ۳. فراهم آوردن داده‌ها ۴. مدل‌سازی ۵. ارزیابی ۶.

گسترش.

در حالت کلی این گام‌ها به‌صورت زیر قابل خلاصه می‌باشند:

۱. پیش‌پردازش ۲. داده‌کاوی ۳. اعتبار سنجی نتایج

اقدامات شباهت احکام نقش مهمی در تحقیقات و برنامه‌های مربوط به متن در زمینه‌هایی مانند متن کاوی، بازیابی صفحه وب و سیستم گفتگو ایفا می‌کند [۴].

در بازیابی صفحه وب، شباهت جملات به‌عنوان یکی از بهترین فن‌های بهبود اثربخشی بازیابی، اثبات‌شده است که در آن از عناوین برای ارائه اسناد در صفحه پیدا کردن صفحه استفاده می‌شود. در بازیابی تصویر از وب، استفاده از متن کوتاه پیرامون تصاویر می‌تواند به‌دقت بازیابی بالاتری نسبت به استفاده از کل سندی که در آن تصویر تعبیه‌شده است، دست یابد. در متن کاوی، از شباهت جمله به‌عنوان معیاری برای کشف دانش غیب از پایگاه داده‌های متنی استفاده می‌شود. علاوه بر این، ترکیب شباهت متن کوتاه برای برنامه‌هایی مانند جمع‌بندی متن، طبقه‌بندی متن و ترجمه ماشینی تفکیک مفهوم کلمه و برای خلاصه‌سازی استخراج، خوشه‌بندی متن، همچنین تشخیص ناهنجاری‌ها مفید است. این کاربردهای نمونه نشان می‌دهد که محاسبه شباهت جمله به یک مؤلفه عمومی برای جامعه تحقیقاتی درگیر در بازنمایی و کشف دانش مرتبط با متن تبدیل‌شده است. اقدامات شباهت متن همچنین برای ارزیابی انسجام متن مفید بوده است [۲۵].

یکی از اولین کاربردهای شباهت متن شاید مدل برداری در بازیابی اطلاعات باشد، جایی که سندی که بیشتر مربوط به یک پرس‌وجو ورودی است با رتبه‌بندی اسناد در یک مجموعه به ترتیب معکوس از شباهت آن‌ها با پرس‌وجو داده‌شده مشخص می‌شود. یک چالش بزرگ در زبان‌شناسی محاسباتی، ساختن مدل‌های نمایندگی برای فعال کردن درک طبیعی است. به‌منظور آموزش و ارزیابی این روش‌ها، جامعه چندین چالش و مجموعه داده از جمله ماشین‌آلات درک، پاسخگویی به سؤالات خودکار پیشنهاد کرده است. پاسخگویی کوتاه درجه‌بندی، استنباط زبان طبیعی^۱ و تشابه متن معنایی. در این کارها، سیستم درسی زبان طبیعی باید دو قطعه متن را جفت کند و سپس خروجی مانند ارتباط بین یک سؤال و یک متن را (درک دستگاه)، یک سؤال و یک پاسخ (پاسخ به سؤال) ارائه دهد [۹].

شباهت معنایی همچنین در بسیاری از برنامه‌های وب معنایی مانند استخراج جامعه، تولید هستی‌شناسی و تفکیک موجودیت نقش دارد. همچنین برای جستجوی تویتر مفید است، جایی که لازم است با دقت اندازه‌گیری ارتباط معنایی بین مفاهیم یا موجودات انجام شود. در بازیابی اطلاعات یکی از مشکلات اصلی بازیابی مجموعه‌ای از اسناد و بازیابی تصاویر توسط زیرنویس‌ها است که از لحاظ معنایی با یک پرس‌وجو کاربر مشخص‌شده در موتور جستجوگر وب مرتبط است. در زمینه بانک اطلاعاتی نیز می‌توان از شباهت متن برای تطبیق طرح‌واره برای حل ناهمگونی معنایی برای سیستم به اشتراک‌گذاری داده‌ها، سیستم ادغام داده‌ها، سیستم ارسال پیام و سیستم مدیریت داده هم‌تا به نظیر استفاده کرد. همچنین برای عملیات پیوند رابطه در دیتابیس که صفات پیوند متنی به یکدیگر شبیه هستند مفید است. همچنین دارای صحت دامنه برنامه است از جمله ادغام و جستجوی داده‌ها از منابع ناهمگن، تمیز کردن داده‌ها و استخراج داده‌ها. سنجش تشابه معنایی متن، بین کلمات، اصطلاحات، جملات، پاراگراف و سند نقش مهمی در علم کامپیوتر و زبان محاسباتی ایفا می‌کند. این برنامه همچنین در چندین زمینه کاربردهای زیادی دارد، مانند

¹ Natural language inference (NLI)

داده‌ورزی زیست پزشکی و ژئوی اطلاعات. علاوه بر زبان‌شناسی، محاسبات معنایی نیز در زمینه‌های خاص حوزه، مانند زیست پزشکی ظاهر می‌شود. شباهت معنایی در برنامه‌های دیگر مانند ترجمه و تفسیر و ترجمه حافظه نیز نقش اساسی دارد. با این حال، محاسبه تشابه معنایی بین جملات یک کار پیچیده و دشوار است. باگذشت سال‌ها، کارهای مشترک SemEval برای تنظیم دقیق و کامل کردن این اقدامات تشابه و کشف ماهیت معنا در زبان کارکرده است. از دیگر کاربردهای شباهت معنایی:

- اصلاح اشتباهات نوشتاری که از آن با عنوان غلطیاب املایی یاد می‌شود.
- شناسایی متون تقلبی یا دزدی علمی و ادبی که در مقالات لاتین با عنوان سرقت ادبی^۱ مطرح است.
- یافتن عبارات جایگزین در موتورهای جست‌وجو جهت افزایش دقت و صحت در بازیابی اسناد.
- تشخیص تقلب در کد نرم‌افزار.
- استخراج پاسخ در سامانه‌های پرسش و پاسخ.
- طبقه‌بندی اسناد بر اساس موضوع یا اندازه متن در فایل متن.

۴.۱ لزوم استفاده از تخمین شباهت معنایی متون

با توجه به کاربردهای که در بخش قبل برای تخمین شباهت معنایی ذکر شد و همچنین رشد نمایی تعداد متون، لزوم استفاده از یک روش مکانیزه برای به دست آوردن تشابه بین متون ضروری به نظر می‌رسد. روزانه هزاران خبر به یک خبرگزاری ارسال می‌شود، میلیون‌ها ایمیل ردوبدل می‌شود، متون تقلبی یا دزدی علمی و سرقت ادبی زیادی صورت می‌گیرد. همه موارد برشمرده لزوم استفاده از روش‌های مکانیزه را بیش‌ازپیش نمایان می‌نماید.

امروزه رشد اطلاعات در رسانه‌های دیجیتالی تجزیه و تحلیل مقادیر زیادی از داده‌های متنی را تشویق می‌کند. این مورد توجه محققان علوم داده و هوش مصنوعی و همچنین صنعت اینترنت را به خود جلب می‌کند. کاربران اینترنت مسئولیت بخش قابل توجهی از این رشد را بر عهده دارند. آن‌ها اطلاعاتی را در شبکه وارد می‌کنند که برای به اشتراک گذاری دانش نیز استفاده می‌شود. بخش مهمی از این دانش در مخازن از جمله انجمن‌های پرسش و پاسخ، روزنامه‌های دیجیتالی و دایره‌المعارف‌های دیجیتالی یافت می‌شود. با توجه به تکثیر بی‌شمار اطلاعات در این مخازن، نگرانی‌های زیادی در مورد نحوه تغذیه و مصرف دانش کاربران وجود دارد. برخی از این نگرانی‌ها شامل حذف موارد اضافی در انجمن‌های پاسخگو به سؤال یا بهره‌برداری از افزونگی‌ها برای ارزیابی اطمینان از اخبار در رسانه‌ها یا به‌سادگی فشرده‌سازی اندازه متن است. دستیابی به این پردازش گسترده اطلاعات برای

¹ plagiarism

داوران انسانی به وضوح غیرممکن است. در این سناریو، روش‌های آماری پردازش زبان طبیعی کمک قابل توجهی است [۲۶].

۵.۱ اهداف تحقیق

با توجه به کاربردهای شباهت یابی و لزوم به کارگیری آن که در بخش‌های قبلی به آن اشاره شد، هدف این تحقیق تخمین شباهت معنایی مجموعه‌ای از متون با روش‌های آماری و دسته‌بندی با ماشین بردار پشتیبان است. هدف اصلی روش‌های شباهت متن بهبود توانایی رایانه‌ها در درک زبان بشری است که می‌تواند از طریق روش‌های مختلف برای بیان معنایی مشابه باشد. یکی از اهداف اصلی کار شباهت معنایی متن ایجاد یک چارچوب یکپارچه با ترکیب چندین مؤلفه معنایی مستقل به منظور یافتن تأثیر آن‌ها در چندین کار پردازش زبان طبیعی است تدوین چنین چارچوبی یک مسئله مهم تحقیقاتی است. یکی از چالش‌های پیش‌رو در روش‌های بیان‌شده در این زمینه رسیدن به دقت حداکثری نتایج به دست آمده از تخمین شباهت است. همچنین ضریب همبستگی نیز یکی دیگر از معیارهایی است که باید به بهبود آن توجه کرد. بنابراین ما درصدد هستیم تا این معیارها را بهبود بدهیم.

۶.۱ ساختار تحقیق

در این فصل به تعریف شباهت معنایی پرداختیم و کاربردهای آن و هدف از انجام این پژوهش را بیان نمودیم. در فصل دوم این تحقیق به معرفی انواع مختلف روش‌های شباهت یابی اسناد متنی و ارزیابی کلی الگوریتم‌های معرفی شده به همراه اشکال یا فرمول‌ها خواهیم پرداخت. این روش‌ها غالباً در سه دسته کلی روش‌های مبتنی بر پیکره و روش‌های مبتنی بر معنا و مبتنی بر رشته قرار می‌گیرند که هر یک مزایا و معایب خاص خود را دارا می‌باشند.

در فصل سوم به بررسی روش پیشنهادی خود در این زمینه می‌پردازیم و مدل خود را ارائه می‌دهیم. همچنین به شرح جزئیات تحقیق پرداخته و روش خود را به تفصیل بیان خواهیم کرد.

در فصل چهارم نتایج به دست آمده از روش ارائه شده و جزئیات پیاده‌سازی بیان خواهد شد و با روش‌های قبلی بیان شده در این زمینه مقایسه خواهیم کرد.

در فصل پنجم در خصوص نتایج به دست آمده در فصل چهارم بحث و بررسی خواهیم کرد و اینکه آیا نتایج حاصل یک جواب عمومی می‌باشد یا اینکه در حوزه خاصی قابل استفاده است. همچنین پیشنهادی را که برای بهبود آن مفید است در آینده بررسی گردد، ارائه خواهیم کرد.

فصل ۲

پیشینه تحقیق

۱.۲ مقدمه

این فصل به مرور کارهای انجام شده در ارتباط با شباهت معنایی جملات و متون کوتاه و طولانی می پردازد. روش های مختلفی در این زمینه وجود دارد که هرکدام مزایا و معایب خاص خود را دارند. به طور کلی روش های شباهت یابی می تواند بر اساس ساختار متن یا معنای آن انجام بگیرد. ابتدا با مبانی نظری در زمینه شباهت یابی آشنا می شویم، سپس مروری بر کارهایی که در این زمینه انجام شده است خواهیم داشت. در انتها به جمع بندی کارهای بیان شده در این فصل می پردازیم و مقایسه ای بین آن ها انجام خواهیم داد.

۲.۲ تاریخچه شباهت معنایی متن

شناسایی شباهت معنایی متن در متن های کوتاه در سال ۲۰۰۶ گزارش شد که در [۲۷]؛ [۲۸] ارائه گردید. پس از آن، تمرکز روی اسناد بزرگ یا کلمات فردی تغییر یافت. از سال ۲۰۱۲، وظیفه تشابه معنایی نه تنها برای یافتن شباهت بین دو متن محدود نیست، بلکه برای ایجاد نمره تشابه توسط کارهای مختلف [۲۹]؛ [۳۰] نیز ایجاد شد.

در سال ۲۰۱۲ اولین کار آزمایشی در SemEval 2012، به عنوان بخشی از اجلاس SEM 2012، با موفقیت بزرگ برگزار شد: ۳۵ گروه با ۸۸ اجرای سیستم شرکت کردند. علاوه بر این، یک کارگاه پشتیبانی DARPA در دانشگاه کلمبیا برگزار گردید. در سال ۲۰۱۳، شباهت معنایی متن به عنوان کار مشترک رسمی اجلاس SEM 2013 انتخاب شد. بر این اساس در STS 2013، دو وظیفه تنظیم گردید: وظیفه اصلی هسته، که مشابه کار ۲۰۱۲ است. و یک کار آزمایشی در نوع شباهت تایپ شده بین سوابق نیمه ساختاریافته [۳۱]؛ سیستم UKP در SemEval 2012 بهترین کار را در شباهت معنایی متن انجام داده است.

در سال ۲۰۱۲ استفاده از شبکه های عصبی برای پردازش زبان طبیعی مورد توجه بسیاری قرار گرفت، یکی از مقالات اصلی که از شبکه های معلولیت در کارهای سنتی پردازش زبان طبیعی استفاده کرده است (کولبرت و همکاران، ۲۰۱۱) می باشد. در سال های گذشته مدل های شبکه عصبی برای کارهای مختلفی از ترجمه ماشین (زو و همکاران، ۲۰۱۳؛ چو و همکاران، ۲۰۱۴) برای پاسخگویی به سؤال استفاده شده اند (وستون و همکاران، ۲۰۱۵) [۳۲].

وظایف مشترک STS از سال ۲۰۱۲ به طور سالیانه برگزار می شود و محل برگزاری ارزیابی الگوریتم ها و مدل های پیشرفته ای است. در طول این دوره، مجموعه ای متنوع از ژانرها و منابع داده شده مورد بررسی قرار گرفته است (سرفصل های خبری، توضیحات تصویری و فیلم، واژه نامه هایی از منابع لغوی از جمله شبکه واژگان، FrameNet، OntoNote، انجمن های گفتگوی وب و مجموعه داده های پرسش و پاسخ) [۲].

هرسال STS به یکی از محبوب‌ترین کارها در رقابت SemEval تعلق دارد. بهترین سیستم STS در [۳۳] از شباهت واژگانی و تحلیل معنایی صریح استفاده کرد. در SemEval 2013، بهترین مدل از مدل‌های معنایی مانند تحلیل معنایی، منابع اطلاعات خارجی (شبکه واژگان) و فن‌های تطابق N-گرم استفاده کردند. برای SemEval 2014 و 2015 بهترین سیستم از طریق [۳۴]؛ [۳۵]؛ [۳۶] معرفی شدند. آن‌ها الگوریتم جدیدی را معرفی کردند که کلمات را بین دو جمله قرار می‌دهد. آن‌ها نشان دادند که از این رویکرد می‌توان برای شباهت معنایی متن نیز استفاده کرد.

روش‌های شباهت معنایی متن را می‌توان به‌طور کلی به دودسته تقسیم کرد: روش‌های کنترل نشده و روش‌های نظارت‌شده. روش‌های کنترل نشده لی و همکاران [۳۷] و میهالچه و همکاران [۳۸] شاید نمایندگان آغازگر تحقیقات شباهت معنایی متن بانفوذ قابل توجهی باشند. روش بی‌نظیر سلطان و همکاران [۳۹] تأثیرات عمده‌ای را در کارهای STS SemEval از سال ۲۰۱۴ داشته است. در سال‌های اخیر، سامانه‌های تحت نظارت STS مدل‌های یادگیری عمیق را با فن‌های پردازش زبان طبیعی سنتی ترکیب کرده‌اند و بهترین عملکرد را به دست آورده‌اند.

در حال حاضر الگوریتمی دقیق‌تر است که به شباهت انسانی نزدیک‌تر باشد. با استفاده از مجموعه‌ای از معیارها به‌عنوان شباهت رشته‌ای کلمه و طولانی‌ترین الگوریتم متعارف مشترک برای محاسبه شباهت معنایی بین جملات استفاده می‌شود [۴۰].

۳.۲ پردازش زبان طبیعی

پردازش زبان‌های طبیعی یکی از زیرشاخه‌های بااهمیت در حوزه گسترده علوم رایانه و هوش مصنوعی است که به تعامل بین کامپیوتر و زبان‌های طبیعی (انسانی) می‌پردازد؛ بنابراین پردازش زبان‌های طبیعی بر ارتباط انسان و رایانه متمرکز است. پس چالش اصلی و عمده در این زمینه درک زبان طبیعی و ماشینی کردن فرایند درک و برداشت مفاهیم بیان‌شده با یک‌زبان طبیعی انسانی است. به تعریف دقیق‌تر، پردازش زبان‌های طبیعی عبارت است از استفاده از رایانه برای پردازش زبان گفتاری و زبان نوشتاری. بدین معنی که رایانه‌ها را قادر سازیم که گفتار یا نوشتار تولیدشده در قالب و ساختار یک‌زبان طبیعی را تحلیل و درک نموده یا آن را تولید نمایند. در این صورت، با استفاده از آن می‌توان به ترجمه‌ی زبان‌ها پرداخت، از صفحات وب و بانک‌های اطلاعاتی نوشتاری جهت پاسخ دادن به پرسش‌ها استفاده کرد، یا با دستگاه‌ها، مثلاً برای مشورت گرفتن به گفت‌وگو پرداخت. هدف اصلی در پردازش زبان طبیعی، ایجاد تئوری‌های محاسباتی از زبان، با استفاده از الگوریتم‌ها و ساختارهای داده‌ای موجود در علوم رایانه است. بدیهی است که در راستای تحقق این هدف، نیاز به دانشی وسیع از زبان است و علاوه بر محققان علوم رایانه، نیاز به دانش زبان‌شناسان نیز در این حوزه می‌باشد. با پردازش اطلاعات زبانی می‌توان آمار موردنیاز برای کار با زبان طبیعی را استخراج کرد. کاربردهای پردازش زبان طبیعی به دودسته کلی قابل تقسیم است: کاربردهای نوشتاری و کاربردهای گفتاری. از کاربردهای نوشتاری آن می‌توان به استخراج اطلاعاتی خاص از یک متن، ترجمه یک متن به زبانی دیگر یا یافتن مستندات خاص در یک پایگاه داده نوشتاری (مثلاً یافتن

کتاب‌های مرتبط به هم در یک کتابخانه) اشاره کرد. نمونه‌هایی از کاربردهای گفتاری پردازش زبان عبارت‌اند از: سامانه‌های پرسش و پاسخ انسان با رایانه، سرویس‌های اتوماتیک ارتباط با مشتری از طریق تلفن، سامانه‌های آموزش به فراگیران یا سامانه‌های کنترلی توسط صدا [۴۱]؛ [۴۲].

۴.۲ ویژگی‌های مهم و مبانی نظری در شباهت معنایی

در این بخش؛ برای سنجش شباهت بین دو جمله ابتدا با مبانی نظری در این زمینه آشنا می‌شویم.

۱.۴.۲ شبکه واژگان

شبکه واژگان پایگاه داده‌ای لغوی برای زبان انگلیسی بشمار می‌آید که در سال ۱۹۸۵ میلادی پیاده‌سازی آن شروع شد و مراحل تکامل آن همچنان ادامه دارد. این منبع لغت حاصل تحقیقات انجام‌شده در آزمایشگاه علوم شناختی دانشگاه پرینستون و در ارتباط با مدل‌سازی دانش لغوی انسان است و به‌طور هم‌زمان مزایای یک فرهنگ جامع و فرهنگ لغت را دارا است. شبکه واژگان در واقع شبکه معنایی از تمام کلمات زبان انگلیسی است که در حال حاضر به‌صورت متن‌باز قابل استفاده است.

جملات در شبکه واژگان با توجه به تعاریف کلی که توسط اصطلاحات ارائه می‌شوند، به‌عنوان گروه‌های مترادف یعنی نحو^۱ یا مفهوم^۲ سازماندهی می‌شوند. با توجه به خصوصیات فنی و عناصر محاسباتی، رویکردهای مبتنی بر شبکه واژگان می‌توانند در سه دسته اصلی قرار بگیرند: مبتنی بر مسافت مسیری، مبتنی بر محتوای اطلاعات و مبتنی بر ویژگی‌ها. رویکردهای مبتنی بر فاصله مسیری لبه‌های جداکننده دو مفهوم، همراه با عمق مفهوم، چگالی محلی، نوع پیوند و استحکام پیوند است.

شبکه واژگان، یک پایگاه داده آنلاین واژگان انگلیسی است که به‌طور گسترده به‌عنوان دانش قبلی برای محاسبه شباهت معنایی بین دو مفهوم مورد استفاده قرار می‌گیرد. این واژگان را به اسم، افعال، صفت‌ها و قیدها به‌عنوان قالب مجموعه مترادف تقسیم می‌کند. هر مفهوم شامل گروهی از کلمات است که به معنای مشترک است. همگام‌سازی‌ها در شبکه واژگان به‌طور کلی با ۱۱ نوع روابط معنایی نظیر: هایپونامی / هایپرنیم is-a، مرامی گری / هلمونی (بخشی از عضو-عضو، ماده-ماده) و antonymy (کلمه متضاد) و غیره به ارث رسیده‌اند. تقریباً ۸۰٪ روابط is-a هستند. بنابراین، شبکه واژگان را می‌توان به‌عنوان ساختار سلسله مراتبی مانند درخت در نظر گرفت که در آن کلیت مفهوم به‌صورت یکنواخت با عمق آن‌ها کاهش می‌یابد. روابط ارثی معمولاً برای محاسبه شباهت معنایی مورد استفاده قرار می‌گیرد، در حالی که روابط دیگر همچون کلمه‌شناسی / هلو نامی و آنتونی نیز برای ارزیابی ارتباط معنایی مناسب‌تر هستند. آثار مرتبط با استفاده از ساختار سلسله

¹synset

²concept

مراتبی و خصوصیات معنایی شبکه واژگان برای اندازه‌گیری شباهت معنایی، که شامل طول لبه، عمق و تراکم است. رابطه is-a از شبکه واژگان است که عمدتاً برای محاسبه کوتاه‌ترین طول مسیر در اقدامات شباهت استفاده می‌شود از طرف دیگر، توپولوژی مفاهیم در شبکه واژگان در تعیین مفهوم اطلاعات ذاتی موجود در مفاهیم کمک می‌کند که می‌تواند به‌عنوان اقدامات ارزیابی کمکی برای بهبود دقت مورد استفاده قرار گیرد [۲].

۲.۴.۲ مفهوم اطلاعات

در تئوری اطلاعات، محتوای اطلاعات، سوابق متغیر یا سیگنال تصادفی، میزان اطلاعات به‌دست‌آمده هنگام نمونه‌برداری است. به‌طور رسمی، محتوای اطلاعات متغیر تصادفی است که بدون توجه به اندازه‌گیری متغیر تصادفی برای هر رویداد در تئوری احتمال تعریف شده است. در واقع محتوای اطلاعات در واحد اطلاعات بیان می‌شود. ارزش مورد انتظار سوابق، آنتروپی نظری اطلاعات است، میانگین اطلاعاتی که یک ناظر انتظار دارد در هنگام نمونه‌گیری از متغیر تصادفی، در مورد یک سیستم به دست آورد. با توجه به یک متغیر تصادفی x با عملکرد توده احتمال $PX(x)$ ، اطلاعات خود اندازه‌گیری x به‌عنوان نتیجه x به‌عنوان $I(x) := -\log[PX(x)] = \log\left(\frac{1}{PX(x)}\right)$ تعریف شده است. به‌طور گسترده یک واقعه E داده می‌شود با احتمال P ، محتوای اطلاعات به‌طور مشابه تعریف می‌شود:

$$I(E) := -\log[Pr(E)] = -\log(P) \quad (1.2)$$

به‌طور کلی، پایه لگاریتمی انتخابی برای بسیاری از خصوصیات نظری اطلاعات اهمیت ندارد. با این حال، واحدهای مختلف اطلاعات بر اساس گزینه‌های محبوب پایه تعیین می‌شوند. اگر پایه لگاریتمی ۲ باشد، این واحد به نام شانون نامیده می‌شود اما از بیت نیز استفاده می‌شود. اگر پایه لگاریتم طبیعی باشد، واحدی که به آن "nat" کوتاه گفته می‌شود. اگر قرار باشد که لگاریتم پایه ۱۰ باشد، به واحدها یا نمودارهای رقمی یا رقم اعشار گفته می‌شود. آنتروپی شانون از متغیر تصادفی x این‌گونه تعریف شده است.

$$H(X) = \sum_x -PX(x) \log PX(x) = \sum_x PX(x) I(x) \stackrel{\text{def}}{=} E[I(X)] \quad (2.2)$$

برابر است با تعریف محتوای اطلاعات مورد انتظار اندازه‌گیری x [۴۳]؛ [۴۴].

۳.۴.۲ ویکی‌پدیا^۱

ویکی‌پدیا به بزرگ‌ترین دایره‌المعارف ارائه امکانات مانند چندزبانه، پوشش گسترده و به‌روز تبدیل شده است و به‌صورت رایگان قابل دسترس همگان می‌باشد. واحد اصلی در ویکی‌پدیا مقاله است که

¹Wikipedia

به صورت متن آزاد نوشته شده است که به دنبال یک مجموعه جامع از دستورالعمل‌های سرمقاله و ساختاری جهت ارتقاء سازگاری و انسجام است. هر مقاله یک مفهوم واحد را توصیف می‌کند و یک مقاله واحد برای هر مفهوم وجود دارد. عنوان هر مقاله عبارت مختصری است که در آن اصطلاح‌نامه متداول استفاده شده است و همچنین عنوان مفهوم هم نامیده می‌شود. مقالات با مقالات دیگر پیوند دارند. هایپر لینک، لینک همسایه نامیده می‌شود. لینک‌های ویکی‌پدیا به عنوان یک منبع اضافی مترادف که با تغییر مسیر هدایت می‌شود مفید نیستند. مقالات ذکر شده از دیگر ورودی‌های دانشنامه، به آن‌ها از طریق لینک‌های داخلی اشاره دارد، که این مقالات مرجع هستند [۲۳].

