

ارائه‌ی مدلی برای عقیده‌کاوی در سطح ویژگی سند برای نظرات کاربران هتل‌ها

شهریار محمدی* الهام خلیج**

*دانشیار، دانشکده مهندسی صنایع، گروه مهندسی فناوری اطلاعات، دانشگاه صنعتی خواجه نصیرالدین طوسی
** کارشناس ارشد دانشکده مهندسی صنایع، گروه مهندسی فناوری اطلاعات، دانشگاه صنعتی خواجه نصیرالدین طوسی

تاریخ پذیرش: ۱۴۰۰/۰۳/۲۹

تاریخ دریافت: ۱۳۹۹/۱۱/۰۴

نوع مقاله: پژوهشی

چکیده

امروزه بررسی نظرات و عقاید کاربران در بستر اینترنت بخش مهمی از فرآیند تصمیم‌گیری مردم در رابطه با انتخاب محصول یا استفاده از خدمات را شامل می‌شود. با وجود اینترنت و دسترسی ساده به وبلاگ‌های مربوط به نظرات در زمینه صنعت گردشگری و هتلداری، منابع غنی و عظیمی از عقاید بصورت متن موجود می‌باشد که می‌توان از روش‌های متن کاوی برای کشف دانش نهفته در این متون استفاده کرد. با توجه به اهمیت نظرات و عقاید کاربران در صنایع، به‌ویژه صنعت گردشگری و هتلداری، مباحث عقیده‌کاوی و تحلیل احساسات مورد توجه متصدیان امور قرار گرفته است. در این مقاله یک روش ترکیبی و جدید بر اساس یک رویکرد رایج در تحلیل احساسات، استفاده از واژگان و الگوریتم ژنتیک برای تولید ویژگی‌هایی برای طبقه‌بندی بار احساسی نظرات ارائه شده است. بدین صورت که دو روش ساخت فهرست واژگان یکی با استفاده از روش‌های آماری و دیگری با استفاده از الگوریتم ژنتیک ارائه شده است. واژگان فوق‌الذکر با فرهنگ واژگان احساس عمومی و استاندارد لیو بینگ آمیخته می‌شوند. نتایج نشان می‌دهد روش پیشنهادی از روش‌های پایه براساس واژه‌نامه‌های احساسی روی این مجموعه داده بهتر عمل کرده و معیارهای ارزیابی صحت، دقت، بازخوانی و معیار F با استفاده از روش پیشنهادی، به ترتیب ۹۴٫۶۵، ۹۴٫۵۳، ۹۴٫۸۹ و ۹۳٫۱۷ می‌باشند.

واژه‌های کلیدی: متن کاوی، عقیده‌کاوی، تحلیل احساسات در سطح ویژگی، داده کاوی، الگوریتم ژنتیک، طبقه بندی.

۱- مقدمه

تحلیل احساس و عقیده کاوی^۱ با استفاده از الگوریتم‌های داده کاوی و متن کاوی به صورت سیستماتیک و بدون نیاز به مطالعه تمامی متن‌های موجود، عقاید، احساسات، ارزیابی‌ها، رفتارها و گرایش‌های کاربران را که به صورت داده‌های متن بیان شده‌اند را آنالیز می‌کند. افزایش اهمیت تحلیل احساس با رشد رسانه‌های اجتماعی مانند توئیتر،

تجزیه و تحلیل احساسات^۱ با هدف کشف خودکار نگرش اساسی انسان‌ها نسبت به یک موجودیت انجام می‌شود. در حال حاضر، تجزیه و تحلیل احساسات از داده‌های متنی به طور گسترده‌ای برای ارزیابی رضایت مشتری و تجزیه و تحلیل‌ها استفاده می‌شود [۱] و [۲].

^۱ Sentiment Analysis

نویسنده مسئول: شهریار محمدی mohammadi@kntu.ac.ir

^۲ Opinion Mining

استفاده شده و همچنین با توجه به میزان سودآوری صنعت گردشگری و تاثیر مستقیم تجربه و نظر سایر گردشگران برای انتخاب دوباره خدمات ارائه شده و یا پیشنهاد انتخاب به سایرین، لزوم پرداختن بیشتر به این موضوع احساس می‌شود. بنابراین در این مقاله صنعت گردشگری و هتلداری مورد بررسی قرار گرفته است که در مقاله‌های قبلی کمتر مورد توجه بوده همچنین در این مقاله علاوه بر تعیین مثبت و یا منفی بودن نظرات در داده‌های متنی ویژگی‌های خاص مربوط به هتل مورد توجه قرار گرفته، در این مقاله از روش پیشنهادی جدید ترکیب الگوریتم ژنتیک و فرهنگ واژگان احساسی استفاده شده است. تحلیل احساسات مشتریان در محیط‌هایی مثل خدمات گردشگری و هتلداری، بیمه، موسسات مالی و بانک‌ها، خرده‌فروشی‌ها، شرکت‌های تجارت الکترونیک و فروش آنلاین و... می‌تواند بسیار کاربردی باشد [۵].

۱-۱ هدف از تحقیق

این پژوهش به دنبال راه حلی برای بهبود تحلیل حجم عظیمی از نظرات متنی می‌باشد که معمولا ساختار نیافته یا نیمه‌ساختار هستند و برای انجام این کار سعی می‌شود طبقه‌بندی رتبه‌ی احساسی واژگان دقیق‌تر از روش‌های قبلی محاسبه گردد. سؤالاتی که در این پژوهش به آن‌ها پرداخته می‌شود به شرح زیر هستند:

- ✓ تأثیر بهره‌گیری از الگوریتم ژنتیک بر بهبود پارامترهای ارزیابی طبقه‌بندی واژگان احساسی به چه صورت است؟
- ✓ آیا طبقه‌بندی واژگان احساسی با روش ترکیبی بیان شده که ویژگی‌های خاصی را نیز در نظر می‌گیرد، نتایج بهتری از روش‌های پایه به همراه دارد؟

شبکه‌های اجتماعی، نظرسنجی‌های آنلاین، وبلاگ‌ها و همچنین سهولت بازیابی آنلاین نظرات کاربران همزمان شده است. در پژوهش پیرمحمدیانی و محمدی بیان شده است که امروزه سیستم‌های تحلیل احساس تقریبا در همه ی زمینه‌ها مورد استفاده قرار می‌گیرند، زیرا آرا و عقاید در تمام فعالیت‌های انسانی مهم بوده و تاثیر کلیدی بر فرآیندهای تصمیم‌گیری دارند [۳]. در مطالعه‌ی ژانگ و همکاران^۳ با اشاره به تاثیر توسعه‌ی سریع فناوری‌های وب ۲،۰ و محتوای تولید شده توسط کاربر^۴، به تحلیل و بررسی نظرات آنلاین درباره‌ی سفر در صنعت گردشگری پرداخته شده است. همچنین وبسایت‌هایی نظیر^۵ TripAdvisor، Expedia^۶ که گردشگران نظرات، عقاید و تجربیات خود در استفاده از خدمات ارائه شده را به اشتراک می‌گذارند، معرفی شده است، این نوع وبسایت‌ها محبوب‌ترین منابع اطلاعاتی برای کسب اطلاعات در راستای تهیه برنامه‌ی سفر و نحوه‌ی رزرو بلیط و هتل هستند [۴]. براساس پژوهش پابلوس و همکاران^۷ فن‌آوری تجزیه و تحلیل متن مبتنی بر پردازش زبان طبیعی^۸ می‌تواند بطور خودکار مقادیر زیادی از بررسی‌ها و نظرات مشتری را از منظر مفهوم کلمه، به کارگیرد. این روش به طور گسترده‌ای در شناسایی موضوع و استخراج مفهوم نظر و متن، مورد استفاده قرار می‌گیرد [۵].

در اکثر مقالات مطالعه شده، پایه و اساس تحلیل احساسات مبتنی بر شمارش کلمات احساسی و تعیین بار مثبت و یا منفی کل متن است [۲، ۶، ۷] و [۸، ۹] نقطه ضعف این روش آن است که ممکن است در یک متن چند ویژگی بیان شده باشد و محاسبه‌ی بار معنایی کلی نتایج دقیقی به همراه نداشته باشد. به عنوان مثال جمله‌ی "دسترسی به مراکز خرید از هتل خوب است اما غذاها کیفیت مناسبی ندارند." در این نظر دو ویژگی مکان هتل و غذا مورد نظر بوده پس طبقه‌بندی کلی نظر مناسب نیست و توجه به هر دو ویژگی مناسب‌تر است.

اخیرا این موضوع مورد توجه دانشمندان فناوری اطلاعات قرار گرفته و روش‌های مختلفی را مورد بررسی قرار داده‌اند [۸، ۱۰، ۱۱، ۱۲] اما کمتر از الگوریتم‌های فراابتکاری

^۶ expediagroup.com

^۷ Pablos et al.

^۸ Natural language processing (NLP)

^۳Zhang et al.

^۴ user generated content

^۵ www. TripAdvisor.com

۲- مروری بر تحقیقات پیشین

عقیده‌کاوی یک فن‌آوری است که به طور خودکار با استفاده از ابزار و نرم افزارهای تجزیه و تحلیل متن، از جمله زبان‌رایانه و پردازش زبان طبیعی، دانسته‌های کامنت‌های آنلاین را استخراج می‌کند. این نرم‌افزارها نظرها، ارزیابی‌ها، نگرش‌ها و احساس مردم را نسبت به سازمان‌ها، اشخاص، افراد، موضوع‌ها، اقدام‌ها و ویژگی‌های آن‌ها را تجزیه و تحلیل می‌کند [۱، ۲، ۱۳].

با توجه به گزارش^۹ UNEP در سال ۲۰۱۶، برای هر دو کشور پذیرنده گردشگر و فرستنده گردشگر، صنعت-گردشگری کمک چشمگیری به تولید منافع اقتصادی آنان می‌کند. توسعه‌ی صنعت گردشگری تا حد زیادی به رضایت مشتریان در خدمات هتل وابسته بوده و در نتیجه تخصص در این زمینه یکی از عوامل مهم برای توسعه‌ی این صنعت است. مطالعه‌ی ژانگ و همکاران^{۱۰} نشان می‌دهد اثربخشی کیفیت و خدمات ارائه شده در هتل‌ها در تقویت صنعت گردشگری سهم بسزایی دارد [۶]. طبق پژوهش بالاز و همکاران^{۱۱}، نظرات نقش اساسی را در فرآیند تصمیم‌گیری افراد و سازمان‌ها دارند زیرا تأثیر عمیقی بر روی نگرش و اعتقادات افراد می‌گذارد. عقیده‌کاوی و تحلیل-احساسات باعث می‌شود تا مشاغل تجارت الکترونیکی بتوانند دانش بیشتری از مشتریان و محصولات خود کسب کنند بدون اینکه هزینه نظرسنجی‌ها را پردازند [۷].

کلیه تکنیک‌های مورد استفاده برای استخراج نظر و عقیده می‌توانند به دو طبقه اصلی تقسیم شوند:

- ✓ رویکردهای مبتنی بر واژه‌ها: این روش واژه‌های دارای بار احساسی متن را با تکیه بر یک فرهنگ واژگان احساسی و رویکرد دانش‌زبانی، طبقه‌بندی می‌کند که شامل یک رویکرد مبتنی بر بدنه و یک رویکرد مبتنی بر فرهنگ لغت است.
- ✓ رویکرد یادگیری ماشین: از الگوریتم‌های یادگیری ماشین بهره‌می‌برد و می‌تواند به سه گروه تقسیم شود: یادگیری نظارت شده، یادگیری نیمه نظارت شده و یادگیری بدون نظارت [۱، ۲، ۱۳].

