



کشف گزارش‌های نقص محصول از متن نظرات آنلاین کاربران

* نرگس نعمتی فرد * محرم منصوری‌زاده * مهدی سخایی‌نیا

* گروه مهندسی کامپیوتر، دانشگاه بوعلی سینا، همدان، ایران

تاریخ پذیرش: ۱۳۹۸/۰۱/۰۹

تاریخ دریافت: ۱۳۹۷/۰۳/۱۰

چکیده

با توسعه وب ۲ و شبکه‌های اجتماعی، مشتریان و کاربران نظرهای خود را درباره‌ی محصولات مختلف با یکدیگر به اشتراک می‌گذارند. این نظرها به عنوان یک منبع ارزشمند، جهت تعیین جایگاه کالا و موفقیت در بازاریابی، می‌تواند مورد استفاده قرار گیرد. استخراج نواقص گزارش شده از میان حجم زیاد نظرهایی که توسط کاربران تولید شده از مشکلات عمده این زمینه تحقیقاتی است. مشتریان و مصرف‌کنندگان با مقایسه محصولات تولیدکنندگان مختلف نقاط قوت و ضعف محصولات را در قالب نظرهای مثبت و منفی بیان می‌نمایند. طبقه‌بندی نظرات بر اساس واژگان حسی مثبت و منفی در متن نظر به اسناد حاوی گزارش نقص و فاقد آن نتیجه درست و دقیقی در پی ندارد. چون گزارش نواقص صرفاً در نظرات منفی صورت نمی‌گیرد. ممکن است که مشتری نسبت به یک کالا حس مثبتی داشته باشد و با این حال در نظر خود یک نقص را گزارش نماید. بنابراین چالش دیگر این زمینه تحقیقاتی طبقه‌بندی درست و دقیق نظرات است. برای حل این مشکلات و چالش‌ها، در این مقاله روشی موثر و کارا برای استخراج نظرهای حاوی گزارش نقص محصول از نظرهای آنلاین کاربران ارائه گردیده است. بدین منظور طبقه‌بند جنگل تصادفی برای تشخیص گزارش نقص و تکنیک بدون ناظر مدل‌سازی موضوعی تخصیص پنهان دیریکله را برای ارائه‌ی خلاصه‌ای از گزارش نقص بکار گرفته شدند. برای تحلیل و ارزیابی روش پیشنهادی از داده‌های وبسایت آمازون استفاده شده است. نتایج نشان داد جنگل تصادفی حتی با تعداد کم داده‌های آموزشی عملکرد قابل قبولی برای کشف گزارش نقص دارد. نتایج و خروجی‌های استخراج شده از اسناد حاوی گزارش نقص، شامل خلاصه‌ی گزارش نقص جهت سهولت در تصمیم‌گیری تولیدکنندگان، یافتن الگوهای وجود گزارش نقص در متن به صورت خودکار و کشف جنبه‌هایی از محصول که بیشترین گزارش نقص مربوط به آنها می‌باشد، نشان‌دهنده توانایی روش تخصیص پنهان دیریکله است.

واژه‌های کلیدی: تشخیص گزارش خرابی، نظر کاوی، تحلیل حسی، تحلیل نظر کاربران، متن کاوی.

۱- مقدمه

گسترش وب^۲، کاربران اینترنتی را در تعامل با یکدیگر و همچنین در تشکیل شبکه‌های اجتماعی برای تولید اطلاعات و انتشار دادگان با حجم زیاد و محتوای مفید بر روی وب توانمند ساخته است. یکی از مهمترین محتوایی که کاربران اینترنتی تولید می‌کنند، اظهار نظر^۱ پیرامون یک موضوع، شی، رویداد یا حتی شخص است. این نقد و بررسی‌ها^۲ عامل مهم و اثرگذاری در فرایند تصمیم‌گیری کاربران اینترنتی در زمینه‌های مختلف است [۱].

اخیراً در زمینه‌ی استخراج ریز اطلاعات از جمله استخراج جنبه‌هایی از محصول که مشتری درباره‌ی آنها نظر خود را بیان کرده است و همچنین نرخ امتیاز دهی، کارها و تحقیقات زیادی صورت گرفته است [۱]. این اطلاعات در تصمیم‌گیری مشتریان هنگام خرید و آگاهی تولیدکننده از حس مشتری نسبت به محصول کمک‌کننده هستند. اما در این میان استخراج گزارش نقص و خرابی کالا به رغم اهمیت فراوان آن، کمتر مورد توجه قرار گرفته است.

در بسیاری از بازخوردهای^۳ مشتریان نسبت به محصول، اطلاعاتی وجود دارد که کشف آنها برای اخذ تصمیمات عملی بسیار مؤثر است. این نوع اطلاعات گزارش‌های نقص محصولات هستند که توسط کاربران فضای مجازی بر اساس تجربه‌ها استفاده از محصول، نوشته می‌شوند. مسلماً شرکت‌های تولیدکننده در فرایند تولید، محصول را از جهات مختلف مورد آزمایش قرار می‌دهند، اما تجربه‌های مشتریان در استفاده از محصول برای تصمیم‌گیری و برنامه‌ریزی مدیران شرکت‌ها از اهمیت بالایی برخوردار است. از طرفی

آگاهی از اینگونه اطلاعات می‌تواند برای مشتریان نیز مفید باشد و از تکرار وقوع خرابی محصول در هنگام استفاده از آن جلوگیری شود. بنابراین استخراج گزارش‌های نقص از متن نظر، یعنی نظر کاوی، هم برای تولیدکننده و هم برای مصرف‌کننده از اهمیت فراوانی برخوردار است [۲].

حجم زیاد نظرات و عدم ساخت‌یافتگی آن، تحقیق در زمینه‌ی نظرکاوی را همواره با مشکلاتی روبرو می‌نماید که رسیدن به دقت بالا در استخراج اطلاعات را دشوار می‌سازد [۳]. پژوهش در خصوص استخراج گزارش‌های نقص^۴ نه تنها از این مشکلات مستثنا نیست، بلکه به دلیل نوع اطلاعات استخراجی با مشکلات دیگری نیز مواجه است. اغلب کاربران برای ابراز حس خود از واژگان حسی^۵ استفاده می‌کنند. بنابراین با تشخیص این لغات کلیدی در متن می‌توان نظر کاربر را استخراج کرد. اما برای کشف گزارش خرابی وجود لغات بیان‌کننده‌ی احساس کافی نیست. گاهی دیده می‌شود که مشتری حس مثبتی نسبت به محصول دارد اما بنابر مشکلی که در حین استفاده از محصول مواجه شده است، گزارشی از نقص محصول را نیز مطرح می‌کند. همچنین در بسیاری از نقد و بررسی‌ها با اینکه مشتری احساسات منفی ابراز کرده است، گزارشی از خرابی و نقص محصول در متن نظر وجود ندارد و فقط سلیقه‌ی شخصی خود را مطرح کرده است. بنابراین صرفاً نمی‌توان گفت گزارش نقص فقط در نقد و بررسی‌ها با قطبیت منفی وجود دارد. به همین دلیل کلاس‌بندی نقد و بررسی‌ها به اسناد حاوی گزارش نقص و فاقد آن، با چالش‌های جدی روبرو است.