در ویکی‌پدیا، اقدامات مبتنی بر مسیر و اقدامات مبتنی بر مفهوم اطلاعات از ساختار طبقه‌بندی برای ارزیابی شباهت معنایی استفاده می‌کنند. ساختار طبقه‌بندی به نام نمودار دسته‌بندی ویکی‌پدیا^۱ شامل حدود ۱۵۰۰۰۰۰ مفهوم است که بسیار کمتر از تعداد مفاهیم ویکی‌پدیا است. از خصوصیات سنتی مانند نحو، واژه‌نامه‌ها و همسایگان استفاده می‌شود که اقدامات مبتنی بر ویژگی‌ها به طور جهانی برای ارزیابی شباهت معنایی بین مفاهیم اتخاذ می‌شوند. ویکی‌پدیا، یک دانشنامه مشترک آنلاین است که می‌تواند به عنوان یک دانش خارجی برای تجزیه و تحلیل متن کوتاه استفاده شود. ویکی‌پدیا ساختار پیوندی متراکم و پوشش گسترده‌ای از اشخاص از جمله اشخاص دارای نام، اشخاص خاص دامنه و اشخاص نوظهور دارد. علاوه بر این، داده‌های زباله ویکی‌پدیا را می‌توان آزادانه از وب به دست آورد. این مزایا محققان و توسعه‌دهندگان را ترغیب کرده است که از آن برای برنامه‌های کاربردی استفاده کنند [۴۵].

۱.۳.۴.۲ تحلیل معنایی صریح مبتنی بر ویکی‌پدیا^۲

تحلیل معنایی صریح مبتنی بر ویکی‌پدیا روشی گسترده است که برای اندازه‌گیری شباهت معنایی بین متون با هر طول استفاده می‌شود. تحلیل معنایی صریح برداری از اشخاص مرتبط با ویکی‌پدیا را برای متن معین به عنوان بیان معنایی آن ایجاد می‌کند و از وکتور برای اندازه‌گیری شباهت استفاده می‌کند. یافتن واحدهای مرتبط با ویکی‌پدیا از یک متن به طور کلی شامل چندین زیر مشخص مانند استخراج اصطلاحات کلیدی، یافتن نهادهای مرتبط برای هر اصطلاح کلیدی و تجمع وزن اشخاص مرتبط است. به منظور حل مشکل، تحلیل معنایی صریح بردارهای وزنی اشخاص مرتبط با هر کلمه را فقط بر اساس قانون اکثریت خلاصه می‌کند. رویکرد جمع‌بندی بردارها برای متون کوتاه در دنیای واقعی مناسب نیست که در آن اصطلاحات کلیدی و نامربوط در تعداد بسیار کمی بار مشاهده می‌شوند. به دلیل عدم اطلاع کافی، قانون اکثریت به خوبی کار نمی‌کند. در چنین شرایطی، تمرکز روی اصطلاحات کلیدی در ضمن فیلتر کردن اصطلاحات پرسروصدا اهمیت دارد [۴۵].

^۱Wikipedia kator Graph (WCG)

^۲Explicit semantic analysis

۴.۴.۲ مدل N-گرم^۱

مدل N-گرم نوعی مدل زبانی احتمالاتی است که برای پیش‌بینی قلم بعدی در دنباله‌ای به شکل مرتبه‌ی (n-1) مارکوف به کار می‌رود. امروزه مدل‌های N-گرم در احتمالات، نظریه‌ی ارتباطات، زبان‌شناسی رایانشی (مثلاً پردازش آماری زبان طبیعی) و فشرده‌سازی داده‌ها کاربرد گسترده دارند. یک N-گرم دنباله‌ای از n کلمه است که به صورت متوالی در یک متن ظاهر می‌شوند. به طور معمول جملات با موضوعات یکسان دارای N-گرم‌های مشابه هستند. به N-گرم با اندازه‌های یک، دو، و سه به ترتیب تک گرم، دو گرم، و سه گرم گفته می‌شود و نام‌گذاری به همین منوال ادامه می‌یابد [۴۶].

۵.۴.۲ طولانی‌ترین پیامد مشترک^۲

طولانی‌ترین پیامد مشترک طول طولانی‌ترین پیامد را در هر دو جمله می‌یابد. پیامد دنباله‌ای است که به همان ترتیب نسبی ظاهر می‌شود، اما لزوماً متناقض نیست. به عنوان مثال، bdf، abg، abc، acefg، aeg، و غیره پیامدهای abcdefg هستند. توجه کنید که LCS لزوماً منحصر به فرد نیست؛ به عنوان مثال LCS، ABC و ACB هم AB و AC هستند. در واقع، مسئله LCS اغلب برای یافتن کلیه پیامدهای مشترک با حداکثر طول تعریف می‌شود [۴۷].

الگوریتم طولانی‌ترین پیامد مشترک برای یافتن طولانی‌ترین بستر یک‌رشته استفاده می‌شود. این دو رشته را باهم مقایسه می‌کند و بر اساس طولانی‌ترین زنجیره مشترک کاراکترها، شباهت را می‌یابد. این می‌تواند به شرح زیر اندازه‌گیری شود:

$$LCS_{\text{Substr}}(S_1, S_2) = LCS_{\text{Suff}}(S_1 \dots i, S_2 \dots j), \max, 1 < i < m, 1 < j < n \quad (3.2)$$

که در آن m و n طول دو رشته است و LCS_{Suff} تابعی است که طولانی‌ترین پسوندهای رایج از پیشوندهای ممکن S1 و S2 را می‌یابد [۴].

۶.۴.۲ فرکانس مدت - فرکانس اسناد معکوس^۳

فرکانس مدت - فرکانس اسناد معکوس در بازیابی اطلاعات، یک آمار عددی است که میزان اهمیت یک کلمه نسبت به یک سند در یک مجموعه‌ای از اسناد را نشان می‌دهد. در واقع هدف این سیستم وزن دهی، نشان دادن اهمیت کلمه در متن است. در واقع Tf-idf روشی است که از طریق آن وزن لگاریتمی کلمه کلیدی به دست می‌آید و تمرکز اصلی آن‌ها صرفاً روی تعداد دفعات تکرار کلمه کلیدی است و بنا بر تنظیمات آن‌ها، هر زمان که تعداد تکرار کلمه کلیدی به صورت درصدی قابل قبول از کل متن ظاهر شود، آن متن را مناسب و هدفمند در نظر می‌گیرند در حالی که این امر اشتباه و یا

¹N-gram²Longest Common Substring(LCS)³term frequency - inverse document frequency (Tf-idf)

حداقل روش ناقص است. Tf-idf از دو عبارت Tf به معنای محدوده فرکانس نرمال و idf که فرکانس محتوایی معکوس نامیده می‌شود تشکیل شده است. مقدار Tf-idf به تناسب تعداد تکرار کلمه در سند افزایش می‌یابد و توسط تعداد اسنادی که در مجموعه هستند و شامل کلمه نیز می‌باشند متعادل می‌شود. به این معنی که اگر کلمه‌ای در بسیاری از متون ظاهر شود احتمالاً کلمه‌ای متداول است و ارزش چندانی در ارزیابی متن ندارد.

برای به دست آوردن ضریب Tf-idf می‌بایست هرکدام از این دو عبارت را به صورت جداگانه محاسبه نموده و حاصل دو عبارت را در هم ضرب کنیم تا نتیجه حاصله، فراوانی وزنی کلمه کلیدی را به ما نشان دهد. فرمول دو عبارت فوق به شرح زیر می‌باشد:

$$Tf = \text{عبارت است از تقسیم تعداد تکرار کلمه بر تعداد کل کلمات محتوا}$$

idf = عبارت است از لگاریتم تقسیم تعداد کل محتوا بر محتواهایی که شامل کلمه مورد نظر هستند

علت استفاده از لگاریتم در این فرمول این است که قصد داریم درون کل محتواهای سایت نرمال‌سازی انجام دهیم.

$$W_{x,y} = tf_{x,y} \times \log \frac{N}{df_x} \quad (۴.۲)$$

۷.۴.۲ چارچوب کلمه به بردار^۱

یک روش آماری برای یادگیری کارآمد یک کلمه مستقل تعبیه شده از یک متن است که توسط توماس میکولو و همکاران توسعه یافت. در سال ۲۰۱۳ در گوگل به عنوان پاسخی برای آموزش کارآمد شبکه آموزش مبتنی بر شبکه عصبی و از آن زمان به عنوان استاندارد بالفعل برای توسعه تعبیه کلمه از قبل آموزش دیده مورد استفاده قرار گرفته است. علاوه بر این، کار شامل تجزیه و تحلیل بردارهای آموخته شده و اکتشاف ریاضی بردار در بازنمایی کلمات است.

ما می‌دانیم که این بازنمایی‌ها به طور شگفت‌آور در ضبط منظم نحوی و معنایی در زبان خوب هستند و این که هر رابطه با یک جبهه بردار خاص رابطه مشخص می‌شود که این اجازه را می‌دهد تا استدلال بردار محور مبتنی بر جبران خسارت بین کلمات نتیجه دهد.

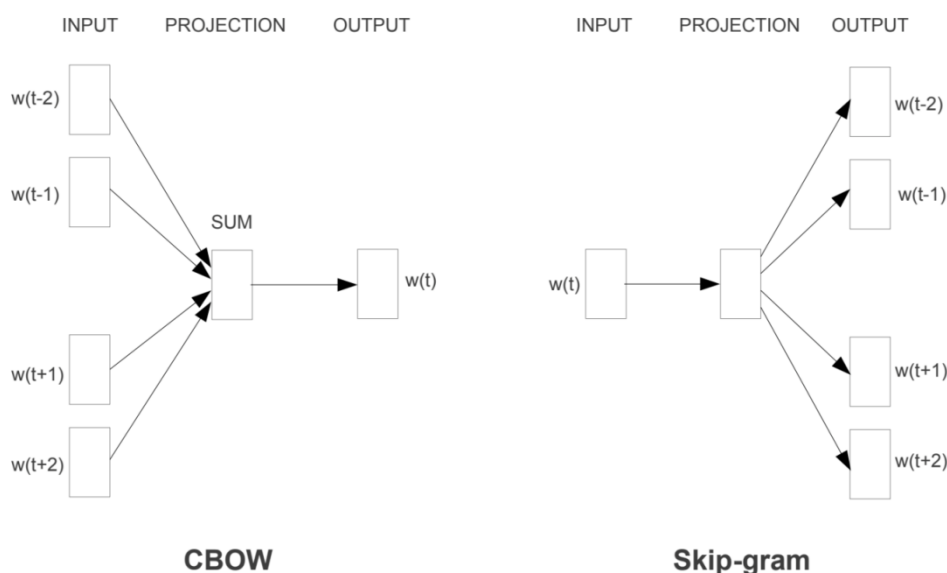
دو مدل یادگیری مختلف در شکل ۱.۲ ارائه شده است که می‌تواند به عنوان بخشی از رویکرد word2vec برای یادگیری کلمه تعبیه استفاده شود.

- بسته‌های کلمات مداوم یا مدل CBOW
- مدل پرس-گرم مداوم^۲

^۱ word2vec

^۲ skip-gram

مدل CBOW تعبیه را با پیش‌بینی کلمه فعلی بر اساس کلمات متن می‌آموزد. مدل پرش-گرم مداوم به پیش‌بینی کلمات همسایه با کمک کلمه فعلی می‌پردازد. مزیت اصلی این رویکرد این است که تعبیه‌های کلمه‌ای با کیفیت بالا می‌توانند به‌طور کارآمد آموخته شوند (فضای کم و پیچیدگی زمانی)، همچنین اجازه می‌دهد تا تعبیه‌های بزرگ‌تر (ابعاد بیشتر) از شرکت‌های بسیار بزرگ‌تر متن (میلیاردها کلمه) آموخته شود.



شکل ۱.۲: برآورد کارآمد از بازنمایی کلمات در فضای بردار

مدل CBOW مدل دیگری است که برای یادگیری بردارهای کلمات مورد استفاده قرار می‌گیرد. این مدل کلمه را از بین کلمات محتوایی منبع پیش‌بینی می‌کند. از آنجایی که چند کلمه محتوایی وجود دارد، بردارهای کلمه متناظر با هر یک را که از طریق ضرب بردار ورودی با ماتریس حاصل‌شده‌اند، میانگین گرفته‌اند. از اینجایی که مرحله میانگین‌گیری بخش زیادی از اطلاعات توزیعی را مسطح می‌کند، بعضی افراد اعتقاد دارند مدل CBOW برای دیتاست‌های کوچک مناسب‌تر است.

۸.۴.۲ مدل کیسه‌ای از کلمات^۱

مدل کیسه‌ای از کلمات نمایشی ساده است که در پردازش زبان طبیعی و بازیابی اطلاعات مورد استفاده قرار می‌گیرد. در این مدل، متنی (مانند یک جمله یا یک سند) به‌عنوان کیسه (چندصفحه‌ای) از کلمات آن نمایش داده می‌شود، بی‌توجه به دستور زبان و حتی دستور کلمه اما با حفظ فراوانی، مدل کیف‌دستی کلمات همچنین برای دید رایانه‌ای استفاده شده است [۴۸].

مدل کیف‌واژه‌ها معمولاً در روش‌های طبقه‌بندی اسناد استفاده می‌شود که در آن فرکانس وقوع

^۱Bag of words

هر کلمه به‌عنوان ویژگی‌ای برای آموزش طبقه‌بندی استفاده می‌شود. در عمل، کیسه‌ای از کلمات عمدتاً به‌عنوان ابزاری برای تولید ویژگی‌ها مورد استفاده قرار می‌گیرد. پس از تبدیل متن به کیسه‌ای از کلمات، می‌توان اقدامات مختلفی را برای توصیف متن محاسبه کرد. متداول‌ترین نوع ویژگی‌ها یا ویژگی‌های محاسبه‌شده از مدل کیسه‌ای از کلمات، فرکانس مدت است، یعنی تعداد دفعاتی که یک اصطلاح در متن ظاهر می‌شود. این نمایه لیست (بردار) ترتیب کلمات موجود در جملات اصلی را حفظ نمی‌کند. این فقط ویژگی اصلی مدل کیسه‌ای از کلمات است. همچنین به‌عنوان مدل فضای برداری شناخته می‌شود. مدل بسته کلمات برای بینایی کامپیوتر مدل بسته کلمات معمولاً در روش‌های دسته‌بندی اسناد مورد استفاده قرار می‌گیرد که در آن وقوع هر کلمه به‌عنوان یک ویژگی برای آموزش طبقه‌بندی آماری استفاده می‌شود.

۹.۴.۲ تخصیص پنهان دیریکلت^۱

تخصیص پنهان دیریکلت، یک مدل تولیدی بدون نظارت است که موضوع را به اسناد اختصاص می‌دهد. کلمات موجود در هر سند به تشخیص موضوع کمک می‌کنند. موضوعات ممکن است پیش‌بینی نشوند، اما تعداد موضوعات باید مشخص شود. ممکن است کلمات بین موضوعات هم‌پوشانی داشته باشند، بنابراین چندین موضوع ممکن است کلمات مشابهی داشته باشند. از کاربردهای LDA به موارد زیر می‌توان اشاره کرده درک موضوعات مختلف در یک مجموعه، شناخت بهتر از نوع اسناد کمی کردن بیشترین کاربرد و مهم‌ترین کلمات در یک مجموعه شباهت اسناد. این موارد مدل متغیرهای نهفته را ایجاد می‌کند:

۱. توزیع موضوعات برای هر سند

۲. توزیع کلمات برای هر موضوع

برخلاف الگوریتم word2vec، که با جملات کاملاً ساختاریافته عملکرد بسیار خوبی دارد، LDA یک کلمه مدل است، به معنی ترتیب کلمات در یک سند، شمارش نمی‌شود. این همچنین بدان معنی است که کلمات متوقف و کلمات نادر باید کنار گذاشته شوند، به طوری که مدل برای کلمات بسیار مکرر و کلمات بسیار کمیاب، که هر دو به موضوعات کلی کمک نمی‌کنند، بیش‌ازحد جبران نمی‌کند.

LDA فرض می‌کند که هر سند توسط یک فرآیند آماری تولید می‌شود. یعنی هر سندی ترکیبی از موضوعات است و هر موضوع ترکیبی از کلمات است.

۱۰.۴.۲ فاصله جنسن شانون^۲

در تئوری و آمار احتمال، واگرایی جنسن شانون روشی برای اندازه‌گیری شباهت بین دو توزیع احتمال است. همچنین به شعاع اطلاعات یا واگرایی کل نسبت به میانگین معروف است. این

^۱Latent Dirichlet Allocation (LDA)

^۲Jensen-Shannon distance

تفاوت بر اساس واگرایی Kullback – Leibler است، با برخی تفاوت‌های قابل توجه از جمله اینکه متقارن است و همیشه دارای ارزش محدود است و برخلاف واگرایی Kullback، این یک معیار را به روش مستقیم تعیین می‌کند. ریشه مربع واگرایی جنسن شانون معیاری است که اغلب به‌عنوان فاصله جنسن و شانون شناخته می‌شود. با استفاده از ماشین‌آلات آزمون‌های آماری، روش‌های مناسبی برای تصمیم‌گیری در مورد مشابه بودن دو سند داریم. نوع آزمون قابل‌استفاده برای شباهت آزمون در ادبیات آمار به‌عنوان آزمون همگنی شناخته شده است.

۱۱.۴.۲ ماشین بردار پشتیبان^۱

ماشین بردار پشتیبان در آزمایشگاه‌های A&T Bell توسط Vapnik و همکارانش توسعه یافته است و مبتنی بر چارچوب‌های آماری است. ماشین بردار پشتیبان یکی از روش‌های یادگیری بانظارت است که از آن برای طبقه‌بندی و رگرسیون استفاده می‌کنند. مبنای کاری دسته‌بندی کننده‌ی ماشین بردار پشتیبان دسته‌بندی خطی داده‌ها است و در تقسیم خطی داده‌ها سعی می‌کنیم خطی را انتخاب کنیم که حاشیه اطمینان بیشتری داشته باشد. حل معادله پیدا کردن خط بهینه برای داده‌ها به‌وسیله روش‌های QP که روش‌های شناخته شده‌ای در حل مسائل محدودیت دار هستند صورت می‌گیرد. قبل از تقسیم خطی برای اینکه ماشین بتواند داده‌های با پیچیدگی بالا را دسته‌بندی کند داده‌ها را به‌وسیله‌ی تابع phi به فضای با ابعاد خیلی بالاتر می‌بریم [۴۹]. برای اینکه بتوانیم مسئله ابعاد خیلی بالا را با استفاده از این روش‌ها حل کنیم از قضیه دوگانگی لاگرانژ برای تبدیل مسئله‌ی مینیمم سازی موردنظر به فرم دوگانگی آن که در آن به‌جای تابع پیچیده‌ی phi که ما را به فضایی با ابعاد بالا می‌برد، تابع ساده‌تری به نام تابع هسته که ضرب برداری تابع phi است ظاهر می‌شود استفاده می‌کنیم.

در یادگیری ماشین، ماشین‌های بردار پشتیبان مدل‌های یادگیری با الگوریتم‌های یادگیری مرتبط هستند که داده‌های مورد استفاده برای طبقه‌بندی و تحلیل رگرسیون را تحلیل می‌کنند. با توجه به مجموعه‌ای از مثال‌های آموزشی، که هر یک به دو دسته دیگر تعلق دارند، یک الگوریتم آموزش ماشین بردار پشتیبان مدلی را ایجاد می‌کند که نمونه‌های جدید را به یک دسته یا گروه دیگر اختصاص می‌دهد و آن را به‌عنوان طبقه‌بندی کننده خطی باینری غیرممکن تبدیل می‌کند. یک مدل ماشین بردار پشتیبان نمایانگر نمونه‌ها به‌عنوان نقاط در فضا است که نقشه برداری شده است به طوری که نمونه‌های دسته‌های جداگانه با یک شکاف واضح که تا حد ممکن گسترده است تقسیم می‌شود. نمونه‌های جدید سپس در همان فضا ترسیم می‌شوند و پیش‌بینی می‌شود که بر اساس طرف شکافی که روی آن قرار می‌گیرند، به یک دسته تعلق دارند.

¹Support vector machines(SVM)

۱۲.۴.۲ تجزیه و تحلیل آماری^۱

تجزیه‌کننده زبان طبیعی برنامه‌ای است که ساختار گرامری جملات را انجام می‌دهد، به‌عنوان مثال کدام گروه از کلمات باهم جمع می‌شوند (به‌عنوان عبارات) و کدام کلمات موضوع یا هدف یک فعل هستند. تجزیه‌کنندگان احتمالی از دانش زبان به‌دست‌آمده از جملات تجزیه‌شده دستی استفاده می‌کنند تا سعی کنند محتمل‌ترین تحلیل جملات جدید را تولید کنند. این تجزیه و تحلیل‌های آماری هنوز هم برخی از اشتباهات را انجام می‌دهند، اما معمولاً بسیار خوب عمل می‌کنند. پیشرفت آن‌ها یکی از بزرگ‌ترین پیشرفت در پردازش زبان طبیعی در دهه ۱۹۹۰ بود. آن‌ها طراحی شده‌اند تا توسط افرادی که می‌خواهند روابط متنی را استخراج کنند، به‌راحتی درک و به کار گرفته شوند.

۱۳.۴.۲ حافظه کوتاه مدت طولانی^۲

حافظه کوتاه‌مدت طولانی یک معماری شبکه عصبی مکرر مصنوعی^۳ است که در زمینه یادگیری عمیق مورد استفاده قرار می‌گیرد. حافظه کوتاه‌مدت طولانی برخلاف شبکه‌های عصبی استاندارد تغذیه‌شده، اتصالات بازخورد دارد. این نه تنها می‌تواند نقاط داده‌های واحد (مانند تصاویر) را پردازش کند، بلکه توالی‌های کل داده‌ها (مانند گفتار یا فیلم) را نیز پردازش می‌کند. به‌عنوان مثال، حافظه کوتاه‌مدت طولانی برای کارهایی از قبیل تشخیص نام‌نویسی دست‌نویس، متصل به آن، تشخیص گفتار و تشخیص ناهنجاری در ترافیک شبکه یا سامانه‌های تشخیص نفوذ^۴ کاربرد دارد.

یک واحد مشترک حافظه کوتاه‌مدت طولانی از یک سلول، یک دروازه ورودی، یک دروازه خروجی و یک دروازه فراموشی تشکیل شده است. سلول مقادیر را در فواصل زمانی دلخواه به یاد می‌آورد و سه دروازه جریان اطلاعات به داخل و خارج از سلول را تنظیم می‌کند.

شبکه‌های حافظه کوتاه‌مدت طولانی به‌خوبی برای طبقه‌بندی، پردازش و پیش‌بینی‌های مبتنی بر داده‌های سری زمانی مناسب هستند، زیرا می‌توان مدت‌زمان ناشناخته‌ای را بین وقایع مهم در یک سری زمانی ایجاد کرد. حافظه کوتاه‌مدت طولانی برای مقابله با مشکلات شیب منفجر و ناپدیدشده که می‌توانند هنگام آموزش با RNN سنتی روبرو شوند، ایجاد شده‌اند. عدم حساسیت نسبی نسبت به طول فاصله، مزیت LSTM نسبت به RNN، مدل‌های پنهان مارکوف و سایر روش‌های یادگیری توالی در کاربردهای بی‌شمار است. چندین واحد معماری از واحدهای LSTM وجود دارد. یک معماری معمول از سلول (قسمت حافظه واحد LSTM) و سه تنظیم‌کننده تشکیل شده است، که معمولاً به‌عنوان دروازه نامیده می‌شوند، جریان اطلاعات در داخل واحد LSTM یک دروازه ورودی، یک دروازه خروجی و یک دروازه فراموشی. برخی از تغییرات واحد LSTM دارای یک یا چند مورد از این دروازه‌ها نیستند یا ممکن است دارای دروازه‌های دیگری باشند. به‌عنوان مثال، واحدهای

¹Statistical Analysis

²Long short-term memory(Lstm)

³Recurrent neural network(RNN)

⁴ Intrusion detection systems(IDS)

عودکننده درج شده (GRU) دروازه خروجی ندارند.

به طور شهودی، سلول مسئول پیگیری وابستگی بین عناصر در دنباله ورودی است. دروازه ورودی میزان ورود مقدار جدید به داخل سلول را کنترل می کند، دروازه فراموشی تا چه حد مقداری در سلول باقی می ماند و دروازه خروجی میزان استفاده از مقدار سلول را برای محاسبه خروجی کنترل می کند. عملکرد فعال سازی دروازه LSTM اغلب عملکرد سیگموئید لجستیک است. اتصالات داخل و خارج از دروازه LSTM وجود دارد، که تعداد کمی از آن ها در حالت عود است. وزن این اتصالات، که باید در طول آموزش یاد بگیرند، نحوه عملکرد دروازه ها را تعیین می کند.

