

ساخت واژگان به صورت خودکار برای

تحلیل نظرات در حوزه بورس

مرتضی آهانگری آهانگرکلایی^۱، علی سبئی^{۲*}، مهدی یعقوبی^۳

^۱ و ^۲ گروه مهندسی کامپیوتر، دانشکده مهندسی گرگان، دانشگاه گلستان، گرگان، ایران



چکیده

با رشد چشم‌گیر رسانه‌های اجتماعی، افراد و سازمان‌ها به‌طور فزاینده‌ای از افکار عمومی در این رسانه‌ها برای تصمیم‌گیری خود استفاده می‌کنند. هدف تحلیل احساسات، استخراج خودکار احساسات افراد از این شبکه‌های اجتماعی است. شبکه‌های اجتماعی مرتبط به بازارهای مالی، از جمله بازارهای سهام، به‌تازگی در مرکز توجه بسیاری از افراد و سازمان‌ها بوده‌اند. افراد در این شبکه‌ها نظرات و عقاید خود را در مورد هر سهم در قالب یک پست یا توییت، به اشتراک می‌گذارند. در واقع تحلیل احساسات در این حوزه، سنجش نگرش افراد به هر سهم است. یکی از رویکردهای پایه‌ای و اصلی در تحلیل خودکار احساسات روش‌های مبتنی بر واژگان است. اغلب واژگان‌های مرسوم به‌صورت دستی استخراج شده‌اند که فرایندی بسیار دشوار و هزینه‌بر است. در این مقاله روشی جدید جهت استخراج واژگان به‌صورت خودکار در حوزه شبکه‌های اجتماعی بورسی ارائه شده است. یک ویژگی خاص این شبکه‌ها، وجود اطلاعات قیمتی هر سهم در هر روز است. با در نظر گرفتن وضعیت قیمتی سهم در روز نظر برای آن سهم، واژگانی برای بهبود کیفیت عقیده‌کاوی در این شبکه‌ها استخراج شد. برای ارزیابی واژگان‌های تولیدشده با استفاده از روش پیشنهادی نیز، با نسخه فارسی واژگان SentiStrength که با هدف استفاده عمومی طراحی شده است، مقایسه شد. نتایج آزمایش‌ها بیست درصد بهبود را در معیار صحت نسبت به استفاده از واژگان عمومی نشان می‌دهد.

واژگان کلیدی: تحلیل احساسات، عقیده‌کاوی، ساخت واژگان، واژگان فارسی.

Automatically generate sentiment lexicon for the Persian stock market

Morteza Ahangari Ahangarkolaei¹, Ali Sebt^{2*}, Mehdi Yaghoubi³

^{1,2,3} Department of Computer Engineering, Faculty of Engineering,
Golestan University, Gorgan, Iran

Abstract

With the significant growth of social media, individuals and organizations are increasingly using public opinion in these media to make their own decisions. The purpose of Sentiment Analysis is to automatically extract people's emotions from these social networks. Social networks related to financial markets, including stock markets, have recently attracted the attention of many individuals and organizations. People on these social networks share their opinions and ideas about each share in the form of a post or tweet. In fact, sentiment analysis in this area is measuring people's attitudes toward each share. One of the basic approaches in automatic analysis of emotions is lexicon-based methods. Most conventional lexicon is manually extracted, which is a very difficult and costly process. In this article, a new method for extracting a lexicon automatically in the field of stock social networks is proposed. A special feature of these networks is the availability of price information per share. Taking into account the price information of the share on the day of tweeting for that share, we extracted lexicon to improve the quality of opinion mining in these social networks. To evaluate the lexicon produced using the proposed method, we compared it with the Persian version of the SentiStrength

* Corresponding author

* نویسنده عهده‌دار مکاتبات

lexicon, which is designed for general purpose. Experimental results show a 20% improvement in accuracy compared to the use of general lexicon.

Keywords: Sentiment Analysis, Opinion Mining, Lexicon Creation, Persian Lexicon.

۱- مقدمه

اینکه «دیگران به چه چیزی فکر می‌کنند»، همیشه در فرایند تصمیم‌گیری برای بیشتر افراد، جزء اطلاعات مهم است. مدت‌ها پیش که شبکه جهانی وب^۱ (www) گسترش نیافته بود، اغلب برای پیدا کردن یک مکانیک خوب جهت تعمیر ماشین خود، از دوستان یا آشنایان کمک می‌گرفتیم، یا با مشورت از مصرف‌کنندگان دیگر تصمیم می‌گرفتیم کدام ماشین ظرفشویی را خریداری کنیم؛ اما اینترنت و شبکه اکنون (در میان امکانات دیگر) این امکان را فراهم آورده‌اند تا از نظرات و تجربیات طیف وسیعی از مردم استفاده کنیم که نه آشنای شخصی ما هستند و نه منتقدان حرفه‌ای شناخته‌شده؛ یعنی افرادی که ما هرگز چیزی از آنها نشنیده‌ایم [1].

با رشد چشم‌گیر رسانه‌های اجتماعی (به‌عنوان مثال، بررسی‌ها^۲، بحث‌های انجمن‌ها، وبلاگ‌ها و شبکه‌های اجتماعی) در وب، افراد و سازمان‌ها به‌طور فزاینده‌ای از افکار عمومی در این رسانه‌ها برای تصمیم‌گیری خود استفاده می‌کنند. با این حال، یافتن و نظارت بر نظرات در سایت‌های نظردهی در وب و بررسی اطلاعات موجود در آن‌ها به دلیل گسترش سایت‌های متنوع، همچنان یک کار دشوار است. هر سایت به‌طور معمول حاوی حجم عظیمی از نظرات است که همیشه در پست‌های طولانی در وبلاگ‌ها به‌راحتی قابل‌فهم نیست. یک خواننده معمولی در شناسایی تارنماهای موردنظر و جمع‌بندی دقیق اطلاعات و نظرات موجود در آن‌ها مشکل خواهد داشت [2].

امروزه مشتریان و صاحبان کسب‌وکار از این نظرات برای شناسایی نقاط قوت و ضعف محصولات استفاده می‌کنند، ولی با توجه به افزایش حجم نظرات، مطالعه موردی و نتیجه‌گیری نهایی از آن‌ها دیگر امکان‌پذیر نبوده و نیازمند سامانه‌ای جهت کسب دانش به‌صورت خودکار تحت عنوان تحلیل احساسات^۳، که به‌عنوان عقیده‌کاوی^۴ نیز شناخته می‌شود، شد. تحلیل احساسات یا عقیده‌کاوی، در سیاست (انتخابات، پیش‌بینی تحولات سیاسی، میزان اتحاد مردم یا جامعه در یک مورد و ...)، علوم اجتماعی و

¹ World Wide Web

² Reviews

³ Sentiment Analysis

⁴ Opinion Mining

روان‌شناسی (تحلیل مسائل اجتماعی و فرهنگی، تحلیل تأثیر اتفاقات مختلف در رفتار مردم و ...)، مدیریت و رهبری (کمک در تصمیم‌سازی و تصمیم‌گیری، آگاهی از میزان رضایت یا مطلوبیت مشتریان، مشترکان یا گروهی از مخاطبان، و ...) کاربردهای فراوانی دارد.

تاکنون پژوهش‌های زیادی در حوزه عقیده‌کاوی و تجزیه و تحلیل احساسات در زبان‌های انگلیسی، چینی و روسی انجام شده‌است. با وجود این که زبان فارسی، زبان اصلی ایران، افغانستان و تاجیکستان است و بیش از ۱۱۰ میلیون نفر در سرتاسر دنیا به این زبان سخن می‌گویند، ولی در متن‌های فارسی پژوهش‌های بسیار کمی در تجزیه و تحلیل احساسات انجام شده‌است و همچنان مشکلات و چالش‌های بسیاری در حوزه تحلیل احساسات و عواطف در زبان فارسی وجود دارد.

امروزه از تجزیه و تحلیل احساسات می‌توان در کاربردهای مختلفی استفاده کرد. بخشی از آن شامل تشخیص احساسات نسبت به موضوعات خاص در حوزه بررسی محصول، مدیریت ارتباط با مشتری، بازارهای مالی و چهره‌های سیاسی است [3].

بازارهای مالی، از جمله بازارهای سهام، بازارهایی هستند که به‌تازگی در کشور ایران، افراد بسیاری به آن توجه کرده‌اند. به همین دلیل، شبکه‌های اجتماعی که مرتبط با این بازارها هستند نیز، بسیار در مرکز توجه هستند. در این شبکه‌های اجتماعی، کاربران نظرات خود را بر اساس وضعیت سهام‌ها و یا کلیت بازار، در قالب یک پست یا توییت منتشر می‌کنند. نظرات کاربران می‌تواند انعکاسی از عملکرد، اخبار یا ارزش قیمتی شرکت‌ها و بر اساس تحلیل‌هایی که از یک سهم ارائه می‌شود، باشد. یا حتی بدون هیچ تحلیل و پشتوانه، به‌طور کامل، احساسی و غیرفنی باشد. به همین دلیل می‌توان گفت که این نظرات با احساسات مثبت یا منفی همراه است.

تحلیل احساسات در این حوزه، سنجش نگرش افراد به هر سهم است، که می‌تواند برای مدیران سطح بالا در سازمان‌های مختلف، شرکت‌هایی که رضایت یا عدم رضایت سهام‌داران برایشان مهم است، یا برای پیش‌بینی روند آتی یک سهم استفاده شود. این نگرش‌ها به‌طور لزوم ممکن است درست نباشد، ولی جمع‌بندی این نظرات می‌تواند یک موج در جهت صعود یا نزول یک سهم ایجاد کند.

خاص از این روش استفاده می‌شود. ساخت واژگان احساسی در این روش را می‌توان یا با استفاده از وقوع کلمات با یکدیگر ایجاد کرد؛ به‌عنوان مثال، اگر یک کلمه همراه کلمه‌ای با قطبیت مثبت بیاید، می‌توان گفت آن کلمه نیز قطبیت مثبت دارد؛ یا می‌توان با استفاده از پیکره در یک حوزه خاص، قطبیت کلمات یک فرهنگ واژگان را که از قبل وجود دارد، با آن پیکره تطبیق داد.

ما در کار خود از روش مبتنی بر پیکره برای ساخت واژگان پیشنهادی خود استفاده می‌کنیم. پیکره مورد استفاده در این روش در حوزه بورس است که ما آن‌ها را از پایگاه سهام‌باب^۴ استخراج کردیم و شامل ۱.۱۰۰.۰۰۰ نظر از سهامداران بورس است. که لغات و قطبیت آن‌ها به صورت خودکار استخراج می‌شود. آزمایش‌های ما نشان می‌دهد که استفاده از این روش برای ساخت واژگان، در مقایسه با واژگانی که به صورت عمومی برای تعیین احساسات در سطح جملات استفاده می‌شود، دقت را در حدود ۲۰ درصد افزایش می‌دهد.

در ادامه این مقاله به موضوعات زیر پرداخته می‌شود. بخش ۲، به برخی از مهم‌ترین پژوهش‌های مربوط به توسعه واژگان احساسی برای زبان فارسی پرداخته می‌شود. در بخش ۳، فرایند توسعه واژگان احساسی که ما تولید کردیم، به طور مفصل شرح داده می‌شود. ارزیابی و تست واژگان تولیدشده در بخش ۴ انجام گرفت. و در بخش ۵، نتیجه‌گیری و پیشنهادهایمان مطرح شد.