در مقاله‌ی سینگ و همکاران^{۱۲} از این رویکرد برای کشف دانسته‌های نظرات متنی مربوط به محصولات مختلف با استفاده از یک روش طبقه‌بندی برای تجزیه و تحلیل، استفاده می‌شود [۱۵].

۲-۱ سطوح عقیده‌کاوی

- استخراج افکار در سطح سند، احساسات غالب و کلی را به‌جای موارد و جزئیات مطرح در مباحث در نظر می‌گیرد [۱۳]. وظیفه‌ی عقیده‌کاوی در سطح سند تعیین قطبیت کلی یک سند است که شامل چند جمله است.
- عقیده‌کاوی در سطح جمله، بطور ویژه متمرکز بر هر جمله است. خواه جمله بیان‌شده گرایش مثبت، منفی یا خنثی داشته باشد. طبقه‌بندی ذهنی یکی دیگر از وظایف در سطح جمله است که بخش‌های ذهنی و عینی اسناد را استخراج می‌کند. مسئله تجزیه و تحلیل مبتنی بر سطح جمله نیز به همین صورت تعریف می‌شود با این تفاوت که نتیجه تشخیص احساس برای هر جمله به‌طور جداگانه بررسی می‌شود [۱۶].
- استخراج ویژگی‌های ارائه‌شده در یک متن یا نظر و بیان گرایش احساس مثبت و یا منفی بر روی آن‌ها، تجزیه و تحلیل احساسات و یا عقیده‌کاوی در سطح ویژگی سند نامیده می‌شود [۹].

۲-۲ رویکردهای متفاوت برای عقیده‌کاوی و

تحلیل احساسات

رویکردهای مبتنی بر واژگان و نظارت نشده^{۱۳}: این رویکرد، در تعیین گرایش مثبت یا منفی متن با استفاده از مجموعه قوانین و اکتشاف‌های حاصل از دانش و قواعد زبان شناسی مورد استفاده است. اقدامات معمول برای اجرای مرحله اول، علامت‌گذاری هر کلمه و تعیین گرایش احساسی مربوط به آن با کمک یک فرهنگ واژگان احساسی و در مرحله دوم،

^{۱۲} Singh et al.

^{۱۳} Unsupervised Lexicon-based Approaches

^۹ United Nations Environment Program

^{۱۰} Zhang et al.

^{۱۱} Balazs et al.

رویکردهای مبتنی بر مفهوم: این رویکرد شامل استفاده از علم هستی‌شناسی برای پشتیبانی از عقیده‌کاوی و تحلیل احساسات است. هستی‌شناسی به عنوان مدلی تعریف می‌شود که دانش یک حوزه‌ی معین را برای کامپیوتر با دستورهایی اگر و آنگاه مفهوم‌سازی می‌کند. معمولاً به صورت نمودارهایی ارائه می‌شود که در آن مفاهیم مدنظر، به گره‌های مرتبط با هم و بصورت متصل کشیده می‌شوند.

ترکیب کلمات و تحلیل احساسی و تاثیر ترکیب کلمات و در آخر، بررسی این که ترکیب‌ها ۱۴ چگونه بر قطبیت و گرایش تأثیر می‌گذارند و این را در نمره احساسات نهایی منعکس می‌کنند. در نهایت مراحل بعدی شامل جمع‌بندی و مصورسازی نظر به کمک نرم افزار می‌باشد.

رویکردهای مبتنی بر یادگیری با نظارت ۱۵: با نام روش‌های مبتنی بر یادگیری ماشین یا روش‌های آماری برای طبقه‌بندی احساسات شناخته می‌شوند و از الگوریتم‌های داده‌کاوی تشکیل شده که الگوهای زیربنایی را از داده‌های آموزش داده شده یا برجسب گذاری شده یاد می‌گیرند، سپس در مرحله بعدی الگوریتم برای طبقه‌بندی داده‌های جدید بدون برجسب کلاس پیش بینی می‌شود، و سپس با استفاده از بازنمایی کلاس‌های پیدا شده توسط الگوریتم به عنوان ورودی برای عقیده‌کاوی استفاده می‌شوند.

جدول ۱. مروری کلی بر مهم‌ترین مقالات مطالعه شده مختص صنعت هتل‌داری

شماره	مقاله	حوزه‌ی مورد مطالعه	الگوریتم و روش مقاله	هدف و نتیجه کلی مقاله
۱	[۹]	داده‌های متنی نظرات از ۲ هتل از سایت Tripadvisor.com مورد مطالعه قرار گرفته‌است.	از شباهت‌های محتوا و احساسات برای تعیین تشابه دو جمله استفاده شد. برای شناسایی تعدادی جملات برگزیده ^۱ ، از الگوریتم خوشه‌بندی مدیود ^۱ استفاده شده است.	ارائه یک تکنیک جهت خلاصه کردن چند متن برای شناسایی جملات کلیدی و مهم از نظرات هتل را ارائه می‌دهد
۲	[۱۲]	داده‌های متنی نظرات هتل‌های اروپایی جمع آوری شده از وبسایت هتل‌ها	استفاده از الگوریتم‌های چند جمله ای بی‌ساده ^۱ و بی‌ساده‌ی برنولی ^۱ جهت طبقه‌بندی نظرات متنی	بهبود نتایج حاصل از طبقه‌بندی نظرات متنی
۳	[۱۳]	داده‌های متنی نظرات هتل‌های اروپایی جمع آوری شده از وبسایت هتل‌ها	انجام دقیق و مرحله به مرحله تکنیک‌های تحلیل احساس را برای پیش پردازش متن و تجزیه و تحلیل سپس تولید مصورسازی نتایج	کشف مفاهیم و واژه‌های کلیدی در تحلیل احساسات
۴	[۲]	استفاده از دو مجموعه داده از دو وبسایت از Slashdot و Epinions	استفاده از روش یادگیری ماشین بدون نظرات	طبقه بندی مجموعه داده‌های مورد بررسی به نظرات + یا -
۵	[۶]	این مجموعه داده شامل ۸۰۰ بررسی صادقانه و ۸۰۰ بررسی فریبنده در بین ۲۰ هتل محبوب شیکاگو توزیع شده‌است.	استفاده از الگوریتم‌های یادگیری ماشین با نظرات	تشخیص نظرات واقعی کاربران از نظرات اسپم یا دروغین توسط ربات یا افراد مغرض
۶	[۲۸]	داده‌های سه آژانس گردشگری آنلاین چین	روش‌های تجزیه و تحلیل معنایی و مصورسازی با نرم افزار گفی ^۱	مقایسه داده‌های متنی سه آژانس گردشگری چین و یافتن نقاط قوت و ارتباط آن‌ها با دیگر عوامل
۷	[۱۵]	داده‌های جمع آوری شده از وبسایت هتل	تجزیه و تحلیل آماری کلمات موضوعی	یافتن کلمه‌های موضوعی مانند اتاق مهمان (مثلاً نمای هتل، راحتی)، مکان (به عنوان مثال، نزدیک فرودگاه) و حمل و نقل (به عنوان مثال، شاتل، پارکینگ)
۸	[۷]	داده‌های جمع آوری شده از وبسایت هتل	تجزیه و تحلیل آماری کلمات موضوعی	باید زمینه را در طبقه بندی کلمه‌های موضوعی در نظر گرفت. جذابیت بدنی، حس شوخ طبعی و باتجربه بودن راهنماهای تور از عوامل مهم مؤثر بر تعامل آن‌ها با گردشگران است.

۳-۲ مراحل اصلی فرآیند

شامل جمع آوری داده ها، پیش پردازش متن، فرآیند اصلی، جمع بندی نتایج و تجسم بوسیله نمودارها و شکل ها است [۵]. مرحله جمع آوری داده ها: در حال حاضر برای دستیابی به این کار دو رویکرد وجود دارد.

اول از طریق رابط برنامه‌نویسی وب سایت^{۱۶} و دوم استفاده از خزنده‌های وب^{۱۷} به منظور دستیابی به داده‌ها از وب سایت‌های مورد نظر است.

مرحله پیش پردازش متن: متداول ترین تکنیک ها عبارتند از:

۱- نشانه گذاری^{۱۸}: که عملکرد آن باعث جدا کردن رشته متن کامل به لیستی از کلمات جداگانه می‌شود.

۲- یافتن ریشه و بن کلمه^{۱۹}: برای مثال واژه‌های شخص، اشخاص، شخصیت به بن آن‌ها یعنی شخص تبدیل می‌شوند.

۳- حذف کلمات بی‌اثر^{۲۰}: عملکرد آن باعث حذف واژگانی که برای ساخت زبان کاربرد داشته اما در محتوای معنایی آن تاثیر ندارند می‌شود. برخی از این کلمات در زبان انگلیسی ، a ، the و هستند.

۴- بخش‌بندی جمله^{۲۱}: عملکرد آن باعث تبدیل پاراگراف‌ها به جملات می‌شود.

۵- برچسب گذاری بخشی از گفتار^{۲۲}: عملکرد آن باعث برچسب خوردن هر کلمه، یک جمله یا بخشی از گفتار است. مانند صفت، اسم، فعل، ضرب‌المثل یا پیشگفتار. کاربرد این روش به عنوان ویژگی فرآیند یادگیری ماشین استفاده می‌شود.

۳-۲ روش پیشنهادی

در این مقاله، دو فرهنگ واژگان نظر آگاه به‌زمینه با روش-های^{۲۳} FBSA و^{۲۴} ALGA بر روی مجموعه داده آموزشی نظرات ساخته می‌شود [۱۷]. سپس، بر روی هر مجموعه داده، با استفاده از هر کدام از این دو فرهنگ واژگان نظر، ویژگی‌هایی محاسبه

می‌شوند. این ویژگی‌ها در کنار ویژگی‌های محاسبه شده با فرهنگ واژگان نظر عام منظوره لیوبینگ^{۲۵} قرار می‌گیرند و به

این ترتیب، برای ویژگی‌های حاصل، انتخاب ویژگی توسط آزمون t^{۲۶} صورت می‌گیرد.

آزمون t برای تعیین اختلاف میانگین یک گروه با یک مقدار پیش فرض و یا میانگین‌های دو گروه به کار می‌روند. در واقع، آزمون t یک نوع آمار استنباطی است برای تعیین اینکه آیا بین میانگین دو گروه اختلاف معنی داری وجود دارد یا خیر، چرا که

ممکن است در ویژگی‌های خاصی مرتبط باشند. در انتها، مدلی برای دسته‌بندی از روی مجموعه داده‌های آموزشی ساخته می‌شود و بر روی مجموعه داده آزمایشی اعمال می‌شود.

روش FBSA مبتنی بر ایجاد فرهنگ واژگان احساسی با استفاده از کامنت‌ها و نظرات است که با تناوب و میزان تکرار کلمات، بار احساسی آن‌ها مشخص می‌شود و روش ALGA

فرهنگ واژگان احساسی نظرات و عقاید بیان شده‌ی متنی را با استفاده از الگوریتم ژنتیک می‌سازد [۱۲].

۳-۱ تولید واژگان توسط FBSA

در روش پیشنهادی، در این بخش که از تولید واژگان توسط FBSA استفاده می‌شود، از روش تحلیل عبارات مبتنی بر فرکانس به دلیل تخمین دقیق رتبه‌های کلمات مثبت و منفی استفاده شده است. در روش FBSA برای

^{۲۲} Part-of-Speech (POS)

^{۲۳} Frequency Based Sentiment Analysis (FBSA)

^{۲۴} Adaptive Lexicon learning using a Genetic Algorithm (ALGA)

^{۲۵} Bing Liu's Opinion Lexicon English.

^{۲۶} t-test

^{۱۶} API

^{۱۷} Web crawlers

^{۱۸} Tokenization

^{۱۹} Stemming or lemmization

^{۲۰} Stopword Removal

^{۲۱} Sentence Segmentation

در رابطه‌ی (۲) اگر کلاس $R_{I,k}$ منفی باشد، از این رابطه استفاده می‌شود.