۵.۲ پیشینه تحقیق

شباهت معنایی مبتنی بر معنا که به معنی جملات بستگی دارد، رویکردهای مختلفی دارد. این روش ها فن های مختلفی را برای مقایسه دو جمله از نظر معنایی به کار می گیرند. رویکرد اول، مبتنی بر پیکره است که شباهت کلمات را بر اساس تجزیه و تحلیل آماری پیکره می یابد. همچنین، از یادگیری عمیق برای تجزیه و تحلیل جملات استفاده می کند. رویکرد دوم، مبتنی بر دانش است که معنی کلمات و روابط بین کلمات در یک شبکه معنایی گنجانده شده است. رویکرد سوم مبتنی بر ساختار می باشد که از اطلاعات ساختار یک جمله برای به دست آوردن معنی آن جمله استفاده می کند.

۱.۵.۲ روش های مبتنی بر پیکره

دو روش مختلف برای تجزیه و تحلیل آماری یک متن وجود دارد. مورد اول استفاده از آنالیز آماری عادی مانند LSA و دوم استفاده از یادگیری عمیق است. در وهله اول، یک قشر بزرگ از طریق شمارش کلمات در پیکره و اسناد به صورت آماری تحلیل می شود. محاسبه Tf-idf که به منظور وزن دادن به کلمه استفاده می شود، یک هدف مهم در تجزیه و تحلیل بدنه است. علاوه بر این، استفاده از یادگیری عمیق نتایج خوبی را در نمایش کلمات از نظر معنایی نشان داده است [۱].

۱.۱.۵.۲ روش های آماری

شباهت آماری آموخته شده از داده ها یعنی مجموعه نوشته ها، که مجموعه ای از متن نوشتاری یا گفتاری است. در این روش ابتدا یک مدل آماری ساخته شده و سپس شباهت تخمین زده می شود. در ارزیابی [۵۰] چترجی (۲۰۰۱) نتایج اولیه یک طرح اندازه گیری شباهت را که مبتنی بر یک مدل خطی است و ضرایب آن با استفاده از فن رگرسیون چندگانه تعیین شده است، گزارش می دهد. داده های تجزیه و تحلیل از نظر سنجی تعدادی از پاسخ دهندگان جمع آوری شده است. سه جنبه عمده شباهت، یعنی عمل گرا، نحوی و معنایی در نظر گرفته شده است. از هر مخاطب خواسته شده است تا شباهت بین جفت های مختلف جملاتی را که با دقت طراحی شده اند تا یکی از انواع فوق شباهت را

ارزیابی کند. تجزیه و تحلیل آماری از این ارزیابی‌ها، درک عمومی انسان را در مورد شباهت اصولی نشان می‌دهد، که به طراحی یک برنامه ارزیابی مناسب کمک می‌کند. ترجمه ماشینی مبتنی بر مثال^۱ مبتنی بر ایده انجام ترجمه با تقلید از مثال‌های ترجمه گذشته است. در این نوع سیستم ترجمه، مقدار زیادی از نمونه‌های ترجمه بین دو زبان در یک پایگاه داده متنی ذخیره می‌شوند. ترجمه جملات ارزیابی شده پس از آن تغییر یافته برای به دست آوردن ترجمه جمله ورودی داده شده مناسب است. بدیهی است، این طرح بستگی به این دارد که طرح ارزیابی چقدر خوب و مؤثر است. هرچه جمله ارزیابی شده به آن نزدیک‌تر باشد، تطبیق آن با نیاز ترجمه فعلی آسان‌تر خواهد بود و به تبع آن، کیفیت کلی ترجمه بهبود می‌یابد. اولین کار در جهت هدف فوق، یافتن مؤلفه‌های کلیدی است که در اندازه‌گیری شباهت بین جملات نقش دارند. یک لیست جامع باید تعیین شود. برای هر یک از این مؤلفه‌ها مقادیر ممکن مشخص می‌شوند و مسافت نسبی آن‌ها تعیین می‌شود.

در ارزیابی [۴۵] شیراکاوا و همکاران (۲۰۱۵) یک روش اندازه‌گیری شباهت معنایی مبتنی بر ویکی‌پدیا ارائه داده‌اند که برای متون کوتاه در دنیای واقعی در نظر گرفته شده است. این مقاله روشی را ارائه می‌دهد که نوعی تحلیل معنایی صریح (ESA) است که یک کیف از اشخاص ویکی‌پدیا (صفحات ویکی‌پدیا) را به عنوان متن معنایی آن به متن اضافه می‌کند و از وکتور موجودات برای محاسبه شباهت معنایی استفاده می‌کند. اضافه کردن نهادهای مرتبط به یک متن، نه یک کلمه یا عبارت، یک مشکل عملی چالش برانگیز است زیرا معمولاً از چندین زیر مشخص تشکیل می‌شود، روش پیشنهادی در این کار با استفاده از نوی بی‌زین، یک مکانیزم وزن احتمالی مبتنی بر قضیه بی‌ز، این مشکل جمع را حل می‌کند. همچنین در این مقاله به این نتیجه رسیده است که کاهش ابعاد بردار به اشخاص نماینده ویکی‌پدیا به ندرت عملکرد را تحت تأثیر قرار می‌دهد در حالی که اندازه بردار را کاهش می‌دهد و از این رو فضای ذخیره‌سازی و زمان پردازش محاسبه شباهت کسینوس را نشان می‌دهد.

هدف پینه‌پرو و همکاران (۲۰۱۷) در [۵۱] ارائه یک رویکرد جدید برای اندازه‌گیری شباهت بین جملات نوشته شده به زبان پرتغالی با استفاده از ویژگی‌های آماری و یادگیری عمیق برای غلبه بر مشکلات معنا بوده است. این مسئله با استفاده از یک رویکرد یادگیری عمیق برای اندازه‌گیری شباهت بین کلمات، بر مشکل معنا غلبه می‌کند. ایده اصلی این است که چهار ویژگی مختلف را بر اساس تطبیق مستقیم^۲ Tf-idf، چارچوب کلمات به بردار و اندازه جملات ترکیب کند تا شباهت نهایی را به دست آورد. در واقع این چهار ویژگی مختلف را برای اندازه‌گیری شباهت جملات برای پرتغالی‌ها ترکیب می‌کند، همچنین از چارچوب کلمات به بردار و یک روش مبتنی بر ماتریس برای مقابله با مشکل معنی استفاده می‌کند. علاوه بر این، از یک سیستم رگرسیون برای به دست آوردن مقدار تشابه با ترکیب ویژگی‌های کار شباهت معنایی متن و طبقه‌بندی‌های مختلف یعنی ماشین بردار پشتیبان، نوی بی‌زین و شبکه عصبی به کار RTE استفاده شده است. این الگوریتم ویژگی‌های استخراج شده را به عنوان ورودی با شباهت پیشنهادی در مجموعه آموزش دریافت می‌کند و خروجی

^۱ Example-Based Machine Translation (EBMT)

^۲ Direct Matching

تابعی است برای ترکیب این ویژگی‌ها که باعث ایجاد شباهت نهایی در مجموعه آزمون می‌شود. برای تعیین کمیت شباهت نهایی از یک الگوریتم رگرسیون استفاده شده است.

در [۵۲] ماجومدر و همکاران (۲۰۱۸) روشی ارائه داده‌اند که قادر خواهد بود میزان تفسیر بین دو جمله را با درجه‌بندی میزان هم‌ارزی معنایی در مقیاس درجه‌بندی شده از ۰ (تراز نشده) به ۵ (معادل معنایی) اندازه‌گیری کند. هدف مقاله اضافه کردن یک لایه قابل تفسیر نیست بلکه روشی را تدوین کرده است که می‌تواند شباهت‌ها و تفاوت‌های بین این دو جمله را توضیح دهد. آن‌ها برای داده‌های استاندارد به دقت ۰/۶۴ برای نوع و نمره تراز رسیده‌اند.

برای توسعه این روش ابتدا تکه‌ها / عبارت را در بین جملات شناسایی کرده، سپس آن‌ها را در بین جملات قرار داده است. در طول تراز عبارت روابط میان آن‌ها را مشخص کرده و پس از آن بر اساس رابطه، نمره شباهت بین عبارت اندازه‌گیری می‌شود. مسئله تشابه متن معنایی قابل تفسیر، به چهار وظیفه فرعی تقسیم شده است. اولین کار تولید درختان پارس برای جملات است (یعنی تکه‌ها را مشخص کرده است). وظیفه بعدی ترسیم تکه‌های جملات است و سوم اختصاص دادن استدلال تراز از فهرستی از برچسب‌های احتمالی. سرانجام، نمره همبستگی / شباهت بین عبارت تراز وسط قرار می‌گیرد. ابتدا جملات ورودی را از قبل پردازش کرده و سپس تکه‌های جملات مشخص شده‌اند. علاوه بر این، تکه‌ها وارد روند تراز شده‌اند، که می‌تواند آن‌ها را بر اساس نمره مرتبط بودن تنظیم کند. سرانجام، نمره کلی اندازه‌گیری شده است و بر اساس آن یک برچسب رابطه اختصاص داده شده است.

در [۲۶] فراندز و همکاران (۲۰۱۹)، یک روش تعبیه‌کننده ارائه داده‌اند، که هدف آن یادگیری بازنمایی جملات بدون نظارت از متن بدون برچسب است. در حقیقت یک روش بدون نظارت را پیشنهاد کرده که یک جمله را به عنوان یک سری وزن دار تعبیه کلمات مدل می‌کند. وزن این سریال با استفاده از اطلاعات متقابل شانون^۱ در میان کلمات، جملات و پیکره متناسب است. در این بررسی به بازنمایی جمله‌ها پرداخته تا رفتار آماری معنای جمله را مورد مطالعه قرار دهد. آن‌ها با استفاده از داده‌های SICK به همبستگی ۰/۷۲۴ رسیده‌اند.

۲.۱.۵.۲ روش‌های مبتنی بر یادگیری عمیق

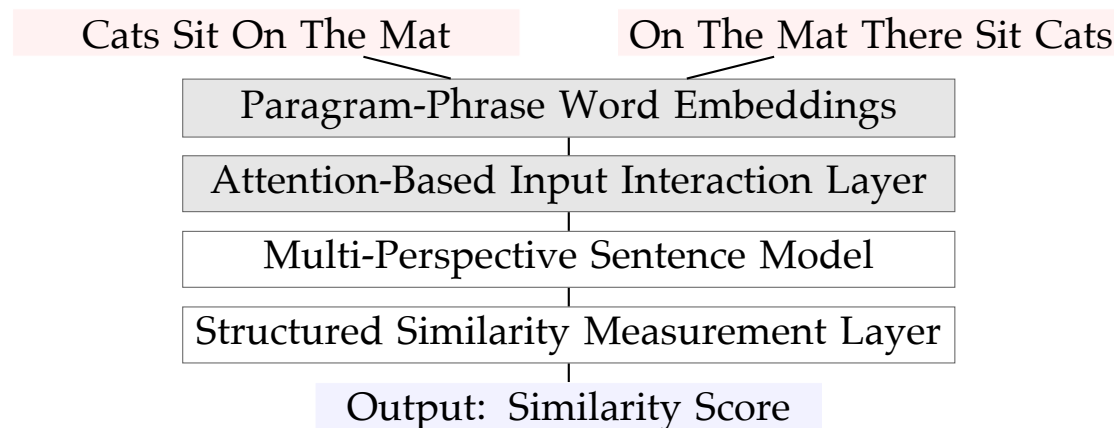
در [۵۳] هی و همکاران (۲۰۱۶) از یک شبکه عصبی کانولوشن مبتنی بر شباهت معنایی متن استفاده کرده‌اند که یک لایه تعامل ورودی مبتنی بر توجه را توسعه داده و آن را با شبکه عصبی کانولوشن چندرسانه‌ای که با استفاده از ضمایم کلمه پاراگرام-عبارت^۲ آموزش داده شده در جفت‌های ترجمه ترکیب کرده است. در اینجا از شبکه عصبی حلقوی چندمنظوره^۳ استفاده شده است که جملات ورودی را به صورت جداگانه تلقی می‌کند و عمدتاً از تعاملات حساس به متن چشم‌پوشی می‌کند. برای برطرف کردن این مشکل لایه تعامل ورودی مبتنی بر توجه را توسعه داده، دو جمله مستقل

¹Mutual Information(MI)

² PARAGRAM-PHRASE

³multi-perspective convolutional neural network(MPCNN)

ورودی را به یک جفت جمله مرتبط می‌کند که به مدل کمک می‌کند تا کلمات ورودی مهم را برای اندازه‌گیری بیشتر تشابه انتخاب کند و همچنین از سازگاری کلمات قوی پاراگرام- عبارت استفاده می‌کند. ساختار کلی مدل در شکل ۲.۲ نشان داده شده است.



شکل ۲.۲: نمای کلی مدل

مدل پایه: شبکه عصبی مصنوعی چندبعدی: از مدل چندبعدی شبکه عصبی کانولوشن استفاده کرده است که دو جزء اصلی دارد:

۱. یک مدل چندبعدی جمله برای تبدیل یک جمله به نمایندگی: یک شبکه عصبی کانولوشن اطلاعات مختلف در هر جمله را با استفاده از انواع مختلف فیلترهای کانولوتور^۱، انواع توزیع و اندازه پنجره‌ها را در برمی‌گیرد.

۲. یک لایه اندازه‌گیری شباهت ساختاریافته: از این لایه برای مقایسه مناطق محلی بازنویسی جمله‌ها از متریک‌های شباهت چندگانه استفاده می‌کند. در این بررسی با استفاده از مجموعه داده‌های آزمون رقابت STS انگلیسی SemEval 2016 که شامل پنج مجموعه داده از حوزه‌های مختلف می‌باشد به ضریب پیرسون ۰/۸۰۴ رسیده است.

مولر (۲۰۱۶) از یک شبکه سیامس که در واقع شبکه طولانی مدت را برای داده‌های برچسب گذاری شده از جفت توالی‌های متغیر طول ارائه داده استفاده کرده است. برای این برنامه‌های کاربردی، ارائه دهنده‌های تعبیه کلمات^۲ با اطلاعات مترادف به LSTM‌ها ارائه داده که از یک بردار اندازه ثابت برای رمزگذاری معنی اساسی بیان شده در جمله (بدون در نظر گرفتن اصطلاح خاص / نحو خاص) استفاده می‌کند. با محدود کردن عملیات بعدی با تکیه بر متریک ساده منهتن، مجازات جمله را که توسط این مدل آموخته شده است را مجبور کرده تا یک فضای بسیار ساختاریافته را شکل دهد که هندسه آن روابط معنایی پیچیده را نشان می‌دهد. در این بررسی با استفاده از مجموعه داده‌های SICK به دقت ۸۴/۲ و ضریب همبستگی پیرسون ۰/۸۸۲۲ دست‌یافته‌اند [۵۴].

در ارزیابی [۵۵] وایتینگ و همکاران (۲۰۱۶) شش معماری ترکیبی را باهم مقایسه می‌کنند،

¹ convolutional

² wordembedding

آن‌ها را در مجموعه داده‌های شباهت متنی حاشیه‌نویسی که هر دو از همان توزیع داده‌های آموزش و همچنین از طیف گسترده‌ای از حوزه‌های دیگر ترسیم شده است، ارزیابی می‌کند. پیچیده‌ترین معماری‌ها، مانند شبکه‌های عصبی مکرر حافظه کوتاه‌مدت، بهترین عملکرد را در داده‌های درون دامنه دارند. با این حال، در سناریوهای خارج از حوزه، معماری‌های ساده مانند کلمه با میانگین بسیار زیاد از LSTM‌ها بهتر است. ساده‌ترین مدل میانگین ما حتی با سامانه‌هایی که برای کارهای خاص تنظیم شده‌اند بسیار رقابتی است و در عین حال بسیار کارآمد و آسان برای استفاده است. در این بررسی بر روی سه کار تحت نظارت NLP آزمایش صورت گرفته است: شباهت جملات، بازپرداخت و طبقه‌بندی احساسات و با استفاده از مجموعه داده‌های SICK به ضریب همبستگی ۰/۷۰ درصد رسیده است.

زو و همکاران یک مدل D-LSTM مبتنی بر LSTM را پیشنهاد داده‌اند و بر روی داده‌های SICK به ضریب همبستگی پیرسون ۰/۸۳۰۵ و ضریب همبستگی اسپرمن ۰/۷۷۲۹ و خطای MSE، ۰/۳۴۴۲ دست‌یافته‌اند [۵۶].

پونتز و همکاران (۲۰۱۸) سامانه‌ای را برای اندازه‌گیری شباهت معنایی جملات، ترکیبی از شبکه عصبی راجعه و عودکننده را طراحی کردند. این سیستم از یک شبکه همخوانی برای در نظر گرفتن زمینه محلی کلمات و یک LSTM برای در نظر گرفتن زمینه جهانی جملات استفاده می‌کند. این ترکیبی از شبکه‌ها به حفظ اطلاعات مربوط به جملات کمک می‌کند و محاسبه شباهت بین جملات را بهبود می‌بخشد. در این بررسی از فن NN‌ها استفاده شده است. ابتدا، از CNN سیامس استفاده کرده تا متن محلی کلمات را در یک جمله تحلیل کند و بازنمایی از اهمیت یک کلمه و همسایگی آن ایجاد نماید. سپس، از LSTM سیامس استفاده کرده تا کل جمله را بر اساس کلمات و متن‌های محلی آنالیز کند. در آخر، شباهت معنایی جفت جملات را با استفاده از فاصله منهتن پیش‌بینی نموده است. در انتها به ضریب همبستگی پیرسون ۰/۸۵۴۹ و ضریب اسپرمن ۰/۷۹۳۳ و خطای ۰/۲۸۹۸ دست‌یافته است [۵۷].

در [۵۸] راناسینگه و همکاران از معماری‌های مکرر سیامس برای اندازه‌گیری شباهت معنایی متن استفاده کرده‌اند. در این روش شبکه‌ها دارای بردارهای نمایشی با همان معنای مشابه هستند و باعث مقایسه آسان‌تر آن‌ها می‌شوند. با توجه به اینکه وزن‌ها در شبکه‌های فرعی به اشتراک گذاشته می‌شود، پارامترهای کمتری برای آموزش وجود دارد که به نوبه خود به معنای نیاز به داده‌های آموزشی کمتری و تمایل به بیش‌ازحد مناسب بودن است. با توجه به میزان نیروی انسانی مورد نیاز برای تولید مجموعه داده برای STS شبکه‌های عصبی سیامی می‌توانند راه‌حل ایده آل برای کار شباهت معنایی متن باشند. راناسینگه و همکاران با کار بر روی مجموعه داده‌ی SICK به ضریب همبستگی ۰/۸۹۰۱ در روش GRU دست‌یافته‌اند. این روش از بقیه روش‌ها نتیجه بهتری داشته است.

۲.۵.۲ روش‌های مبتنی بر معنا

در [۵۹] لی و همکاران (۲۰۰۶) الگوریتمی را ارائه داده‌اند که اطلاعات معنایی و اطلاعات مربوط به ترتیب کلمه را در جملات ذکر کرده، سپس شباهت معنایی دو جمله با استفاده از اطلاعاتی که از یک پایگاه داده ساختاری واژگانی ساخت یافته به دست آورده و از آمار بدنه محاسبه شده است استفاده می‌کند. استفاده از دیتابیس واژگانی باعث می‌شود که از مفهوم دانش مشترک استفاده کند و درج آمار اطلاعات جسمی، باعث می‌شود تا با حوزه‌های مختلف سازگار باشد. متن به‌عنوان دنباله‌ای از کلمات در نظر گرفته می‌شود که هر یک از آن‌ها اطلاعات مفیدی دارند. کلمات به همراه ساختار ترکیبی آن‌ها باعث می‌شود که متن معنای خاصی را منتقل کند. این روش به‌صورت پویا یک کلمه مشترک را با استفاده از تمام کلمات مجزا در جفت جملات تشکیل می‌دهد. برای هر جمله، یک ویکتور معنایی خام با کمک یک بانک اطلاعاتی واژگانی به دست می‌آورد. برای هر جمله یک ویکتور مرتبه کلمه تشکیل می‌دهد و مجدداً از اطلاعات پایگاه داده واژگان استفاده می‌کند، از آنجا که هر کلمه در یک جمله متفاوت به معنای کل جمله است، اهمیت یک کلمه با استفاده از محتوای اطلاعاتی که از یک جسم به دست می‌آید، وزن می‌گیرد. با ترکیب بردار معنایی خام با محتوای اطلاعاتی از بدنه، یک ویکتور معنایی برای هر یک از دو جمله به دست می‌آید. شباهت معنایی بر اساس دو بردار معنایی محاسبه می‌شود. شباهت سفارش با استفاده از دو بردار مرتبه محاسبه می‌شود. سرانجام شباهت جملات با ترکیب شباهت معنایی و تشابه نظم حاصل می‌شود. اندازه‌گیری شباهت الگوریتم این بررسی ضریب همبستگی پیرسون به‌طور معقول خوب 0.816 با رتبه‌های انسانی، در سطح 0.01 معنی‌دار به دست آورد.

در [۴۰] امیر و همکاران (۲۰۱۶) روشی ارائه داده‌اند که روشی کارآمد برای تشخیص پاراگراف‌ها در متون است که با یک عبارت خاص مطابقت دارد. این مقاله از یک معیار شباهت بر اساس هسته‌های معنایی استفاده کرده است. ابتدا از Stanford Parser به‌عنوان یک ابزار NLP برای تفسیر جملات استفاده کرده است. سپس بر اساس روابط دستوری هسته‌های معنایی را تشکیل داده که شامل موضوعات سه‌گانه، فعل‌ها و اشیا است. در واقع هسته‌ها معنای عمومی جملات را بیان می‌کند که انعطاف‌پذیری خاصی برای مقابله با ساختارهای ناهمگن دارند و سرعت محاسبات شباهت معنایی را زیاد می‌کند. بر اساس هسته‌های استخراج شده، یک تجزیه و تحلیل برای استخراج اطلاعات مرتبط انجام شده است. برای انجام این کار از یک روش یادگیری ماشین برای یافتن بهترین راه برای جمع کردن ویژگی‌ها در یک مدل معنایی استفاده شده است. در این بررسی خطای استاندارد باقی‌مانده 0.16 تعیین شد.

ویکی‌پدیا یک مخزن دائره‌المعارف مستقل از دامنه و یک شبکه معنایی برای محاسبه شباهت معنایی بین مفاهیم را فراهم می‌کند. اقدامات مبتنی بر ویژگی‌های سنتی به ترکیب خطی از خصوصیات مختلف با دو محدودیت اصلی، اطلاعات ناکافی و از بین رفتن اطلاعات معنایی متکی است. در این بررسی، چندین روش اندازه‌گیری شباهت معنایی هیبریدی را با استفاده از محتوای اطلاعات و

¹Residual standard error(RSE)

ویژگی‌های مفاهیم ارائه داده است، که از محدودیت‌های معرفی شده در بالا جلوگیری می‌کنند. با توجه به ادغام خصوصیات گسسته در یک جزء، دو مدل بازنمایی معنایی به نام‌های CORM و CARM را ارائه داده، سپس، شباهت معنایی را بر اساس این مدل‌ها محاسبه نموده است و IC مقوله‌ها را به‌عنوان مکمل اندازه‌گیری شباهت معنایی در نظر گرفته است. به‌طور خلاصه، رویکردهای مورد استفاده در این بررسی در تعیین شباهت معنایی بین مفاهیم کارآمدتر بوده و از همبستگی انسانی بهتری نسبت به روش‌های قبلی مانند word2vec و NASARI برخوردار است. این روش ابتدا از همسایگان به‌عنوان یک ویژگی اضافی برای تقویت نمایش مفاهیم استفاده می‌کند. ثانیاً، چند جفت دسته مشابه از دو مفهوم را انتخاب کرده و برای ایجاد ویژگی‌های معنایی اضافی با کمک اقدامات مبتنی بر IC، دو زیر گراف را بنا کرده است. سوم، ویژگی‌های همسایگان و دسته‌ها را در یک بازنمایی مفهوم ادغام نموده است. سرانجام، چندین فرمول را برای اندازه‌گیری شباهت بین دو مفهوم بر اساس مدل‌های نمایش معنایی پیشنهاد کرده که باعث شده است، اطلاعات گمشده در یک ویژگی با استفاده از دیگر خصوصیات مرتبط تکمیل شود و اطلاعات معنایی نهفته در بین خصوصیات مختلف، استخراج و در اندازه‌گیری شباهت معنایی استفاده شود [۲۳].

شباهت معنایی متن می‌خواهد درجه همگرایی بین دو جمله را ارزیابی کند. بیشتر معنای سطحی کلمات را در نظر می‌گیرد ولی همبستگی و ارتباط مفهومی دو کلمه را حساب نمی‌کند. برای حل این مشکل از مفهوم انتقال فضا استفاده می‌شود که از مفاهیم اسم با روابط is-a در شبکه واژگان استفاده می‌کند که این مدل هم فقط اسم‌ها را در نظر می‌گیرد و فعل و صفت و عبارات دیگر را نادیده می‌گیرد در نتیجه اطلاعات از بین می‌رود. راه‌حل ارائه شده برای حل این مشکل باید از پیوندهای مشتق شده در شبکه واژگان برای مرتبط کردن فعل‌ها و صفت‌ها و غیره با مفاهیم مربوط به اسم آن‌ها استفاده کرد و همچنین برای بهبود دقت مدل مبتنی بر مفهوم انتقال فضا اهمیت انواع مختلف کلمات را با توجه به تعیین وزن‌های آن‌ها تعیین می‌کند.