۲- پیشینه پژوهش

مهم‌ترین و اصلی‌ترین مرحله در فرایند جلب رضایت مشتری، شناسایی توقعات، انتظارات و گاهی الزامات طرح‌شده از طرف مصرف‌کننده است. باتوجه‌به این که تمام مشتریان به صورت حضوری در دسترس نیستند، یا افرادی که نظرات مهمی در مورد محصول دارند، آن را با تولیدکنندگان در میان نمی‌گذارند، می‌توان با استفاده از شبکه‌های اجتماعی و با در نظر گرفتن این نکته که افراد در این شبکه‌ها نظرات خود را با دیگران به اشتراک می‌گذارند، اطلاعات مورد نیاز خود را بدون نیاز به حضور مشتری به دست آورد.

کاربردهای تحلیل احساسات در سه سطح سند^۵، جمله و وجه^۶ بررسی می‌شود [6]. تمرکز کار ما بر روی

یک ویژگی منحصربه‌فرد در این شبکه‌ها این است که اطلاعات قیمتی هر سهم در روز درج نظر برای آن سهم موجود است. از آنجاکه سهم‌ها، همه‌روزه در حال رشد یا ریزش هستند، عقاید یا نظراتی که سهامداران یا تحلیل‌گران، روزانه منتشر می‌کنند، از این رشد و ریزش‌های سهم تأثیر می‌پذیرند. ما در این پژوهش ثابت کردیم که احساس نهفته در نظرات یا توییت‌های درج‌شده در روزهایی که سهم در حالت ریزش یا رشد است، با اطلاعات قیمتی آن سهم هم‌بستگی دارند. در سامانه پیشنهادی از این موضوع جهت استخراج واژگان به صورت خودکار برای بهبود کیفیت عقیده‌کاوی در این شبکه‌ها استفاده شد.

به‌طور کلی دو روش اصلی برای تحلیل احساسات وجود دارد. روش مبتنی بر واژگان^۱ و روش مبتنی بر یادگیری ماشین [4]. روش‌های یادگیری ماشین، برای پیش‌بینی قطبیت احساسات بر اساس مجموعه داده‌های آموزشی و تست عمل می‌کنند. و می‌توانند از روش‌های یادگیری با ناظر (از داده‌های برچسب‌گذاری‌شده برای دسته‌بندی متن استفاده می‌کنند) و بدون ناظر (از داده‌های خام برای دسته‌بندی متن استفاده می‌کنند) استفاده کنند. درحالی‌که رویکرد مبتنی بر واژگان نیازی به آموزش اولیه برای استخراج داده‌ها ندارد. از یک فهرست کلمات که از پیش تعریف شده است، استفاده می‌کند. به‌طوری‌که هر کلمه با یک احساس خاص همراه است. ما برای کار خود از روش مبتنی بر واژگان استفاده کردیم.

دو نوع واژگان احساسی برای تحلیل احساسات وجود دارد؛ واژگان احساسی عمومی که هر کلمه در آن با یک قطبیت همراه است و واژگان احساسی برای حوزه خاص^۲، که قطبیت هر کلمه باتوجه‌به آن حوزه مشخص می‌شود. برای مثال کلمه «مشخص» در حوزه فیلم قطبیت منفی دارد، اما در حوزه سیاست قطبیت مثبت خواهد داشت.

برای ساخت واژگان احساسی، سه روش کلی وجود دارد. روش دستی، که هر کلمه توسط یک خبره باتوجه‌به قطبیت آن کلمه برچسب‌گذاری می‌شود، و کاری زمان‌بر و پرهزینه است. ساخت واژگان احساسی بر اساس فرهنگ واژگان، مانند WordNet [5] که مترادف و متضاد کلمات در آن قرار دارد. و ساخت واژگان احساسی با استفاده از پیکره^۳. اغلب برای ایجاد واژگان احساسی برای حوزه‌های

¹ Lexicon based

² Domain Specific

³ Corpus

⁴ <https://www.sahamyab.com/>

⁵ Document Level

⁶ Aspect Level



پیش‌بینی احساسات در سطح سند است؛ که ما در آن به هر سند یک برچسب مثبت یا منفی اختصاص می‌دهیم. به‌طور کلی اطلاعات موجود در اسناد متنی را می‌توان به دو دسته تقسیم کرد: (۱) عینی^۱ (۲) ذهنی^۲. منظور از عینی، واقعیت‌ها، دستورهای واقعی و قابل مشاهده درباره موجودیت‌های مستقل و اتفاقاتی است که در جهان می‌افتد. اما اطلاعات ذهنی، بازتاب عواطف انسانی یا مشاهداتی است که مردم نسبت به دنیای خارج و اتفاقات آن دارند [7]. نظرات کاربران جزء اسناد متنی ذهنی هستند که به‌سرعت در دنیای مجازی در حال تولید هستند؛ بنابراین، حجم زیادی دارند و امکان بررسی نظرات به‌صورت دستی نیست. از این‌رو، نیازمند به‌کارگیری و انتخاب روش‌های مناسب جهت بررسی نظرات به‌صورت خودکار هستیم. به‌صورت کلی روش‌های موجود برای تحلیل احساسات به دو دسته تقسیم می‌شوند؛ روش‌های بر اساس یادگیری ماشین و روش‌های بر اساس واژگان. تاکنون بیشتر کارهایی که در حوزه تحلیل احساسات صورت گرفته، از روش‌های یادگیری ماشین بهره برده‌اند [8, 9]. بیشتر این کارها از الگوریتم ماشین بردار پشتیبان^۳ و نایو بیز^۴ برای کار خود استفاده کرده‌اند. آنها می‌توانند از روش‌های یادگیری با ناظر^۵ یا بدون ناظر استفاده کنند. روش‌های با ناظر، برای طبقه‌بندی متن از داده‌های دارای برچسب استفاده می‌کنند، درحالی‌که روش‌های بدون ناظر فقط از داده‌های خام استفاده می‌کنند [10]. روش‌های مبتنی بر واژگان برای تشخیص احساسات یک متن، از یک واژگان استفاده و با محاسبات آماری بر روی کلمات مثبت و منفی، قطبیت متن داده‌شده را محاسبه می‌کنند. از مهم‌ترین مزایای این روش، سرعت بالا و عدم نیاز به داده‌های آموزش است. و عیب اصلی این روش نیز عدم مقیاس‌پذیری آن است [11]. در ادامه این بخش به بررسی برخی از روش‌های ارائه‌شده مبتنی بر یادگیری ماشین و مبتنی بر واژگان خواهیم پرداخت.

۱-۲- روش‌های مبتنی بر یادگیری ماشین

ترنی^۶ [12]، یک الگوریتم یادگیری ساده بدون نظارت برای دسته‌بندی بررسی‌ها به‌عنوان توصیه‌شده (انگشت

شست بالا)^۷ یا توصیه نشده (انگشت شست پایین)^۸ ارائه داده است. که در آن دسته‌بندی یک بررسی با میانگین گرایش معنایی عبارات موجود در بررسی، که حاوی صفت یا قید است، پیش‌بینی می‌شود. الگوریتم پیشنهادی او هنگام ارزیابی در ۴۱۰ بررسی از پایگاه Epinions، نمونه‌برداری از چهار دامنه مختلف (بررسی اتومبیل، بانک، فیلم، و مقصد سفر) به میانگین صحت ۷۴٪ می‌رسد.

رانی^۹ و کومار^{۱۰} [13]، تحلیل احساسات را بر روی داده‌های بررسی فیلم، که به زبان هندی است، با استفاده از تنظیمات مختلف پیکره‌بندی شبکه عصبی پیچیده^{۱۱} انجام دادند. آن‌ها نتایج به‌دست‌آمده را توسط الگوی شبکه عصبی پیچیده خود با پیشرفته‌ترین نتایج الگوریتم‌های سنتی یادگیری ماشین مقایسه کردند. نتایج کار آن‌ها نشان داد که الگوی پیشنهادی آن‌ها قادر به دستیابی به عملکرد بهتری نسبت به رویکردهای سنتی یادگیری ماشین است و به دقت ۹۵٪ رسید.

برای غلبه بر کاستی‌ها در روش‌های تحلیل احساسات فعلی، یک روش تحلیل احساسات بر اساس شبکه‌های عصبی بازگشتی^{۱۲} به نام حافظه طولانی کوتاه‌مدت دوجهته^{۱۳} ارائه شده است [14]. آن‌ها داده‌های خود را که شامل ۱۵۰۰۰ متن است، بعد از تعبیه کلمات^{۱۴} به‌وسیله الگوی Word2vec و محاسبه وزن کلمات به‌وسیله الگوریتم TF-IDF، به برداری وزن‌دار تبدیل و سپس به شبکه حافظه طولانی کوتاه‌مدت دوجهته اعمال کردند. آزمایش‌های آن‌ها نشان می‌دهد که الگوی پیشنهادیشان دقت^{۱۵} و F1-Score بالاتری نسبت به شبکه عصبی پیچیده، شبکه عصبی بازگشتی، حافظه طولانی کوتاه‌مدت و الگوریتم نایو بیز دارد. دقت و F1-Score به‌دست‌آمده در کار آن‌ها به ترتیب در حدود ۹۱ و ۹۲ درصد است. در کاری مشابه [15]، برای تعبیه کلمات، به‌جای استفاده از Word2vec از الگوی Doc2vec استفاده کردند؛ همچنین برای استخراج هرچه بهتر ویژگی‌ها از شبکه عصبی پیچیده نیز در کنار شبکه حافظه طولانی کوتاه‌مدت استفاده کردند. آن‌ها با آزمایش‌ها نشان دادند که روش پیشنهادی آن‌ها در مقایسه با شبکه عصبی پیچیده، حافظه طولانی کوتاه‌مدت، حافظه طولانی

⁷ Thumbs Up

⁸ Thumbs Down

⁹ Rani

¹⁰ Kumar

¹¹ Convolutional Neural Network (CNN)

¹² Recurrent Neural Network (RNN)

¹³ Bidirectional Long Short-Term Memory (BiLSTM)

¹⁴ Word Embedding

¹⁵ Precision

¹ Objective

² Subjective

³ SVM

⁴ Naïve Bayes

⁵ Supervised

⁶ Turney

ارائه شد. برای آزمایش مجموعه‌دادگان، از شش الگویی که به‌تازگی برای تحلیل احساسات بر اساس وجه، روی حوزه‌های مختلف به زبان انگلیسی، کار شده‌است و تمرکزشان بر روی روش‌های یادگیری عمیق است، استفاده شد. در میان نتایج تمامی الگوها روی کار پیشنهادی آن‌ها، الگوی TD-LSTM نتایج شگفت‌انگیزی داشته‌است، زیرا این الگو نسبت به الگوهای دیگر در مجموعه‌دادگان‌های انگلیسی نتیجه ضعیف‌تری داشت، اما در کار آن‌ها نتیجه بهتری نسبت به سایر الگوها دارد.

در کاری دیگر [25]، از دو الگوی یادگیری عمیق (رمزگذارهای خودکار^۸ و شبکه‌های عصبی پیچیده) در مجموعه داده‌های بررسی فیلم فارسی استفاده شد. نتایج به‌دست‌آمده از این دو الگو با پرسپترون چند لایه^۹ مقایسه شد، نتایج نشان می‌دهد که رمزکننده‌های خودکار از پرسپترون چندلایه، دقت بالاتری دارند و الگوی شبکه عصبی پیچیده پیشنهادی نیز از رمزکننده‌های خودکار با دقت ۸۲/۶ درصد، عملکرد بهتری دارد.