در این روابط، n_i تعداد رکوردها در D_i ، $R_{I,k}$ شماره رکورد k در مجموعه داده D_i و $TF(w_i, R_{I,k})$ تعداد رخ داده‌ها در w_i در $R_{I,k}$ است. در دیگر کلمات، $ferq_+(w_i, D_i)$ و $ferq_-(w_i, D_i)$ تعداد رخ داده‌ی w_i در رکوردهای مثبت و منفی در D_i مجموعه داده هستند.

در این روش، D_i رکوردهای مجموعه داده‌های آموزشی را شامل می‌شود و رکوردهای تست را در بر نمی‌گیرد. مقایسه $freq_+$ و $freq_-$ تنها زمانی معنی دار است که تعداد سوابق مثبت و منفی برابر است زیرا رکوردهای تست در نظر گرفته نشده و بنابراین، داده‌های آموزشی و تست برابر نیست.

به همین دلیل، از ضریب مبتنی بر سوابق در کلاس‌های مثبت و منفی استفاده شده است و سپس فرکانس نرمال با استفاده از رابطه‌ی (۳) محاسبه می‌شود:

$$ferq_-(w_i, D_i) = \frac{n_p(i)}{n_N(i)} \cdot ferq_-(w_i, D_i) \quad (3)$$

بنابراین، $n_N(i)$ و $n_p(i)$ تعداد رکوردهای مثبت و منفی را در D_i مجموعه داده نشان می‌دهند. در رابطه‌ی (۴) از رتبه بندی عبارت برای هر کلمه در واژگان استفاده می‌شود که این مقدار برای هر کلمه محاسبه می‌شود:

$$Score(w_i, D_i) = \frac{ferq_+(w_i, D_i) - ferq_-(w_i, D_i)}{ferq_+(w_i, D_i) + ferq_-(w_i, D_i)} \quad (4)$$

در رابطه‌ی فوق مقدار رتبه بین -1 تا $+1$ است. هر چقدر این عدد به 0 نزدیکتر باشد، این کلمه به فاعل یا همان کلمه‌ی ویژگی، نزدیکتر است. رتبه‌ی نزدیک به $+1$ مثبت بودن کلمه و رتبه‌ی نزدیک به -1 منفی بودن کلمه را نشان می‌دهد. بطور مثال اگر کلمه‌ی Love ۲۸ بار در عبارت بصورت جداگانه ظاهر شود، و ۳ بار بصورت منفی که ۲ بار در یک نظر و ۱ بار در نظر دیگر، در مجموع بصورت $ferq_+(love, D_i) = 28$ و $ferq_-(love, D_i) = 3$ نمایش داده می‌شود. پس از این مرحله، همان طور که در رابطه‌ی (۴) نشان داده شد،

یافتن بار احساسی واژگان تمام کلمات در نظر گرفته می‌شود و یکی از مزیت‌های این روش عدم حذف کلمات توقف^{۲۷} است زیرا می‌تواند در نظرات هتل‌ها تعیین کننده باشد. بدین صورت که رتبه‌ی عبارت برای هر کلمه براساس فرکانسی از کلمات در مجموعه داده‌های آموزشی محاسبه می‌شود. اگر فرض شود که p مجموعه داده داشته باشیم، D_1 تا D_p شامل نظرات و برجسب‌ها هستند که نیمی از آنها شامل نظرات مثبت و نیمی از آنها دارای نظرات منفی هستند. بنابراین، مجموعه داده‌های آموزش و تست با مدل اعتبارسنجی k -دسته^{۲۸} تقسیم می‌شوند. در این روش نمونه اصلی به‌طور تصادفی به زیرنمونه‌هایی با اندازه k تقسیم شده و در هر مرحله یک زیر نمونه مورد تحلیل قرار می‌گیرد.

از زیرنمونه‌های k ، که در هر مرحله بصورت تصادفی بدست آمده‌اند، یک زیرنمونه منفرد به‌عنوان داده‌های آزمایشی برای اعتبارسنجی الگوریتم و مدل ذخیره شده و زیرنمونه‌های دیگر که در واقع تعداد آن‌ها حالا $k-1$ شده است نقش داده‌های آموزشی را خواهند داشت. فرایند اعتبارسنجی، که k بار تکرار می‌شود، هر بار بصورت تصادفی مجموعه داده‌ای را انتخاب می‌کند، با هر یک از این نمونه‌های k دقیقاً یک بار داده‌ها اعتبارسنجی می‌شوند. نتایج k می‌تواند برای برآورد میانگین مورد استفاده قرار بگیرد. در این روش همه مشاهدات برای آموزش و اعتبار مورد استفاده قرار می‌گیرند، و هر مشاهده برای اعتبارسنجی به‌طور دقیق استفاده می‌شود و این مورد برتری این روش نسبت به نمونه‌گیری تصادفی تکراری است. لازم به ذکر است در این مقاله K برابر با ۱۰ در نظر گرفته شده است. بنابراین، برای هر کلمه w_i در مجموعه D_i داده‌های آموزشی، دو مقدار تجمعی تعریف می‌شود: فرکانس‌های مثبت و منفی^{۲۹}.

$$ferq_+(w_i, D_i) = \sum_{k=1}^{n_i} TF(w_i, R_{I,k}) \quad (1)$$

در رابطه‌ی (۱) اگر کلاس $R_{I,k}$ مثبت باشد، از این رابطه استفاده می‌شود.

$$ferq_-(w_i, D_i) = \sum_{k=1}^{n_i} TF(w_i, R_{I,k}) \quad (2)$$

^{۲۹} positive and negative frequencies

^{۲۷} Stop-Word

^{۲۸} K-Fold

واژگان نیز براساس داده‌های آموزشی ایجاد می‌شود و برای محاسبه‌ی ویژگی‌ها براساس داده آموزش و تست بکار می‌رود. در این روش تفاوت بین کلمات بدون هشتگ (#) و با هشتگ نیز در نظر گرفته می‌شود، کلماتی که با هشتگ در نظر گرفته می‌شوند، نشان دهنده تاکید و مهم بودن است. لذا، در برخورد با چنین کلماتی به دو صورت رفتار می‌شود و این کلمات در مجموعه آموزشی بسته به مثبت یا منفی بودن نظر ۲ بار شمارش می‌شوند .
 در این مرحله فهرست واژگان احساسی مبتنی بر تحلیل فرکانس ساخته شد، در ادامه به تولید واژگان براساس ABALGA پرداخته می‌شود تا در نهایت این واژگان باهم ادغام شوند.

۲-۳ نحوه‌ی استخراج ویژگی‌ها و جنبه‌های مختلف مطرح شده در متن نظرات

در این مقاله از روش IOB- encoding که در کتابخانه-ی nltk.corpus.reader با دستور import ConllChunkCorpusReader فراخوانی می‌شود که، برای استخراج جنبه‌های مختلف و صریح بیان شده در هر توثیت استفاده شده است. در این روش کلمات به کار رفته در توثیت‌ها برچسب گذاری می‌شوند که در آن B-POS نشانه‌ی جنبه‌ای است که در ابتدای توثیت شناسایی شده B، نشانه‌ی شروع جنبه‌ی جدید و O نشانه‌ی عدم شناسایی به عنوان جنبه و ویژگی می‌باشد. نمونه‌ای از خروجی در جدول (۴) نمایش داده شده است. جدول ۴. توثیت های برچسب گذاری شده با IOB-encoding

Words: Bathroom was clean , but bed is not comfort
Labels: B-POS O O O O B O O O

در این مقاله از مجموعه داده‌ی مورد نظر ۱۰۰۰ توثیت بررسی شد و پنج ویژگی با جنبه‌ی اصلی مختص هتل شناسایی شد و به صورتی که در جدول (۵) نمایش داده شده، دسته‌بندی گردید تا بتوان جنبه‌های مهم را شناسایی و بار احساسی کلمات مطرح شده را در نظرات حول این ویژگی‌ها شناسایی کرد .

رتبه‌ی کلمه پس از نرمال شدن محاسبه می‌شود. در ادامه پنج ویژگی برای طبقه بندی مجموعه داده‌ها بصورت زیر معرفی شده اند:

- Fpos: تعداد رتبه‌های کلمات مثبت در نظرات
- Fneg: تعداد رتبه‌های کلمات منفی در نظرات
- Pwords: تعداد کلمات مثبت در رکورد براساس رتبه
- Nword: تعداد کلمات منفی در رکورد براساس رتبه
- Score: مجموع همه‌ی رتبه‌ها در نظرات

در این روش، تولید واژگان مبتنی بر تحلیل عبارات فرکانسی به معنی این است که واژگان تولید شده در عبارت برای کل متن استفاده می‌شود. از آنجا که استثنائی برای مواردی است که نیاز است طبقه بندی شود، واژگان براساس مجموعه داده‌های آموزشی تولید می‌شوند.

هر رکورد در مجموعه داده (هر نظر) باید به یک بردار ویژگی تبدیل شود. فرض شود که یک رکورد شامل "It is good" است، بردار ویژگی توسط جدول (۲) محاسبه می‌شود.

جدول ۲. بردار ویژگی It is good [۱۷]

FPos	FNeg	PWords	NWords	Score
۰٫۸۴۹	-۰٫۱۸۶	۲	۱	۰٫۶۶۳

در روش پیشنهادی هر رکورد به چندین ویژگی براساس طول جمله تبدیل می‌شود که در جدول (۲) به پنج ویژگی تقسیم شده است. سپس این مدل برای مجموعه آموزشی ایجاد و برای مجموعه تست بکار می‌رود.

هسته بردار ویژگی از ویژگی‌های فوق الذکر ساخته شده که به صورت زیر است و در جدول (۳) نشان داده شده است. این جدول در واقع، نمونه‌ای از محاسبه ویژگی‌ها برای یک رکورد خاص است که با روش FBSA ایجاد شده است.

جدول ۳. واژگان ساده ایجاد شده توسط روش FBSA [۱۷]

It	For	Is	Good	Of
+۰٫۲۲۸	+۰٫۳۰۷	-۰٫۱۸۶	+۰٫۰۶۲۱	-۰٫۰۷۴

جدول ۵. دسته بندی جنبه‌های استخراج شده

Value	Location	Service	Meal	Room
Price	Railway		Breakfast	Bed
Amount	View	Check-in	Lunch	Bathroom
Rate	Airport	Check-out	Dinner	View
Cheap	Mall	Staff	Coffee	Shower
Worth	Far	Ticket	Tea	Air
Low	Close	Transport	Drink	condition
Money	Near		Restaurant	Bedsheets
Economic	Metro		Bar	Tv
Fee	distance			Furniture
expensive	market			

۳-۳-۱ ادبیات تحقیق (الگوریتم ژنتیک)

امروزه الگوریتم ژنتیک جایگاه ویژه‌ای در میان الگوریتم‌های بهینه سازی برای حل مسائل پیچیده دارد زیرا از لحاظ محاسباتی ساده، در عین حال قدرتمند است، همچنین در هر مرحله فضای جستجو در مجموعه‌ی داده محدود نمی‌شود [۱۹].

الگوریتم‌های فراابتکاری همچون الگوریتم ژنتیک، یکی از الگوریتم‌های جستجو به حساب می‌آید و از طریق تعامل با اعضا، در پی یافتن جواب بهینه‌ی سراسری هستند. در همه این الگوریتم‌ها، جواب‌های بهتر، شانس بیشتری برای حضور در تکرارهای بعدی الگوریتم و تولید نسل بعد دارند که این ویژگی خاص " حیات مناسب‌ترین^{۳۰}" موجب یافتن نتایج بهتر است [۱۸].

شرط پایان الگوریتم، رسیدن به حداکثر تعداد تکرار از پیش تعیین شده، عدم بهبود جواب در چند تکرار پیاپی می‌باشد [۲۰].