¹ Opinion

² Reviews

³ feedback

⁴ defect

⁵ Opinion words

هزینه‌ی زیاد و زمانبر بودن برچسب زنی حجم عظیم نظرها، به عنوان بخشی از فرایند پیش پردازش از تعداد کمی داده‌ی آموزشی جهت کلاس‌بندی اسناد حاوی گزارش نقص استفاده شده است. سپس خروجی‌های مثبت آن را برای خلاصه‌سازی و ارائه اطلاعات کاربردی از اسناد حاوی گزارش نقص، به کمک تکنیک بدون ناظر، تخصیص پنهان دیریکله، به کار برده شده است. نکته حائز اهمیت دیگر اینکه روش پیشنهادی در این مقاله مستقل از دامنه می‌باشد.

ساختار مقاله به این شرح است. پس مقدمه در بخش دوم کارهای مرتبط و مشابه بررسی خواهد گردید و سپس در بخش سوم روش پیشنهادی با تاکید بر انگیزه‌ها و کاربردهای آن ارائه می‌گردد. در بخش چهارم نتایج آزمون و ارزیابی روش پیشنهادی را روی مجموعه دادگان وسیع گزارش شده است و با مقاله با بخش نتیجه‌گیری به پایان خواهد رسید.

۲- کارهای مرتبط

در طول دهه‌ی گذشته، تعداد زیادی از تحقیقات روی نظر-کاوی به صورت عام و همچنین تشخیص و استخراج جنبه متمرکز بوده‌اند که اطلاعات مفیدی هم از متن نظر مشتری‌ها استخراج کردند [۴]. با این حال مطالعات بسیار کمی (تنها یک مورد [۵]) برای استخراج گزارش نقص محصول از متن نظر آنلاین مشتری صورت گرفته است.

از سوی دیگر پایگاه داده‌ای به همراه برچسب حاوی گزارش نقص و غیر آن نیز در دسترس نمی‌باشد، بنابراین استخراج گزارش نقص یک مسئله‌ی جوان است. در ادامه به کارهایی اشاره می‌کنیم که اهداف آنها تا حدودی به مساله مورد تمرکز ما نزدیک است. عمدتاً وجه مشترک این کارها با کار ما این است که آنها نیز به استخراج ریز اطلاعات می‌پردازند.

از اولین تحقیقات انجام شده در این زمینه روش‌های مبتنی بر فرکانس است که به کمک تکنیک‌های متفاوت فیلتر کردن، عبارات اسمی را که فرکانس بالاتری داشته باشند به عنوان جنبه استخراج می‌کنند [۶] [۷] مسئله‌ی تولید خلاصه از نظرات کاربران براساس جنبه‌ی محصول، در [۸] مورد مطالعه قرار گرفت. برای این کار به تشخیص تعداد تکرار جنبه‌ها می‌پردازد و گروه‌های اسمی پر تکرار را به

استخراج گزارش نقص از متن نظرات کاربران یک موضوع جدید است که راه حل‌های خیلی زیادی در این زمینه ارائه نشده است. یکی از راه حل‌هایی که وجود دارد استفاده از کلماتی که به اصطلاح رنگی^۶ گفته می‌شود، است. کلمات رنگی می‌توانند مجموعه‌ای از جنبه‌های یک محصول باشند که ممکن است دچار خرابی و نقص شده‌اند. اما استفاده از این لغات ما را در تشخیص اینکه آیا نظر حاوی گزارش نقص هست یا خیر کمک نمی‌کند، بلکه فقط امکان استخراج نقص در نظری که حاوی گزارش نقص است را فراهم می‌نماید. همچنین این روش وابسته به دامنه است که کاربرد عمومی آن را محدود می‌سازد.

روشی دیگر جهت تشخیص گزارش نقص در متن نظر مشتری استفاده از ایده‌ای است که در تکنیک نظارت از راه دور وجود دارد. تکنیک نظارت از راه دور روشی برای تولید مجموعه‌ی داده‌های آموزشی است. این تکنیک از نشانه‌های تقریبی^۷ به عنوان برچسب‌های مثبت در متن برای آموزش کلاس‌بندی استفاده می‌کند. عباراتی که مشتریان معمولاً برای بیان نقص محصول به کار می‌برند مانند "not allow, not let, no ability, bug, crash, ..." می‌تواند در تشخیص اسناد حاوی گزارش نقص کمک کننده باشد. اما این نشانه‌های تقریبی باعث افزایش پاسخ مثبت کاذب هم می‌شوند.

روش پیشنهادی در این مقاله با استفاده از طبقه‌بند جنگل تصادفی، گزارش‌های نقص را از مجموعه نظرات به صورت خودکار استخراج نموده و خلاصه‌سازی می‌کند. با توجه به