۲-۲- روش‌های مبتنی بر واژگان

بیشتر افرادی که در حوزه تحلیل احساسات کار می‌کنند، تمرکز خود را بر روی روش‌های مبتنی بر یادگیری ماشین قرار داده‌اند و تعداد کمی از آن‌ها توجه خود را به روش‌های مبتنی بر واژگان معطوف کردند. تا به حال فهرست واژگان زیاد و خوبی برای کارهای انگلیسی‌زبان تولید شده‌است [26-29]. اما به تولید واژگان برای تحلیل احساسات به زبان فارسی زیاد توجه نشده‌است.

یکی از واژگان‌های قطبیت شناخته‌شده برای انگلیسی، SentiWordNet است [30]. در SentiWordNet، سه امتیاز به هر یک از مجموعه‌های هم‌معنی WordNet اختصاص داده شده‌است که نشان می‌دهد این مجموعه‌ها چقدر مثبت، منفی و خنثی هستند. این منبع شامل بیش از ۱۱۷۰۰۰ مجموعه هم‌معنی است. پیشنهاد اصلی ساخت SentiWordNet طبقه‌بندی مجموعه‌های هم‌معنی WordNet با استفاده از واژه‌نامه‌های مجموعه‌های هم‌معنی است. نکته مهمی که باید در مورد امتیازات اختصاص‌یافته به هر مجموعه هم‌معنی در نظر گرفته‌شود، این است که این امتیازات قدرت قطبیت را نشان نمی‌دهد. آن‌ها فقط نشان می‌دهند که یک مجموعه هم‌معنی چقدر مثبت، منفی

کوتاه‌مدت دوجهته و CNN-LSTM صحت^۱ بالاتری دارد که صحت به‌دست‌آمده برای روش پیشنهادی آن‌ها، بر روی مجموعه‌دادگان خود که از مقالات فرانسوی به‌دست‌آمده از روزنامه‌های ملی و بین‌المللی است، در حدود ۹۱ درصد است.

تاکنون کارهای زیادی با استفاده از الگوریتم ماشین بردار پشتیبان و نایو بیز برای تحلیل احساسات انجام شده‌است [16-20]. اولین کاری که برای مسئله دسته‌بندی در تحلیل احساسات به زبان فارسی صورت گرفت [21]، از دو روش استاندارد ماشین بردار پشتیبان و نایو بیز، در حوزه بررسی فیلم استفاده شده‌است. همچنین شش ویژگی حضور^۲ و تکرار^۳ Unigram، Bigramها و Trigramها برای نمایش اسناد، با هم مقایسه شدند. باتوجه به ارزیابی‌هایی که انجام دادند، متوجه شدند که در کار آن‌ها الگوریتم ماشین بردار پشتیبان نسبت به نایو بیز عملکرد بهتری داشته‌است و همچنین، استفاده از ویژگی Unigram در مقایسه با Bigram و Trigram کارایی دسته‌بندی کننده را بهبود می‌بخشد و همچنین، در نظر گرفتن فقط حضور یک ویژگی نتیجه بهتری از تکرار آن دارد. در کاری مشابه [22]، با بررسی چهار معیار اطلاعاتی مختلف، شامل فراوانی اسناد^۴، واریانس تکرار اصطلاح^۵، اطلاعات متقابل^۶ و اطلاعات متقابل اصلاح‌شده^۷ که توسط آن‌ها پیشنهاد شد، متوجه شدند که روش پیشنهادی عملکرد به‌نسبت بهتری از رویکردهای فراوانی اسناد، واریانس تکرار اصطلاح و اطلاعات متقابل دارد. اطلاعات متقابل اصلاح‌شده به‌طور کلی می‌تواند به ۸۵٪ از معیار F-Score برسد.

نسخه دیگری از ماشین بردار پشتیبان که با الگوریتم بهینه‌سازی ازدحام ذرات برای داده‌های بررسی فیلم در تویتر استفاده شده [23]، نشان می‌دهد که استفاده از الگوریتم بهینه‌سازی ازدحام ذرات برای تعیین شاخص‌های ماشین بردار پشتیبان، و همچنین استفاده از ویژگی‌های n-grams و به‌ویژه Unigrams می‌تواند صحت را تا ۴ درصد بهبود بخشد. این بهبود با پاک‌سازی داده‌ها، می‌تواند تا حدود ۲ درصد بیشتر نیز افزایش یابد.

نخستین مجموعه‌دادگان به نام Pars-ABSA که به‌طور کامل مبتنی بر وجه در زبان فارسی است، در [24]

¹ Accuracy

² Presence

³ Frequency

⁴ Document Frequency (DF)

⁵ Term Frequency Variance (TFV)

⁶ Mutual Information (MI)

⁷ Modified Mutual Information (MMI)

⁸ Autoencoders

⁹ Multilayer Perceptron (MLP)



یا خنثی است. در نتیجه، جمع این سه نمره برای هر مجموعه هم‌معنی برابر با یک است. یکی دیگر از واژگان قطبیت عمومی شناخته‌شده برای زبان انگلیسی، واژگان NRC است [31]. این واژگان، به هر کلمه برجسب‌های احساسی، عصبانیت، ترس، انتظار، انزجار، ... و همچنین برجسب‌های مثبت و منفی اختصاص می‌دهد که حدود ۱۵۰۰۰ کلمه است و به‌صورت دستی از طریق Amazon's Mechanical Turk ایجاد شده‌است.

نخستین رویکرد تحلیل احساسات به زبان فارسی [32] که مبتنی بر واژگان است، با ارائه یک چارچوب برای تحلیل احساسات انجام شد. آن‌ها دو منبع برای تحلیل احساسات در زبان فارسی معرفی کردند: ۱- یک واژگان فارسی که مرتبط است با کلمات احساسی فارسی به همراه قطبیت آن. ۲- یک مجموعه‌دادگان که به‌صورت دستی جمع‌آوری و توسط یک شخص خبره برجسب‌گذاری شده‌است. آن‌ها از نظریهٔ دمپستر شافر^۱ برای تعیین قطبیت هر سند استفاده کردند. نتایج کار آن‌ها نشان می‌دهد که روش پیشنهادی آن‌ها نسبت به روش‌های مبتنی بر یادگیری ماشین، رتبهٔ F-Score بالاتری در حدود ۹۰ درصد کسب می‌کند. آن‌ها در کاری دیگر [33]، برای کمک به رفع مشکل کمبود منابع برای تحلیل احساسات در زبان فارسی، دو منبع جدید به نام‌های SPerSent و CNRC را برای تحلیل احساسات در زبان فارسی را ارائه می‌دهند. SPerSent، یک مجموعه‌دادگان در سطح جمله است که هر جمله با دو برجسب همراه است، یک برجسب دودویی برای تعیین قطبیت و یک برجسب رتبه‌بندی پنج ستاره. CNRC، یک واژگان فارسی است که برای ایجاد آن از لغت‌نامهٔ NRC [34] همراه با سه مرحله پردازش استفاده شده‌است. آن‌ها برای ارزیابی واژگان CNRC، آن را با نسخهٔ فارسی واژگان‌های NRC و Senti_Str [35] بر روی مجموعه‌دادگان SPerSent توسط الگوریتم یادگیری ماشین نایو بیز مقایسه کردند، نتایج نشان داد که واژگان CNRC کارایی و دقت بالاتری را نسبت به دو واژگان دیگر دارد. آن‌ها در کاری مشابه با این روش [36]، برای اثبات این که ترجمهٔ مستقیم واژگان انگلیسی به فارسی در کار تحلیل احساسات کیفیت مناسب را ندارد، چهار واژگان را با هم مقایسه کردند که شامل نسخهٔ فارسی واژگان‌های NRC، SentiStrength، CNRC و Adjectives است. نتایج نشان می‌دهد که ترجمهٔ مستقیم استفاده‌شده در NRC، ضعیف‌ترین

¹ Dempster-Shafer

عملکرد را به همراه دارد؛ درحالی‌که پیش‌پردازش و پالایش واژگان مورد استفاده در SentiStrength و CNRC باعث بهبود عملکرد می‌شود. همچنین، نتایج نشان می‌دهد که استفاده از فقط صفت در مقایسه با استفاده از NRC، منجر به نتایج بهتری می‌شود.

ثابتی و همکاران [37]، یک روش جدید مبتنی بر گراف، برای انتخاب و گسترش مجموعهٔ اولیه^۲ برای تولید واژگان قطبیت عمومی، به نام LexiPers پیشنهاد داده‌اند، که این واژگان قطبیت شامل بیش از ۶۰۰۰ کلمه است. آن‌ها از رویکردهای مبتنی بر واژگان و همچنین، پیکره بهره‌مند شدند.

امیری و همکارانش روش دیگری را برای تحلیل احساسات مبتنی بر واژگان، ارائه دادند [38]. آنها واژگانی فارسی را که متشکل از صفات، واژه‌ها و اصطلاحات نظریافته در دو دستهٔ رسمی و غیررسمی است، جمع‌آوری کردند. همچنین واژگان جمع‌شده با فارسی استاندارد یا فارسی منسوخ‌شده که توسط عدهٔ خاصی از بومی‌زبانان استفاده می‌شود، مطابقت دارد. آن‌ها همچنین، یک رابط شبکه ایجاد کردند که افراد بومی را قادر می‌سازد به‌صورت دستی یک امتیاز به کلمات واژگان اختصاص دهند. به‌دنبال ایجاد یک واژگان احساسات فارسی حاشیه‌نویسی‌شده، آن‌ها یک پایپ‌لاین^۳ زبانی بر اساس چارچوب GATE طراحی و توسعه دادند. که اجزای آن شامل یک توکن‌ساز^۴ فارسی، تقسیم‌کنندهٔ جمله^۵، برجسب‌گذار POS^۶ و فرهنگ لغت^۷ بود. در نتیجه، آن‌ها صحت حدود ۶۵٪ را گزارش کردند که در مقایسه با رویکردهای مبتنی بر واژگان مشابه، به‌عنوان یک پیشرفت در نظر گرفته‌شد.

آقای دهخرفانی [39]، یک روش مبتنی بر ترجمهٔ جدید برای ایجاد واژگان قطبیت در زبان‌هایی که منابع واژگانی کمی برای تحلیل احساسات وجود دارد، پیشنهاد می‌دهد و آن را برای زبان فارسی به‌کارمی‌برد. برای ساخت این واژگان قطبی که SentiFars نام دارد، از چندین منبع واژگان قطبیت انگلیسی به نام‌های NRC، SentiWordNet، SenticNet و واژگان قطبیت Liu's کمک گرفت. روش پیشنهادی آن‌ها در چهار مرحله انجام می‌شود. نخست، کلمات به فارسی ترجمه می‌شوند. سپس کلمات ترجمه‌شده به‌صورت دستی برجسب مثبت، منفی

² Seed

³ Pipeline

⁴ Tokenizer

⁵ Sentence Splitter

⁶ Part of Speech Tags

⁷ Gazetteer

داده‌ها یا توییت‌هایی که برای این کار در نظر گرفته شده‌است، و همچنین، برای ارزیابی و آزمون استفاده می‌شود، از پایگاه سهام‌یاب استخراج و در پایگاه داده خود ذخیره کردیم. همچنین، ارزش قیمتی هر سهم را در زمان درج آن، از سایت رسمی بورس ایران^۲ استخراج و در یک جدول جداگانه از پایگاه داده ذخیره کردیم. جدول (۱) اطلاعات سهم‌های ذخیره‌شده را در پایگاه داده نمایش می‌دهد.