۳-۳-۲ پیش پردازش

مراحل پیش پردازش در ABALGA شامل موارد زیر است:

- ✓ جداسازی کلمات در هر بررسی^{۳۱}
 - ✓ حذف کلمات توقف
 - ✓ واژه‌های فیلتر شده براساس برچسب‌های مثبت یا منفی
 - ✓ لمس کردن کلمات باقی مانده (به جز اصطلاح جنبه)
 - ✓ ساخت مجموعه‌ای از کلمات ریشه^{۳۲} از مرحله قبل
 - ✓ حذف کلمات با فرکانس کمتر از ۳
 - ✓ گرفتن پنجره با اندازه ثابت در حدود اصطلاحات
- بعنوان نمونه‌ای از پیش پردازش که در موارد بالا ذکر شد، منظور از جداسازی، تشخیص مرز کلمات در متون است،

۳-۳ تولید واژگان توسط ABALGA

در این مقاله یادگیری فهرست واژگان تطبیق شده با ویژگی‌ها، با استفاده از الگوریتم ژنتیک و در سطح توییت طراحی شده‌است. هر توییت ممکن است یک جمله کوتاه یا یک جمله طولانی از ۱۴۰ کاراکتر تا ۲۸۰ کاراکتر باشد. روش پیشنهادی برای نظرات کاربران صنعت هتلداری استفاده می‌شود و در این مقاله متن نظرهای کوتاه‌تر و یا بلندتر مدنظر قرار نگرفته است.

در الگوریتم پیشنهادی، یک توالی ژنی تعریف می‌شود تا زمانی که کلمات موجود در مجموعه آموزش، مشغول یادگیری واژگان هستند، ادامه می‌یابد. هر ژن در توالی ژن ذکر شده دارای نمره‌ای برای کلمه مربوطه است. به عبارت دیگر، درصد بهینه سازی رتبه‌ی احساسی واژگان هستیم که این بهینه سازی با عملیات کراس آور و جهش که در الگوریتم ژنتیک در جهت ارائه بهترین راه حل صورت می‌گیرد، انجام می‌شود که در ادامه تشریح شده است. در ALGA، رتبه‌های کلمات، ژن‌ها هستند. کمبود ALGA در مشکلات مبتنی بر جنبه از ساختار نظرات ناشی می‌شود.

بطور مثال در تحلیل نظرات، گاهی به هر توییت یک برچسب نسبت داده می‌شود که این برچسب برای کل کلمات در نظر گرفته می‌شود. از طرفی دیگر هر جمله ممکن است جنبه‌های مختلفی داشته باشد که ناشی از بخش‌های مختلف نظر مطرح شده باشد که نیاز به تجدید نظر دارد لذا همانطور که در بخش قبلی توضیح داده شد، جنبه‌های مختلف بیان شده در هر توییت استخراج می‌شود.

^{۳۲} lemmatized

^{۳۰} Survival of the fittest

^{۳۱} tokenization of words in each review

مجموعه قرار دارد. هر ژن برای یک کلمه ریشه‌ی مربوطه در ورودی دارای یک رتبه‌ی شناور در محدوده ۱- تا ۱ است. در مرحله‌ی اول رتبه‌ها به صورت تصادفی به هر ژن بصورت عددی در بازه‌ی ۱ - تا ۱ داده می‌شود. تابع برازندگی برای از بین بردن محاسبات زائد و سرعت بخشیدن به کار می‌رود و در کروموزوم‌ها ذخیره می‌شود. در روش پیشنهادی برای هر توالی ژن، تابع برازندگی فقط یک بار در جهت سرعت بخشیدن به روند، محاسبه می‌شود مگر اینکه با جهش یا کراس آور تغییر کند [۱۹]. پارامتر سن تعداد دفعاتی را که توالی فعلی پس از انجام جهش یا کراس آور شکست خورده است، را شمارش می‌کند. پس از رسیدن به حداکثر سنی، الگوریتم، کروموزوم‌ها را از مخزن پدر حذف می‌شوند و کروموزوم فرزند ایجاد می‌شود. این روند بعد از عمل جهش و کراس آور بوجود می‌آید [۱۹].

کروموزوم همچنین اطلاعات استراتژی را که یکی از توابع ایجاد، جهش یا کراس آور است، ذخیره می‌کند و نشان می‌دهد که کدام عملکرد منجر به کروموزوم فعلی شده است.

۳-۳-۴ تابع ایجاد

تابع ایجاد بعنوان ورودی طولی از کلمات را می‌گیرد و لیستی با همان طول را تولید می‌کند که حاوی مقادیر شناور تصادفی در دامنه ۱- تا ۱+ است که در واقع رتبه کلمات در مجموعه کلمات است. بطور مثال مقدار ژن مربوطه واژه "polite" در کروموزوم ۱ می‌تواند ۰,۳۱+ باشد که نشانگر درجه احساسات نسبتاً مثبت برای واژه‌ی مذکور است، در حالی که در کروموزوم ۲، مقدار ژن مربوطه می‌تواند ۰,۴۶- باشد، که نشان دهنده رتبه احساسات منفی است. نمونه‌ای از این توالی در جدول (۷) نشان داده شده است.

جدول ۷. نمونه‌ای از توالی‌های ژن‌ها در الگوریتم ژنتیک

کلمات	near	Far	clean	warm	danger	fair	Polite
توالی	۰,۴۵+	۰,۷۸-	۰,۵۲+	۰,۲۵+	۰,۴۲-	۰,۲۳+	۰,۸۰+
۱							

بدین صورت که متن را به دنباله‌ای از کلمات تبدیل می‌کند. در مورد دوم، حذف کلمات توقف، برخی از کلمات همانند the, is و... که ارزش احساساتی ندارند، از جمله حذف می‌شوند. در مرحله‌ی فیلتر کردن حذف یا نگهداری برخی از کلمات می‌تواند در نظر گرفته شود. همچنین در مراحل بعد بازگردان شکل کلمه به حالت ریشه و بن، حذف کلمات کمتر از سه حرفی که ارزش محاسباتی ندارند. در مرحله‌ی آخر، گرفتن پنجره با اندازه پنج در حدود جنبه یا اصطلاح room می‌تواند در توثیت "The room was clean and I satisfied the view is good." بصورت زیر در نظر گرفته می‌شود. در این مرحله تعداد کلمات توقف هم شمارش می‌شود و جدول (۶) نشان داده شده است.

جدول ۶. پنجره‌ی کلمات در حدود جنبه

[5words before]Aspect		[5words after]				
The	room	was	clean	and	i	satisfied

در ادامه الگوریتم ABALGA شامل ساختار کروموزوم‌ها^{۳۳}، تابع ایجاد^{۳۴}، برازندگی^{۳۵} و تابع جهش^{۳۶} و تابع کراس آور^{۳۷} و انتخاب والد‌ها تشریح خواهند شد.

۳-۳-۳ ساختار کروموزوم‌ها

کروموزوم‌ها در الگوریتم ABALGA اطلاعات زیر را نگهداری می‌کنند:

- ژن‌ها
- برازندگی
- سن
- استراتژی

موازی سازی با مجموعه‌ای از کلمات ریشه ذکر شده در بخش پیش پردازش، به عنوان یکی از پارامترهای ورودی در ABALGA، یک توالی ژنی به طول یکسانی در این

^{۳۶} Mutate function

^{۳۷} Crossover function

^{۳۳} Chromosome structure

^{۳۴} Create

^{۳۵} Fitness

مبتنی بر ABALGA در پنجره به شرح زیر محاسبه می‌شود:

$$Label(D_m, W_l, L_k) = \begin{cases} Positive & if \ ABALGA(D_m, W_l, L_k) > 0 \\ Negative & if \ ABALGA(D_m, W_l, L_k) \leq 0 \end{cases} \quad (7)$$

در روش پیشنهادی از تابع برازندگی برای قیاس کردن کروموزم‌ها استفاده می‌شود. الگوریتم ABALGA دارای دوهدف است: هدف اصلی این است که تعداد تطبیق‌ها را حداکثر کند. درواقع کروموزومی که تعداد تطبیق بالاتری داشته باشد نسبت به دیگر کروموزم‌ها ارجح خواهد بود. اگر حالت تطبیق‌ها محاسبه شده یکسان باشد، در تابع برازندگی، مقدار قدرمطلق تفاضل برای بدست آوردن رتبه استفاده می‌شود. دیگر هدف این است که به تعداد تطبیق na به صفر برسد.

این عملکرد پاداش و مجازات تجدید نظر شده ABALGA، در روش ما، شبیه به ALGA [۱۸] است.

۳-۳-۶ تابع جهش

عملکرد جهش یک توالی ژنی و تابع برازندگی و همچنین پنجره‌های اطراف اصطلاحات را به عنوان پارامترهای ورودی در نظر می‌گیرد و برای بهبود تابع برازندگی خود سعی در تغییر تابع برازندگی قبلی دارد. واژگان کاندیدا در پنجره عبارت‌هایی هستند که برچسب محاسبه شده از رابطه‌ی (۷) با برچسب ابعاد واقعی آن مطابقت نداشته است.

۳-۳-۷ تابع کراس‌آور

توالی ژن به نام‌های والدین و اهداکننده‌ها و مقادیر تابع برازندگی آن‌ها، پارامترهای ورودی تابع کراس‌آور الگوریتم ABALGA است. هنگامی که ژن‌های والدین و اهداکننده یکسان هستند، از تابع ایجاد برای جایگزینی ژن‌های اهداکننده با توالی ژن جدید استفاده می‌شود زیرا توالی ژن‌های یکسان در مخزن والدین مطلوب نیستند. اگر این محدودیت ارضا نشود، سپس، به تعداد تصادفی از دفعات، در دامنه ۱ تا ۱۰، تابع در ژن‌های اهداکننده که ژن‌های مشابه آن‌ها در ژن‌های والدینی که یکسان نیستند،

توالی ۲	+۰.۳۲	-۰.۵۱	+۰.۳۲	+۰.۴۱	-۰.۱۷	-۰.۳۶	-۰.۱۴
توالی ۳	-۰.۲۱	-۰.۶۵	+۰.۴۱	-۰.۳۶	-۰.۴۵	+۰.۴۸	+۰.۱۲
توالی ۴	+۰.۶۳	+۰.۱۲	+۰.۲۳	-۰.۵۲	+۰.۲۱	-۰.۲۲	+۰.۷۹

۳-۳-۵ برازندگی

در روش پیشنهادی تابع پاداش و مجازات^{۳۸} الگوریتم ABALGA متفاوت از ALGA است. عملکرد پاداش و مجازات دوباره در مقایسه با تابع پاداش و مجازات منفرد در ALGA طراحی شده است. با دادن مجموعه داده و با استفاده از Lk واژگان برای هر کلمه‌ی ریشه‌یابی شده، $lem(w_i)$ در پنجره اطراف واژه‌ی جنبه‌ی مربوطه در الگوریتم ABALGA از رابطه‌ی (۵) محاسبه می‌شود:

$$ABALGA(D_m, W_l, L_k) = \sum_{W_j} S_k(lem(w_i)W_j) \quad (5)$$

بطوریکه $S_k(lem(w_i)W_j)$ بصورت رابطه‌ی (۶) محاسبه می‌شود:

$$S_k(lem(w_i)W_j) = \begin{cases} -S_k(lem(w_i)W_j) & if \ w_{j-1} \in in \ Negations \\ S_k(lem(w_i)W_j) & otherwise \end{cases} \quad (6)$$

در رابطه‌ی (۵)، $S_k(lem(w_i)W_j)$ رتبه‌ای از هر کلمه‌ی lemmatized، w_j در پنجره اطراف واژه‌ی جنبه‌ی مربوطه است. در رابطه‌ی (۶) رتبه‌ی احساسات یک کلمه در یک پنجره را بیان می‌کند، اگر کلمه‌ای با بار منفی قبل از کلمه‌ی جاری باشد، در محاسبه رتبه‌ی پنجره معکوس می‌شود. لیستی از کلمات با بار منفی در جدول (۸) نشان داده شده است.