در لغت به معنی دود کننده است. از آنجایی که دود برای Smoky^۶

نشان دادن و اطلاع رسانی یک وضعیت استفاده می‌شود، از

کلمه «رنگی» به جای آن استفاده کردیم

^۷ noisy signals

مدل‌های مختلف خودرو ارائه می‌دهد و نقص‌هایی که در شکایات آمده است را خلاصه و سازمان‌دهی می‌کند. به منظور استخراج نقص‌ها از پایگاه‌داده، چهار موجودیت کلیدی تعریف می‌کند از جمله مدل موتور و سال، اجزای موتور، علائم و تاریخ تصادف که این چهار موجودیت توسط ماژول‌های استخراج موجودیت بدست آورده است. در مدل احتمالی نقص فرض می‌کند شکایت ثبت شده در مجموعه‌ی شکایات از توزیع نقص‌ها، تولید شده است که می‌توان روابط بین آنها را توسط روش احتمالی مولد، مدل کرد. سپس کار بعدی را به صورت یک مدل دامنه‌گرا با تکنیک تخصیص پنهان دیریکله پیشنهاد دادند. برای تعریف یک نقص حوزه‌گرا، مدل تخصیص پنهان دیریکله استاندارد را به تخصیص پنهان دیریکله دوبعدی برای خلاصه کردن نقص‌ها از نظرهای مشتریان توسعه دادند. روش پیشنهاد شده بر مشکل خوشه‌بندی‌های بدون ناظر با استفاده از تعداد زیادی ویژگی مخصوص یک حوزه که در شناسایی نقص شرکت دارند، غلبه می‌کند. این مدل به عنوان یک فرایند قابل توسعه ابتدا اجزای محصول سپس توصیفی از نقص را تولید می‌کند [۱۶]. هردو روش دامنه‌گرا هستند و از طرفی فقط شکایات را در نظر می‌گیرند که همگی حاوی گزارش خرابی هستند و خلاصه‌ای از گزارش خرابی را به صورت ساخت‌یافته ارائه می‌دهند ولی در روشی پیشنهاد ما اسناد با حس و عقیده‌های متفاوت مورد بررسی قرار می‌گیرند. زیرا ممکن است سند با اینکه عقیده‌ی مثبتی درباره‌ی کالا داشته باشد نقصی را نیز گزارش کند و سندی که حاوی حس منفی نسبت به محصول است گزارشی از خرابی محصول مطرح نکرده باشد. هدف ما این است که اگر سندی حاوی گزارش نقص است بتوان آن را کشف کرد. از سوی دیگر کار ما وابسته به دامنه نیست هرچند کارهای دامنه‌گرا نتایج دقیق‌تر خواهند داشت.

در [۱۷] یک سیستم پیشنهاددهنده با هدف استخراج پیشنهادهای مشتریان برای بهبود محصول توسط طراحی شده است؛ با این بینش که پیشنهادات با استفاده از کلمه‌های "wishes" یا "regret" در نظرات کاربران ظاهر می‌شوند. بنابراین تشخیص پیشنهاد، متکی به الگوهای نحوی- معنایی جهت بدست آوردن اینگونه عبارات است.

عنوان جنبه در نظر می‌گیرد. نقطه قوت این روش‌ها این است که در عین سادگی بسیار موثر عمل می‌کنند. اما نقطه‌ی ضعف آنها در تولید تعداد زیادی غیرجنبه است و جنبه‌هایی که تعداد کمتری تکرار شده‌اند، از دست می‌روند. همچنین نیاز به تنظیمات دستی دارد که برای هر پایگاه داده‌ای متفاوت خواهد بود.

کارهایی که اخیراً انجام شده از تکنیک‌هایی بر پایه‌ی مدل‌سازی استفاده شده است. بعضی روش‌های با ناظر مدل‌سازی آماری مانند (HMM, CRF) هستند [۹] [۱۰] و بعضی تکنیک‌های بدون ناظر مدل‌سازی موضوعی مانند تخصیص پنهان دیریکله (LDA) می‌باشند که جنبه‌های محصول و نرخ آنها را استخراج می‌کنند و اطلاعات مفیدی در اختیار مشتریان هنگام تصمیم‌گیری خرید محصول قرار می‌دهند [۱۱] [۱۲] [۱۳]. مدل تخصیص پنهان دیریکله که در [۱۴]. مطرح شد یک ابزار مفید برای خلاصه‌سازی متن است. خانم مقدم و همکاران نیز به کشف گزارش نقص محصول از متن نظرات مشتریان آنلاین پرداخته‌اند که جهت خلاصه‌سازی گزارش‌های نقص از LDA با مجموعه ویژگی‌های کیسه واژگانی^۸، اسم‌ها، فعل‌ها، عبارات اسمی، عبارات فعلی و دو-گرمی استفاده کرده‌اند. خانم مقدم تکنیک نظارت از راه دور را برای ساخت داده‌های آموزشی به کار گرفت و برای این‌کار الگوهایی به صورت دستی به عنوان نشانه‌هایی از وجود گزارش نقص در متن نظرها، استخراج کرد او کار خود را روی بازخوردهای " eBay App Reviews " انجام داده است [۱۴].

یک روش احتمالی برای تشخیص نقص از شکایات مردم در [۱۵] پیشنهاد شده است که هدفش فرموله کردن شکایات است. این کار وابسته به دامنه است و کارش را درباره‌ی

⁸ Bag of Words

۳- روش پیشنهادی

۳-۱- انگیزش

از آنجایی که حجم اسناد تولید شده بر بستر وب بسیار عظیم و به صورت پویا در حال رشد است، جنگل تصادفی را جهت تشخیص اسناد حاوی گزارش نقص در نظر گرفتیم. جنگل تصادفی روی داده‌های بسیار بزرگ قابل اجرا است و می‌تواند هزاران متغیر را بدون حذف آنها مدیریت نماید. از سوی دیگر جنگل تصادفی یک طبقه‌بند با ناظر است که نیاز به داده‌های برچسب خورده دارد.

جهت خلاصه سازی و ارائه اطلاعات کاربردی از اسناد حاوی گزارش نقص، تخصیص پنهان دیریکله (LDA) را که یک روش مدل‌سازی موضوعی است استفاده کردیم. این روش بدون ناظر است و نیازی به داده‌های برچسب خورده ندارد. با توجه به اینکه اسناد تحت بررسی گزارش‌های نقص هستند؛ انتظار داریم LDA در موضوع‌بندی این اسناد نوع نقص گزارش شده در آنها را به عنوان متغیر پنهان در نظر گرفته و اسناد را بر اساس آن دسته‌بندی نماید.

۳-۲- نمادها، مفاهیم و اصطلاحات

مجموعه $P = \{P_1, P_2, P_3, \dots, P_M\}$ شامل محصولات که توسط شرکت‌ها تولید می‌شوند. $R_p = \{r_1, r_2, r_3, \dots, r_n\}$ نیز برای هر محصول، مجموعه‌ای از نقد و بررسی‌هایی است که توسط مشتریان در بستر وب قرار گرفته است. در مجموعه R برخی از نقد و بررسی‌ها حاوی گزارش نقص محصول هستند که این نظرها را با $D(\text{Defect})$ و سایر بازخوردها را با $O(\text{Others})$ نشان می‌دهیم.

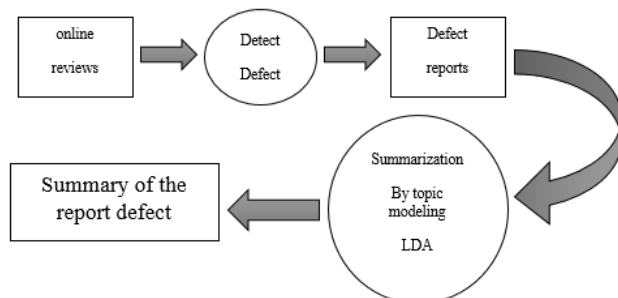
سند: منظور از سند در این پژوهش متن کامل یک نظر است. مثلا شکل (۲) و همچنین شکل (۳) نمونه‌هایی از نظرهای موجود در مجموعه تحت بررسی هستند. یک نظر می‌تواند کوتاه (در حد یکی دو جمله) یا بلند (در حد ده‌ها جمله) باشد. نظرات بلند معمولا حاوی اطلاعات متنوعی مانند تجربه خرید، معرفی محصولات مشابه، بیان نقاط قوت و ضعف محصول و حتی گاهی مطالب بی‌ربط به محصول مورد بحث هستند.

