

کد مقاله ۵۶

تجزیه و تحلیل احساسات متن با استفاده از نظریه ساختار بیانی

شیماسیار^۱، ایوب باقری^۲، احمد یوسفان^۳

^۱ دانشجوی کارشناسی ارشد، دانشکده مهندسی برق و کامپیوتر، دانشگاه کاشان، کاشان،
Sayyarshima@grad.kashanu.ac.ir

^۲ استادیار، دانشکده علوم اجتماعی، دانشگاه اوترخت، هلند
a.bagheri@uu.nl

^۳ مربی، دانشکده مهندسی برق و کامپیوتر، دانشگاه کاشان، کاشان
yoosofan@kashanu.ac.ir

چکیده

امروزه تجزیه و تحلیل احساسات یکی از برجسته ترین زمینه های تحقیقاتی در حوزه ی پردازش زبان طبیعی است و هدف این مقاله طبقه بندی احساس متن به دو دسته مثبت و منفی است. نادیده گرفته شدن اطلاعات ساختاری گفتار همواره چالش اصلی در روش های سنتی طبقه بندی متن است. نظریه ساختار بیانی محتویات یک متن را به بخش هایی با سطوح مختلف اهمیت تقسیم می کند و ارتباط سلسله مراتبی آن را در نظر می گیرد.

در این پژوهش نظرات کاربران در مورد محصولات یک وب گاه بررسی می شود. متن هر نظر ابتدا به درخت گفتمان تجزیه می شود سپس یک طرح وزنی تعریف می شود و احساس متن محاسبه می شود. هدف این پژوهش ارائه روشی است که بر اساس جایگاه معنایی کلمات کار کند و احساس متون را در سطح سند با استفاده از نظریه ساختار بیانی و با تکیه بر قواعد وزنی بدست آورد. هر چند این قواعد نیاز به درک و مطالعه زبان شناسی دارد، با این وجود سعی شده است با در نظر گیری پارامتری جدید در قواعد وزنی، معیارهای ارزیابی (صحت و F) بهبود داده شود. با توجه به آزمایش های انجام شده بهترین نتیجه مقدار صحت ۶۷/۸ درصد و معیار F ۷۱/۹ درصد است. این نتایج نسبت به دو مقاله مشابه پیشین مقدار صحت را ۸ درصد بهبود داده است.

کلمات کلیدی

تجزیه و تحلیل احساسات، مبتنی بر قاعده، نظریه گفتمان بیانی، واحد گفتمان ابتدایی.

دغدغه های مردم و افکارشان پی برد. وب یکی از اصلی ترین فضاها برای بیان نظرات و عقاید درباره انواع موضوع ها از جمله محصول ها و سرویس های دریافتی است. شناخت احساسات نویسندگان و دانستن تجربیات آنها در مورد مسائل مختلف، می تواند برای فعالیت های تجاری و دیگر زمینه ها بسیار ارزشمند باشد. این تجربیات، بخش مهمی از اطلاعات را برای فرآیند تصمیم

۱- مقدمه

امروزه با توسعه شبکه جهانی وب و پیدایش وب گاه هایی نظیر آمازون، توئیتر و همانند آن و دیگر ابزارهای اینترنتی و شبکه های اجتماعی می توان به

تغییر فرهنگ لغت تجزیه و تحلیل می‌کند. در واقع هدف این وزن‌دهی ایجاد تمایز بین عبارات مهم‌تر و کم‌اهمیت‌تر متن و ارائه طبقه‌بندی دقیق‌تر است. نتایج نشان می‌دهد که روش‌هایی که این نظریه را بکار می‌برند نسبت مدل‌های پایه‌ای که از روش ساختار بیانی استفاده نمی‌کنند و به تنهایی روی احساسات واژگان کار می‌کنند بهتر عمل می‌کند اما تحقیقات قبلی در مورد تحلیل احساسات متن با استفاده از این نظریه کم است.

ساختار این مقاله به این صورت تنظیم شده است که بخش ۲ مروری بر تحقیقات گذشته است. بخش ۳ بیان می‌کند که چگونه نظریه ساختار بیانی معنای یک سند را آشکار می‌کند. پس از آن در بخش ۴ چارچوب کلی روش پیشنهادی ارائه می‌شود. بخش ۵ نتایج حاصل از پیاده‌سازی و مقایسه با سایر روش‌ها بیان می‌شود و نهایتاً در بخش ۶ نتیجه‌گیری انجام می‌شود.

۲- مروری بر تحقیقات پیشین

پردازش زبان طبیعی در کاوش متن به شدت مورد بهره‌برداری قرار گرفته است. اما پردازش گفتمان هنوز به صورت پیشرفته فراگیر نیست [۹]، لذا تحقیقات در تجزیه و تحلیل احساسات بر مبنای نظریه ساختار بیانی کم است [۷]. تحلیل گفتمان می‌تواند به عنوان حوزه‌ای از زبان‌شناسی تعریف می‌شود که با چگونگی ساختن معنا در متن یا پاراگراف به جای یک جمله واحد سروکار دارد [۱۰]. گفتمان یکی از سخت‌ترین سطح‌های زبان برای پردازش خودکار با توجه به پیچیدگی آن است. در حقیقت، این دشواری تحلیل گفتمان خودکار را به چالش تبدیل کرده و نیازمند توسعه ابزارهای پردازش زبان طبیعی است.

تجزیه و تحلیل احساسات با توجه به عمق مساله و توجه به جزئیات در سه سطح سند، جمله و جنبه بررسی می‌شوند. در این مقاله تجزیه و تحلیل احساسات در سطح سند انجام شده است و منظور از تحلیل در سطح سند، بررسی این موضوع است که کل سند دارای احساس مثبت یا منفی است.

برای بیش از یک دهه روش‌های بسیاری برای تجزیه و تحلیل متن بکار می‌روند. این روش‌ها عمدتاً به دو دسته تقسیم می‌شوند.

دسته اول روش‌های مبتنی بر یادگیری ماشین هستند که از آنها به عنوان روش‌های نظارتی نیز یاد می‌شود. این روش‌ها از الگوهای نمایش برداری متن استفاده می‌کنند و تلاش می‌کنند تا یک طبقه‌بندی احساس در متن براساس فرکانس وقوع کلمات مختلف آموزش دهند [۱۱]. روش‌های ماشین بردار پشتیبان (SVM)، آنتروپی، k نزدیکترین همسایه در این دسته قرار می‌گیرند. Zim و همکارانش از اطلاعات معنایی و روابط گفتمان برای تجزیه و تحلیل احساسات در سطح جمله استفاده کرده‌اند و شبکه‌های منطق مارکوف را برای طبقه‌بندی بکار برده‌اند [۱۲]. یک مطالعه با بهره‌گیری نظریه گفتمان بیانی شبکه عصبی LSTM را روی درخت وابستگی RST مدل کرده و طبقه‌بندی نظرات را در سطح سند انجام داده است [۱۳].