جدول ۸. لیست کلمات منفی [۱۷]

Not	n't	No	Barely	rarely	Never	Hardly
-----	-----	----	--------	--------	-------	--------

برای تمام پنجره‌هایی که حاوی اصطلاح جنبه هستند، محاسبه می‌شود. اگر رتبه برای یک پنجره مثبت باشد، برچسب جنبه مربوط به عنوان مثبت و در غیر این صورت منفی شمارش می‌شود. بنابراین، پیش بینی برچسب

^{۳۸} Reward and Penalty

NegCount: تعداد کلمات، در پنجره فعلی، با توجه به واژگان در حال استفاده، دارای رتبه‌ی منفی است.
CountSum: تعداد کلمات با نمره مثبت منهای تعداد کلمات با نمره منفی، در پنجره فعلی، مطابق واژگان در حال استفاده.

PosSum: مجموع رتبه‌های مثبت کلمات در پنجره فعلی برای واژگان جاری

NegSum: مجموع رتبه‌های منفی کلمات در پنجره فعلی برای واژگان جاری

طبقه بندهای متعددی شامل درخت تصمیم، درخت تصادفی، جنگل تصادفی، ماشین بردار پشتیبان، بیزین و طبقه بندی نزدیکترین همسایگی برای روش پیشنهادی در نظر گرفته می‌شود. در روش پیشنهادی، واژگان لیو بینگ بعنوان پشتیبان اصلی در الگوریتم FBSA یا ABALGA استفاده می‌شود.

بنابراین؛ واژگان SentiWordNet بعنوان آخرین لایه برای رتبه دهی کلمه بصورت رابطه‌ی (۸)، (۹) و (۱۰) محاسبه می‌شود:

$$Pos_{score}(w_i) = \frac{\sum_i^n p(i)}{totalSyn(w_i)} \quad (8)$$

$$Neg_{score}(w_i) = \frac{\sum_i^n N(i)}{totalSyn(w_i)} \quad (9)$$

$$Obj_{score}(w_i) = \frac{\sum_i^n o(i)}{totalSyn(w_i)} \quad (10)$$

در روابط فوق، $p(i)$ ، $N(i)$ و $O(i)$ رتبه‌های مفعولی هستند. $totalSyn(w_i)$ در واقع تعداد Synset برای w_i است.

در نهایت، رتبه نهایی برای هر کلمه با استفاده از واژگان SentiWordNet با استفاده از رابطه‌ی (۱۱) محاسبه می‌شود:

$$word\ score(w_i) = \quad (11)$$

$$\begin{cases} Pos_{score} & \text{if } Pos_{score} > Neg_{score} \text{ and } Pos_{score} > Obj_{score} \\ Neg_{score} & \text{if } Neg_{score} > Pos_{score} \text{ and } Neg_{score} > Obj_{score} \\ Obj_{score} & \text{otherwise} \end{cases}$$

ABALGA واژگان احساسی را به صورت پویا و FBSA احساسات درون کلمات را بصورت استاتیک (ثابت) ضبط می‌کنند، رتبه کلمات احساسی بصورت پویا، ممکن است نادرست نشان داده شود و از سوی دیگر، واژگان استاتیک رتبه دقیق تری از احساسات ارائه می‌دهند

قرار می‌دهد. سپس ژن‌های مستقر در نسخه‌ای از ژن‌های والدین استخراج و جایگزین می‌شوند. در صورت پیشرفت پس از این جایگزینی، تابع ژن‌های فعلی را برمی‌گردانند و در صورت عدم پیشرفت، تابع دوباره سعی می‌کند تا به یک آستانه حداکثر برسد، سپس در آن مرحله آخرین نسخه کپی شده ژن‌های والدین را برمی‌گرداند.

۳-۳-۸ انتخاب والدین

هدف اصلی سوق دادن جستجو به بخشهایی از فضا که امکان یافتن جوابهای با کیفیت بالاتر وجود دارد. در هر نسل تعدادی از عناصر جمعیت این فرصت را پیدا می‌کنند که تولید مثل کنند. به این عناصر که از میان جمعیت انتخاب می‌شوند، والدین می‌گویند.

به این منظور در ابتدا همه‌ی واژه‌ها به عنوان والد شناخته می‌شوند، سپس در هر مرحله با استفاده از تابع برازندگی که در قسمت قبلی توضیح داده شد هر کروموزومی که تعداد تطبیق بالاتری داشته باشد به عنوان والد انتخاب می‌شود [۱۹].

۳-۴ ادغام واژگان و طبقه بندی

به منظور طبقه بندی نظرات، واژگان FBSA، ABALGA، فرهنگ واژگان SentiWordNet، و فرهنگ واژگان احساسی لیو بینگ با یکدیگر ترکیب شده و برای رتبه دهی کلمات استفاده می‌شوند. بدین صورت که اندازه پنجره ۵ در نظر گرفته می‌شود و کلمات بیان شده در توثیت به عنوان دامنه‌ی ورودی به این پنجره‌ها و امتیازات شمارش شده‌ی زیر بعنوان ویژگی‌ها برای هر جنبه استخراج می‌شوند. بنابراین، ویژگی‌های استخراج شده برای هر پنجره بصورت زیر هستند:

ScoreSum: مجموع رتبه‌های کلمه در پنجره بر طبق واژگان

NormalizedScoreSum: مجموع رتبه کلمات در یک پنجره با توجه به واژگان استفاده شده طول، نشانه‌ها، از پنجره تقسیم می‌شود.

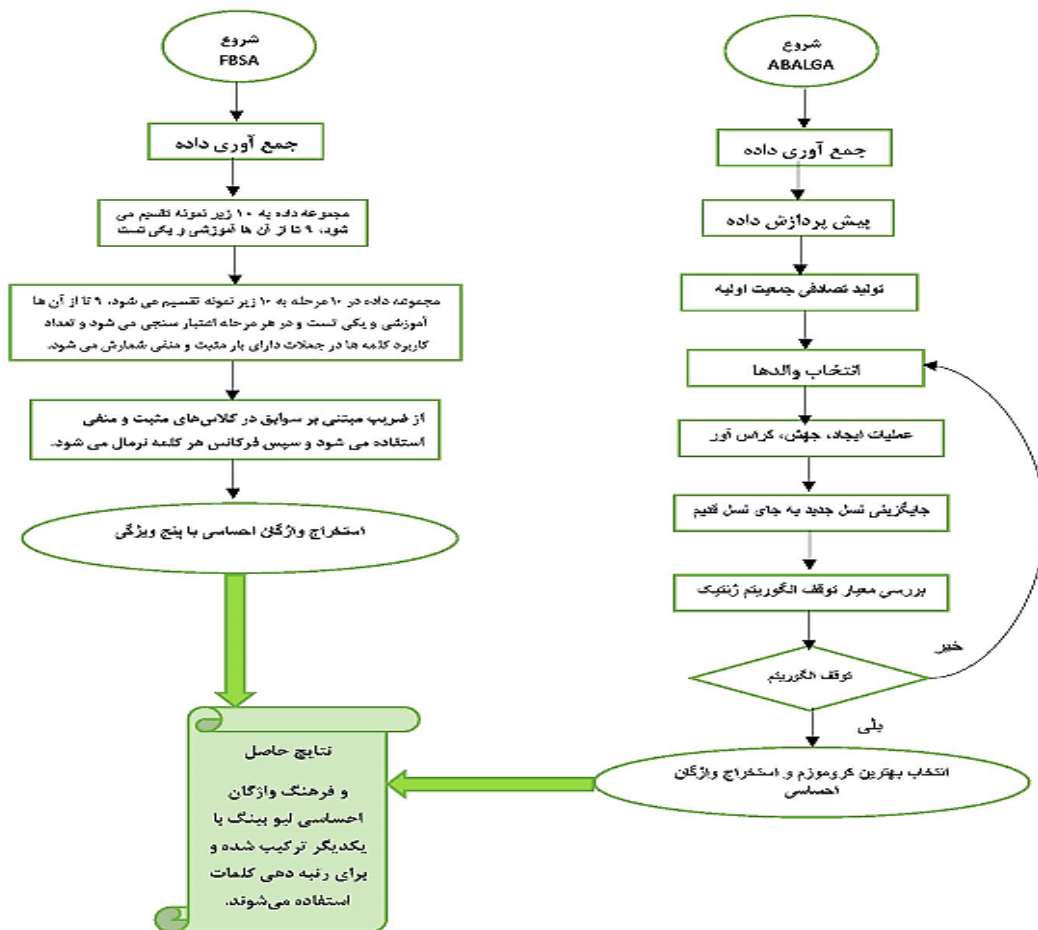
DistSum: مجموع، هر کلمه، از نمرات تقسیم شده با فاصله آن از اصطلاحات.

PosCount: تعداد کلمات، در پنجره فعلی، با توجه به واژگان در حال استفاده، دارای رتبه‌ی مثبت است.

فرهنگ واژگان احساس نظر را با استفاده از الگوریتم ژنتیک می‌سازد استفاده و در مرحله‌ی نهایی به طبقه بندی مجموعه داده‌ی نظرات صنعت هتلداری پرداخته شد. در شکل (۱) مراحل کلی بیان شده در روش پیشنهادی را نشان می‌دهد.

اما قادر به در نظر گرفتن رتبه‌ی احساسی در یک متن نیستند. از این‌رو، ادغام واژگان پویا و استاتیک معیارهای عملکرد را در مقایسه با استفاده از تنها یک نوع واژگان افزایش می‌دهد لذا این الگوریتم رتبه‌دهی احساسی واژگان در این مقاله برای استخراج ویژگی‌های کلمات در نظر گرفته شده است.

در روش پیشنهادی، پس از پیش پردازش مجموعه داده‌ی نظرات کاربران، به استخراج ویژگی‌های متا پرداخته شد. در این روش، از یک روش ترکیبی و جدید با استفاده از روش FBSA که مبتنی بر ایجاد فرهنگ واژگان نظر با استفاده از میزان تکرار و تناوب کلمات است و ABALGA که



شکل ۱. مراحل کلی مدل پیشنهادی

(Negative_Review,۸)Review_Total_Negative_Word_Counts,۹)Total_Number_of_Reviews,۱۰)Positive_Review,۱۱)Review_Total_Positive_Word_Counts,۱۲)Total_Number_of_Reviews_Reviewer_Has_Given,۱۳)Reviewer_Score,Tags,۱۴)days_since_review,۱۵)lat,۱۶)lng, ۱۷)Sty

ما در این پژوهش تنها ستون نظرات مثبت و نظرات منفی کلاس‌بندی شده را در نظر می‌گیریم. که نمونه‌ی آن در شکل زیر نمایش داده شده‌است.

Negative_Review	Positive_Review
I am so angry I made this post.... No negative	Location of hotel was good Very nice and amazing hotel
Rooms are nice but elderly a bit difficult My room was dirty	Very clean and staff are polite Great location
Im sad because of food..	Foods are delicious

شکل ۲. نمونه‌ای از مجموعه داده

۲-۴ معیارهای ارزیابی

در مطالعات قبلی صورت‌گرفته روی داده‌های متنی از جمله روی داده‌های مربوط به صنعت گردشگری و هتل‌داری، از معیارهای ارزیابی صحت^{۴۱}، بازخوانی^{۴۲}، دقت^{۴۳} و معیار^{۴۴} F استفاده شده است، لذا برای مقایسه مدل پیشنهادی این مقاله با سایر روش‌ها از معیارهای فوق که از معیارهای اصلی سنجش میزان دقت مدل‌ها می‌باشد، استفاده گردیده است.