گزارش نقص: بازخوردهایی که به طور واضح به سختی در استفاده، خطا، اشکال و ناتوانی محصول اشاره دارند به عبارت دیگر مشتری جنبه‌هایی از محصول را که درست کار

نمی‌کنند یا نیاز به ترمیم دارند گزارش می‌کند. گزارش نقص معمولا در قالب چند جمله متوالی یک نظر ارائه می‌شود. مثلا در شکل (۱) جملات و عبارات *"It's heavy, hard to push, a 1-day battery life, freeze up and crashes all the time"* نقص موجود در کتابخوان الکترونیکی مد نظر را گزارش می‌کنند.

I am still waiting for the perfect ebook reader. I bought the Nook for these reasons: 1) It reads industry-standard ePub-format ebooks, 2) it's tightly integrated with the Barnes & Noble ebook store, 3) the ebooks are encrypted in a well-documented easily-understood format that is portable across multiple devices so they can be decrypted and read in, say, your iPad's Nook reader software, or even in a Sony Reader (with the very latest firmware), without having to be re-purchased because of DRM nonsense. The problem is that the Nook simply doesn't live up to its promise. The "paper-white" display is more an off-beige, and reflects light in a way that makes it hard to read with a reading light (necessary because it has no backlight, as is true for all ePaper devices). It's heavy and the buttons to change pages are hard to push, especially with gloved hands as you might have while reading outdoors on a cool day. The "5 day battery life" in reality for me has been a 1-day battery life, read a book, and it needs to be recharged, and be darn sure to turn it off. The thing freezes up and crashes all the time even with the very latest software, and is excruciatingly slow even with the very latest software. The latest software added classifications for the ebooks so you could sort them into pseudo-folders, which is necessary given how excruciatingly slow the Nook is to scroll through its book list (get about 50 books on the list and you're in for major pain), but the clunky way they implemented this makes those of us who've gotten used to modern user interfaces frown and shake our heads. Sad to say, I really can't recommend any current eBook reader. Either they're too clumsy to use (Nook), have no books available for them (Sony), or have a proprietary eBook format that locks you into a single vendor (Kindle). I'm seriously considering buying an iPad, yes, it will only work for 9 hours or so on a battery charge, but that's true of the Nook too in real actual use and the iPad is usable for a lot of other things too. It's just disappointing that I can't get an ePaper-based reader that meets my criteria (non-proprietary ebook format, long battery life, compact, decent user interface), and instead have things either crippled by bad design decisions or crippled by having a proprietary ebook format that locks you into a single vendor. Well, I don't like

استفاده گردید. مرحله‌ای که برای ارائه‌ی گزارش نقص باید انجام شود در شکل (۲) آمده است.



شکل ۱- گام‌های اصلی کشف گزارش نقص

۴- آزمون و ارزیابی

۴-۱- معرفی روش آزمون

دقت^۹، بازنمایی^{۱۰} و اندازه F^{۱۱} معیارهای کاربردی در حوزه بازیابی اطلاعات هستند که میزان تناسب اسناد بازیابی شده توسط سیستم را با نیاز کاربر تعیین می‌کنند. این سه معیار به صورت زیر تعریف می‌شوند.

$Precision = \frac{Q_{related}}{Q_{retrieved}}$	رابطه (۱)
$Recall = \frac{Q_{related}}{N_{related}}$	رابطه (۲)
$F_Score = 2 \frac{Precision \times Recall}{Precision + Recall}$	رابطه (۳)

crippled, so I'll look elsewhere, thank you very much...

شکل ۱- نمونه یک نظر درباره یک کتابخوان الکترونیکی

موضوع: منظور از موضوع نظر نوع نقص گزارش شده در آن است. مثلاً سختی یافتن یک منو یا کلید خاص در واسط کاربری برنامه یک نوع نقص است و کاهش سریع باتری نقصی از نوع دیگر است. گزارش نقص‌هایی که مشابه دارند در یک دسته قرار می‌گیرند. این دسته بندی کار مطالعه و ارزیابی نظرات کاربران را خیلی آسان می‌سازد. مثلاً اگر تعداد قابل توجهی از نظرات کاربران در ذیل موضوع «کاهش سریع توان باتری» قرار بگیرند، می‌توان نتیجه گرفت که مساله محصول مورد بحث جدی است.

۳-۳- تشریح روش پیشنهادی

کلاس‌بندی اسناد به گزارش نقص و سایر با روش جنگل تصادفی و داده‌هایی که به صورت دستی برچسب زده شده، با مجموعه ویژگی کیسه واژگانی انجام گردید. جنگل تصادفی یک روش باناظر است که برای طبقه‌بندی دو کلاسی عملکرد خوبی دارد [۱۸]. اما استخراج گزارش نقص از متن نظر کاربران اینترنتی نوعاً مسئله‌ای است که نمی‌توان آن را به صورت باناظر حل کرد، به دلیل اینکه حجم نظرها بسیار زیاد است و برای تکنیک‌های باناظر برچسب زنی این حجم عظیمی از متن‌ها زمان‌گیر، هزینه‌بر و مستعد خطاست. به جهت بهره‌گیری از مزایای کلاس بندی باناظر، در این پژوهش از تعداد داده‌های آموزشی کمی جهت کلاس‌بندی استفاده شد و در واقع یک کلاس‌بندی ضعیف روی اسناد انجام گردید. بعد از کلاس‌بندی اسناد و مشخص شدن اسناد حاوی گزارش نقص توسط جنگل تصادفی از بین حجم عظیمی از اسناد، مهمترین مرحله نحوه‌ی ارائه‌ی گزارش نقص محصول است، زیرا مطالعه‌ی کل اسناد حاوی گزارش نقص که اطلاعات اضافی دیگری نیز دارند برای مدیران و تولیدکنندگان خسته‌کننده و زمان‌بر است. هدف، ارائه‌ی خلاصه‌ای کاربردی از گزارش‌های نقص می‌باشد. تکنیک تخصیص پنهان دیریکله (LDA) با مجموعه ویژگی دوگرمی را به منظور خلاصه‌سازی بعد از کشف اسناد حاوی گزارش نقص

⁹ Precision

¹⁰ Recall

¹¹ F-Score/ Measure

تصادفی، β ، معادل با نسبت تعداد اسنادی که درست طبقه‌بندی شدند؛ بدست آمد. مشابه قبل، می‌توان این دقت را بر روی کل پیکره هم تعمیم داد. یعنی اگر جنگل تصادفی M سند از کل پیکره را گزارش نقص تشخیص دهد؛ می‌توان انتظار داشت که حدود $M \cdot \beta$ تای آنها واقعا گزارش نقص باشند. با توجه به مقادیر α و β ، معیارهای دقت و بازنمایی و اندازه F محاسبه گردیده است.