نتایج تجربی نشان می‌دهد که روش‌های مبتنی بر یادگیری اغلب صحت بالاتری نسبت به روش‌های مبتنی بر واژگان دارند (این روش‌ها در

گیری تشکیل می‌دهد و می‌تواند مزایای رقابتی را برای کسب و کار به ارمغان بیاورد، بنابراین تجزیه و تحلیل نظرات و فهم احساس آن اهمیت بسزایی دارد [۱]. تجزیه و تحلیل متن در حوزه‌های مختلف از جمله بازارهای مالی [۲]، سیستم‌های اقتصادی و مدیریت برندها در سازمان‌ها کاربرد دارد [۳].

در حوزه‌ی متن‌کاوی معمولاً تجزیه و تحلیل احساسات آنرا نظر‌کاوی^۳ می‌نامند. منظور از تجزیه و تحلیل احساسات، روش‌ها و برنامه‌های کاربردی از زبان‌شناسی محاسباتی است که به شناسایی و استخراج مفهوم ذهنی در متون می‌پردازد [۴]. پیشرفت در تجزیه و تحلیل احساسات خودکار اهمیت حیاتی برای سیستم‌های پشتیبانی و تصمیم‌گیری امروز دارد. به کمک این سیستم‌ها، سازمان‌ها می‌توانند به شیوه‌ای به موقع و موثر با راهبردهای مناسب بازاریابی و روابط عمومی پاسخ دهند.

رایج‌ترین روش برای تجزیه و تحلیل متن، مدل‌های به اصطلاح سبد کلمات^۴ است که فراوانی کلمات یا ترکیب‌های n تایی‌شان را در نظر می‌گیرد [۵]. این روش‌ها معمولاً به طور انحصاری به تعداد تکرار کلمات یا عبارات خاص مربوط تکیه می‌کنند. بدین ترتیب، روابط معنایی واقعی و ساختار متن را بین بخش‌های سند، جملات یا حتی عبارت‌ها نادیده می‌گیرند [۶،۷].

همچنین یک سند متنی شامل بخش‌ها، پاراگراف‌ها یا جملاتی است که گاهی اوقات این بخش‌ها برچسب‌هایی متفاوت با یکدیگر دارند. برای نمونه: «در حالی که جان همیشه شکایت می‌کند که از این نوع فیلم‌ها متنفر است، ولی او اعتراف کرد که از این فیلم لذت برده است.» در نظر فوق تعداد کلمات منفی بیشتر از تعداد کلمات مثبت است، در حالی که احساس آن مثبت است.

در متن «این فیلم فوق‌العاده بود، هرچند پف‌فیل مانده بود.» واضح است که احساس متن نسبت به فیلم مثبت است، در حالی که نظر در مورد پف‌فیل ذکر شده منفی است و جمله‌ی دوم نظری نامرتبط به نظردهی درباره‌ی فیلم بیان کرده است [۸].

بنابراین جملاتی مانند نمونه‌های بالا ایجاب می‌کند که ساختار متن را بدانیم و به چنین عبارات زایدی نقش کمتری در تعیین احساس کلی متن بدهیم، همچنین مدل‌سازی روابط بین بخش‌های سند و استفاده صحیح از روابط بین این بخش‌ها منجر به طبقه‌بندی بهتر و برچسب زدن دقیق‌تر سند نیز می‌شود. برای غلبه بر محدودیت‌ها و چالش‌های قبلی، یک روش جالب جدید بکارگیری نظریه ساختار بیانی برای درک معنی و مفهوم متن نظرات است.

بدین ترتیب این مقاله، رویکردی جدید برای نفوذ به اطلاعات و فهم ساختار سند با استفاده از نظریه ساختار بیانی پیشنهاد کرده است و یک طرح وزن‌دهی ارائه داده است که با توجه به هر سه ویژگی ساختار گفتمان یعنی عمق، نوع گره و نوع ارتباط بین واحدها، اقدام به محاسبه احساس متن می‌کند. این پژوهش، احساسات را با ترکیب ویژگی عمق به روابط گفتمان و

پارامتر عمق را نیز در نظر گرفته، در آن ترکیب کرده و تاثیر آن را مورد مطالعه قرار دادیم. علاوه بر این ما سعی در توسعه روابط بیانی در این تحقیق کرده و از فرهنگ لغت کامل تری نسبت به آنچه در تحقیق گفته شده است، استفاده کردیم.

۳- نظریه ساختار بیانی^۸

نظریه ساختار بیانی (RST) یک نظریه گفتار محبوب در پردازش زبان طبیعی است که توسط Mann و همکارانش در سال ۱۹۸۸ پیشنهاد شده است. این نظریه یکی از ابزارهای مفید برای توصیف متن و مشخص کردن ساختار آن است که به وضوح می‌تواند اهمیت بخش‌های مختلف متن را تعیین کند [۲۰]. بر اساس این نظریه، متن به واحدهای کوچکتری تبدیل می‌شود که واحد گفتمان ابتدایی^۹ نامیده می‌شود. این واحدها، کوچکترین واحد اطلاعاتی گفتمان و بخش تجزیه‌ناپذیر جمله را شکل می‌دهد. واحدهای گفتمان با توجه به نوع ارتباطی که با یکدیگر دارند، برچسب گذاری می‌شوند و به دو قطعه متنی غیرهم‌پوشان، هسته و پوسته محدود می‌شوند. هسته حامل اطلاعات اصلی در مورد موضوع است و پوسته اطلاعاتی اضافی و پشتیبان بسته به نوع رابطه‌ای که با هسته یا پوسته دیگر دارد، حمل می‌کند. در این نظریه همواره این فرض برقرار است که هسته نسبت به پوسته تاثیر بیشتری در فهم و تفسیر متن دارد [۹].

هر واحد گفتمان در یک رابطه سلسله مراتبی به دیگر واحدهای گفتمان، بسته به ساختار سند متصل می‌شود. این سلسله مراتب در نهایت به یک درخت دودویی برای ساده کردن پردازش‌های محاسباتی نگاشت می‌شود که درخت گفتمان نامیده می‌شود [۲۰].

در این درخت، گره‌های داخلی با توجه به گفتمان معنایی خود، به یکی از ۱۸ نوع ارتباطی که در جدول (۱) بیان شده است، برچسب می‌خورند. برای مثال ردیف ۳ رابطه‌ی علت بین دو گره، نشان دهنده این است که متن یک گره، دلیل آورده شده برای رویدادی است که در گره دوم اتفاق افتاده است. روابط بین واحدهای گفتمان برای درک متن بسیار مهم هستند، زیرا آنها اطلاعاتی را در مورد چگونگی ارتباط بخش‌های متنی با یکدیگر برای تشکیل یک گفتمان منسجم می‌دهند. البته درخت گفتمان بسته به نوع تجزیه‌گر و تعداد ارتباطاتی که بین واحدها تعریف می‌شود، متفاوت است.

جدول (۱): نوع ارتباط بین واحدهای گفتمان [۲۰]

نوع ارتباط	توضیحات
پیش زمینه	پوسته اطلاعاتی برای درک هسته فراهم می‌کند.
نسبت دادن	پوسته توضیحاتی با پیش فرض‌هایی برای شناخت هسته است.
علت	پوسته منجر به نتیجه‌ای می‌شود که در هسته بیان شده است.
مقایسه	به شباهت‌ها و تفاوت‌ها بین هسته و پوسته اشاره دارد.
شرایط	موقعیت‌های غیرقابل درک یا فرضی، آینده که تحقق آنها بر تحقق رویدادهای ارائه شده در هسته تاثیر می‌گذارد.
پایملا	پوسته اطلاعاتی در مورد اثرات رویدادهای ارائه شده در هسته دارد.
تضاد	پوسته مقایسه یا تضادها را توصیف می‌کند.
جزئیات	اطلاعات اضافی موضوعات موجود در هسته در پوسته بیان می‌شود.