۲-۴-۱ دقت، بازخوانی و ماتریس درهم ریختگی^{۴۵}

ماتریس درهم ریختگی یک ماتریس مربعی $N \times N$ در می-باشد که N همان تعداد برجسب‌ها و کلاس‌های مشخص شده در دسته بند تعریف شده است. پارامترهای مرتبط با دقت و بازخوانی توسط الگوریتم‌های داده‌کاوی بطور کلی و بخصوص در موضوع مقاله تحلیل احساسات و عقیده‌کاوی بصورت زیر تعریف می‌شوند:

۴- نتایج و تفسیر داده‌ها

در این بخش مجموعه داده‌ی مورد نظر مورد تحلیل قرار گرفت و نتایج حاصل از تحلیل‌هایی که توسط الگوریتم‌های مربوطه در مدل پیشنهادی بکار رفته بود با نتایج دیگر مقاله‌ها مقایسه شد. نتایج حاصل نشان داد که مدل پیشنهادی نسبت به مقاله‌های قبلی عملکرد مناسبتری دارد.

۴-۱ مجموعه داده‌ها

در این مقاله مجموعه داده‌ای متنی مرتبط با نظرات کاربرانی که از خدمات هتل استفاده کرده‌اند، فراهم شده که شامل لیستی از نام هتل‌های اروپایی و نظرات و عقاید متنی به زبان انگلیسی توسط کاربرانی^{۳۹} که قبلاً از خدمات هتل مربوطه استفاده کرده‌اند می‌باشد که از سایت کگل^{۴۰} فراهم آورده شده است.

سایت کگل به پژوهشگران این امکان را می‌دهد تا مجموعه داده‌های مناسب مقاله در زمینه علم داده را پیدا کنند همچنین اگر مجموعه داده‌ای دارند که نیاز به تحلیل دارد را منتشر کنند و با کمک افراد حرفه‌ای و متخصص در یک محیط دانش مبتنی بر وب، دانش پنهان در مجموعه داده را کشف و مدل‌های مناسب و کاربردی بسازند.

مجموعه داده منتخب این مقاله شامل داده‌های متنی ۱۴۹۳ هتل با مقادیر رتبه‌بندی شده منحصر بفرد در اروپا می‌باشد. این مجموعه داده‌ها شامل ۵۱۵۰۰۰ نظر مشتری است که همگی به زبان انگلیسی بیان شده‌اند. برای تحلیل بیشتر، مکان‌های جغرافیایی هتل‌های مختلف در این مجموعه داده در نظر گرفته شده که در قالب فایل‌های CSV تهیه و تنظیم شده است.

مجموعه داده‌ی فوق‌الذکر دارای ۱۷ ستون است، که عناوین آن عبارت است از:

۱)Hotel_Address, ۲)Additional_Number_of_Scoring, ۳)Review_Date, ۴)Average_Score, ۵)Hotel_Name, ۶)Reviewer_Nationality, ۷

^{۴۳} Precision

^{۴۴} F-Measure

^{۴۵} Confusion matrix

^{۳۹} ۵۱۵K hotel reviews in Europe

^{۴۰} WWW.Kaggle.com

^{۴۱} Accuracy

^{۴۲} Recall

۴-۲-۳ صحت

صحت معیاری است که در روش‌های ارایه شده در داده کاوی برای طبقه بندی‌ها کاربرد داشته و به میزان نزدیکی پیش بینی‌های مدل پیشنهادی با میزان اندازه گیری شده با مقدار واقعی اشاره دارد و بوسیله‌ی رابطه‌ی (۱۵) به صورت تقسیم تعداد کلماتی که به درستی طبقه‌بندی شده اند به تعداد کل کلمات در مجموعه داده مورد نظر محاسبه می‌شود.

$$Accuracy = \frac{T_p + T_N}{N} \quad (15)$$

۴-۳ نتایج آزمایش‌ها

۴-۳-۱ محیط آزمایش و شبیه سازی مدل پیشنهادی

الگوریتم مدل پیشنهادی این مقاله در محیط برنامه نویسی آنا کوندا^{۵۱} که یک توزیع متن باز برای زبان‌های برنامه نویسی پایتون و R می‌باشد در غالب نرم افزار اسپایدر^{۵۲} و نرم افزار R پیاده‌سازی شده است.

۴-۳-۲ پارامترهای استفاده شده برای ارزیابی مدل پیشنهادی

برای ارزیابی نتایج حاصل شده از روش و مدل پیشنهاد شده با سایر روش‌ها از سه معیار بیان شده در قسمت قبلی یعنی دقت، صحت و معیار بازخوانی، مورد استفاده قرار گرفته اند.

۴-۳-۳ نتایج تحلیل روش پیشنهادی

در این مقاله، نظرات کاربران به خوب (مشتری راضی) و بد (مشتری ناراضی) تقسیم می‌شود. بدین شکل که رتبه کلی نظرات بد بصورت $ratings < 5$ و نظرات خوب بصورت $ratings \geq 5$ می‌باشند که در شکل (۳) نشان داده شده است.

پارامتر $TP^{۴۶}$ بیان کننده‌ی تعداد ویژگی‌هایی است که به درستی توسط مدل انتخاب شده به عنوان ویژگی بازیابی شده‌اند.

پارامتر $FP^{۴۷}$ مبین تعداد ویژگی‌هایی است که به صورت نادرست توسط مدل انتخاب شده به عنوان ویژگی بازیابی شده‌اند.

پارامتر $TN^{۴۸}$ مربوط به تعداد ویژگی‌هایی است که به درستی توسط مدل انتخاب شده به عنوان ویژگی بازیابی نشده‌اند.

پارامتر $FN^{۴۹}$ مربوط به تعداد ویژگی‌هایی است که به نادرستی توسط مدل انتخاب شده به عنوان ویژگی بازیابی شده‌اند.

با استفاده از پارامترهای تعریف شده در بالا و اطلاعات بازیابی شده توسط الگوریتم های متن کاوی می‌توانیم کارایی مدل های مختلف را با استفاده از رابطه‌های زیر ارزیابی کنیم.

$$Precision = TP / (TP + FP) \quad (12)$$

$$Recall = TP / (TP + FN) \quad (13)$$

معیار Precision در رابطه (۱۲) میزان دقت مدل انتخاب شده را مشخص می‌کند و میزان برچسب‌های درست را نشان میدهد

معیار Recall در رابطه (۱۳) بیان کننده نسبت تعداد داده‌های متنی درست دسته‌بندی شده که توسط ما برچسب زده شده در یک کلاس خاص، به تعداد کل داده‌های متنی است که باید در همان کلاس خاص دسته‌بندی شوند.

۴-۲-۲ معیار ترکیبی $F^{۵۰}$

برای ارزیابی عملکرد دسته‌بندها بسیار مورد استفاده قرار می‌گیرد و از ترکیب دو پارامتر دقت و بازخوانی حاصل می‌شود. این معیار ترکیبی به صورت زیر بدست می‌آید:

$$F\text{-measure} = 2 * ((precision * recall) / (precision + recall)) \quad (14)$$

^{۵۰} Hybrid F-measure

^{۵۱} AnaConda

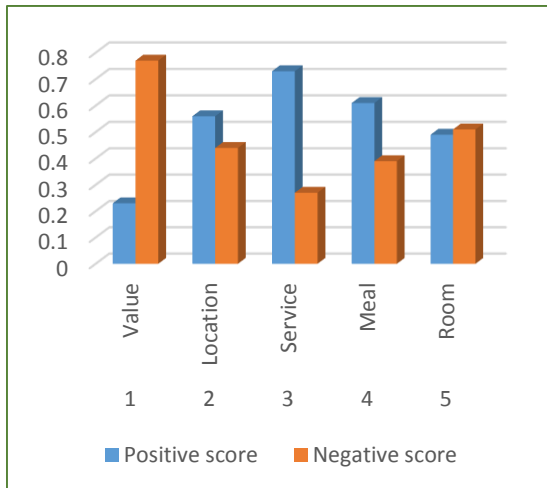
^{۵۲} Spyder

^{۴۶} True Positive

^{۴۷} False Positive

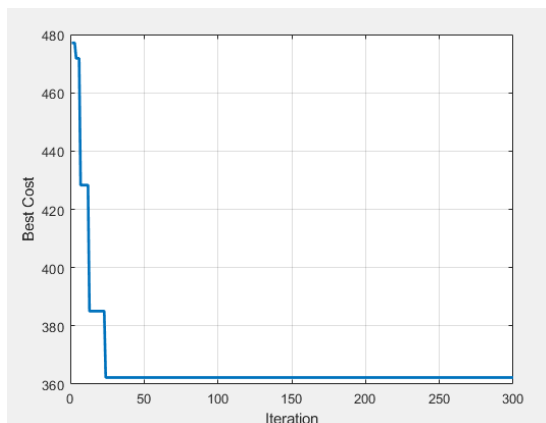
^{۴۸} True Negative

^{۴۹} False Negative



نمودار (۳) امتیاز بار احساسی ویژگی‌های استخراج شده

پس از استخراج ویژگی‌های ترکیبی، الگوریتم ژنتیک جهت بهینه سازی واژگان بکاررفته که در شکل (۳) نشان داده شده است.



شکل (۳) همگرایی بهینه سازی ویژگی‌ها

شکل (۳) نشان می‌دهد الگوریتم ژنتیک پس از ۳۰۰ بار تکرار به بهترین تابع برازندگی خود می‌رسد و همگرا می‌شود.

در جدول (۹) معیارهای دقت، صحت و F در مرحله ی آخر که ادغام و طبقه‌بندی می‌باشد مرحله به مرحله نشان داده شده است.

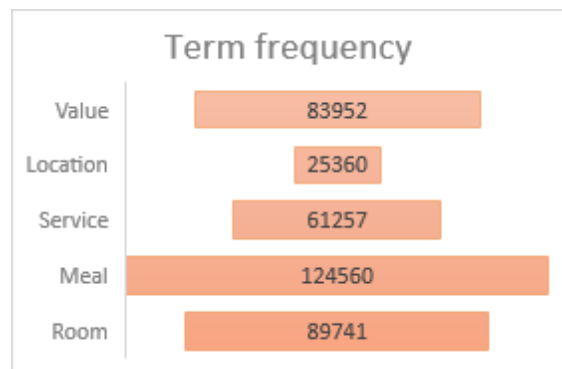
جدول (۹) معیارهای محاسبه شده در مرحله ی ادغام و طبقه بندی

نتایج مراحل ادغام و طبقه بندی	Precision	Recall	F-measure
ScoreSum	۹۰,۴	۹۰,۵	۹۰,۴
NormalizedScoreSum	۹۱,۸	۹۱,۸	۹۱,۸
DistSum	۹۰,۶	۹۰,۷	۹۰,۶
PosCount	۹۱,۹	۹۲,۰	۹۱,۹
NegCount	۹۴,۸	۹۳,۲	۹۳,۱
CountSum	۹۲,۰	۹۲,۰	۹۲,۰
PosSum	۹۴,۱	۹۴,۰	۹۴,۰

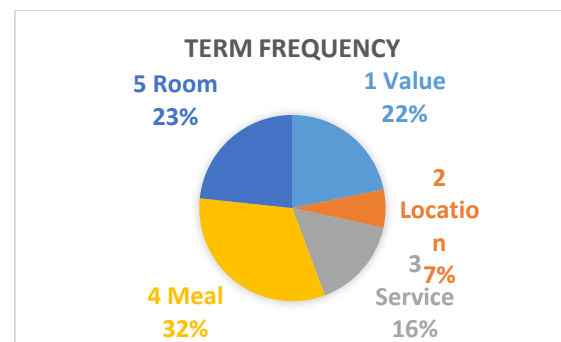
review	is_bad_review
0 I am so angry that i made this post available...	1
1 No Negative No real complaints the hotel was g...	0
2 Rooms are nice but for elderly a bit difficul...	0
3 My room was dirty and I was afraid to walk ba...	1
4 You When I booked with your company on line y...	0

شکل (۳) دسته بندی نظرات خوب و بد

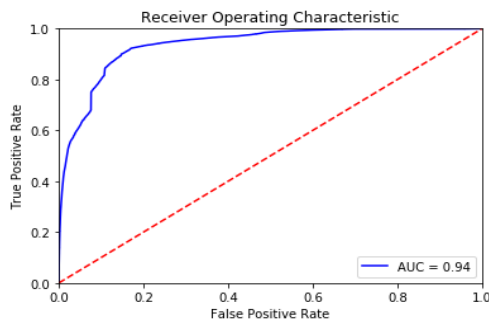
حال در ادامه با توجه به جنبه‌های استخراج شده با توجه به اینکه کلمه‌های دسته‌بندی شده در هر گروه دقیقاً چندبار در سند ظاهر می‌شود و اهمیت نسبی آن‌ها برای متصدیان امور جهت اتخاذ تصمیمات مقتضی در نمودار (۱) و نمودار (۲) به ترتیب تعداد و درصد تکرار آن‌ها و در نمودار (۳) امتیاز بار احساسی استنباط شده از آن‌ها را نشان داده شده است.