یکی از روش‌های ارزیابی و آزمون موضوع‌های بدست آمده از LDA استفاده از داوری خبرگان است. اگر چه به نظر می‌رسد که این یک فرض قوی است که فضای پنهان و نامعلومی که توسط مدل‌سازی موضوعی پیدا شده است معنادار و مفید باشد ولی ارزیابی هر یک از این فرض‌ها کار دشواری است. زیرا پیدا کردن موضوعات یک فرآیند بسیار دشوار و هزینه بر است. یعنی یک لیست استاندارد کامل از موضوعات برای هر متنی وجود ندارد. بنابراین با مطالعه‌ی اسناد در هر گروه ارزیابی صورت خواهد گرفت.

مدل‌سازی موضوعی به این نحو است که واژه‌ها و ترکیب‌هایی به عنوان موضوع اسناد استخراج می‌شوند. سپس خبرگان این واژه‌ها و عبارت‌ها را در قالب موضوعات قابل فهم و استنباط پالایش و معرفی می‌کنند. پس از استخراج واژگان و عبارت اولیه، در هر موضوع واژه یا عبارتی را که با سایر واژه‌ها پیوستگی معنایی ندارد حذف نموده و موضوعات نهایی با توجه به سایر واژه‌ها مشخص می‌گردد. سپس به جهت تشخیص مرتبط بودن یا نبودن اسناد به موضوع، اسناد بررسی می‌گردند.

۴-۲- معرفی دادگان و ابزارها

تعداد ۲۰ هزار نظر مشتریان درباره‌ی محصولات الکتریکی از سایت آمازون گرفته شده است. برچسب زنی اسناد به صورت دستی و تحت نظر خبره صورت گرفته است و به اسناد برچسب حاوی گزارش نقص (D) و سایر (O) زده شده است.

پیش‌پردازش و نرمال‌سازی متن به کمک ابزار متن پردازشی سنتی GnuWin32 انجام گردید. این فرایند شامل حذف کلمات بی‌اثر، حذف کارکترهای غیر الفبایی انگلیسی مانند #، \$ و ...، ریشه‌یابی و کوچک کردن تمام

در این روابط N تعداد کل اسناد پیکره، $N_{related}$ تعداد اسناد مرتبط با پرس و جوی خاص، $Q_{retrieved}$ تعداد اسناد بازیابی شده برای این پرس و جو و $Q_{related}$ تعداد اسناد بازیابی شده مرتبط با پرس و جو هستند. دو معیار دقت و بازنمایی معمولاً در حالت تقابل با همدیگر هستند و افزایش یکی سبب کاهش دیگری می‌شود. از این رو مقایسه دو سیستم بازیابی بر اساس اندازه F که ترکیبی از هر دوی این معیارهاست؛ انجام می‌گیرد. لازمه‌ی محاسبه دقت و بازنمایی داشتن داده‌های برچسب خورده هست و از طرفی کار با حجم عظیمی از اسناد هزینه‌ی برچسب‌زنی بالایی را می‌طلبد. از آنجایی که در این پژوهش برای کلاس‌بندی از تعداد داده‌های آموزشی کم استفاده شده است و روش خلاصه‌سازی یک تکنیک بدون ناظر می‌باشد، از اینرو نیازی به برچسب زنی کل اسناد نیست. برای ارزیابی نتایج کلاس-بندی، ۱۰ درصد از کل اسناد به صورت تصادفی انتخاب و برچسب زنی شدند. بر اساس نسبت تعداد اسناد حاوی گزارش نقص در این مجموعه، تعداد اسناد حاوی گزارش نقص در کل پیکره بر اساس روابط (۴) و (۵) تخمین زده شد.

$N_d = \alpha N_{total}$	رابطه (۴)
$\alpha = \frac{N_{defect}}{N_{random}}$	رابطه (۵)

در روابط فوق N_{random} تعداد اسنادی است که به صورت دستی برچسب خورده‌اند. از این میان N_{defect} تعداد اسنادی است که برچسب مثبت دارند. با توجه به انتخاب تصادفی و یکنواخت اسناد برای برچسب زنی دستی، می‌توان این نسبت را به کل پیکره تعمیم داد. از این رو برآورد می‌شود تعداد N_d سند در پیکره حاوی گزارش نقص باشند. در محاسبات و تصمیم‌گیری‌های بعدی این تعداد لحاظ شده است. همچنین برای ارزیابی دقت جنگل تصادفی ده درصد از اسناد پیکره به صورت تصادفی انتخاب و برچسب زنی شدند. سپس این مجموعه با جنگل تصادفی به دو کلاس گزارش نقص و سایر طبقه‌بندی شدند. نهایتاً دقت جنگل

I've had my NST since last July, and I've been very happy with it. Pros:- Battery life when wifi is off is as good as advertised.- The page turn rate was already fast, but the update last November made it super speedy. Seriously, I've played with the current e-ink Kindles, and the difference in refresh rate might SEEM small on paper, but in practice it's a very noticeable difference.- I never had the wifi problems others did after the November firmware update, but I understand that the latest firmware update should solve it.- The NST feels really nice when you're holding it, and I really like that there are physical page turn buttons in addition to the onscreen touch turning - I use both, depending on whether I'm sitting up or lying down when I read, etc.- Navigation is very intuitive.- I bought a \$5 4gb micro SD card and have had no problems using it with the NST.- Sideloading non-DRMed, non-B&N; content downloaded from places like Project Gutenberg is very easy, and it all goes into the same library as your B&N; downloaded content. Cons:- The user interface is easy to use, but it's also VERY basic, and there are very few features. You can organize your books into "shelves," but only on the device itself, and it's a cumbersome process. You can sort books by title, author, and date added, and... that's pretty much it. The NST isn't a tablet and I don't want it to be, but there are still some pretty simple features that don't seem like it would have been that hard to add, and more flexible organization is one of them.- Once a book's in the library, you can't really see anything about it besides the title and author. Metadata from sideloaded content doesn't show up, and even purchased content requires you to be on wifi to see more information. That's about it. My cons list is really more of a wishlist, and I wouldn't hesitate to recommend the NST to anyone, especially at the lowered \$99 price. I imagine a new touch reader will probably be released within the next few months, but unless the user interface is majorly updated, it's hard to imagine I'd feel the need to upgrade, since the current version meets almost all of my needs perfectly well.