ادامه بحث می‌شود)، اما بسیار وابسته به دامنه و اندازه مجموعه داده هستند که غالباً موجب ناکارآمدی و وابستگی تحلیل‌ها به دامنه در سطح سند می‌شود [۱۴]. روش‌های مبتنی بر یادگیری هنگامی که داده‌های آموزشی ناکافی یا بایاس شده باشند شکست می‌خورد، در حالی که روش‌های مبتنی بر واژگان بسیار رقابتی هستند زیرا آنها نیاز به تلاش کمتری در ایجاد سند برچسب دار انسان دارند و به کمیت و کیفیت مجموعه داده‌های آموزشی حساس نیستند [۱۱]. لذا از روش‌های غیرنظارتی استفاده کردیم.

روش‌های غیر نظارتی یا روش‌های مبتنی بر واژگان به عنوان رویکرد دیگر در تحلیل احساسات هستند. چنین روش‌هایی از یک فرهنگ لغت (واژگان) برای تحلیل متن از منظر احساس کمک می‌گیرند و با جستجو در فرهنگ لغت مربوطه احساس کلمات را ارزیابی و وزن دهی می‌کنند سپس با محاسبه‌ی فراوانی آن، احساس متن را برچسب گذاری می‌کنند. غالباً طبقه‌بندی را با استفاده از قاعده "اگر .. آنگاه .." انجام می‌دهند [۱۵].

نخستین بار در سال ۲۰۰۷، Voll و Taboada ثابت کردند که نظریه گفتمان بیانی می‌تواند نتایج تحلیل احساسات مبتنی بر قاعده را بهبود بخشد. آنها متن را با در نظر گرفتن اینکه صفات در یک متن درجه اهمیت یکسانی ندارند و صفات موجود در یک بخش خاص یا اشاره‌کننده به جنبه خاصی از متن اهمیت بیشتری دارند، وزن دهی کردند و حاشیه‌نویسی متن را به صورت دستی انجام دادند [۱۶] هرچند حاشیه‌نویسی دستی روابط بسیار زمان‌بر و هزینه‌بر بود.

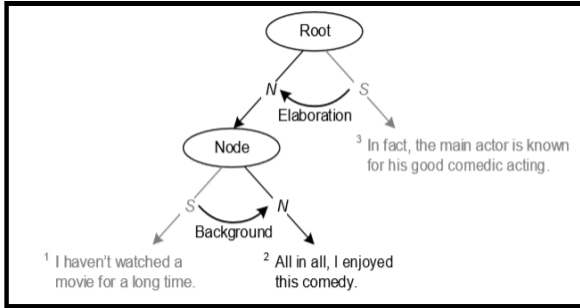
Heerschop و همکارانش چارچوبی به نام Pathos را بر مبنای نظریه ساختار بیانی در سطح سند پیشنهاد کردند [۱۰]. Wang و همکارانش یک چارچوب غیرنظارتی به کمک نظریه گفتمان بیانی ارائه دادند و با در نظر گرفتن نوع روابط معنایی بین واحدها تحلیل احساسات را بهبود دهند اما حاشیه‌نویسی آن به صورت دستی انجام شد و مجموعه داده مورد استفاده کوچک و حدود ۶۰۸ نظر را دربرگرفت [۱۷].

تحقیقات ذکر شده احساس متن را در سطح سند طبقه بندی می‌کنند ولی از نظریه گفتمان بیانی در سطح درون جمله‌ای استفاده می‌کنند.

Bhatia و همکارانش سه روش، با استفاده از درخت گفتمان وابستگی پیشنهاد کردند و نشان دادند که این نظریه می‌تواند تجزیه و تحلیل احساسات را در سطح سند بهبود بخشد، اما آنها تنها نوع واحدهای گفتمان و وابستگی پوسته و هسته را در قاعده پیشنهادی خود در نظر گرفتند و روش‌های مبتنی بر یادگیری ماشینی که آنها پیشنهاد کردند در مقایسه با پایه خود پیشرفت کمی داشت [۱۸].

در این میان کمتر تحقیقی است که هر سه ویژگی درخت گفتمان را در نظر بگیرد و مبتنی بر قاعده باشد، بنابراین هدف ما استفاده از هر سه ویژگی و افزودن پارامتر عمق است.

نزدیکترین و شبیه‌ترین روش به این مقاله رویکرد تحقیق [۱۹] است که DFDS نام دارد، اما طرح قواعد پیشنهادی ما نه تنها براساس نوع واحدهای گفتمان (هسته و پوسته بودن) و نوع ارتباط بین واحدهای گفتمان است، بلکه

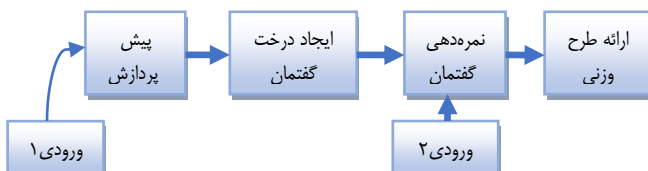


شکل (۲): درخت تجزیه شکل (۱) با ریشه Root [۱۳]

دو روش برای تبدیل متن به صورت درخت گفتمان وجود دارد. روش اول استفاده از تجزیه‌گرها و روش دوم حاشیه‌نویسی به صورت دستی است. همانطور که واضح است حاشیه‌نویسی دستی گرچه دقت بیشتری دارد اما بسیار زمان‌بر و سنگین است، بنابراین از تجزیه‌گرها استفاده می‌کنیم. پژوهش‌های قبلی روش‌های آماری، یادگیری ماشین و الگوریتم‌های مبتنی بر قاعده را برای ایجاد درخت تجزیه خودکار پیشنهاد کرده‌اند. به عنوان مثال SPADE، یک تجزیه‌گر گفتمان احتمالاتی و مبتنی بر قاعده است که در سطح جمله توسعه داده شده است [۲۱]. تجزیه‌گر دیگری که به اصطلاح HILDA نام دارد، یک طبقه‌بندی SVM را برای شناسایی واحدهای گفتمان و برجسب زدن نوع‌شان در رابطه بکار می‌برد [۲۲]. DPLP یکی از بهترین تجزیه‌گرهای گفتمان روی پایگاه داده محک^{۲۳} Treebank است. این تجزیه‌گر یک الگوریتم رمزگشایی حریصانه تجزیه متن، بکار می‌برد و به عنوان یک تجزیه‌گر عمومی، متن باز و قابل دسترس است [۲۳]. در این بین تنها دو تجزیه‌گر DPLP و HILDA در سطح سند به شناسایی گفتمان می‌پردازند. در این مقاله از تجزیه‌گر DPLP استفاده می‌کنیم. این تجزیه‌گر معیار F بالاتری نسبت به HILDA دارد [۲۳].