برای بهره‌برداری از مفهوم جملات محتوای اطلاعات را معرفی کرده، که برای اندازه‌گیری مفاهیم اسم به کار می‌رود. محتوای اطلاعات در مقایسه با روش‌های دیگر دو تا مزیت دارد. (۱) طبقه‌بندی مبتنی بر دانش یعنی رابطه‌ی is-a از مفاهیم در شبکه واژگان می‌تواند شباهت‌های ذهنی و ارزیابی شباهت معنایی را نشان دهد. (۲) تعیین مفهوم اطلاعات با روش آماری با شهود انسانی سازگار است مثلاً هسته کلمه شامل اطلاعات بیشتری نسبت به خود کلمه است.

با این حال روش محتوای اطلاعات هم به دلیل مشکل محاسبه هم‌پوشانی مفهومی در میان مفاهیم چندگانه به‌ندرت استفاده می‌شود. به‌طور خاص اسامی در یک جمله متعلق به مفاهیم مختلف نقشه‌برداری می‌شوند. که یا به مفاهیمی که دارای ارتباط زیرین هستند یا مفاهیم پراهمیت هستند. برای رفع مشکل هم‌پوشانی محتوای اطلاعات پیشنهاد نمایشی از جمله دانش مبتنی بر دانش یعنی مفهوم فضای انتقال داده شده است و برای رفع پیچیدگی محاسباتی متد فوق برای محتوای اطلاعات متنی روش مبتنی بر افزایش اطلاعات برای کاهش پیچیدگی محاسبات محتوای اطلاعات پیشنهاد شده است. سپس شباهت معنایی جملات را می‌توان با اندازه‌گیری نسبت محتوای اطلاعات مشترک بین آن‌ها محاسبه کرد. با این حال بسیاری از روش‌های مبتنی بر دانش فقط مفاهیم

را در یک سطح در نظر می‌گیرند که باعث از بین رفتن اطلاعات زیادی در جملات می‌شود که دو دلیل اساسی دارد:

۱. فقدان محتوای اطلاعات کلمات نامنظم در شبکه واژگان. چون تنها مفاهیم اسم قادر به ایجاد یک ساختار سلسله مراتبی هستند و مفاهیم فعل و قید و صفت و غیره. در زیر شبکه‌های چندگانه وجود ندارند.

۲. فقدان واحدهای نامشخص (oovNEs): تعداد نهادهای نام‌گذاری شده در شبکه واژگان کاملاً محدود است. از این رو مقادیر محتوای اطلاعات تعداد زیادی از کلمات را نمی‌توان مستقیماً به دست آورد.

به منظور حل چالش‌های فوق مدل انتقال فضا را مبتنی بر این جنبه‌ها توسعه داده است. (۱) لینک‌های مشتق شده در شبکه واژگان برای مرتبط کردن فعل‌ها صفت‌ها و واژگان با مفاهیم اسم مربوط به هر کدام استفاده می‌شود. (۲) علاوه بر این oovNEها توسط نرم‌افزارهای مانند stanford corenlp کشف شده است. که برای تشخیص دو کلمه شبیه به هم ولی بامعنای متفاوت مثل کیف و کیف به کار می‌رود که افزایش اطلاعات جملات توسط oovNEها بر اساس گرایش یاد گرفته شده از NEهای شناخته شده برای کاهش بیشتر اطلاعات از دست رفته استفاده می‌کند. علاوه بر این برای بهبود دقت مدل مبتنی بر مفهوم انتقال فضا اهمیت انواع مختلف کلمات را برای ارزیابی شباهت معنایی با توجه به تعیین وزن‌های به آن‌ها اندازه می‌گیرد. روش CTS بر روی داده‌های لی و همکاران به همبستگی ۰/۸۶۷۲ رسیده است [۵].

در ارزیابی [۶۰] پاور (۲۰۱۸) روشی را ارائه داده است که با استفاده از یک الگوریتم قبلاً ایجاد شده و ادغام آن با قشر دامنه با استفاده از آمار دامنه می‌پردازد. دامنه تفکیک شده به عنوان داده یادگیری نظارت شده برای الگوریتم عمل می‌کند. همچنین برای محاسبه شباهت بین افعال عمل در اهداف یادگیری با اشاره به طبقه‌بندی بلوم، شاخص Bloom را معرفی کرده است. طبقه‌بندی بلوم دستور العمل‌های کلیدی کلمه کلیدی را ارائه می‌دهد و یک ساختار سلسله مراتبی را هنگام تعیین نتایج یادگیری مورد استفاده قرار می‌دهد. این روش در ابتدا تشابه معنایی بین کلمات را با توجه به معنی کلمه در متن جمله محاسبه می‌کند. سپس برای شکل دادن یک بردار معنایی برای هر دو جمله به طور جداگانه استفاده می‌شود. بردارهای معنایی شکل گرفته برای محاسبه شباهت معنایی استفاده می‌شوند. وکتور کلمه سفارش با در نظر گرفتن ساختار نحوی جملات، یعنی وقایع کلمات مربوط به یکدیگر ساخته می‌شود. طبقه‌بندی بلوم یک مدل ساختاری مشهور، تثبیت شده و سلسله مراتبی است که شامل افعال اکشن در سطوح مختلف سلسله مراتب است. هرچه سلسله مراتب را بالا می‌بریم، سطح فعل عمل افزایش می‌یابد. سه لایه برتر، آنالیز، سنتز و ارزیابی افعال را با تفکر انتقادی نشان می‌دهد. در اینجا طبقه‌بندی بلوم را به عنوان بخش جداگانه‌ای از روش‌شناسی پیاده‌سازی کرده است و آن را در استفاده از متدلوژی شباهت جمله اصلی محدود نموده است.

در ارزیابی [۲۴] تی و همکاران (۲۰۱۹) هر متن کوتاه را به عنوان دو بردار متراکم نشان داده‌اند: اولی با استفاده از شباهت کلمه به کلمه مبتنی بر بردارهای کلمه از قبل آموزش دیده، دومی با استفاده

از تشابه کلمه به کلمه مبتنی بر منابع خارجی دانش ساخته شده است. همچنین یک الگوریتم پیش‌پردازش ایجاد کرده، که زنجیره‌ای از نام‌های اصلی را با هم متصل کرده و تقسیم کلمه را برای حفظ معنای افعال اصطلاحات و عبارات انجام داده است. روش پیشنهادی از شباهت بین متن‌های کلمه ساخته شده با استفاده از کلمات جاسازی شده و ارتباط معنایی بین مفاهیم مبتنی بر منابع خارجی دانش بهره‌برداری می‌کند. همچنین این دو قطعنامه هسته‌ای اشخاص نام‌گذاری شده را در دو متن کوتاه انجام می‌دهد تا موجودات را به هم پیوسته و تقسیم‌بندی کلمات را برای حفظ معنای افعال اصطلاحات و عبارات حفظ کند. در نتیجه ضریب همبستگی پیرسون 0.8378 را به دست آورده است.

۳.۵.۲ روش‌های مبتنی بر رشته

سیستم ukp از مدل رگرسیون خطی ورودی برای آموزش داده‌ها استفاده می‌کند تا چندین متن را با پیچیدگی‌های مختلف ترکیب کند. این سیستم با استفاده از N -گرم کلمات به تحلیل مقادیر بردار تجزیه و تحلیل صریح معناساختی و تجمیع شباهت کلمه بر اساس منابع واژگانی-معنایی می‌پردازد. از سیستم جایگزینی واژگان و ترجمه ماشین آماری برای اضافه کردن لکسیم‌های اضافی استفاده کرده است تا شکاف‌های واژگانی را کاهش دهد. به منظور مقایسه‌ی شباهت کلمات دوبعدی در سطح متن، از استراتژی ارائه شده توسط میهالکا (Mihalcea) استفاده شده است که با استفاده از مجموع نمرات idf وزن دار شباهت بین هر کلمه با همتای خودش در متن دیگر را به دست آورده است. از جایگزینی واژگانی و ترجمه ماشین آماری و سپس با استفاده از همین ویژگی‌ها اما با تعداد تصادفی برای جفت جملات استفاده کرده است [۶۱].

در بررسی [۶۲] آتوم و همکاران (۲۰۱۶) اقدامات شباهت کلمه قبلی در مجموعه داده‌های معیار مورد بررسی قرار گرفته و سپس از یک شباهت کلمه ترکیبی در یک اندازه‌گیری تشابه متن رمان (TSM) استفاده شده است. TSM پیشنهادی مبتنی بر محتوای اطلاعات و روابط معنایی شبکه واژگان است. TSM شامل تطابق دقیق کلمات، طول هر دو جمله در یک جفت و حداکثر شباهت بین یک کلمه و متن مقایسه شده است. این مقاله یک اندازه‌گیری تشابه متن جدید را بر اساس اندازه‌گیری شباهت کلمه از فاصله مشترک با محتوای مشترک و همچنین محتوای اطلاعات کلمات مقایسه شده ارائه کرده است. شباهت متن پیشنهادی بر اساس طول متن مقایسه شده و کل کلمات دقیق وزن می‌شود. اندازه‌گیری شباهت بیشتر از اقدامات شباهت مقایسه شده را دارد و در سطح 0.05 معنی دار است. دلیل دستیابی زیاد به این روش ناشی از به‌کارگیری اطلاعات تکمیلی (کالبد و محتوای اطلاعات) و اثربخشی اندازه‌گیری شباهت کلمه وام شده است.

در [۶۳] بابا و همکاران (۲۰۱۷) از نمایندگی توزیع کلمات برای تعریف یک سند مورد ارزیابی، استفاده کرده‌اند. این تحقیق یک روش تشخیص سرقت ادبی را بر اساس حداکثر مقدار طول جمله یعنی طولانی‌ترین توالی مشترک با وزن تعریف شده توسط نمایندگی توزیع پیشنهاد کرده است. که این روش در برنامه‌های کاربردی که نیاز به شناسایی دقیق سرقت‌های پیچیده دارند مفید است.

نمایندگی توزیع شده به عنوان یک تابع در نظر گرفته می شود که یک کلمه را به یک بردار با یک بعد کوچک ترسیم می کند و فاصله بین بردارها یک شباهت بین کلمات است که با بردارها مطابقت دارد. نمایش ساده توزیع شده با کاهش ابعاد یک نمایشگر بردار ساده بر اساس فرکانس یک کلمه در دسترس است. در این روش به دقت $83/70$ درصد برای داده های wllcs رسیده است.

شباهت متنی معنایی و استنباط زبان طبیعی دو وظیفه رایج در زمینه زبان طبیعی است که برای محک زدن مدل های بیان جمله درجایی که دو جمله در یک زوج قرار دارند استفاده می شود. در چنین وظایفی جملات به عنوان کیسه کلمات، توالی ها، درختان یا حلقه ها ارائه می شوند، اما مدل مورد نظر به زوج های کلمه ای مبتنی است. نتایج به دست آمده در این بررسی در پنج مجموعه داده با توجه به مدل توجه به کلمه محور، خطا را تا 41% کاهش داده است [۹].

۴.۵.۲ روش های ترکیبی

در [۶۴] اسلام و همکاران (۲۰۰۸) یک روش برای اندازه گیری شباهت معنایی متون با استفاده از یک اندازه گیری مبتنی بر بدنه از شباهت کلمه معنایی و یک نسخه عادی و اصلاح شده از طولانی ترین نتایج مشترک الگوریتم رشته ارائه داده است. همچنین بر محاسبه شباهت بین دو جمله یا دو پاراگراف کوتاه متمرکز می باشد. روش پیشنهادی این مقاله شباهت دو متن را از اطلاعات معنایی و نحوی (از نظر نظم کلمه مشترک) که در آن وجود دارد، تعیین می کند. در واقع سه کارکرد شباهت را در نظر گرفته است، ابتدا شباهت رشته و شباهت کلمه معنایی محاسبه شده و سپس از یک تابع تشابه ترتیب کلمه اختیاری استفاده کرده است تا بتواند اطلاعات نحوی را در روش بیان شده بگنجانند. در انتها، شباهت متن با ترکیب شباهت رشته، شباهت معنایی و تشابه نظم کلمه مشترک با عادی سازی حاصل می شود. آن ها به ضریب همبستگی پیرسون 0.853 و دقت $72/42$ بر روی مجموعه داده ی STS-65 دست یافته اند.

روستایی و همکاران (۲۰۱۴) از پردازش زبان طبیعی و فن های یادگیری ماشین برای استفاده از یک طرح شباهت معنایی جمله استفاده کرده اند. در بخش اول این طرح مجموعه های مختلف ویژگی های زبانی شامل ویژگی های مبتنی بر رشته، مبتنی بر معنا، مبتنی بر اسم و مبتنی بر نحو استخراج شده است. علاوه بر این، شبکه واژگان برای استخراج اطلاعات معنایی استفاده می شود. در بخش دوم الگوریتم های یادگیری ماشین های مانند k -نزدیک ترین همسایه^۱، مدل رگرسیون لجستیک و ماشین های بردار پشتیبان برای ساخت مدل های طبقه بندی بر روی مجموعه ای از ویژگی های استخراج شده استفاده می شوند. نتایج تجربی در بخش اول نشان می دهد که ویژگی های استخراج شده برای شباهت معنایی جمله معتبر است. علاوه بر این با مقایسه عملکرد الگوریتم های مختلف طبقه بندی در بخش دوم الگوریتم KNN موفق تر است. این طرح از نظر دقت برای محاسبه تشابه معنایی بهتر عمل می کند. برای ارزیابی عملکرد روش پیشنهادی، سه مجموعه داده ی استاندارد متفاوت استفاده می شود. مجموعه داده های MSRP، RTE3، TREC9 هستند. MSRP شامل

^۱K-Nearest neighbor (KNN)

۴۰۷۶ جفت آموزش و ۱۷۲۵ جفت آزمون است. RTE3 شامل ۸۰۰ جمله جفت است. TREC9 حاوی ۱۹۳ سؤال پارافرز ساخته شده توسط ارزیابی‌های انسانی با استفاده از سؤالات اصلی استخراج شده از پرسش پرس‌وجو از سؤالات ارسال شده توسط کاربر است. در این بررسی نشان داده شده است که الگوریتم‌های یادگیری ماشین به‌ویژه KNN و رگرسیون لجستیک از نظر دقت از روش‌های دیگر بهتر هستند [۶۵].

در مقاله [۸] ولی و همکاران (۲۰۱۶) کار را با پیش‌پردازش جفت‌های جمله شروع می‌کنند؛ سپس آن را با انتساب نمرات شباهت، نمرات واژگانی، معنایی و نحوی و درنهایت با نمره کل محاسبه با استفاده از یادگیری تحت نظارت به پایان می‌رسانند. آن‌ها بر روی مجموعه داده‌های عربی به ضریب همبستگی ۰/۹۲ رسیدند.

در مقاله [۷] ریچالسکاو همکاران (۲۰۱۶) برای شناسایی تشابه از ترکیبی از رمزگذار خودکار بازگشتی با یک سیستم جریمه - جایزه^۱ که بر اساس شبکه واژگان برای ارتباط معنایی و یک طبقه‌بندی کننده ماشین بردار پشتیبان، که نتیجه نهایی را از ماتریس‌های تشابه تولید می‌کند استفاده کرده‌اند. این راه‌حل بیشتر توسط یک طبقه‌بندی گروه پشتیبانی می‌شود، تراز کننده را با یک شبکه عصبی دوطرفه گیت مکرر و ویژگی‌های اضافی ترکیب نموده‌اند، سپس از رگرسیون بردار خطی پشتیبانی را برای تعیین مجموعه دیگری از امتیازات استفاده کرده‌اند. اولاً رمزگذار خودکار بازگشتی برای انجام یادگیری بر روی درختان تجزیه استفاده می‌شود. سپس شبکه واژگان فاصله بردارهای رمزگذار خودکار بازگشتی را با استفاده از سیستم جریمه - جایزه بر اساس شباهت‌های معنایی کلمات تنظیم می‌کند.

سیستم MayoNLP دو روش برای ارزیابی شباهت معنایی جملات به کار برده است: اولین روش تطبیق متن بر اساس معادلات معنایی با استفاده از یک شبکه معنایی و پیکره آماری است و دومین روش یک مدل یادگیری ماشین است که از یک مدل معنایی یادگیری عمیق استفاده کرده است. این دو روش را به صورت خطی ترکیب کرده‌اند. یک بردار لغت مشترک JWV با جمع‌آوری کلمات منحصر به فرد که در یک جفت جمله وجود دارد تشکیل شده است. شباهت معنایی با مقایسه معنای پنهان در جملات محاسبه می‌شود و شباهت نحوی بر اساس ترکیب کلمه در دو جمله محاسبه می‌شود. این بررسی به ضریب همبستگی ۰/۸۰۱۵ در بهترین اجرا دست یافته است [۶۶].

در [۳۳] دو اجرا برای شباهت معنایی متن تک‌زبان ارسال شده است، یکی مبتنی بر یادگیری نظارت شده و دیگری سیستم بدون ناظر می‌باشد.

● **UWB sup**: سیستم نظارت شده بر پایه رگرسیون SVM با هسته rbf که از تمام فن‌های استفاده شده در sts به عنوان ویژگی‌های رگرسیون استفاده کرده است تا وابستگی بین ویژگی‌های یکسان را بهتر مدل کند.

● **UWB unsup**: سیستم بدون نظارت بر اساس همگرایی کلمات وزن دار با استفاده از انجام کار متقابل زبانی شباهت معنایی متن را با دو زبان اسپانیایی و انگلیسی در دو مرحله انجام داده است. اول جملات اسپانیایی را از طریق مترجم گوگل به زبان انگلیسی ترجمه کرده است. جملات انگلیسی

¹ award-penalty

ثابت باقی مانده، سپس از سیستم شباهت معنایی متن برای کار با تک‌زبانها استفاده کرده است. این بررسی بر روی مجموعه داده‌های STS2016 به ضریب همبستگی پیرسون ۰/۸۱۱۶ دست‌یافته است.

در [۳] کای و همکاران (۲۰۱۷) یک رویکرد ترکیبی مبتنی بر شبکه واژگان CSSM-ICSP برای اندازه‌گیری شباهت معنایی مفهوم پیشنهاد داده‌اند که از مفهوم اطلاعات مفاهیم برای وزن کمترین فاصله بین مفاهیم استفاده کرده است. این مقادیر مفهوم اطلاعات مشترک را در محاسبه کمترین فاصله مسیر بین مفاهیم مقایسه شده ادغام می‌کند. برای تعیین کمیت مفهوم اطلاعات مفاهیم، در درجه اول با در نظر گرفتن خصوصیات معنایی و ساختارهای مختلف شبکه واژگان، یک مدل مفهوم اطلاعات ذاتی ایجاد شده است. در این مقاله، یک روش ترکیبی را برای محاسبه شباهت معنایی بین جفت کلمات بر اساس شبکه واژگان پیشنهاد کرده است. معیار پیشنهادی غیرخطی فاصله مسیر بین مفاهیم را برای اندازه‌گیری شباهت معنایی آن‌ها تغییر داده است. در این روش، وزن مفهوم اطلاعات مسیر به‌عنوان یک عامل محاسبه در نظر گرفته شده است که در آن اختلاف مقادیر مفهوم اطلاعات بین مفاهیم مجاور برای وزن کردن لبه پیونددهنده مفاهیم به‌کاررفته است. با توجه به روش‌های مختلف شناسایی فاصله معنایی بین مفاهیم، این رویکرد به‌عنوان دو استراتژی ارائه می‌شود: استراتژی ۱ از کوتاه‌ترین مسیری عادی استفاده می‌کند و استراتژی ۲ از همپوشانی مسیر بین دو مفهوم مقایسه شده و کمترین مشترک آن‌ها استفاده می‌کند.

رویکردهای مبتنی بر محتوای اطلاعات به مدل‌های آماری ارزش مفهوم اطلاعات مفاهیم به‌عنوان فاکتورهای محاسبه تکیه کرده‌اند، که برای ارزیابی چگونگی یک مفهوم خاص و آموزنده استفاده شده است. بر اساس بردارهای ویژگی، رویکردهای مبتنی بر ویژگی روی ویژگی‌های مشترک یا ویژگی‌های یکسان مفاهیم تطبیقی متمرکز می‌شوند. در کنار ظهور معیارهای مختلف شباهت، ادغام روش‌های مختلف برای غلبه بر محدودیت یک رویکرد واحد امکان‌پذیر است. در اینجا یک روش ترکیبی مبتنی بر شبکه واژگان را برای ارزیابی شباهت معنایی بین مفاهیم ایجاد کرده است. برخلاف رویکردهایی که نیاز به یک مقیاس بزرگ همراه با اصطلاحات حوزه گرا دارند، این رویکرد فقط از شبکه واژگان استفاده می‌کند و به‌جای تراز هستی‌شناسی بر شباهت معنایی مفاهیم در شبکه واژگان تمرکز می‌کنند. این به‌طور عمده به: اولاً تجزیه و تحلیل مقایسه‌ای ویژگی‌های فنی رویکردهای شباهت موجود در مورد عوامل محاسبه، مزایا و اشکالات علاوه بر دقت پرداخته است. ثانیاً ارائه یک مدل جدید از IC ذاتی مفاهیم، که دو ویژگی معنایی شبکه واژگان را در نظر می‌گیرد. ثالثاً به‌صورت غیرخطی فاصله مسیری با وزن IC را به شباهت بین مفاهیم مقایسه شده توسط یک عملکرد نمایی تبدیل کنید. این مدل بر روی داده‌های M&C به ضریب همبستگی پیرسون ۰/۸۹۲ دست‌یافته است.

در بررسی [۶] ورداسی و همکاران (۲۰۱۷) با به کار بردن سه ویژگی بر روی روش هم‌ترازی به نتایج خوبی در تخمین شباهت معنایی متون رسیده‌اند. دو مورد از ویژگی‌های پیشنهادی این مقاله برای تجزیه و تحلیل ویژه‌ای هستند، که برای شباهت یابی واژه‌های معنایی بهتر عمل می‌کنند. با محاسبه شباهت یک جفت جمله، نمره سطح همبستگی بین آن دو جمله به دست می‌آید. برای این کار نیاز است تا فراتر از تشابه واژگانی پرداخته و شباهت معنایی جملات را در نظر بگیرند. یکی از

ابزارهایی که به طور گسترده برای محاسبه شباهت کلمات به کار می رود شبکه واژگان است. مخزن دانش به عنوان یک شبکه معنایی روابط نحو را ذخیره می کند. علاوه بر شبکه واژگان یکی دیگر از روش های مناسب برای تخمین شباهت معنایی کلمه، روش متداول شبکه مبتنی بر word2vec است. در حقیقت word2vec یک کلمه تعبیه شده را در فضای بردار ایجاد می کند و شباهت معنایی بین دو کلمه به سادگی با محاسبه شباهت کوسین از کلمات بازتاب های بردار محاسبه می کند. روش مشخص شده از نمایش word2vec استفاده می کند و بردارهای کلمات داخل یک جمله را برای تولید ماتریس جمله تکرار می کند. به این ترتیب وظیفه اندازه گیری شباهت جمله، به محاسبه شباهت ماتریس ها کاهش می یابد. یکی از ویژگی های به کار برده شده برای ارزیابی شباهت جمله، خطای تخمین خطی^۱ است که نتیجه خوبی در برآورد خطی یک ماتریس دارد. همچنین از یک ترازنامه تک زبانه استفاده کرده تا به طور صحیح کلمه مشابهی از جفت جملات را با هم جمع کند و برای هر شباهت تخمینی ارزشی ارائه دهد.

در این روش از دو منبع شباهت کلمات به عنوان ورودی ترازنامه استفاده شده است:

۱. شباهت کوسین از نمایش برداری کلمات حاصل از word2vec

۲. عناصر متناظر ماتریس شباهت کلمه لغات معنی لغوی

از ترکیب سه ویژگی خطای تخمین خطی، w2vta و spata با یک سیستم Tf-idf نتیجه برآورد شباهت معنایی به دست آمده است.

مناسب ترین سیستم برای مقایسه فن های مختلف بر روی پنج زبان مختلف (انگلیسی، فارسی، آلمانی، ایتالیایی و اسپانیایی) Luminoso بود که یک سیستم ترکیبی است که به طور مؤثر کلمات و اطلاعات از منابع دانش را ادغام می کند [۶۷].

در مقاله [۶۸] پاکلیاردینی و همکاران (۲۰۱۸) یک روش ریاضی، محاسباتی کارآمد الهام گرفته از CBOW را برای آموزش جمله به جای تعبیه کلمات پیشنهاد داده اند. این مدل، یک مدل ساده بدون نظارت است که اجازه می دهد تعبیه کلمات را با استفاده از بردارهای کلمه ای همراه با تعبیه N - گرم ها هم زمان با ترکیب آموزش و خود بردارهای تعبیه شده انجام دهد. پیچیدگی محاسباتی تعبیه های این مدل در عملیات برداری، در هر کلمه پردازش شده، هم در حین آموزش و هم در استنتاج تعبیه های جمله فقط $O(1)$ است. آن ها بر روی مجموعه داده های SICK به ضریب پیرسون ۰/۷۱۵ رسیده اند. در این پژوهش از مدل Sent2Vec نیز بر روی داده های SICK استفاده کرده اند.

در [۶۹] پاوار (۲۰۱۸) متدلوژی ارائه داده شده است که با درج تشابه معنایی و آمار پیکره، به n محاسبه شباهت معنایی بین جملات می پردازد. برای محاسبه شباهت معنایی بین کلمات و جملات، روش پیشنهادی از یک رویکرد مبتنی بر لبه با استفاده از بانک اطلاعاتی واژگان پیروی می کند. این روش در هر دو معیار و میانگین مجموعه داده های شباهت انسان آزمایش شده است. روش پیشنهادی از پایگاه داده شبکه واژگان برای مقایسه معنای مناسب کلمه استفاده می کند. یک وکتور معنایی برای هر جمله تشکیل شده است که شامل وزن اختصاص داده شده به هر کلمه برای هر کلمه دیگر از جمله دوم در مقایسه است. این مرحله همچنین محتوای اطلاعات کلمه را در نظر

¹Linear estimation error(LEE)

می‌گیرد، به‌عنوان مثال فرکانس کلمه از یک کالای استاندارد. شباهت معنایی بر اساس دو بردار معنایی محاسبه می‌شود. برای هر جمله یک بردار مرتبه تشکیل می‌شود که شباهت نحوی بین جملات را در نظر می‌گیرد. سرانجام، شباهت معنایی بر اساس بردارهای معنایی و بردارهای ترتیب محاسبه می‌شود. در این ارزیابی برای تشابه کلمه به ضریب همبستگی پیرسون 0.8753 رسیده‌اند و برای شباهت جمله، همبستگی به‌دست‌آمده 0.8794 است.