شکل ۳- دیدگاه یک مشتری درباره یک محصول الکترونیکی

روش مقدم و همکاران [۵] به دلیل استفاده از نشانه‌های نویری نتایج مثبت کاذب به تعداد زیادی رخ داده است، اما این روش اکثر اسناد حاوی گزارش نقص را یافته است. این

کردن تمام حروف می‌باشد. سپس واژه‌نامه‌ای^{۱۲} از واژه‌های اسناد ایجاد گردید. به این دلیل این واژه‌نامه را ایجاد شد که امروزه کاربران اینترنتی در نوشته‌های خود از واژه‌هایی استفاده می‌کنند که ممکن است حتی در لغت‌نامه‌های مفصل نیز موجود نباشند، مانند mer30، goooooood، 5-star و ... در نهایت جهت خلاصه‌سازی توسط مدل‌سازی موضوعی (LDA) از کتابخانه‌ی gensim^{۱۳} در زبان python استفاده گردید.

۳-۴- نتایج آزمایش‌ها

۳-۴-۱- ارزیابی تشخیص اسناد حاوی گزارش نقص

مقادیر معیارهای دقت، بازنمایی و اندازه F برای کلاس‌بندی اسناد در جدول (۱) آمده است.

جدول ۱- نتایج ارزیابی جنگل تصادفی

روش	معیار F	دقت	بازنمایی
نظارت از راه دور	۰.۵۶	۰.۴۰	۰.۹۱
جنگل تصادفی	۰.۵۴	۰.۷۲	۰.۴۳

ردیف اول این جدول نتایج کار مقدم و همکاران است که با روش نظارت از دور خرابی‌های گزارش شده را استخراج می‌نماید [۵]. در این روش از ۵۰ هزار نظر برچسب خورده (به صورت دستی) استفاده شده است. به دلیل متفاوت بودن اندازه و محتوای مجموعه دادگان نمی‌توان بین این دو مطالعه مقایسه دقیق و قاطعی انجام داد اما می‌توان برخی از جوانب قدرت و ضعف هر دو روش را برشمرد.

¹² Vocabulary

¹³ <https://pypi.python.org/pypi/lda>

بخشی از دادگان را به خوبی طبقه‌بندی نماید. تجمیع نتایج این درخت‌ها در یک قالب ساده اما کارآمد به تولید نتایج طبقه‌بندی دقیقی و کارآمدی منجر می‌شود. اسناد متنی اعم از نظر و غیر آن معمولاً در قالب بردارهای واژگانی بیان می‌شود. از این رو برای هر سند تعداد زیادی ویژگی استخراج می‌شود که طبقه‌بندی‌هایی مثل جنگل تصادفی می‌توانند از آن بهره ببرند. نتایج آزمایش‌های ما نیز نشان می‌دهد جنگل تصادفی تقریباً نیمی از اسناد حاوی گزارش نقص را بازیابی کرده است. کلاس‌بندی جنگل تصادفی با تعداد داده‌های آموزشی کم در کشف گزارش نقص نیز نتیجه‌ی مطلوبی دارد و با وجود تعداد داده‌های آموزشی کم توانسته دقت بالایی در عملکرد خود داشته باشد. بررسی اسنادی که به اشتباه کلاس‌بندی شده‌اند (مانند شکل (۳)) نشان داد این اسناد با وجود اینکه مشتری به نقص کالا اشاره کرده‌است، به محصول علاقه زیادی داشته و از لغاتی که حس مثبت را ابراز می‌کنند استفاده زیادی کرده است. در جدول (۲) تعدادی از جملات این نظر همراه با برچسب‌شان آمده است. می‌توان دید که مشتری علاقه زیادی به محصول داشته و بسیار از آن تمجید کرده است. همچنین در ضمن این تعریف و تمجید نواقصی را هم گزارش نموده است.

جدول ۲- نمونه‌هایی از جملات نظر شکل، حاوی نظر

مثبت و گزارش نقص

گروه	موضوع	مجموعه واژه‌های مرتبط
اول	برنامه‌های نرم‌افزاری	Card, software, app, memory tablet, install, window, android, format, driver, download, comput
دوم	بازگشت محصول	Back, return, got, bought, amazon, didn
سوم	رسانه‌ی ذخیره‌سازی	Tape, disk, record, clean, drive, ver, casset, floppy, cleaner, maxel
چهارم	باتری و شارژر	Batteri, charg, charger, usb, plug, power, recharge, garmin, adapt, Connect, cord, port, fit
پنجم	دستگاه پخش کننده موسیقی و فیلم	Player, dvd, soni, disc, year, rio, skip, mp3, panason, movi, repair
ششم	پخش کننده صدا	Radio, sound, speaker, good, headphone, better, like, even, much Look, volum
هفتم	پانل نگهدارنده‌ی تلویزیون	Case, lock, mount, palm, cover, plastic
هشتم	کتاب خوان الکترونیکی	Nook, book, kindl, read, purchas, barn, custom, nobl, screen, Service

روش جزو دسته عمومی روش‌های بازنمایی زیاد^{۱۴} قرار می‌گیرد. گرچه کشف همه نواقص گزارش شده اهمیت زیادی دارد، اما تعداد زیاد نمونه‌های مثبت کاذب سبب می‌شود که مدل‌سازی موضوعی LDA موضوعات حاشیه‌ای و پس زمینه متعددی تولید کند [۱۳] برای درک بهتر این رویداد تصور کنید که همه نظرات اعم از گزارش نقص و سایرین را تحلیل موضوعی نماییم. با توجه به اینکه تنوع محصولات و اشیا مورد بحث در این نظرها، تفکیک موضوعی به تفکیک نظرات بر اساس نوع محصول مورد بحث متمایل خواهد شد. یعنی بیش از اینکه وجود یا عدم وجود گزارش نقص در یک نظر معیار تخصیص آن نظر به موضوع خاصی باشد، نوع و مدل محصول سبب خواهد شد تا نظرات در گروه‌های مختلف قرار بگیرند. ما به عنوان مطالعه اولیه چنین فعالیتی را انجام دادیم و دریافتیم که استفاده از موضوع‌بندی بدون توجه به ماهیت نظر گمراه کننده است.