۴- روش پیشنهادی

به طور کلی ابتدا پایگاه داده نظرات وارد می‌شود. پس از پیش پردازش یک نظر، به کمک یک تجزیه‌گر، درخت گفتمان مربوطه ایجاد می‌شود و نمره احساس واحدهای گفتمان با استفاده از یک روش تجزیه و تحلیل در سطح جمله بدست می‌آید. در نهایت با توجه به ساختار درخت گفتمان و روابط بین واحدهای گفتمان یک سری قواعد وزن‌دهی مطرح کرده‌ایم که نمره جملات در سطح جمله را به طور بازگشتی با هم جمع کرده و احساس کلی متن را محاسبه می‌کنیم؛ در صورتی که نتیجه محاسبات عدد مثبتی باشد آنگاه احساس متن را مثبت و در غیر اینصورت منفی طبقه‌بندی می‌کنیم. در ادامه گام‌های این روش با توجه به شکل (۳) شرح داده می‌شود.



اطلاعات افزایش دهنده توانایی بالقوه برای انجام کارهای هسته پوسته توضیحاتی برای ارزیابی موضوعات موجود در هسته دارد. دلایلی برای شرايطی که در هسته ارائه شده، در برگرفته می‌گیرد. هیچ سلسله مراتبی خاصی بین واحدهای گفتمان وجود ندارد. شرايطی که تحقق آن از طریق تحقق وضعیت ارائه شده در هسته انجام نمی‌شود. ترکیب مواد متن را توصیف می‌کند. موضوع بین واحدهای گفتمان تغییر کرده است. یک واحد گفتمان بیانی نظریه‌هایی در مورد موضوع را بیان می‌کند. پوسته به احتمال زیاد تمایل به تحقق هسته را دارد.

توانایی^{۱۴}
ارزیابی
توصیف
اتصال^{۱۵}
در غیر این صورت
کنترل-متنی^{۱۶}
تغییر موضوع^{۱۷}
موضوع-نظر^{۱۸}
چگونگی-روش^{۱۹}

روابط بین واحدهای گفتمان نیز با توجه به اهمیت‌شان برای فهم متن به دو دسته تقسیم می‌شوند:

- تک هسته‌ای: دو واحد گفتمان در وضعیت یکسانی نیستند و هسته دارای اهمیت بیشتری است. در این نوع رابطه، یک واحد با اهمیت کمتر (پوسته) به واحد با اهمیت بیشتر (هسته) متصل می‌شود. برای مثال در متن زیر، رابطه علت یک نوع رابطه‌ی تک هسته‌ای است.
 - چند هسته‌ای: دو یا چند واحدی را که اهمیت برابر یا مشابه دارند و همه هسته‌اند، را به هم متصل می‌کند. برای مثال رابطه اتصال (ردیف ۱۲ از جدول (۱)) یک نوع رابطه‌ی چند هسته‌ای است.
- به متن شکل (۱) توجه کنید، شکل (۲) درخت تجزیه آن را نمایش می‌دهد.

I haven't watched a movie for a long time. all in all, I enjoy this comedy. In fact, the main actor is known for his good comedic acting

شکل (۱): نمونه متن ورودی [۱۳]

در این درخت گفتمان دو نوع ارتباط تک هسته‌ای هسته - پوسته وجود دارد. متن به سه بخش تقسیم شده است پس دارای سه واحد گفتمان ابتدایی است که گره‌های برگ این درخت دودویی را تشکیل می‌دهد. نماد N نشان دهنده‌ی هسته و S پوسته است. دو گره داخلی دیگر برای نشان دادن روابط این واحدهای گفتمان مورد نیاز است. برجسب Elaboration در ریشه بیانگر این است که واحد گفتمان ابتدایی ۳ جزئیات بیشتری را در مورد محتوای زیر درخت سمت چپ (به عنوان مثال کمدی) نشان می‌دهد. علاوه بر این، برجسب Background نشان می‌دهد که جمله ۱ قابلیت درک جمله ۳ را افزایش می‌دهد و پیش زمینه‌ای از آن را بیان می‌کند.

شکل (۳): طرح کلی روش پیشنهادی

۱-۴- ورودی ۱:

ورودی ۱ همان پایگاه داده‌ی نظرات است که باید جمع‌آوری شود. این پایگاه داده متنی شامل نظرات کاربران مختلف یک وب‌گاه روی کالا یا محصولی است که خریداری شده است.

۲-۴- پیش پردازش:

در طول آزمایش گام‌های پیش‌پردازش زیر را که در کاوش متن معمول هستند انجام می‌دهیم:

(۱) همه‌ی بخش‌های متن را که به پیام اصلی مرتبط نیستند؛ شامل اطلاعات تماس، شماره صفحات کتاب، لینک‌ها و دستورات HTML را حذف می‌کنیم. (۲) تمیز کردن محتوا؛ حروف تکرار متن و شکلک‌های متن را حذف می‌کنیم. در این راستا هر رخدادی که بیش از دو بار در یک کلمه رخ می‌دهد با دو رخداد جایگزین می‌شود. برای نمونه، در نوشتن کلمه «گرسنگی» لغاتی مانند huuuuungry و huuuuungry، یک حرف به تعداد دلخواه در متن تکرار شده است این کلمات به کلمه hungry تبدیل می‌شوند [۲۴]. این پیش پردازش برای بالا بردن دقت به کار می‌رود.

۳-۴- ایجاد درخت گفتمان

دو نکته‌ی مهم در تحلیل با استفاده از نظریه ساختار بیانی عبارتند از تعیین نوع تجزیه‌گر گفتمان و مجموعه‌ی دقیق روابط بیانی‌ای که در تجزیه و تحلیل وجود دارد. ابتدا متن را با استفاده از تجزیه‌گر DPLP^{۲۴} تجزیه می‌کنیم. این تجزیه‌گر متن را به واحدهای گفتمان، معمولاً عبارات یا جملات، بر اساس تجزیه‌های نحوی ارائه شده توسط ابزار StanfordCoreNLP^{۲۵} تقسیم بندی و پردازش می‌کند سپس ساختار درختی گفتمان را بدست می‌آورد. در این تجزیه‌گر حاشیه‌نویسی و پیش‌پردازش به صورت موازی است.

۴-۴- نمره دهی گفتمان

پردازشی اولیه روی عبارت انجام می‌دهیم برای نمونه کلمات توقف^{۲۶} حذف می‌کنیم، سپس قطبیت کلمات موجود در واحدهای گفتمان را به کمک فرهنگ لغت VADER^{۲۷} بازبایی می‌کنیم. این فرهنگ لغت همان ورودی ۲ در شکل (۳) هست. به موجب آن یک لیست کلمات مثبت یا منفی طبقه‌بندی شده، فراهم می‌شود.

این فرهنگ لغت، یک ابزار تجزیه و تحلیل احساسات است که (۱) در متون سبک رسانه‌های اجتماعی به خوبی کار می‌کند، با این حال به راحتی قابل تعمیم به حوزه‌های دیگر است، (۲) نیاز به هیچ داده آموزشی ندارد، اما از یک واژگان احساسی استاندارد و مبتنی بر ارزش و قابل تعمیم ساخته شده است، (۳) سریع است و می‌تواند برای داده‌های آنلاین هم استفاده شود. علاوه

بر این کلماتی که حاوی اصطلاحاتی مانند LOL^{۲۸} هستند، شکلک‌هایی نظیر D، :، :D، :P، :O، :X و... را در بردارد. ویژگی‌های گرامری مانند علائم نگارشی و نشانه گذاری (!، .، "، ، و...)، بزرگ کوچک بودن کلمه و کلمات تشدید کننده (very, extremely) و منفی را تشخیص می‌دهد. (۴) کلماتی مانند "خوب!!!" را شناسایی می‌کند. این کلمات دارای نشانه گذاری‌هایی هستند که شدت احساس یک کلمه را افزایش می‌دهند.