نمودار (۱) تعداد تکرار جنبه‌های استخراج شده



نمودار (۲) درصد تکرار جنبه‌های استخراج شده



شکل (۶) منحنی ROC و نمرات AUC از رگرسیون لجستیک

همانطور که در اشکال فوق نشان داده شده، سطح زیر نمودار الگوریتم ماشین بردار پشتیبان با وجود نزدیکی به نایوبیزین بیشتر است، بنابراین از عملکرد خوبی نسبت آن و نسبت به رگرسیون لجستیک برخوردار است و در رده دوم قرار دارد. شکل‌های (۷)، (۸) و (۹) صحت در کانفیوژن ماتریس را براساس برجسب‌های واقعی و پیش‌گویی شده نشان می‌دهد.

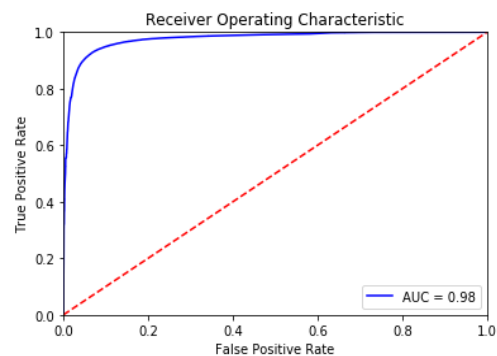
		Predicted Label		
		Negative	Positive	
True Label	Negative	True Neg: 146635 (Num Neg: 154693)	False Pos: 8058	False Pos Rate: 0.05
	Positive	False Neg: 14241	True Pos: 140194 (Num Pos: 154435)	True Pos Rate: 0.91
		Neg Pre Val: 0.91	Pos Pred Val: 0.95	Accuracy: 0.93

شکل (۷) ماتریس کانفیوژن با طبقه بند بیزین

CountSum	۹۲,۰	۹۲,۰
PosSum	۹۴,۱	۹۴,۰
نتیجه نهایی	۹۵,۱۷	۹۳,۸۹

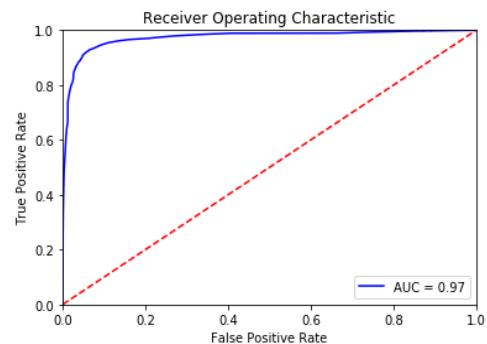
۴-۳-۴ مقایسه طبقه بندی‌ها در روش پیشنهادی

در این بخش به مقایسه طبقه‌بندی‌های مختلف مانند نایوبیزین، ماشین بردار پشتیبان و رگرسیون لجستیک^{۵۳} در روش پیشنهادی پرداخته می‌شود. شکل‌های (۴)، (۵) و (۶) نشان دهنده‌ی نرخ صحت تشخیص نظرات مثبت و منفی می‌باشد.



شکل (۴) منحنی ROC و نمرات AUC از طبقه بندی NB

همانطور که از شکل‌های (۶)، (۷) و (۸) مشخص است، سطح زیر نمودار طبقه‌بند بیزین بیش‌تر بوده و نسبت به دو طبقه‌بند ماشین بردار پشتیبان و رگرسیون عملکرد بهتری دارد. درحالی‌که روش ماشین بردار پشتیبان در رده دوم و روش رگرسیون در رده سوم قرار دارد.



شکل (۵) منحنی ROC و نمرات AUC از ماشین بردار پشتیبان

^{۵۳} Logistic Regression

لجستی				
ک	۸۱,۵۳	۸۲,۱۵	۸۳,۴۶	۸۲,۸۰
رگرسیون				

جدول (۱۱) معیارهای ارزیابی طبقه بندها بعد از استفاده از الگوریتم ژنتیک پیشنهادی

طبقه بند	Accuracy	Precision	Recall	F-measure
بیزین	۹۴,۶۵	۹۵,۱۷	۹۳,۸۹	۹۴,۵۳
SVM	۹۲,۷۹	۹۴,۵۹	۹۰,۷۷	۹۲,۶۴
لجستی				
ک	۸۸,۲۹	۸۴,۴۴	۸۶,۴۶	۸۵,۴۴
رگرسیون				

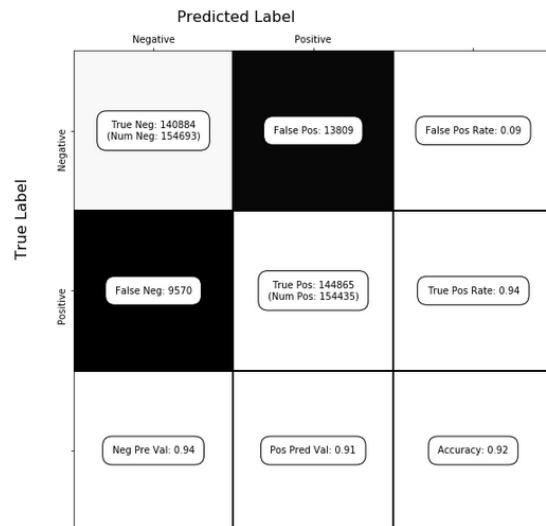
۴-۳-۵ مقایسه واژگان مدل پیشنهادی با روش‌های دیگر

در این قسمت به تحلیل و بررسی و مقایسه معیارهای صحت، دقت و معیار F در مقالات مبتنی بر روش واژگان و روش پیشنهادی با سایر روش‌ها نظیر SentiWordNet [۱۵]، سامها [۱۶]، روش ترکیبی واژه نامه و corpus-based [۱۷] و یادگیری عمیق [۱۸] پرداخته‌ایم.

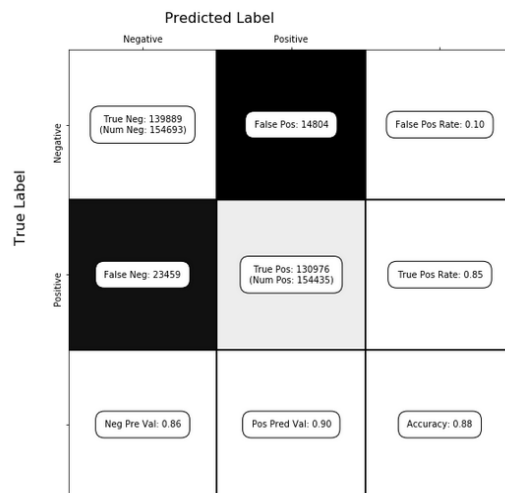
روش پیشنهادی خود را با نتایج چند مقاله که تنها با استفاده از روش محاسبه‌ی واژگان احساسی کارکرده‌اند مقایسه کردیم و نتایج نشان داد روش ترکیبی ما با فرضیات مطرح شده و تعداد ویژگی‌های بیان شده، نسبت به روش این مقالات که بر اساس روش واژگان احساسی هستند، برتری قابل‌وضوحی دارد که در جدول (۱۲) به مقایسه کمی این روش‌ها پرداخته شده است.

جدول (۱۲) مقایسه معیارهای ارزیابی واژگان با روش پیشنهادی

واژگان	Accuracy	Precision	F-measure
SentiWordNet [۲۱]	۷۹,۷۵	۹۰,۱	۷۸,۳
سامها [۲۲]	۵۶	۶۱	۶۰
ترکیب روش واژه‌نامه و corpus-based [۲۳]	۸۵	۷۳	۷۸



شکل (۸) ماتریس کانفیوژن با طبقه بند ماشین بردار پشتیبان



شکل (۹) ماتریس کانفیوژن با طبقه بند رگرسیون

سپس همانطور که در جدول (۱۰) و (۱۱) نشان داده شده، به منظور آزمایش اینکه آیا الگوریتم ژنتیک انتخابی دارای پاسخ بهینه بوده یا خیر، نتایج استفاده از طبقه بندها را قبل و بعد از استفاده از این الگوریتم محاسبه می‌کنیم.

جدول (۱۰) معیارهای ارزیابی طبقه بندها قبل از استفاده از الگوریتم ژنتیک پیشنهادی

طبقه بند	Accuracy	Precision	Recall	F-measure
بیزین	۹۰,۶۵	۹۱,۵۴	۹۲,۷۳	۹۲,۱۳
SVM	۸۹,۳۶	۹۰,۶۹	۸۸,۱۷	۸۹,۴۱

این مقاله بر روی استخراج آرا و عقاید در قالب سند، جمله و تحلیل احساسات در سطح ویژگی بر روی مجموعه داده‌های نظرات کاربران هتل ارایه شد. در روش پیشنهادی، هدف، تعیین گرایش احساسی متن و طبقه‌بندی آن به دسته‌های مثبت و منفی است. در مدل ارایه شده در این مقاله علاوه بر تعیین گرایش کلی احساسی داده‌های متنی مورد آزمایش و تعیین قطبیت کلی مثبت یا منفی، ویژگی‌های مهم بیان شده در آرا مهم کاربران هتل‌ها که بیشتر مدنظرشان بوده نیز استخراج شده است. آگاهی از این ویژگی‌ها، به هتلداران و یا سایر ذینفعان این امکان را می‌دهد که از مهم‌ترین شاخصه‌های مدنظر مشتریان آگاه شوند و از این دانش در راستای سیاست‌گذاری بهتر و در صورت لزوم تغییر رویکرد و در نهایت کسب سود بیشتر استفاده کنند.

در روش پیشنهادی علاوه بر طبقه‌بندی نظرات، مجموع ویژگی‌های بیان شده در نظرات کاربران هتل‌ها برای هر یک از هتل‌ها براساس رتبه‌های مثبت و منفی استخراج شد.

نتایج بدست آمده نشان دهنده‌ی دقت بالای مدل پیشنهادی این مقاله در انتخاب ویژگی‌ها می‌باشد. با تحلیل و بررسی نتایج مشخص شد نظرات متنی کاربرانی که ویژگی‌های شاخص‌تری برای هتل ابراز می‌کنند، تعداد صفات بیشتری برای وصف واژه‌ها و ویژگی‌های مدنظرشان استفاده می‌کنند، نتایج نشان می‌دهد که این روش بسیار بهتر از روش‌های مبتنی بر فرکانس واژه‌ها عمل می‌کند. علاوه بر این با مقایسه این روش با سایر روش‌ها دریافت شد که نتایج نشان دهنده بهبود نتیجه‌ی بدست آمده از طبقه بندی در این پژوهش نسبت به روش‌های پیشین است.