در مقابل روش‌هایی که هدفشان بازنمایی بالاست، روش‌های دارای دقت بالا^{۱۵} بر این حقیقت تاکید دارند که نظراتی که تفکیک و طبقه‌بندی موضوعی می‌شوند اصولاً جزو گروه هدف (در اینجا گزارش نقص) باشند. بر اساس این فرض می‌توان امیدوار بود که تفکیک موضوعی انجام شده روی این نظرات بیانگر تفکیک انواع نقص‌های گزارش شده می‌باشد.

جنگل تصادفی به صورت ذاتی دادگان دارای تنوع زیاد را به خوبی دسته‌بندی می‌کند. دلیل این توانایی این است که این روش با انتخاب تصادفی زیرمجموعه‌های متنوع از ویژگی‌ها، تعداد زیادی درخت می‌سازد که هر یک از آنها می‌تواند

¹⁴ High Recall

¹⁵ High Precision

(cd-player, sound-quality) نیز جزء جنبه‌هایی هستند که زیاد مورد نقد و بررسی مشتریان قرار گرفته‌اند. ترکیب پر تکرار customer service گویای این است که افراد به دلیل وجود نقص به خدمات پس از فروش مراجعه کرده‌اند.

جدول (۴) سه گروه از واژگان قابل توجه در موضوعات را نشان می‌دهد. گروه اول الگوهایی هستند که مستقیماً به نوع نقص اشاره نمی‌کنند اما از حضور آنها می‌توان به وجود گزارش نقص در یک نظر پی برد. این واژگان را نیز می‌توان به عنوان نشانگرهای نویری نقص و جزو موضوع پس‌زمینه برشمرد. وجود این نشانگرها در اغلب موضوعات حاکی از عملکرد مناسب جنگل تصادفی در کلاس‌بندی و استخراج گزارش‌های نقص است.

جدول ۴- نشانگرهای گزارش نقص، نوع نقص و جنبه هدف بدست آمده از تحلیلی موضوعی

واژگان عضو	گروه
a lot-defect, another-one, another-problem, biggest-complaint, big-wast, bung-buck, buy-another, buyer-beware, cant-use, cheap-feel, cheaply-made, common-problem, didn't-expect, disappoint-experience, don't-buy, don't-recommend, dosen't-allow, explain-problem, first-time, frustrate-try, main-problem, money-back, never-able, notic-problem, piece-junk, return product, send-back, s-shame, try-get, try-use, wast-money, wast-time	نشانگر گزارش نقص
slow-type, low-power, crash-often, adapt-fail, soft-reset, permanent-damage, defect-camera, background-hiss, sort-problem, player-stop, difficult-remove, background-noise, sound-horrible, poor-sound, drain-battery, bad-connect, start-freez, radio-faulty, get-hot, badly-written, start-skip, read-bad, give-break, bad-patch, stop-play, really-slow, got-stuck, battery-diy, stop-work, screen-freez, hit-pause, brock-first, whip-antenna, short-antenna, horizont-line, camera-eat, , simply-stop, lot-noise, come-dark	نشانگر نوع نقص

تحلیل موضوعی گزارش‌های نقص با بردارهای دو گرمی قابلیت توصیف و تفسیر بهتری دارد. زیرا اغلب نواقص و نارضایتی‌ها با ترکیب‌ها و اصطلاحات دو کلمه‌ای بیان شده‌اند. مثلاً افعال منفی معمولاً با پیشوندهای don't و can't استفاده شده‌اند. همچنین صفات ساده‌ای مثل bad, worst, low که بیانگر دیدگاه منفی هستند پیش از نام محصول یا جنبه‌ای از آن آمده‌اند.

بعضی از واژگان مختص یک یا دو موضوع هستند. این واژگان کم تکرار معمولاً عنوان یک محصول یا جنبه خاصی از آن هستند. در مقابل، برخی از واژگان پر تکرار هستند. یعنی در سه موضوع یا بیشتر ظاهر شده‌اند. این گروه نشان‌دهنده نوع نقص یا نارضایتی هستند.

پرتکرارترین واژگان افعال کمکی منفی هستند. استفاده از افعال کمکی منفی در بیان نقص و خرابی محصول در زبان انگلیسی رایج و مطابق دستور زبان است. البته لغات isn't و aren't خیلی متمایزکننده نیستند، چون در اکثر جملات وجود دارند. این واژگان به نوعی بیانگر موضوع عمومی موسوم به موضوع پس‌زمینه^{۱۶} هستند. موضوع پس‌زمینه حوزه عمومی مورد بحث همه اسناد تحت بررسی را نشان می‌دهد. به عنوان یک نتیجه جنبه‌ای می‌توان گفت که گزارش نقص محور نظرهای تحت بررسی بوده است. ظهور barn_noble به عنوان واژه پرتکرار نشان می‌دهد که اکثر شکایات از شرکت barn&noble بوده است. همچنین پرتکرارهایی مثل battery power و battery از بین جنبه‌های گوناگون یک محصول بر این حقیقت تاکید می‌کنند که اکثر مشتریان عمر باتری محصول را مورد نقد و بررسی قرار داده‌اند. به صورت مشابه، واژگانی مثل

¹⁶ Background topic

Never-buy: به دلیل عدم رضایت، مشتری اعلام می‌کند که دیگر از این نوع محصول یا محصولات شرکتی هرگز خریداری نکند.

Didn't-expec: مشتری علیرغم تبلیغی که برای محصول شده انتظار چنین نقصی را ندارد.

گروه دوم واژگان به صورت مستقیم نوع نقص را نشان می‌دهند. شاید این گروه را بتوان مهمترین دست‌آورد تحلیل موضوعی قلمداد کرد. در جدول (۴) گروه با عنوان نشانگرهای نقص مشخص شده‌اند.

گروه سوم واژگان جنبه‌هایی از محصولات را نشان می‌دهند که بیشتر مورد بحث بوده‌اند. این گروه از دو جهت اهمیت ویژه دارند. نخست اینکه نشان می‌دهد کاربران در مقایسه و گزینش محصولات به چه ویژگی‌هایی توجه دارند. مثلاً با اینکه اغلب گوشی‌های تلفن همراه و ادوات الکترونیکی پوشیدنی ضدآب یا ضد ضربه نیستند اما این جنبه‌ها کمتر مورد توجه بوده‌اند اما در مقابل باتری، صفحه نمایش و کابل شارژ در مرکز توجه قرار داشته‌اند. دلیل دوم اهمیت این جنبه‌ها این است که با کاوش قواعد انجمنی^{۱۷} بین نشانگرهای نقص و این جنبه‌ها می‌توان نقاط قوت و ضعف محصولات مختلف را به صورت خودکار استخراج و دسته‌بندی کرد.