بدین ترتیب امتیاز هر واحد گفتمان به کمک فرهنگ لغت با جمع بندی نمره کلمات موجود در آن بدست می‌آید.

۵-۴- ارائه طرح وزنی

در طبقه بندی یک سند متنی بعد از بدست آوردن ساختار درختی متن، نیاز است یک طرح وزنی ارائه شود و بدین وسیله امتیاز کلی متن بدست آید.

طرح پیشنهادی اطلاعات را از واحدهای گفتمان دریافت و ترکیب می‌کند سپس آنها را به سطح بالاتر بعدی منتقل می‌کند تا زمانی که به ریشه برسد، در آن نقطه پیش‌بینی ساخته می‌شود. به این ترتیب، این رویکرد صرفاً به ساختار درختی متکی است و هر بخش، امتیازات مثبت و منفی خود را با استفاده از یکی از روش‌های طبقه بندی احساس تخمین می‌زند. در نتیجه نمره احساس یک متن با ترکیب مقادیر وزنی احساس همه‌ی فرزندان در یک حالت تکراری حاصل می‌شود. طرح وزنی پیشنهادی در واقع همان طرح وزنی ارائه شده در تحقیق [۱۹] است که پارامتر عمق با آن ترکیب شده است بدین ترتیب این طرح از هر سه مولفه درخت گفتمان بیانی (یعنی عمق درخت - نوع گره - ارتباط بین واحدهای گفتمان) برای وزنی دهی استفاده می‌کند. علاوه بر آن، این تحقیق یک گره دیگر را به آن ۷ گروه مجموعه قواعد موجود در طرح وزنی اضافه کرده است.

۱-۵-۴- وزن دهی

در طرح وزنی ارائه شده پارامتر، d_i نشان دهنده‌ی عمق است. این پارامتر با الهام از تحقیق [۲] طبق رابطه‌ی زیر بدست می‌آید:

$$d_i = \begin{cases} \frac{i(\alpha-1)}{i_{max}} + 1 & \alpha \geq 0 \\ \frac{i(\alpha+1)}{i_{max}} - \alpha & \alpha < 0 \end{cases} \quad (1)$$

که در آن $\alpha < 1 < -1$ است و α را ضریب عمق می‌نامیم.

در این رابطه، i عمق گره نسبت به گره ریشه و i_{max} حداکثر عمق مربوط به زیر درخت (عمق گره از گره برگ آن زیردرخت) است. می‌توان دو فرض را در نظر گرفت:

○ اگر $\alpha < 1$ مجازاتی را روی قسمت‌های عمیق‌تر متن می‌گذارد.

واحدهای گفتمانی که مجاورت و نزدیکی بیشتری به ریشه دارند در تعیین احساس کلی متن بیشتر دخالت دارد.

○ اگر $\alpha \geq 1$ باشد هر چه متن در عمق دورتری نسبت به ریشه قرار

گیرد، اهمیت بیشتری دارد.

رابطه‌ی زیر را تعریف کرده‌ایم. این رابطه یک نوع رابطه پوسته- هسته است.

$$S_i = d_i \times [s(\text{nucleus}) + s(\text{satellite})] \quad (۸)$$

• گروه هشتم: سایر روابط

$$S_i = d_i \times [w_1 \times s(\text{nucleus}) + w_s \times s(\text{satellite})] \quad (۹)$$

گروه‌های تعریف شده هنگامی تعریف می‌شود که نوع رابطه بین دو واحد گفتمان از نوع هسته-پوسته باشد. در صورتی که نوع روابط چند هسته‌ای باشد، امتیاز آن رابطه مجموعه وزن هسته‌ها ضرب در پارامتر d_i است.

۵- آزمایش

روش پیشنهادی را روی دامنه‌های Book, Electronic, Houseware و DVD [۲۵] انجام دادیم. این دامنه‌ها به ترتیب نظرات سایت آمازون روی محصولات لوازم آشپزخانه، لوازم الکترونیکی، کتاب و دی وی دی در سال ۲۰۰۷ بوده است. هر دامنه شامل ۱۰۰۰ نظر مثبت و ۱۰۰۰ نظر منفی است. هر نظر از ۵ قسمت: نام نظردهنده، مکان، نام محصول، عنوان نظر، متن نظر و تاریخ تشکیل شده است. نویسندگان آن F. و J. Blitzer, M. Dredzr و J. Pereira هستند.

برای ارزیابی روش پیشنهادی پارامترهای صحت و معیار F استفاده شده است. معیار صحت: نسبت متن‌های درست دسته بندی شده بر همه‌ی پیش‌بینی‌هاست که فرمول آن به صورت زیر است:

$$\text{Accuracy} = (\text{TP} + \text{TN}) / (\text{TP} + \text{TN} + \text{FP} + \text{FN}) \quad (۱۰)$$

که در آن:

TP تعداد متن‌هایی است که به درستی برچسب مثبت خورده‌اند

FP تعداد متن‌هایی است که به اشتباه برچسب مثبت خورده‌اند.

TN تعداد متن‌هایی است که به درستی برچسب منفی خورده‌اند.

FN تعداد متن‌هایی است که به اشتباه برچسب منفی خورده‌اند.

معیار F1: میانگین هارمونیک Precision و Recall است.

$$F1 = (2 \times \text{Precision} \times \text{Recall}) / (\text{Precision} + \text{Recall}) \quad (۱۱)$$

که در آن معیار Precision: نشان دهنده‌ی تعداد متونی است که برچسب مثبت خورده‌اند و حقیقتاً دارای برچسب مثبت (مرتبط) هستند. معیار Recall: تعداد متونی که با برچسب مثبت انتخاب شده و درست طبقه بندی شده‌اند.

در ادامه نتایج حاصل از پیاده سازی روش الهام گرفته شده مطرح می‌شود سپس تاثیر نوآوری‌های صورت گرفته در طرح وزنی بیان و پس از تعیین ضریب عمق، به مقایسه روش با روش‌های پیشین می‌پردازیم.

در مجموعه قواعد تعریف شده، $s(\text{nucleus})$ امتیاز هسته و $s(\text{satellite})$ امتیاز پوسته است که چگونگی محاسبه‌ی آن در بخش قبل توضیح داده شد و $W = \{W_s, W_e, W_{vs}, W_{vl}, W_1\}$ مجموعه‌ی وزنی است که اعدادی مثبت صحیح در بازه بین صفر و یک هستند و مقادیر آن در تحقیق [۱۹] آمده است.