این الگوریتم در ابتدا هر دو جمله را متمایز می‌کند و در بخش‌های گفتار آن‌ها را برچسب‌گذاری می‌کند. این رویکرد، تفکیک معنای صحیح کلمه مقایسه را تضمین می‌کند. شباهت بین کلمات بر اساس یک رویکرد مبتنی بر لبه قبلاً محاسبه می‌شود. محتوای اطلاعاتی از یک قشر می‌تواند برای تأثیر شباهت در حوزه خاص استفاده شود. بردارهای معنایی حاوی شباهت‌ها بین کلمات برای جملات تشکیل شده و بیشتر برای محاسبه تشابه جمله استفاده می‌شود. بردارهای ترتیب کلمه نیز برای محاسبه تأثیر ساختار نحوی جملات تشکیل شده‌اند. از آنجا که نظم کلمه بر شباهت کلی کمتر از شباهت معنایی تأثیر می‌گذارد، شباهت ترتیب کلمات به میزان کمتری وزن می‌شود. این متدلوژی بر روی مجموعه داده‌های قبلاً ایجاد شده آزمایش شده است که حاوی نتایج استاندارد و همچنین میانگین نتایج انسانی است.

۵.۵.۲ سایر روش‌ها

در [۷۶] به ازای هر جمله یک بردار شباهت معنایی و یک بردار شباهت چیدمان کلمات تشکیل می‌شود و شباهت بین دو جمله با ترکیب وزن دار این دو بردار را به دست آورده است. الگوریتم این مقاله جهت محاسبه میزان شباهت میان کلمات به دلیل نبودن پیکره کامل در زبان فارسی، از پیکره اینترنتی استفاده کرده است. سپس با کمک این مقادیر و ترکیب دو بردار شباهت معنا و چیدمان به ترکیب نقش نحو معنا پرداخته است. در نهایت نتایج به‌دست‌آمده را در سه سطح فازی دسته‌بندی کرده است این روش بر روی 50 جمله از اخبار برگزیده فارسی پیاده‌سازی شده است.

در [۷۷] با استخراج ویژگی متن کلمات کلیدی را استخراج می‌کند و با استفاده از تجزیه‌کننده‌های متنی درخت‌های تجزیه مربوط به اسناد را به دست آورده و با استفاده از وزن کلمات کلیدی به‌دست‌آمده میزان شباهت بین درخت‌ها را تخمین می‌زند. برای بررسی میزان شباهت بین درختان از الگوریتم به دست آوردن تعداد زیر درختان مشابه استفاده کرده‌اند. در انتها با استفاده از الگوریتم خوشه‌بندی سلسله مراتبی k - نزدیک‌ترین همسایه گروه‌بندی سندهای متنی را انجام داده است.

در [۷۰] جان ویتینگ و همکاران (۲۰۱۶) با استفاده از تعبیه‌های CHARAGRAM، یک رویکرد ساده برای یادگیری مدل‌های ترکیبی شخصیت پایه را برای تعبیه توالی متن ارائه دادند. یک کلمه یا جمله با استفاده از یک بردار شمارش N - گرم کاراکتر نمایش داده می‌شود و به دنبال آن یک دگرگونی غیرخطی برای انجام تعبیه کم بعدی انجام می‌شود. برای ارزیابی از سه روش استفاده کرده است: تشابه کلمه، تشابه جمله و برچسب زدن بخشی از گفتار. در این بررسی به این نتیجه رسیده‌اند که تعبیه‌های CHARAGRAM از معماری‌های پیچیده‌تر بر اساس شبکه‌های عصبی مکرر و

دستیابی به عملکرد جدیدترین روش در چندین کار مشابهت، بهتر است. آن‌ها به ضریب همبستگی پیرسون $۰.۶۳/۳۳$ بر روی مجموعه داده‌ی sl999 و $۰.۵۳/۸۷$ روی مجموعه داده‌ی ws353 دست‌یافته‌اند. رکیب و همکاران (۲۰۱۷) به ترتیب ۲ اقدام مرتبط با کلمه موجود را مبتنی بر Google N-gram و Global Vectors برای نمایندگی کلمه اتخاذ نموده‌اند و آن‌ها را به روش متفاوتی با یک روش مرتبط بودن جمله مبتنی بر عبارت N-gram مبتنی بر Google برای محاسبه مرتبط بودن متن ترکیب کرده‌اند. ترکیب ارتباط Google با کلمه گرام و عبارت Google به‌تنهایی با دستیابی به میانگین وزنی بالاتر پیرسون، یعنی ۰.۶۳۹ و ۰.۶۱۹ منجر شده است. به‌طور مشابه، ترکیبی از ارتباط کلمه مبتنی بر GloVe و مرتبط بودن با عبارت N-gram مبتنی بر Google عملکرد بهتری نسبت به کلمه مبتنی بر GloVe دارد و با دستیابی به میانگین وزنی بالاتر پیرسون، یعنی به ترتیب، برابر ۰.۶۱۹ و ۰.۶۰۵ در همان ۱۴ مجموعه داده رسیده است [۷۱].

در مقاله [۷۲] یک روش بدون ناظر به نام رویکرد منبع نور^۱ برای اندازه‌گیری شباهت معنایی بین متون به زبان‌های مختلف پیشنهاد شده است. این رویکرد برای کار در فضای دوزبانه (یا چندزبانه)، بردارهای مداوم کلمه (یعنی تعبیه کلمه) را از یک زبان به فضای بردار زبان دیگر از طریق مدل ترجمه خطی طراحی کرده است. سپس کلمات را با توجه به شباهت بردارهای آن‌ها در فضای تعبیه دوزبانه تراز کرده و اقدامات مختلف نظارت‌نشده از شباهت معنایی با بهره‌گیری از تعبیه‌های دوزبانه و تراز کلمات را بررسی کرده‌اند. این روش ساده عملکردی نزدیک به روش‌های نظارت‌شده دارد و ثبات در زوج‌های مختلف زبان را نشان می‌دهد. این روش برای استخراج جملات موازی از مجموعه‌ای از نوشتجات^۲ و تشخیص سرقت ادبی قابل مقایسه است. در این مدل به ابزار و منبع مشخصی نیاز نیست و صرفاً برپایه در دسترس بودن مجموعه‌ای نوشتجات بزرگ برای هر دو زبان ورودی استوار است. همچنین به یک مجموعه نسبتاً کوچک از جفت ترجمه کلمات احتیاج دارد.

در [۷۳] از تعبیه جمله BERT استفاده می‌کند. ابتدا ارتباط نظری بین هدف از پیش آموزش داده‌شده مدل زبان پوشیده و وظیفه شباهت معنایی را از لحاظ نظری آشکار کرده و سپس تعبیه‌های جمله BERT را به‌صورت تجربی تجزیه و تحلیل نموده است. این تعبیه جمله همیشه یک فضای معنایی ناهمسانگرد ناهموار برای جملات ایجاد می‌کند، که به عملکرد شباهت معنایی آن آسیب می‌رساند. برای پرداختن به این موضوع در این مقاله توزیع تعبیه‌شده جمله ناهمسانگرد را به یک توزیع گاوسی صاف و همسانگرد از طریق جریان‌های نرمال‌سازی که با یک روش بدون ناظر صورت می‌گیرد تبدیل کرده است. در این روش به نمره همبستگی $۰.۶۵/۶۲$ دست‌یافته‌اند.

۶.۲ جمع‌بندی

در این فصل ویژگی‌های مورد استفاده در شباهت معنایی متن را مورد بررسی قرار دادیم و به بررسی تعدادی از این مقالات در این حوزه پرداختیم. در زیر تعدادی از این مقالات را باهم مقایسه نموده و

^۱cross-lingua

^۲Corpor

نتایج آن‌ها را مورد بحث قرار داده‌ایم. در فصل بعدی روش پیشنهادی خود را برای تخمین شباهت معنایی متون ارائه خواهیم داد.

با توجه به جدول ۱.۲ بر روی داده‌های SICK، راناسینگه و همکاران با روش GRU به ضریب همبستگی بهتری نسبت به بقیه با خطای کمتری در این مجموعه داده رسیده‌اند. بر روی مجموعه داده‌های SemEval نیز افضل و همکاران به ضریب همبستگی بهتری دست پیدا کرده‌اند. پاور و همکاران بر روی مجموعه داده‌های sts-۶۵ که حاوی ۶۵ جفت جمله می‌باشد به همبستگی بهتری دست یافته‌اند.

جدول ۱.۲: مقایسه برخی کارهای صورت گرفته در زمینه شباهت معنایی متن

دقت	مجموع مربع خطا	ضریب همبستگی اسپرمن	ضریب همبستگی پیرسون	مجموعه داده	روش (سال)
۷۲/۴۲	۲۷/۵۸		۰/۸۵۳	STS۶۵	ISLAM and D (2008)
۷۸/۸۸	۰/۲۱۱۲	۰/۸۳۹۶	۰/۸۹۰۱	SICK	GRU (2014)
۸۴/۲	۰/۲۲۸۶	۰/۸۳۴۵	۰/۸۸۲۲	SICK	MaLSTM (2016)
			۰/۸۰۱۵	SemEval-۲۰۱۶	Afzal et al (2016)
		۰/۶۷	۰/۸۹۲	J,C	IC model (2017)
۸۵/۳۵	۰/۳۷		۰/۷۱	AssIN	Proposed Measure (2017)
۷۸/۰۲	۰/۲۸۹۸	۰/۷۹۳۳	۰/۸۵۴۹	SICK	Siamese # local context: 5 + Siamese LSTM(2018)
			۰/۸۷۹۴	STS۶۵	Pawar et al(2018)
			۰/۷۷۹۹	SemEval-۲۰۱۶	CTS(۲۰۱۸)
۶۲	۰/۳۸		۰/۸۶۰	SICK	Gazpio et al(2019)

فصل ۳

روش تحقیق

۱.۳ مقدمه

در فصل قبل مروری بر روش‌های تخمین شباهت معنایی متن و کاربرد آن‌ها در شباهت یابی و همچنین کارهای صورت گرفته در این زمینه و نتایج به دست آمده از آن‌ها داشتیم. در این فصل به رویکرد پیشنهادی خود که مبتنی بر پیش پردازش داده‌های متنی و استفاده از رویکرد تخصیص پنهان دیریکلت به همراه فاصله جنسن شانون برای به دست آوردن نتیجه نهایی شباهت معنایی متون مورد آزمایش خود می‌باشد خواهیم پرداخت. سپس با استفاده از رویکرد ماشین بردار پشتیبان دقت دسته‌بندی را به دست خواهیم آورد. ابتدا به تعریف دقیق مسئله پرداخته و سپس با بررسی دقیق هر یک از ویژگی‌های روش پیشنهادی، مدل خود را به تفصیل بیان خواهیم کرد.

۲.۳ طرح مسئله

در فصل قبل مبانی نظری و روش‌های ارزیابی شباهت معنایی متن مورد مطالعه قرار گرفت. با توجه به کاربردی بودن شباهت یابی در زمینه‌های بسیار از جمله غلطیابی املایی، تشخیص سرقت ادبی، تشخیص تقلب، بازیابی اطلاعات، طبقه‌بندی اسناد و غیره درصدد هستیم تا به دقت بهتری برای شباهت بین جملات برسیم. همچنین میزان همبستگی را افزایش داده و از مجموع مربعات خطای نتایج به دست آمده کم کنیم.

امروزه شیوه‌های بسیاری برای کارهای پردازش زبان طبیعی صورت می‌گیرد. یکی از این شیوه‌ها استفاده از روش‌های آماری در ترکیب با سایر روش‌ها می‌باشد و همچنین نیاز به یک طبقه‌بند سریع و قابل اعتماد برای دسته‌بندی داده‌ها داریم. یکی از روش‌های مورد استفاده در این زمینه استفاده از روش LDA می‌باشد که ما از آن برای به دست آوردن شباهت جملات استفاده می‌کنیم. همچنین از ترکیب این فن با روش آماری جنسن و شانون برای اندازه‌گیری فاصله بین جملات استفاده خواهیم کرد. در نهایت از ماشین بردار پشتیبان برای تخمین دقت فاصله به دست آمده استفاده می‌نماییم.

۳.۳ روش پیشنهادی

برای به دست آوردن شباهت معنایی بین جملات ابتدا باید جملات تبدیل به زبان قابل فهم برای رایانه شوند. به همین منظور باید تبدیل به بردارهای عددی شوند. قبل از تبدیل به بردار عددی به پیش پردازش جملات خواهیم پرداخت. به این صورت که علائم نگارشی اضافه و کلماتی که تأثیر زیادی در معنای متن ندارند را از جمله حذف می‌کنیم و ریشه کلمات را جایگزین خواهیم کرد. سپس به استخراج ویژگی آن‌ها که همان تبدیل به بردار عددی می‌باشد خواهیم پرداخت و فرکانس کلمات

پرکاربرد را برای آموزش LDA به دست خواهیم آورد. در انتها فاصله بین بردارها را محاسبه کرده که شباهت نهایی را به ما خواهد داد. در واقع شباهت معنایی، اندازه‌گیری فاصله مفهومی بین دو موضوع است، هر چه این فاصله به صفر نزدیک‌تر باشد شباهت جفت جمله بیشتر است و هر چه به یک نزدیک‌تر باشد تناقض بین جملات بیشتر می‌باشد. پس از به دست آوردن نمره شباهت هر جفت جمله به دسته‌بندی آن‌ها با استفاده از برچسب نمره شباهت خواهیم پرداخت و دقت دسته‌بندی همچنین همبستگی پیرسون و اسپرمن را به دست خواهیم آورد.

در ادامه به شرح موارد گفته‌شده می‌پردازیم:

۴.۳ ارائه‌ی مدل کلی

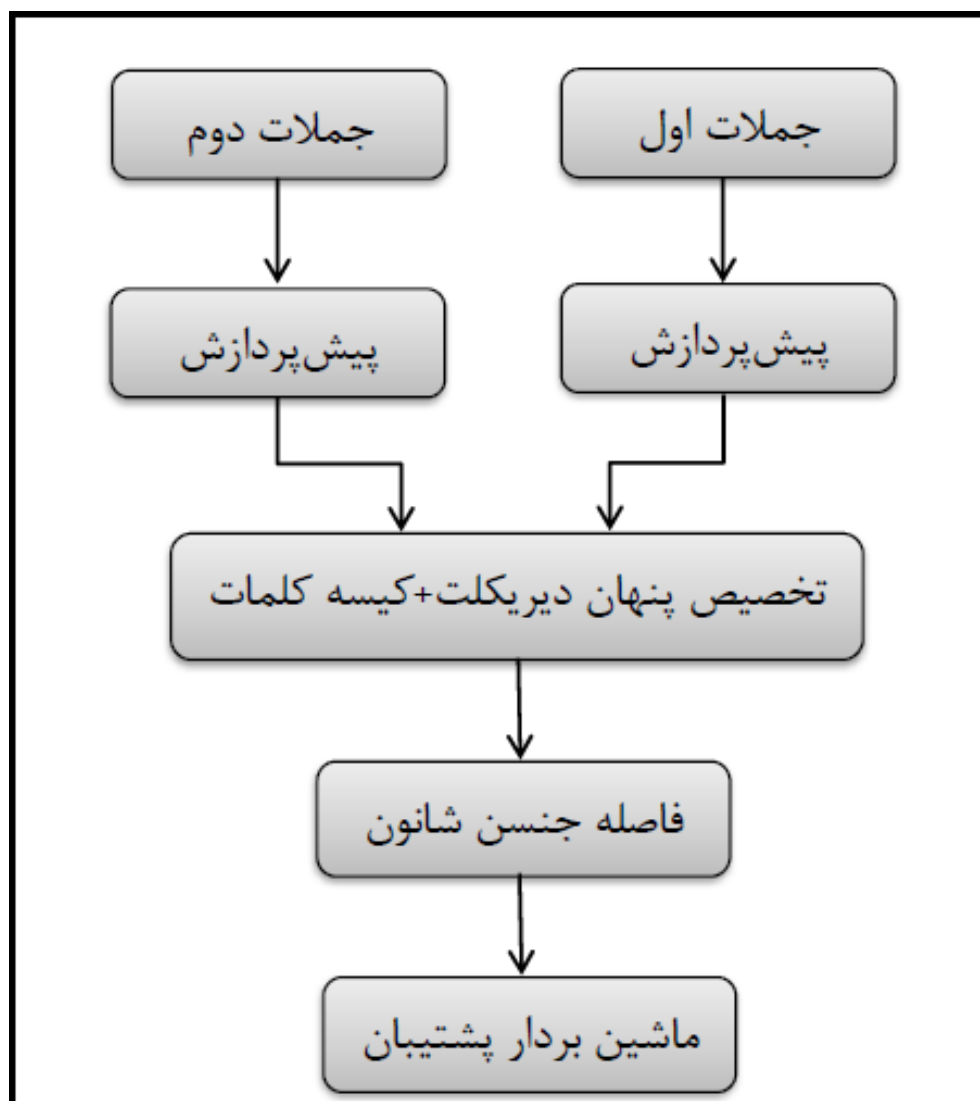
همان‌طور که می‌دانیم شباهت متن باید تعیین کند که نزدیک بودن دو قطعه متن هم از نظر نزدیک بودن سطح (تشابه واژگانی) و هم معنی (تشابه معنایی) وجود دارد. به جای مقایسه کلمه به کلمه، برای گرفتن بیشتر ارتباط معنایی باید به زمینه نیز توجه کنیم. برای در نظر گرفتن شباهت معنایی، باید روی معنای کلی جمله متمرکز شویم. که در آن یک قطعه متن قبل از محاسبه شباهت به یک گروه مرتبط از کلمات مرتبط تقسیم شده است. ما با پردازش جملات علائم اضافی و کلمات بی‌تأثیر را حذف می‌نماییم و جملات را توکن‌بندی می‌کنیم و نقش هر کلمه در جمله را به دست می‌آوریم. ایده اصلی تهیه اسناد به‌عنوان ورودی و دریافت بردارهای ویژگی به‌عنوان خروجی است و در نهایت با اندازه‌گیری فاصله بین این ویژگی‌ها، جملات را مقایسه کنیم.

در شکل ۱.۳ مدل کلی را رسم نموده‌ایم. همان‌طور که در این شکل مشاهده می‌کنیم، مدل کلی دارای ۴ مرحله می‌باشد. در بخش اول، پیش‌پردازش روی داده‌ها انجام می‌گیرد. در مرحله بعد به کمک LDA توزیع بین کلمات برای هر دسته مشخص می‌شود و برای هر جمله بردار عددی را به دست می‌آوریم، در حقیقت این مرحله، مرحله‌ی استخراج ویژگی می‌باشد. سپس این بردارهای ایجادشده و ماتریس به‌دست‌آمده از LDA را به ماژول تخمین شباهت ارائه می‌دهیم تا هر جمله را به یک دسته اختصاص دهد و خروجی را برگرداند و پس‌از آن خروجی ایجادشده را به مرحله بعدی می‌فرستیم تا نتایج به‌دست‌آمده طبقه‌بندی و بررسی شوند و دقت انجام‌شده به دست آید.

۱.۴.۳ پیش‌پردازش متن

همان‌طور که می‌دانیم داده‌های متنی غیر ساخت‌یافته، بزرگ‌ترین منبع داده‌های تولیدشده به‌وسیله انسان محسوب می‌شوند. این داده‌ها باید از زبان انسان به قالب قابل‌خواندن توسط ماشین برای پردازش بیشتر تبدیل شوند.

اولین مرحله در دسته‌بندی متن تبدیل اسناد به‌صورت رشته‌ای از کاراکترها با فرمت‌های مختلف می‌باشد که برای روش‌های یادگیری و طبقه‌بندی نمایش داده می‌شود. روش‌های پیش‌پردازش متن از آن جهت حائز اهمیت هستند که ابزارهای لازم را برای تبدیل متن از زبان طبیعی به فرمت



شکل ۱.۳: ساختار کلی روش پیشنهادی

قابل خواندن توسط ماشین را فراهم می کنند. هر زمان که داده های متنی داشته باشیم، باید چندین مرحله پیش پردازش را برای داده ها انجام دهیم تا کلمات را به ویژگی های عددی تبدیل کنیم که با الگوریتم های یادگیری ماشین کار می کنند. مراحل پیش پردازش یک مسئله عمدتاً به دامنه و خود مسئله بستگی دارد، بنابراین نیازی نیست که همه مراحل را برای هر مشکلی اعمال کنیم. پس از به دست آوردن متن، ما با نرمال سازی متن شروع می کنیم. نرمال سازی متن شامل موارد زیر است:

● تبدیل متن به حروف کوچک

اولین مرحله در پیش پردازش متن تبدیل داده ها به فرم یکنواخت می باشد. بهتر است داده های متنی را به حروف کوچک تبدیل کنیم تا مطمئن شویم به طور مثال کلمه "Nlp" با "nlp" یکسان است.

stem ریشه ممکن است چنین نباشد. از این رو کلمه ریشه به نام lemma نیز نامیده می‌شود که همواره در دیکشنری وجود دارد.

● برچسب‌گذاری کلمات^۱

برچسب‌گذاری کلمات دسته‌بندی واژگانی خاصی هستند که در آن واژه‌ها بر اساس زمینه و نقش معنایشان به بخش‌های مختلف انتساب می‌یابند. به‌طور معمول واژه‌ها می‌توانند در یکی از دسته‌های اصلی زیر قرار گیرند:

اسم^۲: این جزء معمولاً واژه‌هایی را شامل می‌شود که یک شیء یا نهاد را معرفی می‌کنند که می‌تواند زنده یا غیرزنده باشد. برخی نمونه‌ها شامل روباه، سگ، کتاب و غیره هستند. نماد تگ POS برای اسم‌ها به صورت N است.

فعل^۳: افعال همان واژه‌هایی هستند که اعمال، وضعیت‌ها یا رخداد‌های خاصی را توصیف می‌کنند. طیف متنوعی از زیر دسته‌های افعال مانند افعال کمکی، افعال انعکاسی و افعال گذرا و موارد بسیار دیگر هستند. برخی نمونه‌های معمول افعال به صورت دویدن، پریدن، خواندن و نوشتن هستند. نماد تگ POS برای افعال به صورت V است.

صفت^۴: صفت‌ها، اسم‌هایی هستند که برای توصیف یا تعیین کیفیت واژه‌های دیگر که معمولاً اسم یا عبارت‌های اسمی هستند، به کار می‌روند. عبارت گل زیبا یک اسم (N) دارد که با استفاده از یک صفت (ADJ) به صورت زیبا توصیف شده است. نماد تگ POS برای صفت‌ها به صورت ADJ است.

قید^۵: قیدها معمولاً به صورت اصلاح‌کننده^۶ برای واژه‌های دیگر شامل اسم، صفت، فعل و یا دیگر قیدها استفاده می‌شوند. عبارت گل بسیار زیبا دارای قید (ADJ) «بسیار» است که صفت زیبا را تعدیل کرده و میزان زیبایی آن را تعیین می‌کند. نماد تگ POS برای قیدها به صورت ADJ است.

علاوه بر این چهار دسته اصلی اجزای گفتار، دسته‌های دیگری نیز وجود دارند که به‌طور مکرر در زبان انگلیسی ظاهر می‌شوند. این موارد شامل ضمائر، حروف اضافه، حروف ندا، حروف ربط، ضمائر اشاره و موارد بسیار دیگر هستند. به‌علاوه، هر تگ POS مانند اسم (N) را می‌توان به زیر دسته‌هایی مانند اسامی مفرد (NN)، اسامی خاص مفرد (NNP) و اسامی جمع (NNS) تقسیم کرد.

این فرایند طبقه‌بندی و برچسب‌گذاری تگ‌های POS برای کلمات به نام تگ‌گذاری اجزای گفتار نامیده می‌شود. تگ‌های POS برای حاشیه‌نویسی کلمات و تعیین POS آن‌ها استفاده

^۱ part-of-speech tagging(pos)

^۲ Noun

^۳ Verb

^۴ Adjective

^۵ Adverb

^۶ Modifier

می‌شوند که در جهت اجرای آنالیزهای خاص مانند خلاصه‌سازی اسامی و تعیین رایج‌ترین اسم‌ها، ابهام‌زدایی از معنی کلمه و آنالیز گرامری کاملاً مفید هستند. در روش پیشنهادی ما این مراحل انجام‌شده‌اند. در شکل ۲.۳ مراحل این فرآیند نشان داده‌شده است.



شکل ۲.۳: مراحل پیش‌پردازش متن

۲.۴.۳ استخراج ویژگی از متن

همان‌طور که قبلاً بیان کردیم تخمین شباهت معنایی متن در حوزه‌های بسیاری مورد استفاده قرار گرفته است. مشکلی که در دسته‌بندی متن وجود دارد، حجم زیاد ویژگی‌ها است که باعث کاهش دقت نتایج دسته‌بندی می‌شود. برای انتخاب و برای حل این مشکل و کاهش ابعاد ویژگی‌ها از متدهای انتخاب ویژگی استفاده می‌کنند.