باتوجه به سامانه پیشنهادی در شکل (۱)، بعد از جمع‌آوری داده‌ها، پیش‌پردازش‌هایی بر روی آن‌ها صورت می‌گیرد. در ادامه، به بررسی و توضیح هریک از این پیش‌پردازش‌ها می‌پردازیم.

نام سهم	تعداد	سال
دی	85461	1392-1399
فملی	43054	1392-1399
فولاد	35978	1392-1399
حفاری	8917	1392-1399
حکشتی	30907	1392-1399
خودرو	162371	1392-1399
خسپا	157308	1392-1399
ستران	31383	1392-1399
شبندر	147451	1392-1399
شپنا	88118	1391-1399
شتران	22605	1392-1399
تاپیکو	51276	1392-1399
وتجارت	50341	1392-1399
وبملت	77265	1392-1399
وبصادر	107565	1392-1399

(جدول-۱): اطلاعات سهم‌های انتخاب شده
(Table-1): Information of selected shares

۳-۱- حذف توییت‌های غیر ضروری:

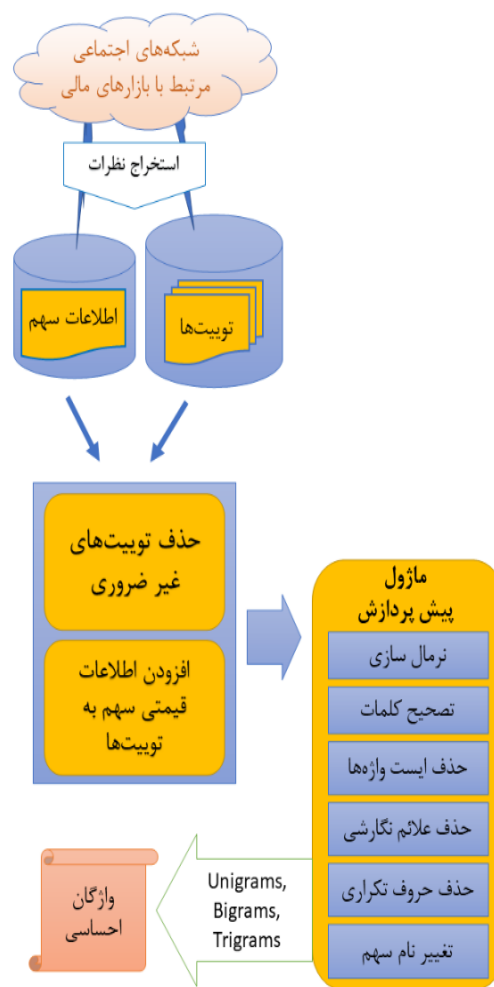
متأسفانه اطلاعات قیمتی روزانه سهم‌ها، برای همه روزهای درج نظر، در اختیار نیست. از آنجاکه ما برای ساخت خودکار واژگان موردنظر به این اطلاعات قیمتی نیاز داشتیم، ناچار به حذف این دسته از توییت‌ها شدیم. تعداد این توییت‌ها بعد از این مرحله به ۹۰۰۰۰۰ توییت کاهش یافت. چند نمونه از این اطلاعات قیمتی به همراه زمان و نام سهم، در جدول (۲) آورده شده‌است.

² <http://www.tsetmc.com/>

و خنثی دریافت می‌کنند. مرحله بعد استخراج ویژگی توسط واژگان قطبیت انگلیسی است. سپس کار دسته‌بندی توسط دسته‌بندی‌کننده منطقی^۱ صورت می‌گیرد که این کار توسط یادگیری نگاشت بین ورودی‌های توصیف‌شده توسط ویژگی‌های استخراج‌شده، و سه برجسب رده (مثبت، منفی و خنثی) انجام می‌شود. در نهایت بعد از آزمایش‌های انجام شده توانست با در نظر گرفتن هر چهار واژگان انگلیسی یادشده برای استخراج ویژگی، به صحت ۹۵/۹۲ درصد برسد.

۳- سامانه پیشنهادی

همان‌طور که پیشتر گفته شد، به دلیل اهمیت بازارهای مالی، ما به صورت خودکار یک واژگان برای تحلیل احساسات در حوزه بورس ایجاد کردیم. این کار در چند مرحله صورت گرفت که در شکل (۱) آمده‌است.



(شکل-۱): روش پیشنهادی برای ایجاد واژگان
(Figure-1): Proposed method for lexicon creation

¹ Logistic Classifier

(جدول ۲-): اطلاعات قیمتی سهم
(Table-2): Share price information

نام سهم	تاریخ	بیشترین قیمت	قیمت پایانی	قیمت دیروز
دی	1399/12/25	19941	19486	18813
خودرو	1399/02/27	1174	1174	587
حکشتی	1397/11/29	4116	4072	4121
حفاری	1393/07/30	5687	5605	5607
فولاد	1388/03/04	1900	1865	1884

۳-۲- افزودن اطلاعات قیمتی سهم به توییت‌ها:

بعد از حذف توییت‌های غیرضروری، اطلاعات قیمتی برای هر توییت به آن اضافه شد. این اطلاعات که قسمتی از آن در جدول (۲) آورده شده، جهت تعیین وزن و قطبیت کلمات واژگان نهایی استفاده شده است. جدول (۳)، بخشی از توییت‌ها را به همراه اطلاعات قیمتی آن نمایش می‌دهد.

(جدول ۳-): توییت‌ها به همراه اطلاعات قیمتی
(Table-3): Tweets with price information

اطلاعات قیمتی	توییت
'1392/08/26: تاریخ', 'خودرو': 'نماد سهم', 'کمترین قیمت', 3150: 'بیشترین قیمت', 'آخرین', 3090: 'قیمت پایانی', 3070: 'قیمت', 3111: 'اولین قیمت', 3075: 'معامله', 25003468: 'حجم', 3075: 'دیروز'	#خودرو #حکشتی امروز هر دو صف می‌شوند. خوشبین
'1393/01/26: تاریخ', 'شبندر': 'نماد سهم', 'کمتری قیمت', 11800: 'بیشترین قیمت', 'آخرین', 11185: 'قیمت پایانی', 11104: 'قیمت', 11104: 'اولین قیمت', 11190: 'معامله', 17345720: 'حجم', 11566: 'قیمت دیروز'	#شبندر سال ۹۴ تا سال ۹۶ شبندر بازم بهتر می‌شود.
'1396/09/27: تاریخ', 'فولاد': 'نماد سهم', 'کمترین قیمت', 2900: 'بیشترین قیمت', 'آخرین', 2871: 'قیمت پایانی', 2854: 'قیمت', 2860: 'اولین قیمت', 2860: 'معامله', 2850: 'دیروز', 32281473: 'حجم'	#فولاد سلام خدمت دوستان فولادی باتوجه به رشد خوب فولاد که نوش همه سهامدارن آن باشد، متأسفانه واگرایی منفی را نشان می‌دهد. احتیاط بیشتر پیشنهاد می‌شود. الهی همه پرسود باشند.

۳-۳- پیش‌پردازش:

یکی از مراحل مهم در تحلیل احساسات، پیش‌پردازش داده‌ها است. انتخاب روش‌های پیش‌پردازش مناسب،

می‌تواند باعث بهبود طبقه‌بندی صحیح داده‌ها شود [40]. به همین دلیل ما پیش‌پردازش داده‌ها را در شش مرحله انجام می‌دهیم.

۳-۳-۱- عادی‌سازی:

برای عادی‌سازی داده‌ها ما از ابزار هضم^۱ که یک کتابخانه برای زبان پایتون^۲ است، استفاده کردیم. این ابزار از ماژول‌های متفاوتی تشکیل شده است. که ما از ماژول عادی‌سازی آن برای کار خود استفاده کردیم. بعضی کلمات در فارسی از چند بخش تشکیل شده‌اند. این کلمات به‌طور معمول، با نیم‌فاصله از هم جدا می‌شوند. اما اغلب، افراد در توییت‌هایی که در شبکه‌های اجتماعی قرار می‌دهند، این نکته را رعایت نمی‌کنند. برای مثال، کلمه «آن‌ها» به‌اشتباه «آن‌ها» نیز نوشته می‌شود. یکی از کارهایی که این ماژول انجام می‌دهد، تبدیل این نوع فاصله‌ها به نیم‌فاصله است. همچنین برخی از حروف در فارسی دارای یونیکد متفاوت هستند. به‌طور خاص، کلمه «ی» و «ک» در بعضی مواقع به‌ترتیب به شکل «ی» و «ك» نیز نوشته می‌شود. این حروف نیز توسط این ماژول به شکل استاندارد خود در می‌آیند. در حقیقت، در این قسمت ما کلمات غیراستاندارد را به شکل استاندارد آوردیم.

۳-۳-۲- تصحیح کلمات:

در بیشتر متون غیررسمی صفحات وب فارسی، کاربران کلمات را همان‌گونه که در مکالمات روزانه استفاده می‌کنند، می‌نویسند. به همین دلیل این متون حاوی تعداد زیادی از کلمات با املاهای غیراستاندارد است. بنابراین، بررسی املاهای کلمات در زبان فارسی چالش برانگیزتر از زبان انگلیسی است [32].

برای حل این مسئله، فهرستی از کلمات به‌همراه تعداد رخداد آن‌ها از توییت‌های ذخیره‌شده در پایگاه داده استخراج شد. سپس اندازه این فهرست در دو مرحله کاهش یافت. در مرحله نخست، از پیکره شناخته‌شده همشهری [41]، به‌منظور تطبیق هر کلمه از فهرست به‌دست‌آمده با آن استفاده کردیم. با این فرض که اگر آن کلمه در پیکره همشهری وجود داشته‌باشد، املاهای آن استاندارد است. در مرحله دوم با فرض اینکه کلماتی که تعداد رخداد آن در کل داده‌ها کمتر از ده بار باشد، تأثیر ناچیزی روی تخمین قطبیت احساسی جملات

¹ <https://www.sobhe.ir/hazm/>

² Python

۳-۳-۶- تغییر نام سهم:

از آنجاکه نام سهم تأثیری بر روند تشخیص احساسات ندارد و باعث افزونگی در واژگان ما می‌شود، در این قسمت نام سهم‌ها را به یک اسم خاص (#سهم) تغییر دادیم. که هم از افزونگی در مجموعه‌دادگان جلوگیری کند و هم روند تطبیق کلمات بهبود یابد.