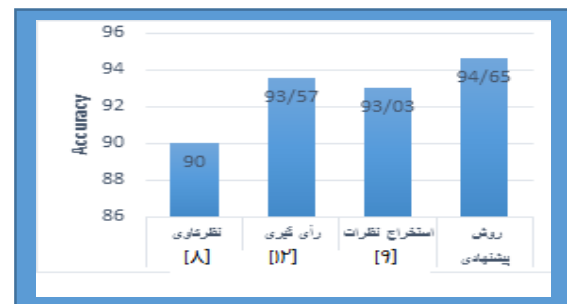
۵-۱- پیشنهادهای آینده

در کارهای آینده می‌توان از روش بیان شده در مجموعه داده‌های بزرگ‌تر و متنوع‌تر و نیز در مورد سایر خدمات و کالاهای مختلف و نیز روش مشابهی را برای داده‌های به زبان فارسی به کاربرد، مدل پیشنهاد شده تنها جنبه‌های صریح مطرح شده در متن نظر را در بررسی می‌کند می‌توان با یافتن جنبه‌های ضمنی مطرح شده مدل را بهبود داد تا تصویر جامعی از کلیه‌ی عقاید مطرح شده را در اختیار متصدیان امور قرار دهیم.

یادگیری عمیق [۲۴]	۸۳	۸۴٫۷	۸۳٫۶
روش پیشنهادی	۹۴٫۶۵	۹۵٫۱۷	۹۴٫۵۳

۴-۳-۶ مقایسه تحلیل نظرات هتل هادرمدل پیشنهادی با روش‌های دیگر

در این بخش به مقایسه معیار صحت بدست آمده از طبقه بندی مدل پیشنهادی با سایر روش‌ها روی مجموعه داده‌های نظرات مربوط به هتل مانند، نظر کاوی [۸] روش رای گیری [۱۸] و روش استخراج نظرات [۹] پرداخته شده است. در شکل (۱۰) این مقایسه نشان داده شده است.



شکل (۱۰) مقایسه تحلیل نظرات هتل روش پیشنهادی با سایر روش‌ها

۵- نتیجه گیری و پیشنهاد برای تحقیقات آتی

عقیده کاوی، تجزیه و تحلیل احساسات نیز نامیده می‌شود که فرآیندی برای کشف عقیده‌ی کاربران درباره موضوع یا محصول یا مسئله‌ی خاصی است. موضوع می‌تواند یک خبر، رویداد، محصول، فیلم، موقعیت هتل، خدمات ارایه شده در هتل و ... باشد. عقیده‌کاوی و تحلیل احساسات موضوعی تحقیقاتی در حوزه متن کاوی، پردازش زبان طبیعی و وب کاوی می‌باشد.

عقیده کاوی و تحلیل احساسات سیستماتیک زیرشاخه‌ای از علم داده کاوی است که برای کسب دانش نهفته و مخفی در مجموعه داده‌های متنی ساختار نیافته با حجم بالا به کار می‌رود. مجموعه داده‌های متنی می‌تواند نظرو عقیده مطرح شده یا کامنت‌های کاربران در شبکه‌های اجتماعی یا وبسایت‌های خاص مرتبط با موضوع مورد نظر باشد، بازخورد و تجربه‌ی مشتری در استفاده از هر محصول یا موضوع خاص یا هر تاپیک دیگری باشد.

hospitality sector using unique attributes and sentiment orientation. *Tourism Management*.

[۷] C. Y. Tsai, M. T. Wang, & H. T. Tseng, "The impact of tour guides' physical attractiveness, sense of humor, and seniority on guide attention and efficiency." *Journal of Travel & Tourism Marketing*, vol. ۳۳, pp ۱۳-۲۰, ۲۰۱۵.

[۸] A. S. Mohammad & M. Al Kadri, "Using Lexicon-Based Opinion Mining to Gauge Customer Satisfaction," *International Journal of Innovative Technology and Exploring Engineering (IJITEE)*, vol. ۹, February ۲۰۲۰.

[۹] A. Ahania, M. Nilashib, O. Ibrahimc, L. Sanzognia & S. Weaven, "Market segmentation and travel choice prediction in Spa hotels through TripAdvisor's online reviews," *International Journal of Hospitality Management*, vol. ۸۰, pp ۷۷-۸۰, July ۲۰۱۹.

[۱۰] T. W. Lui, M. Bartosia, G. Piccoli, & V. Sadhya, "Online review response strategy and its effects on competitive performance," *Tourism Management*, vol. ۶۷, pp ۱۸۰-۱۹۰, August ۲۰۱۸.

[۱۱] M. E. Mowlaei, M. Saniee Abadeh, H. Keshavarz, "Aspect-Based Sentiment Analysis using Adaptive Aspect-Based Lexicons," *Expert System with Applications*, vol. ۱۴۸, June ۲۰۲۰.

[۱۲] L. M. Domingo, J. C. Martín, G. Mandsberg. "Social media as a resource for sentiment analysis of Airport Service Quality (ASQ)," *Journal of Air Transport Management*, vol. ۷۸, pp ۱۰۶-۱۱۵, July ۲۰۱۹.

[۱۳] D. Zhang, J. Tu, L. Zhou & Z. Yu, "Higher tourism specialization, better hotel industry efficiency?," *International Journal of Hospitality Management*, ۲۰۲۰.

[۱۴] T. Chinsha & J. Shibily, "A syntactic approach for aspect based opinion mining," *Proceedings of the ۲۰۱۵ IEEE ۹th International Conference on Semantic Computing*, pp ۸۱-۸۸, March ۲۰۱۵.

در این مقاله جنبه‌های مطرح شده در قالب پنج ویژگی بیان شد، می‌توان با افراد خبره مشورت کرد و جنبه‌های بیشتری را در نظر گرفت، همچنین می‌توان مدل پیشنهادی این مقاله را با سایر مدل‌ها و الگوریتم‌ها ترکیب و نتایج بدست آمده را بررسی کرد.

در نهایت با طراحی شبکه واژگان احساسی جداگانه برای تحلیل هر دسته نظرات در موضوعات مختلف همانند رستوران، نقاط دیدنی، خدمات فرودگاهی و دیگر موارد، تاثیر آن را بر بهبود نتایج بدست آمده بررسی نمود. با توجه به اینکه بار احساسی و معنای ادراک شده از یک واژه در کالاها و خدمات مختلف با یکدیگر متفاوت است، طراحی وساخت فرهنگ واژگان احساسی خاص منظوره برای زمینه های مختلف می‌تواند موجب بهبود نتایج حاصل شده از پژوهش‌های قبلی گردد.

مراجع

[۱] K. Ravi, V. Ravi, "A survey on opinion mining and sentiment analysis: Tasks, approaches and applications," *Knowledge-Based Systems*, vol. ۸۹, pp ۴۶-۱۴, November ۲۰۱۵.

[۲] J. A. Balazs, J. D. Velásquez, "Opinion Mining and Information Fusion: A Survey," *Information Fusion*, vol. ۲۷, pp ۹۵-۱۱۰, ۲۰۱۶.

[۳] ر. پیرمحمدیانی و ش. محمدی "معیارهای ارزیابی ارزش اثرگذاری کاربران رسانه‌های اجتماعی چارچوبی براساس کاوش رسانه‌های اجتماعی"، در دو فصلنامه علمی فناوری اطلاعات و ارتباطات ایران، صفحات ۱۰۹-۱۲۵، بهار و تابستان ۱۳۹۸.

[۴] Z. Zhang, Z. Zhang & Y. Yang, "The power of expert identity: How website recognized expert reviews influence travelers' online rating behavior," *Tourism Management*, vol. ۵۵, pp ۲۴-۱۵, August ۲۰۱۶.

[۵] G. Pablos, M. Cuadros, & M. T. Linaza, "Automatic analysis of textual hotel reviews," *Information Technology & Tourism*, vol. ۱۶(۱), pp ۶۹-۴۵, ۲۰۱۶.

[۶] M.R. Martinez-Torres, S.L. Tora. (۲۰۱۹) A machine learning approach for the identification of the deceptive reviews in the

Cluster Computing, Springe, vol.۲۲, pp۷۱۸۱-۷۱۹۹, August ۲۰۱۹

[۲۴] D. Tang, F. Wei, N. Yang, M. Zhou, T. Liu, & B. Qin, "Learning Sentiment-Specific Word Embedding for Twitter Sentiment Classification," *Learning Sentiment-Specific Word Embedding for Twitter Sentiment Classification*, vol.۱, June ۲۰۱۴.

[۲۵] K. Sailunaz, R. Alhaj, "Emotion and Sentiment Analysis from Twitter Text," *Computational Science*, vol.۳۶, September ۲۰۱۹.

[۲۶] Y. H. Hu, Y. L. Chen, & H. L. Chou, "Opinion mining from online hotel reviews – a text summarization approach," *Information Processing & Management*, vol.۵۳, pp۴۳۶-۴۴۹, March ۲۰۱۷.

[۲۷] J. Li, L. Xu, L. Tang, S. Wang, S., & L. Li, "Big data in tourism research: A literature review," *Tourism Management*, vol.۶۸, pp۳۰۱-۳۲۳, October ۲۰۱۸.

[۲۸] D. Zhang, J. Tu, L. Zhou & Z. Yu, "Higher tourism specialization, better hotel industry efficiency?," *International Journal of Hospitality Management*, ۲۰۲۰.

[۲۹] K. Cheng, J. Li, J. Tang, H. Liu, "Unsupervised Sentiment Analysis with Signed Social Networks," *Proceeding of the Thirty-First AAAI Conference on Artificial Intelligence*, vol.۳۱, pp۳۴۲۹-۳۴۳, ۲۰۱۷.

[۱۶] A. K. Samha, Y. Li & J. Zhang. "Aspect-Based Opinion Extraction from Customerreviews," *Computation and Language*, April ۲۰۱۴.

[۱۷] H. Keshavarz & M. Saniee Abadeh, "Accurate frequency-based lexicon generation for opinion mining," *Journal of Intelligent and Fuzzy System*, September ۲۰۱۷.

[۱۸] H. Keshavarz & M. Saniee Abadeh, "ALGA: Adaptive lexicon learning using genetic algorithm for sentiment analysis of microblogs," *Knowledge Based Systems*, vol.۱۲۲, pp۱-۱۶, April ۲۰۱۷

[۱۹] S. Katoch, S. S. Chauhan & V. Kumar, "A review on genetic algorithm: past, present, and future," *Multimedia Tools and Applications*, vol.۸۰, pp۸۰۹۱-۸۱۲۶, ۲۰۲۱.

[۱۲۰] م. امین طوسی و ه. عزتی "الگوریتم ژنتیک آگاه از بهترین عضو با کاربرد در رنگ‌آمیزی و بعدمتریک گراف،" در دوفصلنامه علمی فناوری اطلاعات و ارتباطات ایران، صفحات ۱۴۳-۱۵۴، بهار و تابستان ۱۳۹۹.

[۲۱] T. Chinsha & J. Shibily, "A syntactic approach for aspect based opinion mining," *Proceedings of the ۲۰۱۵ IEEE 4th International Conference on Semantic Computing*, pp۸۱-۸۸, March ۲۰۱۵.

[۲۲] A. K. Samha, Y. Li & J. Zhang. "Aspect-Based Opinion Extraction from Customerreviews," *Computation and Language*, April ۲۰۱۴.

[۲۳] M. Z. Asghar, A Khan, S. R. Zahra, S. Ahmad & F. M. Kundi, "Aspect-based opinion mining framework using heuristic patterns,"