انتخاب ویژگی متن یک موضوع اساسی و مهم برای استخراج متن و بازیابی اطلاعات است، زیرا در این مرحله واژه‌های کلیدی انتخاب می‌شوند تا به‌عنوان بهترین نمایش‌دهنده برای سند متنی مورد استفاده قرار بگیرند. اگر تعداد واژه‌های کلیدی انتخاب‌شده کم باشد صحت و کارایی سیستم تحت تأثیر قرار می‌گیرد و کاهش می‌یابد و در مقابل اگر تعداد واژه‌های کلیدی انتخاب‌شده زیاد باشد باعث کاهش کارایی سیستم در بعد زمان خواهد شد و سرعت آموزش در فاز آموزش پایین می‌آید. در طول استخراج ویژگی، ویژگی‌های نامربوط یا اضافی حذف می‌شوند. محاسبه وزن را استخراج ویژگی متن می‌نامند. هم‌زمان با استخراج ویژگی با استفاده از تخصیص پنهان دیریکلت به تشخیص اینکه هر کلمه در چه دسته‌ای قرار دارد می‌پردازیم. این دسته‌ها بر اساس نمره شباهت

موجود به دست آمده‌اند. در واقع کلمات را بر اساس نمره شباهت در یک دسته قرار می‌دهیم و این فرآیند را تکرار می‌کنیم تا کل داده‌ها به دسته درست تعلق گیرند. تعداد دسته‌ها را با این فرآیند مشخص می‌کنیم. ممکن است کلمات در میان دسته‌ها هم‌پوشانی داشته باشند، بنابراین چندین دسته کلمات مشابهی دارند.

مدل با اختصاص دادن هر کلمه در هر جمله به یک دسته تصادفی، کار خود را شروع می‌کند. سپس، این مرحله را برای هر کلمه تکرار می‌کند. دسته فعلی را لغو کرده و کلمه را بر اساس احتمال به دسته بعدی اختصاص می‌دهیم. یعنی دسته به جمله فعلی اختصاص دارد و کلمه به دسته فعلی مختص می‌شود. مدل LDA طبقه‌بندی دسته‌ها را انجام نمی‌دهد، بلکه بررسی می‌کند که آیا این کلمه برای این دسته مناسب هست یا نه، یعنی توزیع گسسته بین کلمات برای یک دسته را انجام می‌دهد. می‌توان دسته‌ها را بعد از آموزش نام برد. وقتی هر کلمه را در هر جمله تکرار می‌کنیم دسته فعلی آن را لغو می‌کنیم و کلمه را به دسته جدید اختصاص می‌دهیم. دسته‌ای که کلمه را مجدداً به آن اختصاص می‌دهیم بر اساس احتمالات زیر است:

$$P(\text{document} \text{ "Like" } the \text{ topic}) * (\text{topic} \text{ "Like" } the \text{ word } w) \quad (1.3)$$

$$\Rightarrow \frac{n_{i,k} + \alpha}{N_i - 1 + K_\alpha} * \frac{m_{w,k} + \eta}{\sum_{w \in v} m_{w,k} + V_\eta} \quad (2.3)$$

در این فرمول به ترتیب:

$n_{i,k}$: تعداد واژه‌های اختصاص یافته به عنوان k در سند i

α : پارامتر هموارسازی

N_i : تعداد کلمات در سند i

K : تعداد موضوعات

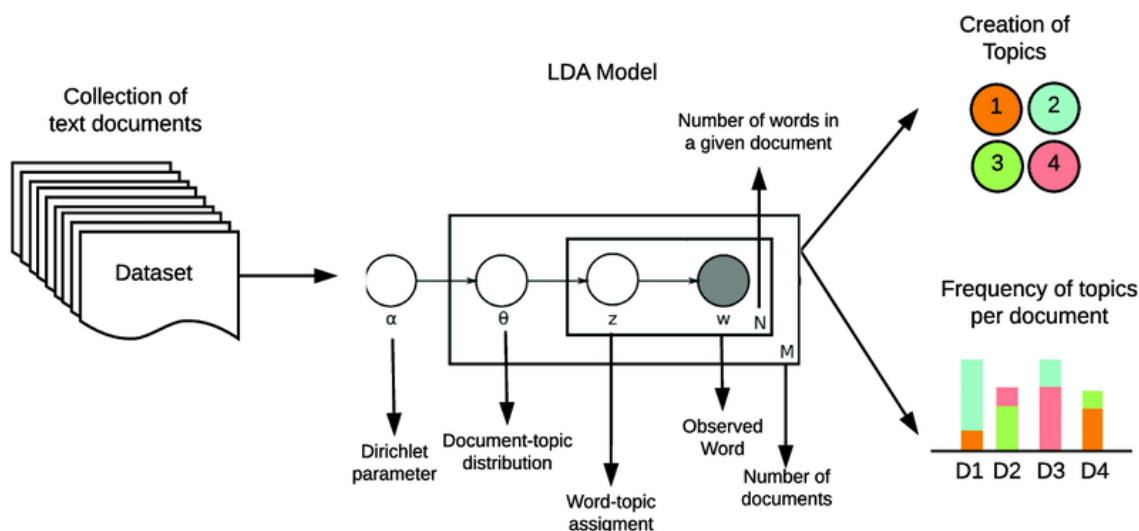
$m_{w,k}$: تعداد واگذاری‌ها، کل مجموعه، کلمه w به موضوع k

η : پارامتر هموارسازی

V : اندازه واژگان یعنی تعداد کلمات مجزا

LDA دارای ۲ ابر پارامتر است: α و η مقدار کم برای α به این معنی است که اسناد فقط تعداد کمی از عناوین را دارند. مقدار زیاد آن نیز به این معنی که اسناد در یک مجموعه شبیه یکدیگر هستند. مقدار کم برای η به این معنی است که موضوعات دارای تعداد کم کلمات کمک‌کننده هستند. مقدار زیادی از η ، به این معنی که موضوعات باهم تداخل دارند و بیشتر شبیه یکدیگر می‌شوند.

الگوریتم LDA به ترتیب کلمات در جمله اهمیت نمی‌دهد. بهتر است از کلماتی که زیاد در پایگاه داده وجود ندارد در این روش حذف شوند. همان‌طور که در مرحله اول ما ایست واژه‌ها را نیز حذف کردیم تا بتوانیم در این روش به نتیجه بهتری برسیم. در شکل ۳.۳ روند کلی LDA نشان داده شده است.



شکل ۳.۳: مدل تخصیص پنهان دیریکلت [۷۴]

در این مرحله هر خروجی از مرحله قبل را به کیسه کلمات تبدیل می‌کنیم. برای به دست آوردن BOW برای هر جمله‌ای که تاکنون آموزش ندیده است فقط کلمات موجود در فرهنگ لغت موجود برای آموزش مدل استفاده می‌کنیم. بنابراین اگر کلمات جدیدی درون جمله وجود داشته باشد هنگام استنباط از توزیع دسته، به آن‌ها توجه نمی‌شود و هیچ خطایی برای کلمات جدید دیده نمی‌شود. اما یک مشکل دارد آن‌هم این است که برخی کلمات ممکن است حذف شوند. برای حل این مشکل مجموعه آموزش را افزایش می‌دهیم و سعی در آموزش بیشتر جملات داریم.

۳.۴.۳ جزئیات روش جنسن شانون و دسته‌بندی

پس از به دست آوردن BOW هر جمله فرکانس‌های جملات را به‌روز می‌کنیم و شناسه‌هایی برای کلمات جدید ایجاد می‌کنیم. هم‌زمان فرکانس‌های جمله را نیز به‌روز می‌کنیم، برای هر کلمه‌ای که در این سند ظاهر می‌شود، فرکانس سند آن را یکی افزایش می‌دهیم.

در نهایت ما یک توزیع دسته برای هر جمله جدید داریم. اگر بخواهیم مشابه‌ترین جملات را در مجموعه بیابیم می‌توانیم این کار را با مقایسه توزیع دسته جمله جدید با تمام توزیع‌های موضوعی جملات در مجموعه انجام دهیم. ما از معیار فاصله جنسن شانون برای یافتن مشابه‌ترین اسناد استفاده می‌کنیم.

واگرایی جنسن شانون یک نسخه متقارن و صاف از واگرایی موازی است. این واگرایی بین دو بردار احتمال p و q به این صورت تعریف شده است:

$$JSD(P \parallel Q) = \frac{1}{2} D(P \parallel M) + \frac{1}{2} D(Q \parallel M) \quad (۳.۳)$$

که M در آن:

$$M = \frac{1}{2}(P + Q)$$

و D واگرایی $Kullback - Leibler$ است:

$$D(P \parallel Q) = \sum_i P(i) \log \frac{P_i}{Q_i} \quad (۴.۳)$$

$$JSD(P \parallel Q) = \frac{1}{2} \sum_i [P_i \log(\frac{P_i}{\frac{1}{2}(P_i) + (Q_i)}) + Q_i \log(\frac{Q_i}{\frac{1}{2}(Q_i) + (P_i)})] \quad (۵.۳)$$

ریشه مربع واگرایی جنسن شانون فاصله جنسن شانون است:

$$\sqrt{JSD(P \parallel Q)}$$

واگرایی هندسی جنسن شانون با در نظر گرفتن میانگین هندسی فرمول بسته‌ای را برای واگرایی بین دو توزیع گوسی ارائه می‌دهد. یک تعریف کلی‌تر، که امکان مقایسه بیش از دو توزیع احتمال را فراهم می‌کند، این است:

$$JSD_{\pi_1, \dots, \pi_n}(P_1, P_2, \dots, P_n) = H(\sum_{i=1}^n \pi_i H(P_i)) \quad (۶.۳)$$

که در آن :

$$p_1 = p_2, p_2 = P_2, \pi_1 = \pi_2 = \frac{1}{2}$$

واگرایی جنسن شانون روشی برای اندازه‌گیری شباهت بین دو توزیع احتمال است که ما در این پژوهش برای اندازه‌گیری شباهت بین جملات استفاده می‌کنیم. برای یافتن مشابه‌ترین جملات در مجموعه، توزیع موضوع سند جدید را با همه توزیع‌های موضوعی اسناد در مجموعه مقایسه می‌کنیم و از معیار فاصله جنسن شانون برای یافتن مشابه‌ترین اسناد استفاده می‌کنیم. آنچه مسافت جنسن و شانون به ما می‌گوید، این است که با مقایسه واگرایی توزیع آن‌ها، کدام جملات از نظر آماری "نزدیک‌تر" و در نتیجه شباهت بیشتری دارند.

هرچه فاصله جنسن شانون کوچک‌تر باشد، دو توزیع مشابه‌تر است و دو جمله شباهت بیشتری به هم دارند. با استفاده از میانگین توزیع احتمالات آن‌ها یا ادغام محتویات هر دو جمله، میانگین سند M را بین دو سند ۱ و ۲ ایجاد می‌کنیم و اندازه هر یک از اسناد ۱ و ۲ با متوسط سند M از طریق $KL(P \parallel M)$ و $KL(Q \parallel M)$ را به دست می‌آوریم. در آخر اختلافات را از نقطه ۲ به طور متوسط انجام می‌دهیم.

پس از این مرحله از ماشین بردار پشتیبان برای به دست آوردن دقت نتایج به دست آمده استفاده می‌شود. علاوه بر انجام طبقه‌بندی خطی، ماشین‌های بردار پشتیبان می‌توانند با استفاده از آنچه

ترفند هسته نامیده می‌شود، یک طبقه‌بندی غیرخطی را انجام دهند. به‌طور ضمنی ورودی‌های آن‌ها را در فضاها دارای ویژگی‌های بسیار بالا نقشه‌برداری می‌کنند. هنگامی که داده‌ها بدون مجوز هستند، یادگیری نظارت‌شده امکان‌پذیر نیست و یک رویکرد یادگیری بدون نظارت مورد نیاز است، که سعی می‌کند خوشه‌بندی طبیعی داده‌ها را به گروه‌ها پیدا کند و سپس داده‌های جدید را به این گروه‌های تشکیل‌شده نقشه ببرد. ماشین بردار پشتیبان با استفاده از برچسب آموزش برای هر گروه، متن جدید را دسته‌بندی کند.

یک الگوریتم یادگیری ماشین معمولی سعی می‌کند مرزی پیدا کند که داده‌ها را به‌گونه‌ای تقسیم کند که خطای طبقه‌بندی غلط به حداقل برسد. SVM بهینه‌ترین مرز تصمیم‌گیری را پیدا می‌کند. بهینه‌ترین مرز تصمیم‌گیری محدوده‌ای است که حداکثر حاشیه را از نزدیک‌ترین نقاط همه کلاس‌ها داشته باشد.

در الگوریتم‌های SVM عملکرد هسته این است که داده‌ها را به‌عنوان ورودی گرفته و آن‌ها را به شکل مورد نیاز تبدیل کند. الگوریتم‌های مختلف SVM از انواع مختلف توابع هسته استفاده می‌کنند. این توابع می‌توانند انواع مختلفی داشته باشند. به‌عنوان مثال خطی^۱، غیرخطی^۲، چندجمله‌ای^۳، تابع پایه شعاعی^۴ و سیگموئید^۵.

هسته چندجمله‌ای یک عملکرد هسته است که با ماشین‌های بردار پشتیبان و سایر مدل‌های هسته‌دار استفاده می‌شود و نشان‌دهنده شباهت بردارها در یک فضای مشخصه نسبت به چند جمله‌ای‌های متغیرهای اصلی است. در شکل ۴.۳ مرز تصمیم‌گیری برای بردار پشتیبان خطی و در شکل ۵.۳ مرز تصمیم‌گیری برای داده‌های غیرخطی قابل تفکیک مشاهده می‌کنیم.

۵.۳ جمع‌بندی

در این فصل روش پیشنهادی خود را به تفصیل و با جزئیات بیان نمودیم. در فصل بعد این روش را با استفاده از سناریوهای ارائه‌شده می‌آزماییم و نتایج را برای هر مرحله ارائه می‌دهیم.

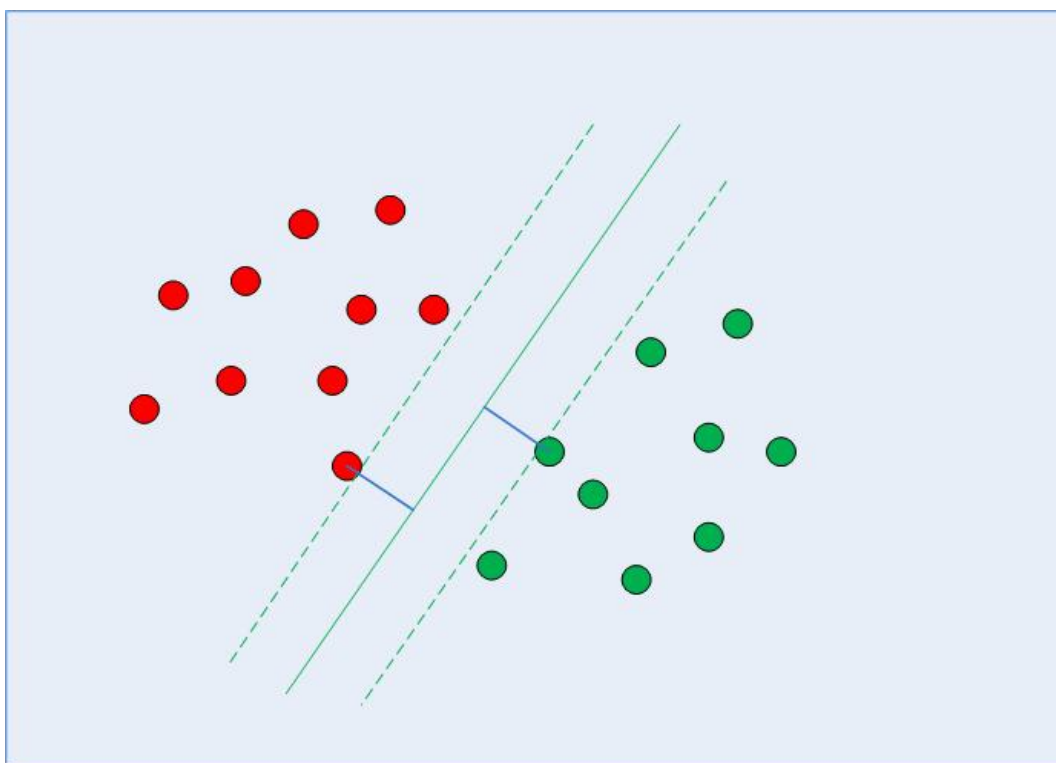
¹linear

²Nonlinear

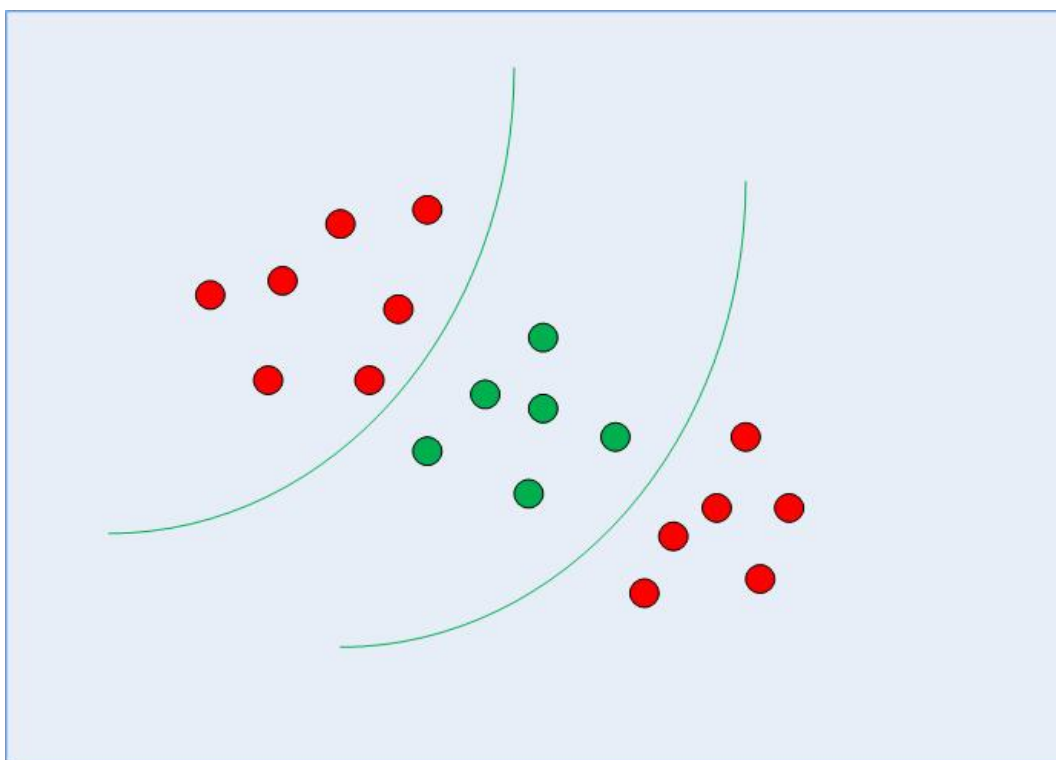
³Polynomial

⁴Radial base function(RBF)

⁵Sigmoid



شکل ۴.۳: مرز تصمیم برای بردار پشتیبان خطی [۷۵]



شکل ۵.۳: مرز تصمیم برای داده‌های غیرخطی [۷۵]

فصل ۴

نتایج و تحلیل آن‌ها

۱.۴ مقدمه

در فصل گذشته به معرفی روش پیشنهادی خودپرداختیم، در این فصل نتایج عملی این روش را شرح خواهیم داد. ابتدا مجموعه دادگان مورد استفاده خود در این پژوهش را معرفی می‌کنیم و پس از آن به جزئیات پیاده‌سازی انجام گرفته می‌پردازیم. در ادامه نتایج رویکرد خود را بررسی کرده و جداول و نمودارهای به دست آمده از تحقیق را توصیف می‌نماییم. در انتها خلاصه‌ای از کارهای صورت گرفته در این فصل را بیان می‌کنیم.

۲.۴ مجموعه دادگان مورد استفاده

بانک اطلاعاتی SICK^۱: جملات داده شامل جملات دانش ترکیبی در وظیفه مشترک SemEval2014 استفاده می‌شود. مجموعه دادگان ما شامل ۱۰۰۰۰ جفت جمله است. ارتباط میان جملات با مقدار بین ۱ تا ۵ نشان داده شده است. این مجموعه داده به عنوان معیار برای ارزیابی فن تشابه جملات استفاده می‌شود. این مجموعه داده حاوی متون نوشته شده به زبان انگلیسی توسط کاربران اینترنت می‌باشد. بخش مهمی از مجموعه داده‌ها از مخازنی مانند اخبار روزانه، مجموعه داده‌های ویکی‌پدیا، توییت‌ها، متن‌ها، عناوین خبری، شرح تصاویر و فیلم‌ها و دیگر موارد است. این مجموعه داده شامل جفت‌های جملات معنایی است که به صورت خودکار انتخاب و دستکاری شده‌اند تا پدیده‌هایی مانند جملات غنی از نظر لغوی، مترادف و متضاد را شامل شوند. هدف از استفاده از این دادگان ارزیابی مدل‌های معنایی توزیع ترکیبی در جملات کامل از طریق ارتباط معنایی و پیام متنی می‌باشد.

مجموعه داده‌های SICK از دو مجموعه پارافراس موجود ساخته شده است: مجموعه داده‌های Image FLickr 8k و مجموعه داده‌های تشابه متن معنایی SemEval 2012 STS MSR-Video Description بین دو عنصر. مجموعه داده‌های قبل از آموزش شامل داده‌های جداگانه جمله‌ای است که برای کارهای قبلی SemEval2013 معنایی متن معنایی ارائه شده است. مجموعه داده‌های قبل از آموزش شامل تقریباً ۱۱۰۰۰ جفت جمله است که دارای برچسب [۱،۵] است.

در مجموعه نهایی همان‌طور که در جدول ۱.۴ نشان داده شده است، نمرات ارتباط به شرح زیر توزیع شد: امتیازات مربوط به ۹۲۳ جفت در محدوده [۱،۲]، ۱۳۷۳ جفت در محدوده [۲،۳]، ۳۸۷۲ جفت در دامنه [۴،۳] و ۳۶۷۲ جفت در محدوده [۴،۵]. در فرآیند برچسب‌گذاری، جملات را به دسته تناقض^۲، مترادف^۳ و خنثی^۴ تقسیم می‌کنند، که منجر به ۵۵۹۵ جفت خنثی، ۱۴۲۴ جفت

^۱ Sentences Involving Compositional Knowledge

^۲ contradiction

^۳ entailment

^۴ neutral

تناقض و ۲۸۲۱ جفت مترادف دلالت شد (طبق جدول ۲.۴)

جدول ۱.۴: طبقه‌بندی مجموعه دادگان مورد استفاده براساس امتیازات

محدوده داده‌ها	تعداد	درصد
محدوده [۱،۲]	۹۲۳	۹/۳
محدوده [۲،۳]	۱۳۷۳	۱۳/۹۵
محدوده [۳،۴]	۳۸۷۲	۳۹/۳۴
محدوده [۴،۵]	۳۶۷۲	۳۷/۳۱

جدول ۲.۴: طبقه‌بندی مجموعه دادگان مورد استفاده براساس برچسب گذاری موجود

برچسب موجود	تعداد
خنثی	۵۵۹۵
متناقض	۱۴۲۴
مترادف	۲۸۲۱

در جدول ۳.۴ نمونه‌ای از مجموعه داده‌ها همراه با ضریب همبستگی موجود مشاهده می‌کنید.

جدول ۳.۴: نمونه‌ای از مجموعه داده‌ها همراه با ضریب همبستگی

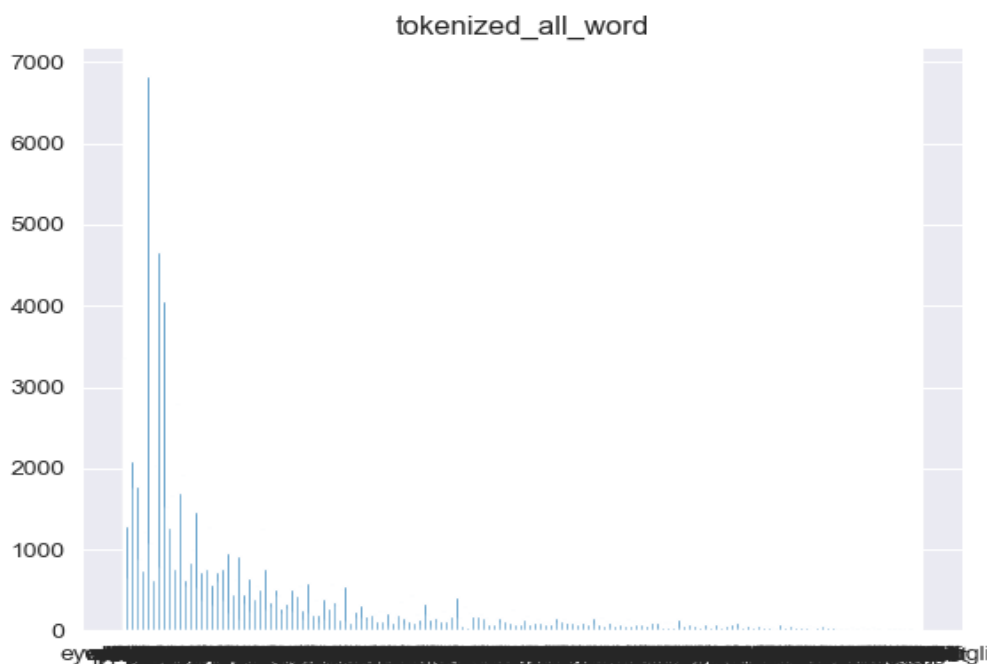
نمره ارتباط	جمله اول	جمله دوم
۴.۹	Two young women are sparring in a kickboxing fight	Two women are sparring in a kickboxing matc
۴	A man with a jersey is dunking the ball at a basketball game	There is no man dunking the ball at a basketball game

علاوه بر داده‌های SICK، مجموعه دادگان SIS65 نیز جهت ارزیابی شباهت معنایی مورد استفاده قرار می‌گیرد. این مجموعه داده شامل ۶۵ جفت جمله می‌باشد. به دلیل بیشتر بودن جملات پایگاه داده SICK و همچنین به دلیل دسترسی رایگان به این مجموعه و استفاده مقالات بسیاری از این پایگاه داده به‌عنوان یک پایگاه داده استاندارد، ما جملات این داده‌ها را مورد ارزیابی قرار می‌دهیم.