۲-۵-۴- مجموعه قواعد

در طرح وزنی ارائه شده واحدهای گفتمان بسته به اینکه خود چه نوع ارتباطی از این گروه‌ها دارند و از نظر اهمیت (هسته/پوسته) در چه جایگاهی هستند، وزن‌های متفاوتی می‌گیرند. بنابراین با توجه به نوع واحد گفتمان می‌توانیم این فرض را رمز گذاری کنیم که گره هسته از اهمیت بیشتری نسبت به پوسته برخوردار است و با توجه به نوع ارتباط می‌توان تاثیر یک واحد گفتمان را در تعیین احساس کلی متن تغییر داد. با توجه پارامتر عمق می‌توان امتیازات احساس واحدهای گفتمان را بسته به عمقشان در درخت گفتمان جریمه کرد. مجموعه روابط در زیر نشان داده شده‌اند که بخشی از آنها در شکل (۱) آمده است.

• گروه اول: روابط *antithesis, concession*

If $s(\text{nucleus}) \times s(\text{satellite}) \leq 0$

$$S_i = d_i \times [w_1 \times s(\text{nucleus}) - w_s \times s(\text{satellite})]$$

else: (۲)

$$S_i = d_i \times [w_1 \times s(\text{nucleus}) + w_s \times s(\text{satellite})]$$

• گروه دوم: روابط *circumstance, background*

$$S_i = d_i \times [w_{vl} \times s(\text{nucleus}) + w_{vs} \times s(\text{satellite})] \quad (۳)$$

• گروه سوم: رابطه *condition*

$$S_i = d_i \times [-w_{vl} \times s(\text{nucleus}) + w_{vs} \times s(\text{satellite})] \quad (۴)$$

• گروه چهارم: روابط *motivation, purpose*

If $s(\text{nucleus}) \times s(\text{satellite}) \leq 0$

$$S_i = d_i \times [w_{vl} \times s(\text{nucleus}) + w_{vs} \times s(\text{satellite})]$$

else: (۵)

$$S_i = d_i \times [w_1 \times s(\text{nucleus}) + w_s \times s(\text{satellite})]$$

• گروه پنجم: روابط

[*evidence, reason, result, restatement, enablement, justification*]:

If $s(\text{nucleus}) \times s(\text{satellite}) \leq 0$

$$S_i = d_i \times [w_{vs} \times s(\text{nucleus}) + w_{vl} \times s(\text{satellite})]$$

else: (۶)

$$S_i = d_i \times [\max(s(\text{nucleus}), s(\text{satellite}))]$$

• گروه ششم: روابط *evaluation, conclusion*

$$S_i = d_i \times [w_e \times s(\text{nucleus}) + w_{vs} \times s(\text{satellite})] \quad (۷)$$

• گروه هفتم: رابطه *Elaboration*

نتایج حاصل از افزودن این رابطه به مجموعه روابط مقادیر دقت، صحت و یادآوری را با استفاده از لغت نامه SWN، به صورت زیر تغییر می‌دهد.

جدول (۴) : نتایج حاصل از افزودن رابطه‌ی جدید به مجموعه قواعد

F	صحت	FN	FP	TN	TP	مجموعه داده
۶۵	۶۱/۴	۲۸۱	۵۰۹	۳۹۱	۷۱۹	لوازم الکترونیکی
۶۶/۵	۶۰/۷۵	۲۲۰	۵۶۵	۴۳۵	۷۸۰	دی وی دی
۶۶/۳	۵۹/۸۵	۲۰۷	۵۹۶	۴۰۴	۷۹۳	کتاب
۶۸/۱	۶۳/۰۵	۲۰۹	۵۳۰	۴۷۰	۷۹۱	لوازم آشپزخانه

با توجه به نتایج به طور کلی معیار صحت در حدود دو درصد افزایش یافته است.

۳-۵- تعیین ضریب عمق

برای محاسبه‌ی پارامتر ضریب عمق (α) موجود در رابطه‌ی عمق که در قواعد وزنی مطرح شد، برای یافتن این ضریب در شبکه جستجو می‌کنیم. برای هر مجموعه داده بازه [۰،۱] به ۲۰ تقسیم مساوی تقسیم کرده و ضریب را جستجو می‌کنیم. بهینه‌ترین و بهترین ضریب عمق، عددی است که به ازای آن معیار ارزیابی صحت بیشترین شود.

در این راستا برای انتخاب روش موثرتر دو فرهنگ لغت مختلف SWN و VADER را برای محاسبه‌ی نمره احساس واحدهای گفتمان بکار گرفتیم. ضرایب بدست آمده برای هر مجموعه داده در جدول (۵) آمده است.

با توجه به جدول (۵) بدست آمدن مقدار منفی برای ضریب عمق در مجموعه داده لوازم الکترونیکی نشان می‌دهد که هر چه به انتهای متون نظرات نزدیک می‌شویم از اهمیت جملات در سند کاسته می‌شود. بدست آمدن مقدار مثبت ضریب عمق در سایر مجموعه داده‌ها حاکی از آن است که هر چه به انتهای متن نزدیک می‌شویم به بخش نتیجه گیری متن نزدیک و اهمیت جملات در سند بیشتر می‌شود.

جدول (۵) : ضریب عمق برای هر مجموعه داده با توجه به فرهنگ لغت

فرهنگ لغت	فرهنگ لغت
VADER	SWN
۰/۹۲۶۴	۰/۸۹۸ کتاب
۰/۹۶۶۹	۰/۹۸۲۰ دی وی دی
-۰/۹۹	۱ لوازم الکترونیکی
۰/۹۸۴	۰/۹۹۵ لوازم آشپزخانه

۴-۵- نتایج روش پیشنهادی

عملکرد روش پیشنهادی را با در نظر گرفتن دو فرهنگ لغت بررسی کردیم. با توجه نتایج بدست آمده در شکل (۴) درمی‌یابیم که استفاده از فرهنگ لغت VADER نسبت به SWN در تعیین احساس متن در هر چهار مجموعه داده موجب کارایی بیشتر می‌شود و پارامتر عمق در تعیین احساسات در متون نقش دارد و موجب افزایش معیار ارزیابی صحت و معیار F می‌شود.

۱-۵- نتایج حاصل از اجرای روش DFDS

همانطور که در بخش ۲ گفته شد روش DFDS به عنوان مبنای روش پیشنهادی قرار گرفته است به طوری که گام‌های پیش پردازش ما را شامل نمی‌شود و پارامتر عمق را نیز در نظر نگرفته است و برای تعیین احساس واحدهای گفتمان از فرهنگ لغتی متفاوت استفاده کرده است.

این تحقیق ابتدا با استفاده از نظریه ساختار بیانی و به کمک تجزیه‌گر DPLP، اسناد را به واحدهای گفتمان ابتدایی تجزیه می‌کند سپس احساسات هر واحد را توسط فرهنگ لغت SentiWordNet (SWN) بدست می‌آورد. سپس طرح وزنی‌ای [۱۹] را ارائه کرده است که احساسات کل متن را پیش بینی می‌کند.

با پیاده سازی این روش به نتایجی متفاوت از آنچه این تحقیق گزارش کرده است، رسیدیم، ممکن است به دلیل باشد که مقاله بخشی از روند کار را ذکر نکرده است. در پیاده سازی صورت گرفته در بهترین حالت مقدار صحت در دو مجموعه داده لوازم الکترونیکی و لوازم آشپزخانه حدود دو درصد کمتر از گزارشات درج شده در مقاله است. این نتایج در جدول (۲) آمده است.