۳-۴- ساخت واژگان:

همان‌طور که در بخش ۱ نیز ذکر شد، سه روش کلی برای ساخت واژگان وجود دارد؛ که کار ما ساخت واژگان بر اساس پیکره است. رویکرد سامانه پیشنهادی برای ساخت واژگان، استفاده از ویژگی‌های n-grams (به‌طور ویژه Unigram، Bigram و Trigram) در داده‌های توییت است. همچنین، اطلاعات قیمتی متناسب با هر توییت برای محاسبه نرخ رشد هر سهم در زمان درج توییت نیز محاسبه می‌شود. دلیل محاسبه نرخ رشد، این است که ما فرض می‌کنیم زمانی که نرخ رشد یک سهم مثبت است، افرادی که توییت‌های خود را در این شبکه‌های اجتماعی برای آن سهم ارسال می‌کنند، به دلیل حس رضایت، کلمات مثبت بیشتری را در نظرات خود به کار می‌برند. همچنین، اگر نرخ رشد سهم منفی باشد، در توییت‌هایی که این افراد منتشر می‌کنند، به دلیل حس ناراحتی، کلمات منفی بیشتری به کار برده می‌شود. با در نظر گرفتن این فرض، برای هر توییت نرخ رشد متناسب با زمان ارسال آن و این که این نظر برای چه سهمی ارسال شده، محاسبه شد. برای محاسبه نرخ رشد، ما از رابطه (۱) - (نسخه اول)، و رابطه (۲) - (نسخه دوم) استفاده کردیم:

نرخ رشد = Growth Rate

قیمت پایانی سهم = P_{Close}

قیمت پایانی سهم در روز قبل = $P_{PrevClose}$

کمترین قیمت سهم در طول روز = P_{Min}

بیشترین قیمت سهم در طول روز = P_{Max}

Growth Rate =

$$\begin{cases} \frac{P_{Close} - \min(P_{Min}, P_{PrevClose})}{\min(P_{Min}, P_{PrevClose})}, & P_{Close} - P_{PrevClose} > 0 \\ \frac{P_{Close} - \max(P_{Max}, P_{PrevClose})}{\max(P_{Max}, P_{PrevClose})}, & P_{Close} - P_{PrevClose} < 0 \\ 0, & P_{Close} - P_{PrevClose} = 0 \end{cases} \quad (1)$$

Growth Rate =

$$\begin{cases} \frac{P_{Close} - P_{Min}}{P_{Min}}, & P_{Close} - P_{PrevClose} > 0 \\ \frac{P_{Close} - P_{Max}}{P_{Max}}, & P_{Close} - P_{PrevClose} < 0 \\ 0, & P_{Close} - P_{PrevClose} = 0 \end{cases} \quad (2)$$

خواهند گذاشت، از فهرست حذف شدند؛ و در نهایت، تعداد ۲۷۸۲ کلمه از فهرست باقی ماند. در جدول (۴) تعدادی از کلمات این فهرست به همراه معادل صحیح آن آورده شده است. در پایان، معادل درست هریک از این کلمات به صورت دستی و بر روی کل دیتابیس اعمال شد.

(جدول-۴): بخشی از فهرست کلمات فارسی غیراستاندارد به

همراه معادل استاندارد آن

(Table-4): Part of the list of informal Persian words with its formal equivalent

معادل صحیح آن کلمات	کلمات با املاي ناصحيح
منفی است	منفی
قیمت‌ها	قیمتا
می‌ریزد	میریزه
می‌گردد	میگرده
نشوند	نشن

۳-۳-۳- حذف ایست‌واژه‌ها:

برای کاهش اندازه داده‌ها و بهبود دقت تحلیل احساسات، ایست‌واژه‌هایی را که تأثیر قابل توجهی در تشخیص احساسات ندارند و پیوسته در مجموعه‌دادگان تکرار شده‌اند، حذف کردیم. برای این کار برخلاف روش‌های دیگر که از فهرست ایست‌واژه‌های شناخته‌شده برای تشخیص استفاده می‌کنند، ما با ایجاد یک Unigrams از کلمات درون مجموعه‌دادگان، به همراه تکرار هر کدام از آن‌ها برای شناسایی ایست‌واژه‌ها استفاده کردیم. با این فرض که تکرار کلمات ایست‌واژه در مجموعه‌دادگان، اختلاف چشم‌گیری با دیگر کلمات خواهند داشت. این روش باعث می‌شود کلمات تأثیرگذار در مجموعه‌دادگان که ممکن است توسط فهرست‌های ایست‌واژه موجود، به‌عنوان یک ایست‌واژه شناخته شوند، در مجموعه‌دادگان باقی بمانند.

۳-۳-۴- حذف علائم نگارشی^۱:

تمامی علائم نگارشی بجز «#» و «.» از مجموعه‌دادگان حذف شد.

۳-۳-۵- حذف حروف تکراری:

در بعضی مواقع در میان توییت‌ها، افراد برای تأکید بر یک موضوع بعضی کلمات را با تکرار بر یک حرف آن می‌نویسند. مثلاً برای تأکید بر کلمه «نه»، از «نهههههههه» استفاده می‌کنند. در این ماژول، تکرار این حروف حذف شد.

¹ Punctuations



در رابطه پیشنهادی هنگامی که نرخ رشد برابر با صفر شود، نشان‌دهنده آن است که این کلمه فاقد قطبیت بالاست و تأثیر آن در روند تحلیل احساسات در کار ما بسیار پایین است. در داده‌های اطلاعاتی ما، «قیمت پایانی» برای هر سهم در هر روز موجود است. قیمت پایانی عبارت است از میانگین قیمت‌های معامله‌شده در طول ساعات معاملاتی آن روز برای یک سهم. ما برای بازتاب بهتر تغییرات قیمتی در طول روز، از «بیشترین قیمت» و «کمترین قیمت» به‌منظور محاسبه نرخ رشد استفاده کردیم. دلیل این کار این است که هیجانات هم‌راستای قیمت پایانی سهم را در نظر بگیریم. برای مثال، اگر در آغاز شروع ساعات معاملاتی، یک سهم با ارزش بالایی شروع شود و در انتهای ساعات معاملاتی با کاهش چشم‌گیری مواجه شود، نظرات یا توییت‌هایی که کاربران برای آن سهم منتشر می‌کنند، از نظر احساسات بسیار مهم و قابل‌توجه است. همچنین، برعکس این موضوع نیز صادق است. یعنی اگر یک سهم در آغاز ساعات معاملاتی بازار با ارزش کم شروع به معامله شود و این ارزش در ادامه ساعات معاملاتی افزایش یابد، نظراتی که کاربران منتشر می‌کنند، اغلب با احساسات مثبت همراه است. از این‌رو، استفاده از کمترین قیمت و بیشترین قیمت، باعث تمایز بهتری برای محاسبه نرخ رشد می‌شود. منظور از تمایز بهتر این است که اگر شرایط ذکرشده اتفاق بیفتد، نرخ رشد برای آن توییت‌ها مثبت‌تر، یا منفی‌تر خواهد شد.

در شرایطی که سهمی از آغاز، در صف خرید باشد، کمترین قیمت سهم (P_{Min}) و قیمت پایانی سهم (P_{Close}) در طول روز یکی می‌شود، همچنین برعکس این موضوع نیز صادق است. به این صورت که اگر سهمی از آغاز در صف فروش باشد، بیشترین قیمت سهم (P_{Max}) و قیمت پایانی سهم (P_{Close}) نیز در طول روز یکی می‌شود، از این‌رو، رابطه (۱) - (نسخه اول) که نسخه پایه برای روش پیشنهادی است، در این حالت به‌جای استفاده از کمترین قیمت و بیشترین قیمت از قیمت پایانی سهم در روز قبل ($P_{PreClose}$) برای محاسبه نرخ رشد استفاده می‌کند. این در حالی است که ضابطه اول و دوم رابطه (۲) - (نسخه دوم) در این حالت برابر با صفر می‌شود. رابطه پیشنهادی (۲)، قطبیت بیشتر را به کلمات یا عباراتی اطلاق می‌کند که در روزهای پرنوسان دیده‌شوند و همچنین، در روزهایی که صف فروش یا خرید قفل است، روش پیشنهادی مقداری نزدیک به صفر می‌دهد. در واقعیت نیز روزهایی

که قیمت‌ها دچار نوسان هستند (چه در جهت مثبت و چه در جهت منفی)، بیشتر شاهد توییت کاربران خواهیم بود. یا برعکس، در روزهایی که نوسانات قیمتی کم است، کاربران فعالیت کمتری دارند و نظرات آن‌ها دارای قطبیت قابل‌توجه کمتری است. به عبارت دیگر، فعالیت کاربران با احساسات قوی را بیشتر در روزهای نوسانی بالا شاهد هستیم. رابطه پیشنهادی (۲) - (نسخه دوم) بیانگر این واقعیت است. همچنین، باید به این نکته توجه داشت که محور محاسبه نرخ رشد در روابط پیشنهادی، بر اساس میانگین سهم است و نوساناتی که بسیار لحظه‌ای هستند و یا حجم معاملات در آن پایین است، در قیمت پایانی سهم تأثیر مهمی نمی‌گذارند، یا به عبارت دیگر، این نوسانات معیار احساسات و قضاوت اهالی بازار نیستند. ما در کار خود هر دوی این روش‌ها را استفاده کرده و نتایج به‌دست‌آمده این موضوع را نیز تأیید می‌کند. نتایج در بخش ۳-۴ (ارزیابی نهایی) قابل‌مشاهده است. از این‌رو، در ادامه بیشتر به روش پیشنهادی (۲) - (نسخه دوم) پرداخته خواهد شد.

به منظور محاسبه نرخ رشد، نخست، بررسی می‌کنیم که قیمت پایانی سهم نسبت به روز گذشته افزایش داشته یا کاهش. اگر قیمت پایانی افزایش داشته‌باشد، از ضابطه اول رابطه (۲)، و اگر کاهش داشته‌باشد، از ضابطه دوم رابطه (۲)، نرخ رشد را محاسبه می‌کنیم.

همان‌طور که ذکر شد، استخراج کلمات واژگان، توسط ویژگی n-gramها انجام شد و ما نرخ رشد و رخداد این کلمات را در کنار آن وارد کردیم. محاسبه نرخ رشد هر یک از این کلمات واژگان به این صورت است که نخست، زمانی که این کلمات توسط ویژگی n-gram از توییت‌ها انتخاب می‌شوند نرخ رشد هر یک از توییت‌ها که آن کلمه در آن قرار دارد، در یک فهرست، برای آن کلمه در نظر گرفته می‌شود. پس از پیمایش تمامی توییت‌ها میانگین این نرخ‌رشد‌ها برای هر کلمه به‌عنوان نرخ رشد هر کلمه در نظر گرفته می‌شود. برای مثال اگر کلمه «فروش» در ۱۰۰۰ توییت وجود داشته‌باشد، میانگین نرخ رشد این توییت‌ها به‌عنوان نرخ رشد این کلمه در نظر گرفته می‌شود.

برای درک عمیق‌تر موضوع، رابطه پیشنهادی در قالب یک مثال ساده از چند توییت به‌همراه نحوه محاسبه نرخ رشد و عملکرد روش پیشنهادی برای حالات مختلف آورده شده‌است. برای این کار در جدول (۵) اطلاعات

1-1-1 Unigram

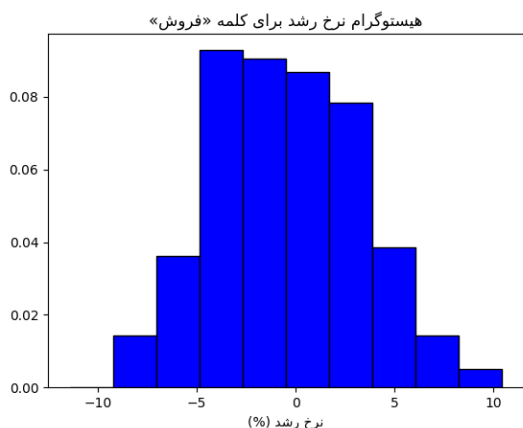
واژگانی که با استفاده از ویژگی Unigram به دست آمد، به نسبت واژگان‌هایی که با استفاده از ویژگی‌های Bigram و Trigram به دست آمد، در کار تحلیل احساسات از کیفیت کمتری برخوردار است. جدول (۶)، تعدادی از کلمات موجود در واژگان unigram را نمایش می‌دهد.