۳.۴ پیاده‌سازی

در این بخش به جزئیات پیاده‌سازی انجام‌شده می‌پردازیم.

مجموعه دادگان با فرمت متنی هستند و طبق فرآیندی که گفته شد طبقه‌بندی شده‌اند. ابتدا ایست‌واژه‌ها از جمله حذف می‌شوند. زیرا ایست‌واژه‌ها تعداد رخداد بالا و بار معنایی کمی دارند، بنابراین حذف این واژه‌ها از پایگاه داده حجم محاسبات را کاهش و کارایی را بهبود می‌دهد. همچنین به کمک یک تابع به توکن‌بندی جملات می‌پردازیم. سپس ریشه کلمات رو با ریشه‌یابی lemma به دست می‌آوریم و pos کلمات را در تابعی دیگر به دست می‌آوریم. نمودار کلمات پیش‌پردازش شده در شکل ۱.۴ نشان داده شده است.



شکل ۱.۴: نمودار داده‌های پیش‌پردازش شده

در ادامه فرکانس کلمات را با استفاده از کتابخانه nltk به دست می‌آوریم. دادگان را به مجموعه آموزش و آزمون تقسیم می‌کنیم. مجموعه آموزش برای آموزش مدل LDA استفاده می‌شود، در حالی که مجموعه آزمون برای بازیابی جملات جدید ورودی به مدل ما استفاده می‌شود. برای اجرای LDA روش‌های متعددی وجود دارد که ما از نمونه‌برداری gibbs استفاده خواهیم کرد، زیرا روشی ساده برای درک موضوع است. کتابخانه‌هایی که برای اجرای LDA مورد استفاده قرار می‌گیرند Genism و Scikit-learn هستند. در میان این دو بسته Genism بهتر عمل می‌کند زیرا مدل‌سازی موضوع با Genism مرجع خوبی برای یادگیری LDA است. ما هم از همین کتابخانه برای پیاده‌سازی LDA استفاده کردیم.

برای دستیابی به نتایج دقیق‌تر است مدل را مجدداً آموزش دهیم. در این مرحله جملات را با استفاده از ماژول استخراج‌کننده ویژگی‌ها به کیسه کلمات تبدیل می‌کنیم. برای به دست آوردن BOW از کتابخانه Genism استفاده کردیم و از کلمات موجود در فرهنگ لغت برای آموزش استفاده

کردیم. تمام توزیع‌های دسته‌های به‌دست‌آمده از LDA را در یک ماتریس قرار می‌دهیم. این کار امکان محاسبه سریع و کارآمد را فراهم می‌کند. هدف ما از استفاده کردن از LDA یافتن میزان شباهت جملات باهم است. برای این منظور ما از پیاده‌سازی آنروپی استفاده کردیم. ما فاصله ماتریس به‌دست‌آمده از دسته‌های LDA و ماتریس موجود از کل مجموعه موجود را با کمک فن جنسن شانون که کمترین فواصل مربوط به آن‌ها را در آرایه به دست می‌آورد محاسبه می‌کنیم با این کار می‌توانیم مشابه‌ترین جملات را به دست آوریم. فاصله جنسن شانون یک روش آماری برای به دست آوردن شباهت بین دو توزیع احتمال می‌باشد. پیاده‌سازی این روش با کتابخانه Scipy انجام گرفت.

۱.۳.۴ طبقه‌بندی نمرات شباهت به‌دست‌آمده

برای طبقه‌بندی داده‌ها نیز از ماشین بردار پشتیبان خطی، هسته چندجمله‌ای و هسته تابع پایه شعاعی و دسته‌بندی‌کننده بیز ساده استفاده کردیم. به‌طور ساده روش بیز روشی برای دسته‌بندی پدیده‌ها، بر پایه احتمال وقوع یا عدم وقوع یک پدیده است. شیوه یادگیری در روش بیز ساده از نوع یادگیری با نظارت است. مرحله آخر هر الگوریتم یادگیری ماشین ایجاد ارزیابی برای هسته چندجمله‌ای است. برای این منظور نیز از کتابخانه sklearn کمک گرفته‌ایم. درست است محاسبات پشت SVM پیچیده است، ولی ما با کمک این کتابخانه توانسته‌ایم به راحتی به پیاده‌سازی این طبقه‌بندی بپردازیم. روش بیز ساده دقت کمتری نسبت به ماشین بردار پشتیبان دارد. بنابراین، ماشین بردار پشتیبان خطی روشی مؤثر در این زمینه می‌باشد.

۴.۴ آزمون و پارامترهای مورد نظر

در این بخش به مقایسه روش پیشنهادی خود با سایر روش‌های انجام‌شده در این زمینه می‌پردازیم. جملات موجود در پایگاه داده خود را با روش‌های از جمله فاصله کوسین، فاصله اقلیدسی، فاصله منتهن بررسی می‌کنیم. همچنین روش خود را با بعضی از مقالات فصل دوم مقایسه خواهیم نمود. از جمله روش‌های استفاده‌شده در فصل دوم که از پایگاه داده موردنظر ما استفاده کرده‌اند:

● روش Sent2Vec

● روش PARAGRAM-PHRASE

● روش MaLSTM features با دسته‌بند SVM

● روش DLSTM features با دسته‌بند SVM

همچنین در این فصل مشاهده نمودیم که برای استخراج ویژگی‌ها روش‌های N-gram، Tf-idf، مدل کیسه‌ای از کلمات و word2vec مورد استفاده قرار گرفته‌اند. برای تخمین شباهت به کمک این ویژگی‌ها از روش‌های شبکه‌های عصبی CNN و LSTM، رگرسیون خطی، بیزین، طبقه‌بندی Bloom و K-نزدیک‌ترین همسایه و ترکیبی از این روش‌ها استفاده شده است. هدف ما از انجام این پژوهش مقایسه این روش با روش‌هایی که از پایگاه داده SICK استفاده کرده‌اند می‌باشد. همچنین به مقایسه روش خود با روش‌های زیر پرداختیم.

- ویژگی word2vec با فاصله کوسین

- ویژگی Tf-idf با فاصله منتهن

- ویژگی Tf-idf با فاصله اقلیدسی

دو جفت جمله انتخاب شده را به عنوان نمونه‌ای از پایگاه داده خود انتخاب کردیم، دقت به دست آمده از هر یک از روش‌های بیان شده را بررسی نمودیم. به این منظور دو جفت جمله ۴۱۴ و ۷۳۰ را انتخاب کرده و روش‌های ذکر شده را روی آن‌ها اعمال کردیم. در جدول ۴.۴ داده مورد آزمایش و در جدول ۵.۴ نتایج به دست آمده از این روش‌ها را مشاهده می‌نماییم. روش اول از فاصله کوسین استفاده کرده‌ایم و نتیجه بین ۱ و ۱ می‌باشد. هر چه دو جمله شباهت بیشتری داشته باشند نتیجه به یک نزدیک‌تر است. در روش دوم فاصله منتهن بین ۰ و ۱ به دست آمده است. فاصله جنسن شانون نیز نتیجه‌ای بین ۰ و ۱ در خواهد گرداند. ما تمام نتایج را به معیار فاصله مورد نظر یعنی بین ۱ تا ۵ تبدیل کردیم. همان‌طور که مشاهده می‌کنید روش ما نتایج صحیح‌تری نسبت به بقیه روش‌ها داشته است. دقت به دست آمده فقط برای همین دو جمله آزمایشی به دست آمده است.

جدول ۴.۴: جملات مورد مقایسه ۴۱۴ و ۷۳۰ مجموعه دادگان

نمره ارتباط	جمله B	جمله A	
۴/۷	The man in the purple hat is operating a camera that	The man in the purple hat isn't operating a camera that smakes videos	DATA ۴۱۴
۳/۴	The man in the purple hat is operating a camera that	Two men are taking a break from a trip on a snowy road	DATA ۷۳۰

طبق آنچه که در فصل گذشته توضیح دادیم روش پیشنهادی ما برای تخمین شباهت جفت جمله پس از بررسی ساختارهای متداول موجود در پردازش زبان طبیعی، به ویژه تخمین شباهت متون و جملات، از روش آماری جنسن شانون با ترکیب تخصیص پنهان دیریکلت استفاده کرده است و به وسیله بردار ماشین پشتیبان به طبقه‌بندی نمرات شباهت به دست آمده پرداختیم تا دقت مدل خود را به دست آوریم.

جدول ۵.۴: روش‌هاى مقايسه شده با روش پژوهش

روش مورد استفاده	نمره شباهت روى جملات ۴۱۴	نمره شباهت روى جملات ۷۳۰	دقت
word2vec-featurese + cosian similarity	۴/۰۸	۲/۷۲	۸۳/۴
Tfidf-featurese + manhattan-distance	۳/۷۳	۳/۲۸	۸۷/۹۱
Tfidf-featurese + euclidean-distance	۳/۴۲	۲/۷۱	۷۶/۲۳
LDA + jensen-shanon	۴/۷	۳/۳	۹۸/۵

داده‌هاى ورودى را به آموزش و تست تقسيم نموديم. مجموعه آموزش براى آموزش مدل LDA استفاده مى‌شود و مجموعه تست براى بازيايى جملات مشابه استفاده شدند. طول کوتاه‌ترين جمله در مجموعه دادگان ما ۴ و طول بيشتريين جمله ۳۰ مى‌باشد. ما ۹۸۴۰ جمله داريم که به صورت تصادفى به دادگان آموزش و تست تقسيم کرده‌ايم. در اکثر اجراها تعداد دادگان آموزش بيشتري از ۶۰۰۰ و تقريباً دو برابر داده‌هاى تست مى‌باشد. ما توانسته‌ايم با اين روش زمان پيش‌پردازش داده‌ها را کاهش داده و به ۱۲٪ ثانيه برسانيم. وقتى داده‌ها را با LDA آموزش مى‌دهيم در ۷۸۹٪ ثانيه آموزش انجام مى‌شود.

۵.۴ معيار ارزيايى

قدرت ارتباط بين دو متغير به عنوان آزمون همبستگى شناخته مى‌شود. براى ارزيايى روش پيشنهاده شده و مقايسه با ساير روش‌ها، از معيار ضريب همبستگى پيرسون، ضريب همبستگى اسپيرمن و دقت طبقه‌بندي استفاده مى‌کنيم.

همبستگى پارامترى - همبستگى پيرسون^۱(r): اين يک وابستگى خطى بين دو متغير (x,y) را اندازه‌گيرى مى‌کند که به عنوان يک آزمون همبستگى پارامترى شناخته مى‌شود زيرا به توزيع داده‌ها بستگى دارد. هنگامى که پيرسون به حوزه تشابه معنائى اعمال مى‌شود، درجه همبستگى بين نمرات آزمون‌هاى جفت مجموعه داده‌ها، انسان و نمرات شباهت مشاهده‌شده را نشان مى‌دهد. با اين حال، همبستگى پيرسون نسبت به افراد خارج از مجموعه داده‌هاى معيار حساس است. اگرچه بسيارى از کارها مشکل اصلى را برطرف کرده‌اند، اما تحقيقات کمى در مورد توزيع داخلى سطل‌هاى مجموعه داده‌هاى معيار متمرکز شده است. مجموعه داده‌هاى معيار متن نماينده و توزيع شده،

^۱pearson rank correlation

طیف گسترده‌ای از مقادیر نمرات شباهت را در خود گنجانده است. بنابراین، مجموعه داده‌های معیار می‌تواند یک مجموعه داده مقطعی در نظر گرفته شود. اگرچه یک روش شباهت متن کامل می‌تواند یک همبستگی پیرسون را گزارش کند، اما همبستگی استاندارد پیرسون از مقادیر همبسته جداگانه در سطح مقطع داده آگاه نیست و دلیل آن این است که از محیط خارج می‌شود. بنابراین، با استفاده از روش خطای میانگین مربع مقیاس یافته عادی، که از خطای مقیاس استاندارد استنباط می‌شود، برای از بین بردن محیط‌های دور استفاده می‌کند.

در مباحث آماری، ضریب همبستگی پیرسون یا ضریب همبستگی حاصل ضرب - گشتاور پیرسون میزان همبستگی خطی بین دو متغیر تصادفی را می‌سنجد. مقدار این ضریب بین ۱- تا ۱ تغییر می‌کند که «۱» به معنای همبستگی مثبت کامل، «۰» به معنی نبود همبستگی، و «-۱» به معنی همبستگی منفی کامل است. این ضریب که کاربرد فراوانی در آمار دارد، توسط کارل پیرسون بر اساس ایده‌ی اولیه‌ی فرانسویس گالتون تدوین شد و با فرمول زیر محاسبه می‌گردد:

$$r_{xy} = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{\sqrt{n \sum x_i^2 - (\sum x_i)^2} \sqrt{n \sum y_i^2 - (\sum y_i)^2}} \quad (1.4)$$

همبستگی غیر پارامتری - کندال (tau) و اسپیرمن^۱ (rho): آن‌ها ضرایب همبستگی مبتنی بر رتبه هستند که به‌عنوان همبستگی غیر پارامتری شناخته می‌شوند. در حقیقت برای سنجش ضریب همبستگی بین دو متغیر تصادفی است. مقدار این ضریب نشان‌دهنده‌ی قابلیت بیان یک متغیر به‌صورت تابعی یکنوا از متغیر دیگر است. همبستگی کامل پیرسون (+۱ یا -۱) زمانی اتفاق می‌افتد که متغیری تابعی یکنوا از متغیر دیگر باشد. ضریب همبستگی اسپیرمن به‌صورت ضریب همبستگی پیرسون بین داده‌های رتبه‌بندی شده تعریف می‌شود. به‌عنوان مثال، اگر n زوج داده به‌صورت داده‌شده، ابتدا رتبه هر داده را به‌صورت حساب کرده و سپس ضریب همبستگی اسپیرمن را به‌صورت زیر حساب می‌کنیم:

$$r_s = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)} \quad (2.4)$$

معیارهای صحت^۲، حساسیت^۳، معیار-F^۴ و دقت^۵ از مهم‌ترین معیارهای مورد استفاده برای ارزیابی هستند. دقت معیاری برای نزدیک بودن اندازه‌گیری‌ها به یک مقدار خاص است، درحالی‌که صحت نزدیکی اندازه‌گیری‌ها به یکدیگر است، یعنی لزوماً با یک مقدار خاص نیست. به‌عبارت‌دیگر: اگر مجموعه‌ای از نقاط داده را از اندازه‌گیری‌های مکرر همان مقدار داشته باشیم، گفته می‌شود مجموعه دقیق است اگر میانگین آن‌ها نزدیک به مقدار واقعی اندازه‌گیری شده باشد. از طرف دیگر، اگر مقادیر نزدیک به هم باشند، مجموعه صحت دارد. این دو مفهوم مستقل از یکدیگر هستند، به

¹spearman rank correlation

²Precision

³Recall

⁴F-measure

⁵Accuracy

این معنی که مجموعه داده‌ها می‌توانند دقیق باشند یا صحت داشته باشند یا هر دو باشند یا هیچ‌یک. هرکدام از آن‌ها طبق فرمول‌های زیر محاسبه می‌شوند:

$$Precision = \frac{t_p}{t_p + f_p} \quad (۳.۴)$$

$$Recall = \frac{t_p}{t_p + f_n} \quad (۴.۴)$$

$$F - measure = ۲ * \frac{Precision * Recall}{Precision + Recall} \quad (۵.۴)$$

$$Accuracy = \frac{t_p + t_n}{t_p + t_n + f_p + f_n} \quad (۶.۴)$$

در این روابط، t_p تعداد پاسخ مثبت درست، t_n تعداد پاسخ منفی درست، f_p تعداد پاسخ مثبت غلط و f_n تعداد پاسخ منفی غلط می‌باشد.

۶.۴ نتایج

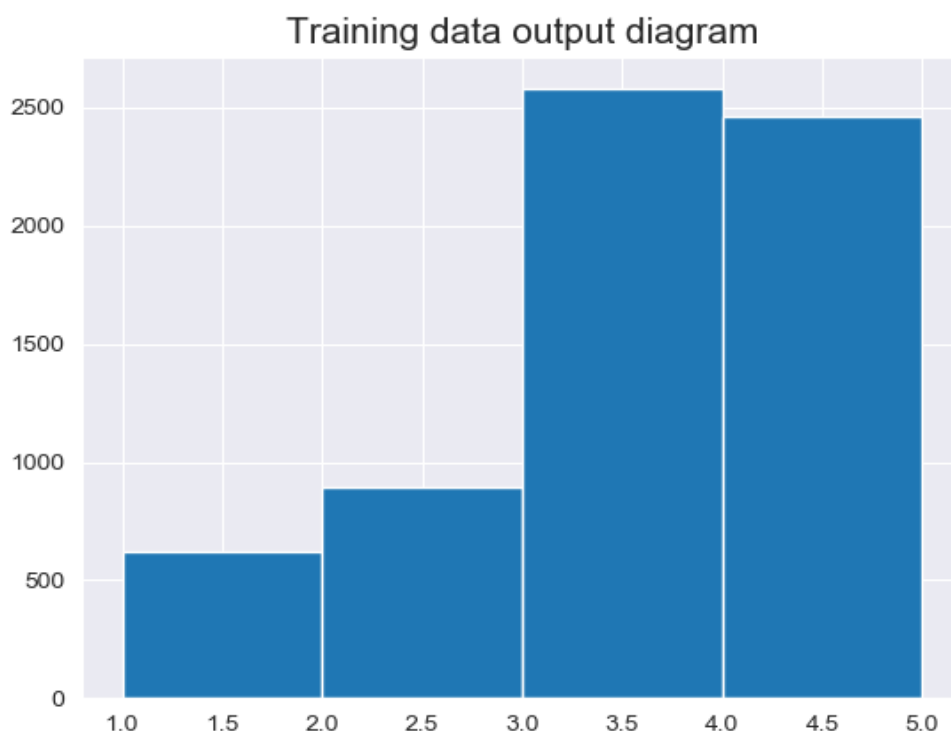
ما آموزش داده‌ها را چندین بار تکرار کردیم و توانستیم مدل خود را به‌خوبی آموزش دهیم. در جدول ۶.۴ نتایج به‌دست‌آمده بر اساس امتیاز ارتباط بین جملات روی داده‌های آموزش، به درصد بیان شده است. همان‌طور که مشاهده می‌کنید ما توانسته‌ایم با روش خود مدل را به‌خوبی آموزش داده و به خطای بسیار کمی برای به دست آوردن نمره ارتباط بین جملات دست پیدا کنیم. بعدازآن با آزمایش مدل نتیجه کار خود را بر روی داده‌های این مجموعه مشاهده می‌کنیم.

جدول ۶.۴: نتایج طبقه‌بندی مجموعه دادگان آموزش براساس امتیازات

محدوده داده‌ها	درصد
محدوده [۱،۲]	۸/۸۹
محدوده [۲،۳]	۱۵/۳۳
محدوده [۳،۴]	۳۹/۰۷
محدوده [۴،۵]	۳۶/۷۱

در شکل ۲.۴ مشاهده می‌کنیم که داده‌های آموزش بر اساس نمرات به ۴ دسته بین ۱ تا ۲ و دسته بعدی بین ۲ تا ۳ و به همین صورت دسته‌بندی شده‌اند.

شکل ۲.۴: نمودار دسته‌بندی داده‌های آموزش



در انتها با ماشین بردار پشتیبان خطی، کرنل چندجمله‌ای و کرنل تابع پایه شعاعی به دسته‌بندی فاصله‌های به‌دست‌آمده پرداختیم. در جدول ۷.۴ با ماشین بردار پشتیبان خطی به‌دقت ۸۹ درصد بر روی داده‌ها دست‌یافته‌ایم. همچنین معیار صحت و حساسیت و معیار F-در هر دسته مشخص شده است.

جدول ۷.۴: ماشین بردار پشتیبان خطی

	Percision	recall	f1-score	Support
۱	۸۹.۰	۷۸.۰	۸۳.۰	۹۱۲
۲	۸۹.۰	۷۸.۰	۸۳.۰	۱۳۵۹
۳	۸۹.۰	۷۸.۰	۸۳.۰	۳۸۳۲
۴	۸۹.۰	۷۸.۰	۸۳.۰	۳۳۵۷
۵	۰۰.۰	۰۰.۰	۰۰.۰	۲۸۲
Accuracy				
avg	۸۹.۰	۸۹.۰	۸۹.۰	۹۷۴۲

در جدول ۸.۴ نیز با ماشین بردار پشتیبان کرنل چندجمله‌ای به دقت ۷۶ درصد بر روی داده‌ها دست یافته‌ایم.

جدول ۸.۴: ماشین بردار پشتیبان کرنل چند جمله‌ای

	Percision	recall	f1-score	Support
۱	۹۰.۰	۹۳.۰	۹۲.۰	۹۱۲
۲	۶۵.۰	۳۰.۰	۴۱.۰	۱۳۵۹
۳	۶۸.۰	۹۵.۰	۸۰.۰	۳۸۳۲
۴	۸۹.۰	۷۵.۰	۸۲.۰	۳۳۵۷
۵	۰۰.۰	۰۰.۰	۰۰.۰	۲۸۲
Accuracy				
avg	۷۶.۰	۷۶.۰	۷۶.۰	۹۷۴۲

جدول ۹.۴ هم با استفاده از ماشین بردار پشتیبان مدل کرنل تابع پایه شعاعی به دقت دسته‌بندی ۸۷ درصد دست یافته‌ایم .

جدول ۹.۴: ماشین بردار پشتیبان کرنل تابع پایه شعاعی

	Percision	recall	f1-score	Support
۱	۹۰.۰	۸۸.۰	۸۹.۰	۹۱۲
۲	۸۰.۰	۹۲.۰	۸۶.۰	۱۳۵۹
۳	۸۷.۰	۹۵.۰	۹۱.۰	۳۸۳۲
۴	۹۰.۰	۸۴.۰	۸۷.۰	۳۳۵۷
۵	۰۰.۰	۰۰.۰	۰۰.۰	۲۸۲
Accuracy				
avg	۸۷.۰	۸۷.۰	۸۷.۰	۹۷۴۲

همانطور که در جدول ۱۰.۴ مشاهده می‌نمایید روش ماشین بردار پشتیبان خطی نسبت به سایر روش‌ها توانسته به دقت بالاتری دست یابد. به همین جهت برای دسته‌بندی از این روش استفاده کردیم.

در زیر ضریب همبستگی پیرسون و اسپرمن را نیز مشاهده می‌کنیم:

$$r_s = 0/925 \quad r_{xy} = 0/891$$

جدول ۱۰.۴: مقایسه دقت بدست آمده روش‌های دسته‌بندی

دقت	روش‌های دسته‌بندی
۰/۸۹	<i>LinerSVM</i>
۰/۷۶	<i>PolynomialSVM</i>
۰/۸	<i>RBF SVM</i>
۰/۷۱	<i>NaiveBayesclassifier</i>

همچنین ما روش خود را با روش‌های موجود در فصل دوم که از مجموعه دادگان SICK استفاده کرده‌اند را مورد مقایسه قرار دادیم. نتایج این مقایسه را در جدول ۱۱.۴ مشاهده می‌نماییم. همانطور که در جدول می‌بینیم روش ما توانسته‌است نسبت *MaLSTM*، به بهبود دقت ۳/۸ درصدی و نسبت به روش *DLSTM*، به بهبود ضریب همبستگی پیرسون دست یابد.