جدول (۲) : الگوریتم DFDS

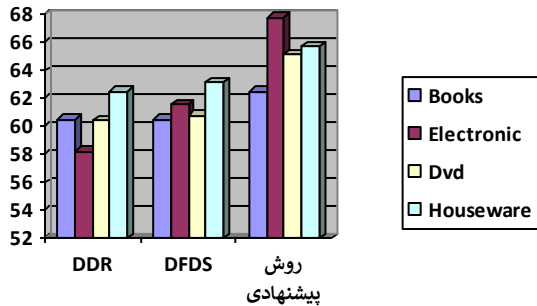
F	صحت	FN	FP	TN	TP	مجموعه داده
۶۳/۸	۵۹/۶	۲۸۷	۵۲۰	۴۸۰	۷۱۳	لوازم الکترونیکی
۶۵/۷	۶۰/۵۵	۲۴۳	۵۴۵	۴۵۵	۷۵۷	دی وی دی
۶۴/۳	۵۸/۶۵	۲۵۲	۵۷۵	۴۲۵	۷۴۸	کتاب
۶۷/۲	۶۱/۵	۲۰۹	۵۶۱	۴۳۹	۷۹۱	لوازم آشپزخانه

۲-۵- تاثیر افزودن رابطه جدید و پیش پردازش

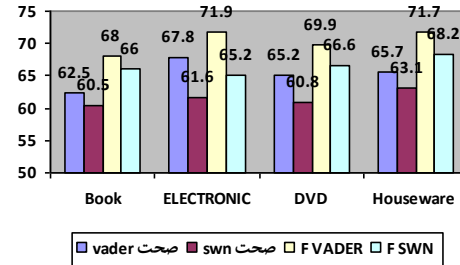
جدول (۳) تعداد روابط گفتمانی موجود در هر یک از گروه‌هایی که در قاعده وزنی بخش ۴-۵ تعریف کرده‌ایم را برای چهار مجموعه داده نشان می‌دهد. همانطور که مشاهده می‌کنیم، از آنجایی که تعداد روابط گروه هفتم زیاد است؛ یک رابطه‌ی جدید تعریف کردیم متمایز از آنچه روش DFDS (ادغام رابطه‌ی گروه هفتم و سایر) ارائه کرده است.

جدول (۳) : تعداد روابط گفتمانی هر دسته در چهار مجموعه داده

مجموعه داده گروه	لوازم آشپزخانه	کتاب	وی دی دی	لوازم الکترونیکی
گروه اول	۲۱۷	۳۰۷	۳۱۶	۲۳۴
گروه دوم	۶۸۳	۷۲۹	۷۷۳	۶۸۶
گروه سوم	۶۲۹	۸۵۹	۸۰۳	۸۳۱
گروه چهارم	۱۵۲۸	۱۹۹۶	۱۷۱۰	۱۷۵۲
گروه پنجم	۱۹۵	۲۱۳	۲۲۶	۱۷۹
گروه ششم	۱	۰	۰	۱
گروه هفتم	۷۲۴۳	۱۳۴۵۳	۱۲۲۲۶	۸۲۹۸
سایر	۶۴۱۴	۱۱۹۳۹	۱۱۳۷۳	۸۰۹۸



شکل (۵): مقایسه صحت مقاله [۱۹] و مقاله [۱۸] با روش این تحقیق



شکل (۴): ارزیابی صحت و F روش پیشنهادی با استفاده از فرهنگ لغت VADER و SWN

۶- نتیجه گیری

نظریه ساختار بیانی یک روش تجزیه و تحلیل احساس است که با تقسیم متن به بخش‌های معنادار، بخش‌هایی را شناسایی می‌کند که نقش بیشتری در تعیین احساس کلی متن دارند. در این مقاله یک روش غیرنظارتی برای تجزیه و تحلیل احساس نظرات سایت آمازون در قالب روش‌های مبتنی بر واژگان توصیف شد. اگرچه در این زمینه و در سطح سند تحقیقات کمی صورت گرفته، ما یک چارچوب بر این مبنا ارائه کردیم که هر سه مولفه گفتمان بیانی به منظور تحلیل احساس متون استفاده شده است.

در آزمایش‌ها ثابت کردیم که با ترکیب پارامتر عمق در مجموعه قواعد، تغییر قواعد و اضافه کردن قاعده جدید می‌توان مقدار صحت و معیار F را بهبود داد. همچنین تغییر در فرهنگ لغت نیز کارایی نسبتاً خوبی در طبقه‌بندی نظرات از منظر قطبیت (مثبت، منفی) بدست می‌آورد و برای مسائلی که در آنها مقدار TP اهمیت بیشتری دارد، روش پیشنهادی کارا تر است.

با توجه به نتایج به دست آمده دریافتیم که صحت در روش پیشنهادی نسبت به دو روشی که مبتنی بر واژگان بود در بهترین حالت حدود ۸ درصد افزایش یافته است، هرچند نسبت به روش‌های نظارتی دارای مقدار صحت کمتری است. اما این روش نسبت به روش‌های نظارتی مزایایی هم دارد، مهمترین آن، این است که نیاز به مقدار قابل توجهی داده برچسب زده برای اهداف آموزشی ندارد و این موضوع در دنیایی که حجم داده‌ها و تعداد دامنه‌ها در حال افزایش است و عملاً جمع‌آوری اطلاعات امکان‌پذیر نیست، بسیار حائز اهمیت است. همچنین در بین چهار مجموعه داده مورد بررسی، مجموعه داده لوازم الکترونیکی دارای بیشترین مقدار صحت است. از آنجایی که نظرات آن دارای متن‌های طولانی‌تری بوده است و همچنین فرهنگ لغت، لغات جامع‌تری در زمینه‌ی لوازم الکترونیکی در بر گرفته است؛ می‌توان گفت افزودن پارامتر عمق و تغییر لغت نامه در این مجموعه داده بیشتر موثر است.

نوع تجزیه‌گر یک عامل مهم و کلیدی در تعیین احساسات و ساختار گفتمان است لذا استفاده از سایر تجزیه‌گرها می‌تواند تغییری ایجاد کند. در کارهای آتی به منظور در نظر گرفتن بخش‌هایی قطب مشخصی ندارند،

۵-۵- مقایسه روش پیشنهادی با روش‌های پیشین

از آنجایی که تعداد تحقیقات کمی در رابطه با تحلیل احساس در سطح سند صورت گرفته است، ما روش پیشنهادی را با دو تحقیق مقایسه کردیم: تحقیق اول، روش مقاله [۱۹] که در بخش ۱-۵- بحث کردیم و نتایج آن را آوردیم همچنین با تحقیق [۱۸] که روش خود را DDR نامیده است و از پارامتر عمق استفاده کرده است؛ نتیجه‌ی این مقایسه‌ها در ارزیابی معیار صحت در شکل (۵) آمده است.