(جدول-۶): قسمتی از واژگان به دست آمده با Unigram (Table-6): Part of the Lexicon obtained with Unigram feature

همان‌طور که پیشتر ذکر شد، واژگان‌های ما به صورت خودکار استخراج شده‌است. برای بررسی بیشتر و مشاهده پراکندگی کلمات به دست آمده با Unigram، نسبت به رشد قیمتی سهم در روز درج نظر، در قالب نمودار هیستوگرام چند نمونه را بررسی کردیم که یکی از

لغت	رخداد در واژگان	میانگین نرخ رشد (%)
فروش	89179	-0.533
خرید	128991	0.051
خوشبین	161192	0.266
حمایت	30432	-0.604

آن‌ها در شکل (۲) آورده شده‌است.



(شکل-۲): نمودار هیستوگرام برای یک نمونه کلمه

به دست آمده با Unigram (Figure-2): Histogram for a sample word obtained with unigram

باتوجه به شکل (۲)، مشاهده می‌کنیم که توزیع آماری کلماتی که با استفاده از ویژگی Unigram استخراج شده‌اند، اغلب مختلط است. در ادامه، دلیل این مسئله را بررسی می‌کنیم. برای مثال، توییت «امروز روز خوبی برای سرمایه‌گذاری است و برای این سهم از فردا صف خرید تشکیل می‌شود»، زمانی که ارزش آن سهم نسبت به روز قبلش در حالت صعودی بود، منتشر شد. همچنین، توییت «این هفته، سهم‌ها وضعیت خوبی ندارند، و اصلاً توصیه به

قیمتی دو سهم به همراه توییتی که برای هر کدام گذاشته شده‌بود، آورده شده‌است. در ادامه روند محاسبه نرخ رشد در رابطه پیشنهادی با استفاده از این اطلاعات قیمتی شرح داده می‌شود.

نمونه اول:

نخست، نیاز است که نرخ رشد سهم نسبت به روز پیش محاسبه شود، تا از این طریق بفهمیم از کدام ضابطه باید استفاده کرد.

$$P_{Close} = 10981, P_{PrevClose} = 10222$$

$$P_{Min} = 10165, P_{Max} = 11210$$

$$P_{Close} - P_{PrevClose} = 10981 - 10222 > 0$$

$$\Rightarrow \text{Growth Rate} = \frac{P_{Close} - P_{Min}}{P_{Min}} = \frac{10981 - 10165}{10165} \approx 8\%$$

نمونه دوم:

$$P_{Close} = 24256, P_{PrevClose} = 25804$$

$$P_{Min} = 23289, P_{Max} = 26578$$

$$P_{Close} - P_{PrevClose} = 24256 - 25804 < 0$$

$$\Rightarrow \text{Growth Rate} = \frac{P_{Close} - P_{Max}}{P_{Max}} = \frac{24256 - 26578}{26578} \approx -8.7\%$$

مقادیر به دست آمده برای نرخ رشد برای هر دو روش پیشنهادی (نسخه اول و دوم) یکسان است. به این ترتیب نرخ رشد برای تمامی توییت‌ها محاسبه می‌شود، و سپس میانگین نرخ رشد هر n-gram بر اساس رخداد آن محاسبه می‌شود.

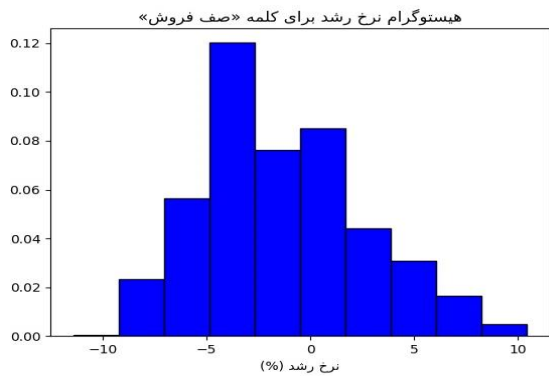
(جدول-۵): نمونه‌ای از توییت‌ها به همراه اطلاعات قیمتی

(Table-5): Sample tweets with price information

اطلاعات قیمتی	توییت
نماد سهم: ستران، بیشترین قیمت: 11210 (+9.66%)، کمترین قیمت: 10165 (-0.55%)، پایانی: 10981	#سهم این سهم پایین نداره کلی خبر خوب در شرکت نهفته است تغییر کاربری زمین و افزایش سرمایه و فروش زیر مجموعه و صادراتی بودن و افزایش نرخ سیمان و خرید حتی در صف خرید هم لذت خودش رو داره در این سهم چرا که فرداش هم صف خرید می‌شود و فردای فرداش هم همینطور
نماد سهم: حکشتی، حداکثر قیمت: 26578 (+3%)، حداقل قیمت: 23289 (-9.74%)، قیمت پایانی: 24256 (-6%)، آخرین معامله: 23431 (-9.2%)، اولین قیمت: 24800 (-3.9%)، قیمت پایانی دیروز: 37087720	#سهم بایداز حقوقی سهم به وزارت اطلاعات و قوه محترم قضایی شکایت بشه

همان‌طور که ذکر شد، ما از ویژگی‌های Unigram، Bigram و Trigram به عنوان ویژگی‌های n-gram استفاده کردیم. در ادامه، مشخصات هریک از واژگان‌های به دست آمده ذکر می‌شود.





(شکل-۳): نمودار هیستوگرام برای یک نمونه

به دست آمده با Bigram

(Figure-3): Histogram for a sample obtained with bigram

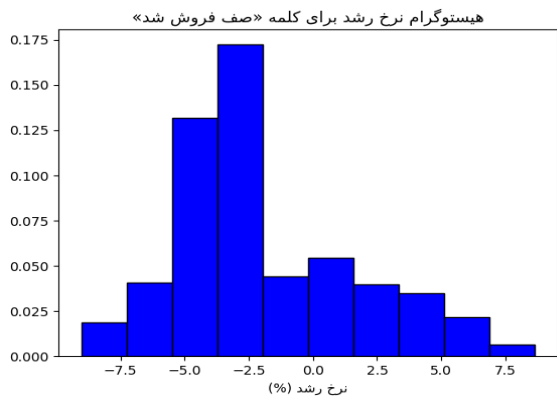
۴-۳-۲: Trigram

جدول (۸)، بخشی از واژگان به دست آمده با Trigram را نمایش می دهد.

(جدول-۸): قسمتی از واژگان به دست آمده با Trigram
(Table-8): Part of the lexicon obtained with Trigram feature

لغت	رخداد در واژگان	میانگین نرخ رشد (%)
صف فروش شد	705	-1.952
صف خرید مبارک	324	2.432
پر سود باشید	2574	-0.385
بازار منفی است	489	-1.064

انتخاب کلمات واژگان با استفاده از ویژگی Trigram، به نسبت واژگان به دست آمده با Bigram، باعث بهبود در کیفیت تحلیل احساسات بر روی توییت های بورسی شد. دلیل این بهبود همان طور که پیشتر ذکر شد، احتمال کمتر دیده شدن ترکیبات سه تایی کلمات، هم در شرایط مثبت بازار و هم در شرایط منفی بازار است. شکل (۴)، نمودار هیستوگرام کلمه «صف فروش شد» را نمایش می دهد.



(شکل-۴): نمودار هیستوگرام برای یک نمونه به دست آمده با

Trigram

(Figure-4): Histogram for a sample obtained with trigram

سرمایه گذاری نمی کنم.» در زمانی منتشر شد که سهم نسبت به روز قبلیش در وضعیت نزولی قرار داشت. اگر به این دو نظر، کمی دقت کنیم، متوجه می شویم که کلمه «خوبی»، هم در روزهای منفی سهم و هم در روزهای مثبت سهم رؤیت می شود. در صورتی که اگر کلمات بعد یا قبل را نیز در نظر می گرفتیم این دو، متفاوت از یکدیگر می شدند. بنابراین، می توان نتیجه گرفت که استفاده از واژگان مبتنی بر Unigram، برای تحلیل احساسات در این حوزه از کارآمدی مناسبی برخوردار نباشد.

۴-۳-۱: Bigram

جدول (۷)، بخشی از واژگان به دست آمده توسط Bigram را نمایش می دهد.

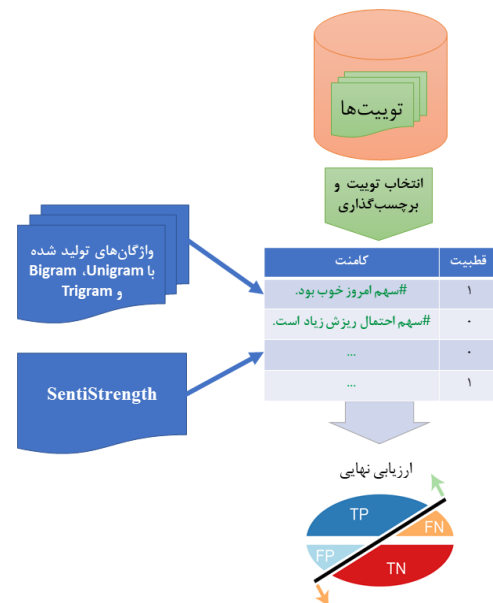
(جدول-۷): قسمتی از واژگان به دست آمده با Bigram
(Table-7): Part of the lexicon obtained with Bigram feature

لغت	رخداد در واژگان	میانگین نرخ رشد (%)
صف فروش	29904	-1.374
صف خرید	40820	0.588
با ضرر	2047	-0.866
خرید خوشبین	2159	0.508

یکی از مزایای استفاده از ویژگی Bigram برای انتخاب کلمات واژگان این است که احتمال این که دو کلمه متوالی که احساس منفی یا مثبتی را انعکاس می دهند، هم در شرایط منفی بازار و هم در شرایط مثبت بازار رؤیت شوند، بسیار کمتر از کلمات واحد (که توسط Unigram انتخاب شده اند)، است. به عبارت دیگر، می توان گفت که توزیع آماری کلماتی که با استفاده از Bigram استخراج می شوند، بیشتر به سمت مثبت یا منفی تجمع می یابند. ولی توزیع کلماتی که با استفاده از ویژگی Unigram استخراج می شوند، بیشتر نزدیک به صفر است. به همین دلیل استفاده از واژگان استخراج شده با ویژگی Bigram، نسبت به واژگان استخراج شده با Unigram، کیفیت تحلیل احساسات را روی کار ما بهبود بخشید. در شکل (۳) توزیع آماری کلمه «صف فروش» برای واژگان به دست آمده با Bigram نمایش داده شده است. نتایج استفاده از این واژگان در کار تحلیل احساسات بر روی توییت های بورسی، در بخش (۴) آورده شده است.

۴- ارزیابی و آزمایش‌ها

به‌طور کلی، موضوعات تحلیل زبان طبیعی برای هر زبان مستقلند و قابل قیاس با زبان‌های دیگر نیستند. همچنین، در بازار بورس ایران اصطلاحات یا کلمات خاصی به کار می‌رود که نمونه مشابه آن را نمی‌توان در بازارهای مالی دیگر یا در زبان‌های دیگر مشاهده کرد. به همین دلیل ما برای ارزیابی یک‌سری آزمایش‌ها و همچنین، یک مقایسه، بر روی واژگان تولیدشده به روش پیشنهادی و واژگان عمومی که در کار آقای بصیری و همکارانش [32] استفاده شده‌است، انجام دادیم. شکل (۵) روند ارزیابی را نمایش می‌دهد. به‌منظور مقایسه تعداد ۱۰۰۰ توییت به‌صورت تصادفی از دیتابیس خود انتخاب و قطبیت هر یک از آن‌ها را به‌صورت دستی مشخص کردیم. به‌منظور انتخاب قطبیت ما از برچسب‌گذاری دودویی (مثبت و منفی) استفاده کردیم. جدول (۹)، مشخصات مجموعه‌دادگان ما را نشان می‌دهد.