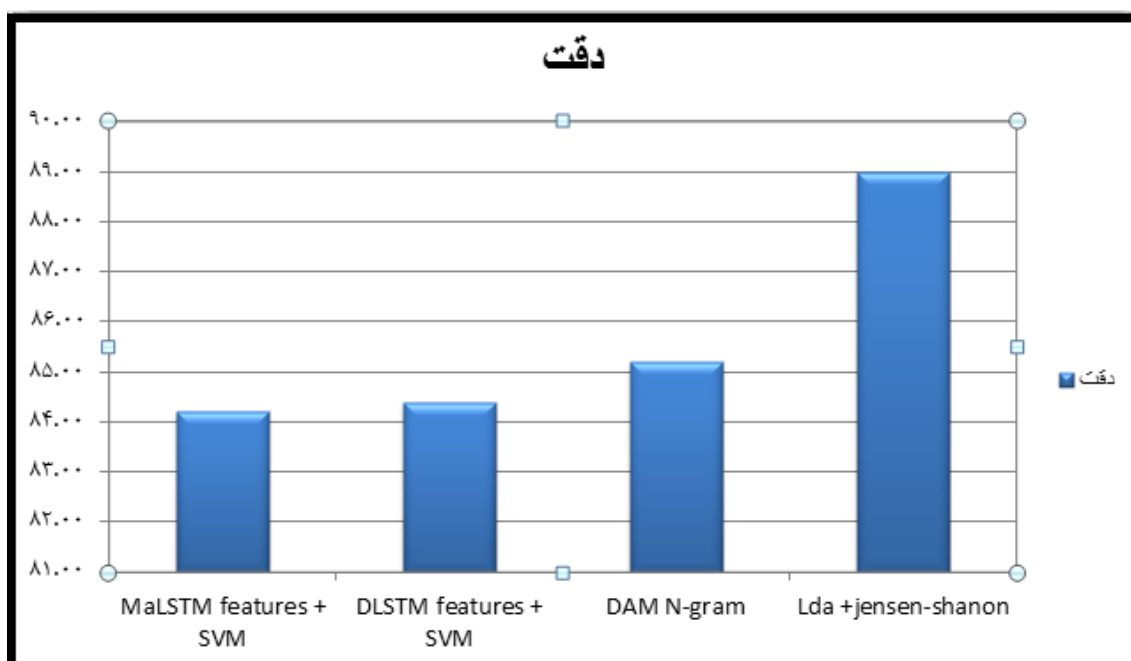
جدول ۱۱.۴: مقایسه روش ارائه شده با سایر روش‌ها

دقت	ضریب همبستگی	روش مورد استفاده
	۰/۷۲	<i>PARAGRAM – PHRASE</i>
۸۱/۴	۰/۷۰	<i>Sent2Vec</i>
۸۵/۲	۰/۸۸۲۲	<i>MaLSTM features + SVM</i>
۸۴/۲	۰/۸۳۰۵	<i>DLSTM features + SVM</i>
۸۹	۰/۸۹۱	<i>LDA + jensen – shanon</i>

در شکل ۳.۴ نمودار مقایسه دقت روش پیشنهادی ما با سایر روش‌هایی که از مجموعه دادگان SICK استفاده کرده‌اند نشان داده شده است. همانطور که می‌بینید دقت این پژوهش نسبت به دیگر روش‌ها بهتر عمل کرده است.

۷.۴ جمع‌بندی

در این فصل به بررسی نحوه پیاده‌سازی روش پیشنهادی خود پرداخته، نتایج به‌دست‌آمده را نیز بیان کردیم. همچنین روش خود را با سایر روش‌ها مقایسه نمودیم. دیدیم که از لحاظ دقت و ضریب همبستگی روش مایکی از بهترین روش‌های موجود میان روش‌های ذکر شده می‌باشد. این روش در به دست آوردن نمره جفت جملات و برچسب کلاس‌ها می‌تواند مؤثر باشد و در کارهای آینده مورد استفاده قرار گیرد.



شکل ۳.۴: نمودار دقت روش‌های بیان شده

فصل ۵

نتیجه گیری

۱.۵ نتیجه‌گیری

در این پژوهش دقت دسته‌بندی داده‌ها اهمیت داشت. در فصل اول به تعریف موضوع پرداختیم. در مورد تاریخچه این مطلب و همچنین کاربردهای شباهت معنایی و لزوم استفاده از آن‌ها صحبت کردیم.

در فصل دوم پیشینه‌ی تحقیق مورد بررسی قرار گرفت. مبانی نظری استفاده‌شده در پژوهش‌های پیشین تعریف گردید. روش‌های زیادی برای استخراج ویژگی‌ها بیان شد از جمله N-گرم‌ها، فرکانس مدت-فرکانس اسناد معکوس، چارچوب کلمه به بردار و دیگر موارد. سپس روش‌های طبقه‌بندی متن و اندازه‌گیری شباهت بیان گردید از جمله روش‌های فاصله کسینوسی، فاصله اقلیدسی، شباهت ژاکارد و فاصله همینگ بودند. همچنین به توضیح روش‌های طبقه‌بندی از جمله روش‌های آماری، شبکه‌های عصبی و روش‌های یادگیری عمیق پرداختیم و دیدیم روش‌های بسیاری برای تخمین شباهت متن وجود دارند از جمله روش‌های آماری، روش‌های معنایی، مبتنی بر رشته و روش‌های ترکیبی که به مواردی از هرکدام اشاره شد و باهم مقایسه شدند.

در فصل سوم به توضیح کار خود پرداخته که ابتدا داده‌های متنی خود را پیش‌پردازش کردیم و از روش آماری تخصیص پنهان دیریکلت و استخراج ویژگی داده‌های پردازش‌شده از فن doc2bow استفاده کردیم که هر جمله را به مجموعه عددی تبدیل می‌کند. سپس از فاصله جنسن شانون برای تخمین شباهت جملات استفاده کردیم که یکی از روش‌های مؤثر در این زمینه می‌باشد. در انتها به دسته‌بندی نتایج با ۳ نوع ماشین بردار پشتیبان خطی، چندجمله‌ای و تابع پایه شعاعی پرداختیم.

در فصل چهارم دادگان مورد استفاده در پژوهش و پیاده‌سازی انجام‌شده را توضیح دادیم. در انتها نتایج به دست آمده یعنی نمره شباهت برای هر ۱۰۰۰۰ جمله را طبقه‌بندی نمودیم تا دقت راهکار خود را بررسی کنیم و به مقایسه روش خود با سایر روش‌هایی که از مجموعه دادگان مورد نظر ما استفاده کرده بودند پرداختیم. دقت اندازه‌گیری یکی از مسائل مهم در این پژوهش است.

در دنیای امروزی داده‌های اینترنتی بیش از پیش مورد استفاده قرار گرفته‌اند و استفاده از کتابخانه‌های دیجیتالی گسترش پیدا کرده است. ممکن است داده‌های مشابه به دلایل مختلفی شکل گرفته باشند. برای بهره‌برداری صحیح از این داده‌ها استفاده از فن‌های یافتن مقدار این شباهت مهم قلمداد می‌شود.

همچنین یکی از روش‌هایی که برای تخمین شباهت معنایی متون مورد استفاده قرار گرفته است روش‌های آماری در ترکیب با سایر روش‌ها می‌باشد. ما نیز از این چارچوب برای پژوهش خود استفاده کرده‌ایم.

مورد دیگری که در این پژوهش وجود دارد، استفاده از پایگاه داده استاندارد می‌باشد که در پژوهش‌های مشابه دیگری نیز مورد استفاده قرار گرفته‌اند و ما به مقایسه روش خود با این روش‌ها پرداخته‌ایم.

۲.۵ کارهای آینده

پژوهش‌های بسیاری می‌تواند در رابطه با این موضوع صورت بگیرد که در زیر تعدادی از این روش‌ها را پیشنهاد می‌دهیم:

- استفاده از داده‌های فارسی به‌جای دادگان انگلیسی و ارائه روشی مناسب برای اینکه دقت تخمین شباهت معنایی داده‌های فارسی را بالا ببریم.
- استفاده از روش‌های یادگیری عمیق برای استخراج ویژگی‌ها از جمله شبکه RBM و شبکه باور عمیق یعنی DBN که می‌تواند نتیجه مطلوبی داشته باشد و یکی از فن‌های مناسب استخراج ویژگی است.
- استفاده از روش‌های یادگیری عمیق و شبکه‌های پیچیده‌تر برای طبقه‌بندی به‌طور مثال شبکه‌های هیبرید عمیق و شبکه باور عمیق و دیگر شبکه‌های عمیق که نتایج مناسبی را برای این کار در پی دارند.
- استفاده از روش‌های ترکیبی برای اندازه‌گیری فاصله بین بردارها مثلاً ترکیب فاصله جنسن شانون و فاصله همینگ .
- استفاده از روش fasttext و glove برای استخراج ویژگی جملات پردازش شده.

مراجع

- [1] Farouk, M. (2019). Measuring sentences similarity: a survey. arXiv preprint arXiv:1910.03940.
- [2] Agirre, E., Banea, C., Cer, D., Diab, M., Gonzalez Agirre, A., Mihalcea, R., Wiebe, J. (2016). Semeval-2016 task 1: Semantic textual similarity, monolingual and cross-lingual evaluation. In SemEval-2016. 10th International Workshop on Semantic Evaluation; 2016 Jun 16-17; San Diego, CA. Stroudsburg (PA): ACL; 2016. p. 497-511.. ACL (Association for Computational Linguistics).
- [3] Cai, Y., Zhang, Q., Lu, W., Che, X. (2018). A hybrid approach for measuring semantic similarity based on IC-weighted path distance in WordNet. Journal of Intelligent Information Systems, 51(1), 23-47.
- [4] Majumder, G., Pakray, P., Gelbukh, A., Pinto, D. (2016). Semantic textual similarity methods, tools, and applications: A survey. Computación y Sistemas, 20(4), 647-665.
- [5] Huang, H., Wu, H., Wei, X., Gao, Y., Shi, S. (2019). Mapping sentences to concept transferred space for semantic textual similarity. Knowledge and Information Systems, 60(3), 1353-1376.
- [6] Vardasbi, A., Faili, H., Asadpour, M. (2017, October). Eigenvalue based features for semantic sentence similarity. In 2017 Artificial Intelligence and Signal Processing Conference (AISP) (pp. 184-189). IEEE.
- [7] Rychalska, B., Pakulska, K., Chodorowska, K., Walczak, W., Andruszkiewicz, P. (2016, June). Samsung Poland NLP Team at SemEval-2016 Task 1: Necessity for diversity; combining recursive autoencoders, WordNet and ensemble methods to measure semantic similarity. In Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016) (pp. 602-608).
- [8] Wali, W., Gargouri, B., Hamadou, A. B. (2017). Enhancing the sentence similarity measure by semantic and syntactico-semantic knowledge. Vietnam Journal of Computer Science, 4(1), 51-60.
- [9] Lopez-Gazpio, I., Maritxalar, M., Lapata, M., Agirre, E. (2019). Word n-gram attention models for sentence similarity and inference. Expert Systems with Applications, 132, .1-11

-
- [10] Leacock, C., Chodorow, M. (1998). Combining local context and WordNet similarity for word sense identification. *WordNet: An electronic lexical database*, 49(2), 265-283.
- [11] Rada, R., Mili, H., Bicknell, E., Blettner, M. (1989). Development and application of a metric on semantic nets. *IEEE transactions on systems, man, and cybernetics*, 19(1), 17-30.
- [12] Wu, Z., Palmer, M. (1994, June). Verbs semantics and lexical selection. In *Proceedings of the 32nd annual meeting on Association for Computational Linguistics* (pp. 133-138). Association for Computational Linguistics
- [13] Jiang, Y., Bai, W., Zhang, X., Hu, J. (2017). Wikipedia-based information content and semantic similarity computation. *Information Processing Management*, 53(1), 248-265.
- [14] Meng, L., Gu, J., Zhou, Z. (2012). A new model of information content based on concept's topology for measuring semantic similarity in WordNet. *International Journal of Grid and Distributed Computing*, 5(3), 81-94.
- [15] Sánchez, D., Batet, M., Isern, D. (2011). Ontology-based information content computation. *Knowledge-Based Systems*, 24(2), 297-303.
- [16] Taieb, M. A. H., Aouicha, M. B., Hamadou, A. B. (2013). Computing semantic relatedness using Wikipedia features. *Knowledge-Based Systems*, 50, 260-278.
- [17] Taieb, M. A. H., Aouicha, M. B., Tmar, M., Hamadou, A. B. (2012, November). Wikipedia category graph and new intrinsic information content metric for word semantic relatedness measuring. In *International Conference on Data and Knowledge Engineering* (pp. 128-140). Springer, Berlin, Heidelberg.
- [18] Jiang, Y., Zhang, X., Tang, Y., Nie, R. (2015). Feature-based approaches to semantic similarity assessment of concepts using Wikipedia. *Information Processing Management*, 51(3), 215-234.
- [19] Petrakis, E. G., Varelas, G., Hliaoutakis, A., Raftopoulou, P. (2006). X-similarity: Computing semantic similarity between concepts from different ontologies. *Journal of Digital Information Management*, 4(4).
- [20] Resnik, P. (1995). Using information content to evaluate semantic similarity in a taxonomy. *arXiv preprint cmp-lg/9511007*
- [21] Rodríguez, M. A., Egenhofer, M. J. (2003). Determining semantic similarity among entity classes from different ontologies. *IEEE transactions on knowledge and data engineering*, 15(2), 442-456.

-
- [22] Tversky, A. (1977). Features of similarity. *Psychological review*,84(4),327.
- [23] Qu, R., Fang, Y., Bai, W., Jiang, Y. (2018). Computing semantic similarity based on novel models of semantic representation using Wikipedia. *Information Processing Management*,54(6)1002-1021
- [24] Nguyen, H. T., Duong, P. H.,Cambria, E. (2019). Learning short-text semantic similarity with word embeddings and external knowledge sources. *Knowledge-Based Systems*, 182, 104842.
- [25] Li, Y., McLean, D., Bandar, Z. A., O'shea, J. D., Crockett, K. (2006). Sentence similarity based on semantic nets and corpus statistics. *IEEE transactions on knowledge and data engineering*, 18(8), .1138-1150
- [26] Arroyo-Fernández, I., Méndez-Cruz, C. F., Sierra, G., Torres-Moreno, J. M., Sidorov, G. (2019). Unsupervised sentence representations as word information series: Revisiting TF-IDF. *Computer Speech Language*,56,107-129.
- [27] Li, Y., McLean, D., Bandar, Z. A., O'shea, J. D., Crockett, K. (2006). Sentence similarity based on semantic nets and corpus statistics. *IEEE transactions on knowledge and data engineering*,18(8),1138-1150.
- [28] Mihalcea, R., Corley, C., Strapparava, C. (2006, July). Corpus-based and knowledge-based measures of text semantic similarity. In *Aaai* (Vol. 6, No.2006,pp.775-780)
- [29] Agirre, E., Cer, D., Diab, M., Gonzalez Agirre, A., Guo, W. (2013). Semantic textual similarity, including a pilot on typed similarity. In *The Second Joint Conference on Lexical and Computational Semantics*.
- [30] Agirre, E., Diab, M., Cer, D., Gonzalez-Agirre, A. (2012, June). Semeval-2012 task 6: A pilot on semantic textual similarity. In *Proceedings of the First Joint Conference on Lexical and Computational Semantics-Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation* (pp. 385-393). Association for Computational Linguistics.
- [31] Agirre, E., Cer, D., Diab, M., Gonzalez-Agirre, A., Guo, W. (2013, June). * SEM 2013 shared task: Semantic textual similarity. In *Second Joint Conference on Lexical and Computational Semantics (SEM), Volume 1: Proceedings of the Main Conference and the Shared Task: Semantic Textual Similarity*(pp32-43).

-
- [32] Neculoiu, P., Versteegh, M., Rotaru, M. (2016, August). Learning text similarity with siamese recurrent networks. In Proceedings of the 1st Workshop on Representation Learning for NLP(pp.148-157).
- [33] Brychcín, T., Svoboda, L. (2016, June). UWB at SemEval-2016 Task 1: Semantic textual similarity using lexical, syntactic, and semantic information. In Proceedings of the 10th International Workshop on Semantic Evaluation(SemEval-2016)(pp.588-594).
- [34] Sultan, M. A., Bethard, S., Sumner, T. (2014). Back to basics for monolingual alignment: Exploiting word similarity and contextual evidence. Transactions of the Association for Computational Linguistics,2,219-230.
- [35] Sultan, M. A., Bethard, S., Sumner, T. (2014, August). Dls@ cu: Sentence similarity from word alignment. In Proceedings of the 8th International Workshop on Semantic Evaluation (SemEval 2014)(pp.241-246).
- [36] Sultan, M. A., Bethard, S., Sumner, T. (2015, June). Dls@ cu: Sentence similarity from word alignment and semantic vector composition. In Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015)(pp.148-153).
- [37] Li Y, McLean D, Bandar ZA, O'shea JD, Crockett K (2006) Sentence similarity based on semantic nets and corpus statistics. IEEE Trans Knowl Data Eng 18(8):1138–1150.
- [38] Mihalcea R, Corley C, Strapparava C (2006) Corpus-based and knowledge-based measures of text semantic similarity. In: Proceedings of the twenty-first national conference on artificial intelligence and the eighteenth innovative applications of artificial intelligence conference, Boston, Massachusetts, USA,pp 775–780.
- [39] Sultan, M. A., Bethard, S., Sumner, T. (2015, June). Dls@ cu: Sentence similarity from word alignment and semantic vector composition. In Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015) (pp. .(148-153
- [40] Amir, S., Tanasescu, A., Zighed, D. A. (2017). Sentence similarity based on semantic kernels for intelligent text retrieval. Journal of Intelligent Information Systems,48(3),675-689.
- [41] Crystal D. , The Cambridge Encyclopedia of Language, 2nd edition, Cambridge University Press, 1996. ISBN 0-521-55967-7 Manning,C. D.
- [42] Schutze, H. , Foundations of Statistical Natural Language Processing, 5th edition, The MIT Press, 2002.ISBN 0-262-13360-1

- [43] Borda, Monica (2011). *Fundamentals in Information Theory and Coding*. Springer. ISBN 978-3-642-20346-6.
- [44] Han, Te Sun Kobayashi, Kingo (2002). *Mathematics of Information and Coding*. American Mathematical Society. ISBN .978-0-8218-4256-0
- [45] Shirakawa, M., Nakayama, K., Hara, T., Nishio, S. (2015). Wikipedia-based semantic similarity measurements for noisy short texts using extended naive bayes. *IEEE Transactions on Emerging Topics in Computing*,3(2),205-219.
- [46] Broder, A. Z. (1997). Syntactic clustering of the Web, *J of Computer Networks and ISDN Systems*, Volume 29,Issue .8-13
- [47] Greenberg, R. I. (2003). Bounds on the number of longest common subsequences. arXiv preprint cs/0301030.
- [48] Sivic, J., Zisserman, A. (2008). Efficient visual search of videos cast as text retrieval. *IEEE transactions on pattern analysis and machine intelligence*, 31(4), 591-606.
- [49] Schölkopf, B., Smola, A. J. (2002). *Learning with kernels*. MIT Press, Cambridge, MA. Learning with kernels. MIT Press, Cambridge, MA.
- [50] Chatterjee, N. (2001, February). A statistical approach for similarity measurement between sentences for EBMT. In *Proceedings of Symposium on Translation Support Systems STRANS-2001*.
- [51] Cavalcanti, A. P., de Mello, R. F. L., Ferreira, M. A. D., Rolim, V. B., Tenório, J. V. S. (2017, October). Statistical and semantic features to measure sentence similarity in portuguese. In *2017 Brazilian Conference on Intelligent Systems (BRACIS)* (pp. 342-347).IEEE.
- [52] Majumder, G., Pakray, P., Avendaño, D. E. P. (2018, January). Semantic Textual Similarity Using Lexical and Cosine Similarity. In *Annual Convention of the Computer Society of India* (pp. 717-732). Springer, Singapore.
- [53] He, H., Wieting, J., Gimpel, K., Rao, J., Lin, J. (2016, June). UMD-TTIC-UW at SemEval-2016 Task 1: Attention-based multi-perspective convolutional neural networks for textual similarity measurement. In *Proceedings of the 10th International Workshop on Semantic Evaluation(SemEval-2016)*(pp.1103-1108).
- [54] Mueller, J.,Thyagarajan, A. (2016, March). Siamese recurrent architectures for learning sentence similarity. In *Thirtieth AAAI Conference on Artificial Intelligence*.

-
- [55] Wieting, J., Bansal, M., Gimpel, K., Livescu, K. (2016). Charagram: Embedding words and sentences via character n-grams. arXiv preprint arXiv:1607.02789.
- [56] Zhu, W., Yao, T., Ni, J., Wei, B., Lu, Z. (2018). Dependency-based Siamese long short-term memory network for learning sentence representations. PloS one, 13(3),e0193919.
- [57] Pontes, E. L., Huet, S., Linhares, A. C., Torres-Moreno, J. M. (2018). Predicting the semantic textual similarity with siamese CNN and LSTM. arXiv preprint arXiv:1810.10641.
- [58] Orasan, C., Mitkov, R., Ranasinghe, T. (2019, September). Semantic textual similarity with siamese neural networks. RANLP.
- [59] Li Y, McLean D, Bandar ZA, O'shea JD, Crockett K (2006) Sentence similarity based on semantic nets and corpus statistics. IEEE Trans Knowl Data Eng 1150–1138:(8)18
- [60] Pawar, A., Mago, V. (2018). Similarity between learning outcomes from course objectives using semantic analysis, blooms taxonomy and corpus statistics. arXiv preprint arXiv:1804.06333.
- [61] Bär, D., Biemann, C., Gurevych, I., Zesch, T. (2012, June). Ukp: Computing semantic textual similarity by combining multiple content similarity measures. In Proceedings of the First Joint Conference on Lexical and Computational Semantics-Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (pp. 435-440). Association for Computational Linguistics.
- [62] Atoum, I., Otoom, A. (2016). Efficient hybrid semantic text similarity using wordnet and a corpus. Int. J. Adv. Comput. Sci. Appl, 7(9), 124-130.
- [63] Baba, K., Nakatoh, T., Minami, T. (2017). Plagiarism detection using document similarity based on distributed representation. Procedia computer science, 111, 382-387
- [64] Islam, A., Inkpen, D. (2008). Semantic text similarity using corpus-based word similarity and string similarity. ACM Transactions on Knowledge Discovery from Data (TKDD), 2(2), 10
- [65] Roostaei, M., Fakhrahmad, S. M., Sadreddini, M. H., Khalili, A. (2014). Efficient calculation of sentence semantic similarity: a proposed scheme based on machine learning approaches and NLP techniques. Scientific Journal of Review, 3(3), 94-106.
- [66] Afzal, N., Wang, Y., Liu, H. (2016). Mayonlp at semeval-2016 task 1: Semantic textual similarity based on lexical semantic net and deep learning semantic model. In Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016) (pp. 674-679).

- [67] Camacho-Collados, J., Pilehvar, M. T., Collier, N., Navigli, R. (2017, August). Semeval-2017 task 2: Multilingual and cross-lingual semantic word similarity. In Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017) (pp.15-26).
- [68] Pagliardini, M., Gupta, P., Jaggi, M. (2017). Unsupervised learning of sentence embeddings using compositional n-gram features. arXiv preprint arXiv:1703.02507.
- [69] Pawar, A., Mago, V. (2018). Calculating the similarity between words and sentences using a lexical database and corpus statistics. arXiv preprint arXiv:1802.05667.
- [70] Wieting, J., Bansal, M., Gimpel, K., Livescu, K. (2016). Charagram: Embedding words and sentences via character n-grams. arXiv preprint arXiv:1607.02789.
- [71] Rakib, R. H., Islam, A., Milios, E. (2018). Improving text relatedness by incorporating phrase relatedness with word relatedness. *Computational Intelligence*, 34(3), 939-966.
- [72] Glavaš, G., Franco-Salvador, M., Ponzetto, S. P., Rosso, P. (2018). A resource-light method for cross-lingual semantic textual similarity. *Knowledge-Based Systems*, 143, 1-9.
- [73] Plaza-del-Arco, F. M., Molina-González, M. D., Ureña-López, L. A., Martín-Valdivia, M. T. (2021). Comparing pre-trained language models for Spanish hate speech detection. *Expert Systems with Applications*, 166, 114120.
- [74] Buenaño-Fernandez, D., González, M., Gil, D., Luján-Mora, S. (2020). Text Mining of Open-Ended Questions in Self-Assessment of University Teachers: An LDA Topic Modeling Approach. *IEEE Access*, 8, 35318-35330.
- [75] Schölkopf, B., Smola, A. J. (2002). *Learning with kernels*. MIT Press, Cambridge, MA. Learning with kernels. MIT Press, Cambridge, MA
- [۷۶] شیرازی حسین، (۱۳۹۱)، پنجمین کنفرانس داده کاوی ایران، ”ارائه رویکرد جدید در محاسبه شباهت معنایی جملات فارسی”، ص ۵۷، تهران
- [۷۷] ملکوتی ن، (۱۳۹۱)، ”ارائه یک روش جدید برای بررسی میزان شباهت اسناد متنی”، دانشکده مهندسی برق و کامپیوتر، دانشگاه شیراز

Abstract

In today's world, due to the increasing volume of internet information and the multiplicity of different digital cultures, the existence of similarities in textual data has increased for various reasons. For this reason, estimating the degree of similarity between them is necessary. Finding similarities between textual data is also used in areas such as information retrieval systems, plagiarism, data mining, document classification, and much more.

By entering a sentence or document or text into the system, its similarity with existing documents is checked and Used in the field in question. So far, various methods of natural language processing and machine learning have been introduced to calculate the similarity of textual data. These methods have achieved different degrees of accuracy, and better research is needed to improve this criterion. Among the statistical methods, which is one of the corpus-based methods, and by combining this method with other methods, interesting results can be achieved.

In this research, data is applied to the preprocessing machine before converting the text to readable format to reduce the features and increase the measurement accuracy. Each word is then assigned a correlation score using Latent Dirichlet Allocation, and this process is repeated based on the probabilities to improve the assignment of the correlation score to the words and the sentence is placed in the desired category. Doc2bow is also used to representation sentences. Then, to calculate the most similar sentences, the Jensen Shannon distance is used, which is obtained by comparing the divergence of the desired label distribution. Finally, the machine uses the linear support vector and the Radial base function to classify similar sentences. The proposed method achieved a classification accuracy of 89 percent and a correlation score of 0.92 percent on the SICK database and was able to be 3.8 percent better accurate than the MaLSTM research.

Keywords: text semantic similarity, Latent Dirichlet Allocation, document to Bag of words, Jensen-Shannon distance, Support vector machines



Shahrood University of Technology

Faculty Of Computer Engineering

MSc Thesis in: Artificial Intelligence Engineering

**Estimating semantic based text similarity by
using statistical methods**

By: Fereshteh Riahi

Supervisor

Morteza Zahedi

Advisor

AliReza Tajari

February 2021