با توجه به نتایج معیار صحت در بهترین حالت یعنی در مجموعه داده لوازم الکترونیکی از ۵۹/۶ درصد به مقدار ۶۷/۸ درصد بهبود یافته و بطوری که صحت حدود ۸ درصد در مجموعه داده لوازم الکترونیکی افزایش یافته است. علاوه بر این در هنگام مقایسه روش پیشنهادی با روش DDR در بهترین حالت مقدار صحت را از ۵۸/۲ درصد به ۶۷/۸ درصد بهبود یافته است. از آنجایی که فرهنگ لغت SWN دارای نویز فراوان است و اکثریت مجموعه‌ی هم‌معنی آن، هیچ قطب مثبت یا منفی ندارند، از فرهنگ لغت VADER برای روش پیشنهادی استفاده کردیم. با توجه به شکل (۵) مقدار صحت در این روش نسبت به این دو روش بهبود دارد.

- [13] Kraus, M. and Feuerriegel, S., "Sentiment analysis based on rhetorical structure theory: Learning deep neural networks from discourse trees", Expert Systems with Applications., Vol. 118, pp. 65–79, 2019.
- [14] Hailong, Z., Wenyan, G. and Bo, J., "Machine learning and lexicon based methods for sentiment classification: A survey", 11th Web Information System and Application Conference. IEEE, pp. 262–265, 2014.
- [15] Asghar, M.Z., Khan, A., Khan, F. and Kundi, F.M., "RIFT: A Rule Induction Framework for Twitter Sentiment Analysis", Arab. J. Sci. Eng., Vol. 43, 2018.
- [16] Voll, K. and Taboada M., "Not All Words Are Created Equal: Extracting Semantic Orientation as a Function of Adjective Relevance", Australasian Joint Conference on Artificial Intelligence., Springer, pp. 337–346, 2007.
- [17] Wang, F., Wu, Y. and Qiu, L. "Exploiting Discourse Relations for Sentiment Analysis", Proceedings of COLING, Vol. 2, pp. 1311–1320, 2012.
- [18] Bhatia, P., Ji, Y. and Eisenstein, J., "Better Document-level Sentiment Analysis from RST Discourse Parsing", arXiv preprint arXiv:1509.01599, pp. 2212–2218, 2015.
- [19] Zhao, Z., Rao, G. and Feng, Z. "DFDS: A Domain-Independent Framework for Document-Level Sentiment Analysis Based on RST", SPRINGER, pp. 297–310, ۲۰۱۷.
- [20] Mann, W.C. and Thompson, S.A., *Rhetorical structure theory: A theory of text organization*. california, 1987.
- [21] Soricut, R. and Marcu, D., "Sentence Level Discourse Parsing using Syntactic and Lexical Information", Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics, pp. 149–156, 2003.
- [22] Hernault, H., Prendinger, H. and Ishizuka, M., "HILDA: A discourse parser using Support Vector Machine classification", Dialogue & Discourse, Vol. 1, No. 3, pp. ۱–۳۳, ۲۰۱۰.
- [23] Ji, Y. and Eisenstein, J., "Representation Learning for Text-level Discourse Parsing", Proc. 52nd Annu. Meet. Assoc. Comput. Linguist., pp. 13–24, 2014.
- [24] Go, A., Bhayani, R. and Huang, L., "Twitter Sentiment Classification using Distant Supervision", CS224N Project Report, Stanford Vol. 1, No. 12, pp. 1–6, 2009.
- [25] www.cs.jhu.edu/~mdredze/datasets/sentiment/
- [1] Pang, B., Lee, L., "Opinion Mining and Sentiment Analysis", Information Retrieval, Vol. 2, No. 1–2, 2008.
- [2] Märkle-Huß, J., Feuerriegel, S., Prendinger, H., "Improving Sentiment Analysis with Document-Level Semantic Relationships from Rhetoric Discourse Structures", Proc. 50th Hawaii Int. Conf. Syst. Sci., pp. 1142–1151, 2017.
- [3] Sun, S., Luo, C., and Chen, J., "A review of natural language processing techniques for opinion mining systems", Information fusion, Vol. 36, pp. 10–2۵, ۲۰۱۷.
- [4] Yue, L., Chen, W., Li, X., Zuo, W. and Yin, M., "A survey of sentiment analysis in social media", Knowledge and Information Systems, pp. 1–47, 2018.
- [5] Pang, B., Lee, L. and Vaithyanathan, S., *Thumbs up? Sentiment Classification using Machine Learning Techniques*, 2002.
- [6] Hogenboom, A., Frasinca, F., De Jong, F. and Kaymak, U., "Polarity classification using structure-based vector representations of text", Decision support systems., Vol. ۷۴, pp. ۴۶–۵۶, ۲۰۱۵.
- [7] Fu, X., Liu, W., Xu, Y., Yu, C. and Wang, T., "Long Short-term Memory Network over Rhetorical Structure Theory for Sentence-level Sentiment Analysis", Asian Conference on Machine Learning, pp. 17–32, 2016
- [8] Taboada, M., Voll, K. and Brooke, J., "Extracting sentiment as a function of discourse structure and topicality", Simon Fraser University School of Computing Science Technical Report, pp. 302–306, ۲۰۰۸.
- [9] Chenlo, J.M., Hogenboom, A. and Losada, D.E., "Rhetorical Structure Theory for polarity estimation: An experimental study", Data and Knowledge Engineering, Vol. 94, pp. 135–147, 2014.
- [10] Heerschop, B., Goossen, F., Hogenboom, A., Frasinca, F., Kaymak, U. and de Jong, F., "Polarity analysis of texts using discourse structure", Proc. 20th ACM Int. Conf. Inf. Knowl. Manag. - CIKM '11, pp. 1061–1070, 2011.
- [11] Tan, S., Wang, Y. and Cheng, X., "Combining learn-based and lexicon-based techniques for sentiment detection without using labeled examples", Proceedings of the 31st annual international ACM SIGIR conference on Research and development in information retrieval, pp. 743–744, 2008.
- [12] Zirn, C., Niepert, M., Stuckenschmidt, H. and Strube, M., "Fine-Grained Sentiment Analysis with Structural Features", Proceedings of 5th International Joint Conference on Natural Language Processing, Chiang Mai, Thailand, pp. 336–344, 2011.

مراجع

پانویس ها

- ^۱ Text mining
^۲ Sentiment analysis
^۳ Opinion mining
^۴ Bag of words
^۵ Support Vector Machines
^۶ Unsupervised Techniques
^۷ lexicon-based method
^۸ Rhetorical Structure Theory (RST)
^۹ Elementary Discourse Units(EDU)



-
- ^۱ Nucleus
 - ^{۱۱} Satellite
 - ^{۱۲} Attribution
 - ^{۱۳} Consequence
 - ^{۱۴} Enablement
 - ^{۱۵} Joint
 - ^{۱۶} Textual-organization
 - ^{۱۷} Topic-Change
 - ^{۱۸} Topic-comment
 - ^{۱۹} Manner-means
 - ^{۲۰} Mononuclear
 - ^{۲۱} Multinuclear
 - ^{۲۲} High -Level Discourse Analyzer
 - ^{۲۳} benchmark
 - ^{۲۴} <https://github.com/jiyfeng/DPLP>
 - ^{۲۵} <http://nlp.stanford.edu/sentiment>
 - ^{۲۶} Stop words
 - ^{۲۷} Valence Aware Dictionary for sEntiment Reasoning
 - ^{۲۸} Laughing Out Loud
 - ^{۲۹} Accuracy