(شکل-۵): مراحل ارزیابی
(Figure-5): Evaluation steps

۴-۱- واژگان SentiStrength:

یک کتابخانه در دسترس است که برای تشخیص قطبیت و قدرت متون اجتماعی کوتاه و غیررسمی استفاده می‌شود [42]. آقای بصیری و همکارانش [32]، از این کتابخانه برای تحلیل احساسات روی کار خود استفاده کردند. از آنجاکه این نرم‌افزار برای زبان انگلیسی طراحی و ایجاد شده‌است، بصیری و همکارانش نخست، فهرست اصلی واژه‌های آن را به‌صورت دستی به زبان فارسی ترجمه، و بعد از حذف کلمات تکراری در کار خود استفاده کردند.

(جدول-۹): مشخصات مجموعه‌دادگان
(Table-9): Dataset specifications

تعداد	برچسب
668	مثبت
332	منفی

۴-۲- تجمیع^۱:

به‌منظور محاسبه قطبیت هر توییت در مجموعه‌دادگان، ما از روش‌های آماری در کار خود استفاده و برای تعیین قطبیت هر توییت از دو روش معروف تجمیع استفاده کردیم؛ روش جمع امتیازها و روش دمپستر شافر.

۴-۱-۱- جمع امتیازها:

این روش امتیازهای مثبت و منفی هر جمله را با هم جمع و عدد حاصل، قطبیت جمله را مشخص می‌کند. محاسبه قطبیت هر توییت به این صورت است که هر کلمه از توییت یک امتیاز در واژگان دارد. به‌عنوان مثال، در واژگان تولیدشده با Unigram، برای هر کلمه یک میانگین نرخ رشد در نظر گرفته شده که همان امتیاز قطبیت کلمه است. بنابراین برای هر توییت، امتیاز هر یک از کلمات آن را از واژگان موردنظر محاسبه و با یکدیگر جمع می‌کنیم و عدد حاصل را به‌عنوان قطبیت نهایی آن توییت در نظر می‌گیریم. از آنجاکه حد آستانه را در این کار صفر در نظر گرفتیم، اگر قطبیت به‌دست آمده بیشتر از صفر شود، آن توییت از نظر احساسی مثبت، و اگر کمتر از صفر شود آن توییت از نظر احساسی منفی در نظر گرفته می‌شود. همچنین، به‌منظور مقایسه با کار آقای بصیری ما از روش آن‌ها (دمپستر شافر) نیز استفاده و نتایج را با هم مقایسه کردیم.

۴-۲-۲- دمپستر شافر:

این راهبرد مبتنی بر نظریه شواهد دمپستر شافر است [43]. این روش برای اولین بار در کار تحلیل احساسات به زبان فارسی توسط بصیری و همکارانش [32]، پیشنهاد شد. آن‌ها در کار خود نمره هر جمله از یک بررسی را به‌عنوان یک رویداد برای نمره کلی آن بررسی در نظر گرفتند. همچنین، بعد از تعریف رویداد در کار خود برای محاسبه تابع جرم^۲، از معادله (۳) استفاده کردند.

¹ Aggregation

² Mass Function

$$(3) \quad \frac{\text{کمترین نمره} - \text{نمره جمله}}{\text{کمترین نمره} - \text{بیشترین نمره}} = (\text{جمله}) \text{ تابع جرم}$$

(جدول ۱۱): نتایج ارزیابی واژگان‌های به‌دست‌آمده با روش پیشنهادی (۲) - (نسخه دوم) و SentiStrength با استفاده از

جمع نمرات

(Table-11): Results of the evaluation of the lexicons based on the proposed method 2 and SentiStrength using sum of scores

واژگان	صحت	دقت	بازیابی	F-Score
واژگان استخراج شده با Unigram	0.66	0.9	0.41	0.56
واژگان استخراج شده با Bigram	0.66	0.85	0.51	0.64
واژگان استخراج شده با Trigram	0.66	0.83	0.55	0.66
واژگان SentiStrength	0.46	0.52	0.08	0.13

همان‌طور که در جدول (۱۱) مشاهده می‌شود، کار تشخیص قطبیت با استفاده از واژگان SentiStrength، که با هدف استفاده عمومی ایجاد شده‌است، در داده‌های بورسی کمترین دقت را دارد. دلیل این دقت پایین این است که کلماتی که در داده‌های حوزه بورس وجود دارد، مختص بورس است و در حوزه‌های دیگر کاربرد زیادی ندارد. از این رو واژگان SentiStrength در میان واژگان‌های دیگر، کمترین دقت را به‌دست‌آورد. همچنین بیشترین رتبه F-Score نیز مربوط به واژگان به‌دست‌آمده با ویژگی Trigram است. دلیل برتری واژگان به‌دست‌آمده با Trigram به نسبت واژگان به‌دست‌آمده با Unigram و Bigram، همان‌طور که در بخش ۳-۴-۳ نیز ذکر شد، همبستگی بیشتر کلمات به‌دست‌آمده با Trigram در کل مجموعه‌دادگان نسبت به اطلاعات قیمتی است. در واقع کلمات به‌دست‌آمده با این ویژگی، امکان این‌که هم در روزهای مثبت و هم در روزهای منفی بازار رؤیت شوند، به‌مراتب کمتر از کلمات به‌دست‌آمده با Unigram و Bigram دارند.

جدول (۱۲) نتایج ارزیابی با استفاده از واژگان‌های پیشنهادی و واژگان SentiStrength را با کمک راهبرد دمپستر شافر، نمایش می‌دهد.

همان‌طور که در جدول (۱۲) نیز مشاهده می‌شود، تشخیص قطبیت با واژگان SentiStrength و جمع نمرات با استفاده از نظریه دمپستر شافر کمترین صحت را دارد و بیشترین صحت را نیز واژگان به‌دست‌آمده با ویژگی Bigram دارد. البته گفتنی است که استفاده از این روش برای جمع نمرات، نسبت به جمع نمرات با استفاده از جمع ساده، باعث بهبود تشخیص قطبیت با استفاده از واژگان SentiStrength شد. اما همان‌طور که مشاهده

در معادله ۳ «نمره جمله»، خروجی واژگان SentiStrength برای آن جمله است. «کمترین نمره» و «بیشترین نمره» به ترتیب کوچک‌ترین نمره منفی و بزرگ‌ترین نمره مثبت در جملات منفی و مثبت است. و در نهایت، برای محاسبه نمره نهایی نمرات رویدادهای به‌دست‌آمده برای هر بررسی را از طریق معادله (۴) جمع می‌کنند.

$$(4) \quad m(A) = \frac{\sum_{X \cap Y=A} m_n(X)m_o(Y)}{1 - \sum_{X \cap Y=\emptyset} m_n(X)m_o(Y)}$$

مقدار جمع کلی یک عدد حقیقی در بازه [۰، ۱] است و گرد می‌شود تا به‌عنوان یک مقدار دودویی (مثبت/منفی) استفاده شود. ما نیز به‌جهت مقایسه با کار آن‌ها، از این روش برای جمع امتیازات استفاده و نتایج را با هم مقایسه کردیم.

۳-۴- ارزیابی نهایی:

برای ارزیابی از چهار معیار ارزیابی معروف استفاده شد. این معیارها شامل صحت، دقت، بازیابی و F-Score است. جدول (۱۰)، نتایج ارزیابی با استفاده از واژگان‌های به‌دست‌آمده با روش پیشنهادی (۱) - (نسخه اول) و همچنین جدول (۱۱)، نتایج ارزیابی با استفاده از واژگان‌های به‌دست‌آمده با روش پیشنهادی (۲) - (نسخه دوم) و واژگان SentiStrength را با کمک جمع امتیازها، نمایش می‌دهد.

(جدول ۱۰): نتایج ارزیابی واژگان‌های به‌دست‌آمده با روش

پیشنهادی (۱) - (نسخه اول) با استفاده از جمع نمرات

(Table-10): Results of the evaluation of the lexicons based on the proposed method 1 and SentiStrength using sum of scores

واژگان	صحت	دقت	بازیابی	F-Score
واژگان استخراج شده با Unigram	0.66	0.62	0.64	0.63
واژگان استخراج شده با Bigram	0.64	0.67	0.57	0.62
واژگان استخراج شده با Trigram	0.61	0.74	0.36	0.48

می‌شود، تجمیع نمرات با استفاده از تئوری دمپستر شافر به نسبت تجمیع نمرات با استفاده از روش جمع ساده، دقت تشخیص احساسات را برای کار ما علاوه بر اینکه افزایش نمی‌دهد، بدر حدود ۷٪ کاهش نیز می‌دهد. از این رو نتیجه می‌گیریم که استفاده از تئوری دمپستر شافر و همچنین، استفاده از واژگان‌های ساخته شده با اهداف استفاده عمومی، علاوه بر اینکه ممکن است در بعضی از حوزه‌ها کاربرد خوبی نداشته باشند، ممکن است نتایج بسیار ضعیفی ایجاد نمایند.

جدول-۱۲): ارزیابی واژگان‌های پیشنهادی و SentiStrength با استفاده از دمپستر شافر

(جدول-۱۲): ارزیابی واژگان‌های پیشنهادی و SentiStrength با استفاده از دمپستر شافر

(Table-12): Result of the evaluation of the proposed lexicons and SentiStrength using Dempster-Shafer theory

واژگان	صحت	دقت	بازیابی	F-Score
واژگان استخراج شده با Unigram	0.49	0.61	0.05	0.1
واژگان استخراج شده با Bigram	0.59	0.77	0.47	0.58
واژگان استخراج شده با Trigram	0.59	0.82	0.31	0.45
واژگان SentiStrength	0.51	0.68	0.26	0.38

۵- نتیجه‌گیری و پیشنهاد

هدف از تحلیل احساسات خودکار عقاید افراد نسبت به موضوعات مختلف در سطح شبکه است. شبکه‌های اجتماعی محیطی است که هم‌روزه افراد به بحث و تبادل نظر می‌پردازند. یکی از این شبکه‌ها که به‌تازگی، بسیار به آن توجه شده، شبکه‌های اجتماعی مرتبط با بازارهای مالی و یکی از آن‌ها بازار بورس است که افراد خرید و فروش سهم را در آن انجام می‌دهند. در این شبکه‌ها افراد عقاید خود را نسبت به سهم‌ها بیان می‌کنند. ما در کار خود عقیده‌کاوی را بر روی توییت‌هایی که هم‌روزه در بازار بورس منتشر می‌شوند، انجام دادیم. ویژگی منحصر به فرد این شبکه‌ها، وجود اطلاعات قیمتی هر سهم در هر روز است. ما از این موضوع جهت استخراج واژگانی به صورت خودکار از بین ۹۰۰۰۰۰ توییت که در این شبکه‌ها موجود است، استفاده کردیم. در واقع، از این ویژگی‌ها برای محاسبه نرخ رشد هر سهم در زمان درج نظر استفاده و امتیازات کلمات به دست آمده را برای واژگان به کمک این نرخ رشد‌ها محاسبه کردیم. برای ارزیابی روش پیشنهادی برای تولید واژگان، واژگان‌های تولید شده را با نسخه فارسی واژگان SentiStrength که آخرین

6- Refrence

۶- مراجع

- [1] B. Pang and L. Lee, "Opinion mining and sentiment analysis," *Foundations and Trends® in Information Retrieval*, vol. 2, no. 1-2, pp. 1-135, 2008.
- [2] B. Liu and L. Zhang, "A survey of opinion mining and sentiment analysis," in *Mining text data*: Springer, 2012, pp. 415-463.
- [3] I. Chaturvedi, E. Cambria, R. E. Welsch, and F. Herrera, "Distinguishing between facts and opinions for sentiment analysis: Survey and challenges," *Information Fusion*, vol. 44, pp. 65-77, 2018.
- [4] M. Taboada, J. Brooke, M. Tofiloski, K. Voll, and M. Stede, "Lexicon-based methods for sentiment analysis," *Computational linguistics*, vol. 37, no. 2, pp. 267-307, 2011.
- [5] G. A. Miller, *WordNet: An electronic lexical database*. MIT press, 1998.
- [6] B. Liu, "Sentiment analysis and opinion mining," *Synthesis lectures on human language technologies*, vol. 5, no. 1, pp. 1-167, 2012.

- [21] M. S. Hajmohammadi and R. Ibrahim, "A SVM-based method for sentiment analysis in Persian language," International Conference on Graphic and Image Processing (ICGIP 2012), vol. 8768, p. 876838, 2013.
- [22] M. Saraei and A. Bagheri, "Feature selection methods in Persian sentiment analysis," International Conference on Application of Natural Language to Information Systems, pp. 303-308, 2013.
- [23] A. S. H. Basari, B. Hussin, I. G. P. Ananta, and J. Zeniarja, "Opinion mining of movie review using hybrid method of support vector machine and particle swarm optimization," Procedia Engineering, vol. 53, no. 7, pp. 453-462, 2013.
- [24] T. S. Ataei, K. Darvishi, S. Javdan, B. Minaei-Bidgoli, and S. Eetemadi, "Pars-ABSA: an Aspect-based Sentiment Analysis dataset for Persian," arXiv preprint arXiv:1908.01815, 2019.
- [25] K. Dashtipour, M. Gogate, A. Adeel, C. Ieracitano, H. Larijani, and A. Hussain, "Exploiting deep learning for persian sentiment analysis," International Conference on Brain Inspired Cognitive Systems, pp. 597-604, 2018.
- [26] E. Cambria, D. Olsher, and D. Rajagopal, "SenticNet 3: a common and common-sense knowledge base for cognition-driven sentiment analysis," in Proceedings of the AAAI Conference on Artificial Intelligence, 2014, vol. 28, no. 1.
- [27] M. Hu and B. Liu, "Mining and summarizing customer reviews," in Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining, 2004, pp. 168-177.
- [28] L. Deng and J. Wiebe, "Mpqa 3.0: An entity/event-level sentiment corpus," in Proceedings of the 2015 conference of the North American chapter of the association for computational linguistics: human language technologies, 2015, pp. 1323-1328.
- [29] A. Neviarouskaya, H. Prendinger, and M. Ishizuka, "SentiFul: A lexicon for sentiment analysis," IEEE Transactions on Affective Computing, vol. 2, no. 1, pp. 22-36, 2011.
- [30] A. Esuli and F. Sebastiani, "Sentiwordnet: A publicly available lexical resource for opinion mining," in LREC, 2006, vol. 6: Citeseer, pp. 417-422.
- [31] S. M. Mohammad and P. D. Turney, "Crowdsourcing a word-emotion association lexicon," Computational intelligence, vol. 29, no. 3, pp. 436-465, 2013.
- [32] M. E. Basiri, A. R. Naghsh-Nilchi, and N. Ghassem-Aghaee, "A framework for sentiment analysis in persian," Open transactions on information processing, vol. 1, no. 3, pp. 1-14, 2014.
- [33] M. E. Basiri and A. Kabiri, "Sentence-level sentiment analysis in Persian," in 2017 3rd International Conference on Pattern Recognition
- [7] S. Li, "Sentiment classification using subjective and objective views," International Journal of Computer Applications, vol. 80, no. 7, 2013.
- [8] B. Pang, L. Lee, and S. Vaithyanathan, "Thumbs up? Sentiment classification using machine learning techniques," arXiv preprint cs/0205070, 2002.
- [9] Y. Jo and A. H. Oh, "Aspect and sentiment unification model for online review analysis," in Proceedings of the fourth ACM international conference on Web search and data mining, 2011, pp. 815-824.
- [10] D. Maynard and A. Funk, "Automatic detection of political opinions in tweets," in Extended Semantic Web Conference, 2011: Springer, pp. 88-99.
- [11] K. Dashtipour, A. Hussain, Q. Zhou, A. Gelbukh, A. Y. Hawalah, and E. Cambria, "PerSent: A freely available Persian sentiment lexicon," in International Conference on Brain Inspired Cognitive Systems, 2016: Springer, pp. 310-320.
- [12] P. D. Turney, "Thumbs up or thumbs down? Semantic orientation applied to unsupervised classification of reviews," arXiv preprint cs/0212032, 2002.
- [13] S. Rani and P. Kumar, "Deep learning based sentiment analysis using convolution neural network," Arabian Journal for Science and Engineering, vol. 44, no. 4, pp. 3305-3314, 2019.
- [14] G. Xu, Y. Meng, X. Qiu, Z. Yu, and X. Wu, "Sentiment analysis of comment texts based on BiLSTM," Ieee Access, vol. 7, pp. 51522-51532, 2019.
- [15] M. Rhanoui, M. Mikram, S. Yousfi, and S. Barzali, "A CNN-BiLSTM model for document-level sentiment analysis," Machine Learning and Knowledge Extraction, vol. 1, no. 3, pp. 832-847, 2019.
- [16] M. Ahmad, S. Aftab, and I. Ali, "Sentiment analysis of tweets using svm," Int. J. Comput. Appl, vol. 177, no. 5, pp. 25-29, 2017.
- [17] M. Ahmad, S. Aftab, M. S. Bashir, N. Hameed, I. Ali, and Z. Nawaz, "SVM optimization for sentiment analysis," Int. J. Adv. Comput. Sci. Appl, vol. 9, no. 4, pp. 393-398, 2018.
- [18] K. Korovkinas, P. Danėnas, and G. Garšva, "SVM and k-means hybrid method for textual data sentiment analysis," Baltic Journal of Modern Computing, vol. 7, no. 1, pp. 47-60, 2019.
- [19] L. Dey, S. Chakraborty, A. Biswas, B. Bose, and S. Tiwari, "Sentiment analysis of review datasets using naive bayes and k-nn classifier," arXiv preprint arXiv:1610.09982, 2016.
- [20] V. Narayanan, I. Arora, and A. Bhatia, "Fast and accurate sentiment classification using an enhanced Naive Bayes model," in International Conference on Intelligent Data Engineering and Automated Learning, 2013: Springer, pp. 194-201.

کامپیوتر دانشگاه گلستان است. زمینه‌های پژوهشی موردعلاقه ایشان پردازش زبان طبیعی، پردازش تصویر و الگوریتم‌های تکاملی است.

نشانی رایانامه ایشان عبارت است از:

a.sebti@gu.ac.ir



مرتضی آهنگری آهنگرکلایی

کارشناس ارشد مهندسی نرم‌افزار از دانشگاه گلستان است. زمینه‌های پژوهشی وی تحلیل احساسات و پردازش زبان طبیعی است.

نشانی رایانامه ایشان عبارت است از:

m.ahangari98@stu.gu.ac.ir



مهدی یعقوبی مدرک کارشناسی

خود را در رشته مهندسی کامپیوتر گرایش نرم‌افزار در سال ۱۳۸۰ از دانشگاه صنعت نفت تهران و مدرک کارشناسی ارشد خود را در

سال ۱۳۸۳ در رشته مهندسی کامپیوتر گرایش نرم‌افزار از دانشگاه صنعتی امیرکبیر دریافت کرد و دکترای خود را در رشته‌ی مهندسی کامپیوتر گرایش هوش مصنوعی در دانشگاه صنعتی شاهرود به پایان رساند. ایشان در حال حاضر استادیار گروه مهندسی کامپیوتر دانشگاه گلستان است. زمینه‌های پژوهشی موردعلاقه ایشان داده‌کاوی، فرآیندکاوی و سامانه‌های اطلاعاتی مبتنی بر فرآیند است.

نشانی رایانامه ایشان عبارت است از:

m.yaghoubi@gu.ac.ir

and Image Analysis (IPRIA), 2017: IEEE, pp. 84-89.

- [34] S. M. Mohammad, S. Kiritchenko, and X. Zhu, "NRC-Canada: Building the state-of-the-art in sentiment analysis of tweets," arXiv preprint arXiv:1308.6242, 2013.
- [35] M. Thelwall, K. Buckley, and G. Paltoglou, "Sentiment strength detection for the social web," *Journal of the American Society for Information Science and Technology*, vol. 63, no. 1, pp. 163-173, 2012.
- [36] M. E. Basiri and A. Kabiri, "Translation is not enough: comparing lexicon-based methods for sentiment analysis in Persian," in 2017 International Symposium on Computer Science and Software Engineering Conference (CSSE), 2017: IEEE, pp. 36-41.
- [37] B. Sabeti, P. Hosseini, G. Ghassem-Sani, and S. A. Mirroshandel, "LexiPers: An ontology based sentiment lexicon for Persian," arXiv preprint arXiv:1911.05263, 2019.
- [38] F. Amiri, S. Scerri, and M. Khodashahi, "Lexicon-based sentiment analysis for Persian text," in *Proceedings of the International Conference Recent Advances in Natural Language Processing*, 2015, pp. 9-16.
- [39] R. Dehkharghani, "Sentifars: A persian polarity lexicon for sentiment analysis," *ACM Transactions on Asian and Low-Resource Language Information Processing (TALLIP)*, vol. 19, no. 2, pp. 1-12, 2019.
- [40] E. Haddi, X. Liu, and Y. Shi, "The role of text pre-processing in sentiment analysis," *Procedia Computer Science*, vol. 17, pp. 26-32, 2013.
- [41] A. AleAhmad, H. Amiri, E. Darrudi, M. Rahgozar, and F. Oroumchian, "Hamshahri: A standard Persian text collection," *Knowledge-Based Systems*, vol. 22, no. 5, pp. 382-387, 2009.
- [42] M. Thelwall, K. Buckley, G. Paltoglou, D. Cai, and A. Kappas, "Sentiment strength detection in short informal text," *Journal of the American society for information science and technology*, vol. 61, no. 12, pp. 2544-2558, 2010.
- [43] G. Shafer, *A mathematical theory of evidence*. Princeton university press, 1976.



علی سبیطی مدرک کارشناسی خود

را در رشته مهندسی کامپیوتر گرایش نرم‌افزار در سال ۱۳۸۵ از دانشگاه یزد و مدرک کارشناسی ارشد خود را در سال ۱۳۸۸ در رشته مهندسی کامپیوتر گرایش هوش

مصنوعی از دانشگاه صنعتی امیرکبیر اخذ کرد و دکترای خود را در رشته‌ی مهندسی کامپیوتر گرایش هوش مصنوعی در دانشگاه صنعتی شاهرود در سال ۱۳۹۶ به پایان رساند. ایشان در حال حاضر استادیار گروه مهندسی