۵- دسته‌بندی

اخیراً تمرکز پژوهشگران حوزه نظرکاوی بیشتر روی استخراج جنبه‌های محصول و تخمین امتیاز آنها از بازخوردها می‌باشد. گرچه استخراج جنبه و تخمین امتیاز آن می‌تواند به مشتریان جهت تصمیم‌گیری در خرید

Battery-compartment, diamond-rio, extern-antena, flash-card, graphic-card, lcd-screen, memori-card, page-turn, nook-tablet, usb-cable, cell-phone, floppy-disk, image-quality, mp3-s, phone-jack, power-cord, dvd-player, recharge-battery, e-reader, mp3-player, nook-color, touch-screan

جنبه هدف

واژگان گروه اول به نوعی بیانگر عقیده، دیدگاه، تصمیم و توصیه نظردهنده هم می‌باشند. برای روشن شدن این مطلب مفهوم تعدادی از آنها توضیح داده شده است.

Don't-buy: مشتری که از محصولی راضی نباشد، حال به دلیل نقصی که دارد یا اینکه عقیده شخصی وی نسبت به محصول منفی باشد دیگران را نیز از خرید محصول منصرف می‌کند.

buy-another, another-one: مشتری محصولی را خریداری کرده است به دلیل نقصی که داشته مجبور شده است یکی دیگر تهیه کند. در برخی اسناد مشتری از محصول دوم راضی است اما در برخی دیگر به نقص مشابه قبلی اشاره کرده است.

Cant-use: به دلیل وجود نقص آن‌گونه که باید از محصول نتوانسته‌اند استفاده کنند مثلاً محصولی علیرغم تبلیغی که کرده است با دستگاه خاصی سازگاری ندارد و مشتری نتوانسته از محصول بهره‌ی کامل ببرد.

First-time: برخی مشتریان از وجود نقص و خرابی محصول در ابتدای استفاده و یا دریافت محصول شکایت کردند.

Send-back, return product: خیلی از مشتریان محصولی که نقص دارد را بازپس می‌دهند.

Try-use, try-get: این دو-گرمی نیز نشانه‌ادی از وجود نقص است. زیرا اسنادی که این لغات را دارند به این نکته اشاره می‌کنند که مشتری با روش‌های مختلف سعی در استفاده از محصول داشته است اما به دلیل نقصی که دارد موفق نشده است.

Wast-money, wast-time: وقتی مشتری از محصول راضی نیست صرف زمان و هزینه برای آن را هدر رفت می‌داند.

¹⁷ Association rules

دو-گرمی خلاصه‌ای از گزارش‌های نقص پرتکرار و اطلاعاتی نظیر اینکه بیشتر کدام جنبه‌های محصول مورد نقد و بررسی هستند را ارائه داد. همچنین روش پیشنهادی توانست به طور خودکار کشف الگو داشته باشد. فهرست واژگان تشکیل‌دهنده موضوع پس زمینه یکبار دیگر بر موفقیت جنگل تصادفی در تشخیص گزارش نقص تاکید دارد.

نشانگرهای نقص بدست آمده را می‌توان در مطالعات جدید برای پیش پردازش و فیلتر نظرهای حاوی نقص استفاده کرد. همچنین می‌توان این الگوها را در مطالعات مبتنی بر نظارت از راه دور [۱] بکار گرفت.

می‌توان از تکنیک نظارت از راه دور با استفاده از الگوهایی که به صورت خودکار در این پژوهش استخراج شد جهت بهبود نتایج کلاس‌بندی استفاده نمود. با توجه به روشی که در این پژوهش مطرح شد برای استخراج اطلاعات دیگری همچون استخراج نظرهایی که از محصول یا سرویس‌های خدماتی ناراضی هستند، نظراتی که عقیده‌ی مثبتی دارند، استخراج نظرهایی که دو یا چند کالا را مقایسه کرده‌اند و استخراج اطلاعات از مقایسه‌ها و نظرهایی که حاوی پیشنهادات مشتریان است می‌توان استفاده کرد.

C. Havasi, "New avenues in opinion mining and sentiment analysis," *IEEE Intelligent Systems*, vol. 28, pp. 15-21, 2013.

5.S. Moghaddam, "Beyond sentiment analysis: mining defects and improvements from customer feedback," in *European Conference on Information Retrieval*, 2015.

6.L.-W. Ku, Y.-T. Liang and H.-H. Chen, "Opinion extraction, summarization and tracking in news and blog corpora," in *Proceedings of AAI*, 2006.

7.M. Hu and B. Liu, "Mining opinion features in customer reviews," in *AAI*, 2004.

8.W. Jin, H. H. Ho and R. K. Srihari, "OpinionMiner: a novel machine

محصول کمک کند؛ مدیران شرکت‌ها و تولیدکنندگان برای اخذ تصمیمات عملی و برنامه‌ریزی‌های تجاری خود نیاز به اطلاعات دقیق‌تری دارند. کشف نقص محصول نقش موثری در ارائه‌ی سریع راه حل کارا و در نتیجه راضی نگه‌داشتن مشتریان دارد. در این مقاله روشی پیشنهاد گردیده است که بتوان اطلاعات کاربردی از بازخورد مشتریان استخراج کرد. روش پیشنهادی مزایای زیر را داراست:

- به صورت خودکار گزارش‌های نقص را از متن نظرهای مشتریان استخراج می‌کند،
- مستقل از دامنه است،
- برای هر پایگاه داده‌ای قابل اجرا می‌باشد و
- به دلیل اینکه مدیران شرکت‌ها مجبور به خواندن کل متن نظر نباشند خلاصه‌ای از گزارش‌های نقص کشف شده را ارائه می‌دهد.

نتایج روی مجموعه داده‌های واقعی از سایت آمازون نشان داد برای کشف گزارش خرابی بررسی و صرفاً تحلیل لغات حسی مفید نیست. اما طبقه‌بندی مثل جنگل تصادفی با دادگان آموزشی کم نیز می‌تواند کلاس‌بندی قابل قبول داشته باشد. تخصیص پنهان دیریکله با مجموعه ویژگی

منابع

1. B. Liu and L. Zhang, "A survey of opinion mining and sentiment analysis," in *Mining text data*, Springer, 2012, pp. 415-463.
- 2.S. Moghaddam and M. Ester, "Opinion digger: an unsupervised opinion miner from unstructured product reviews," in *Proceedings of the 19th ACM international conference on Information and knowledge management*, 2010.
- 3.B. Liu, M. Hu and J. Cheng, "Opinion observer: analyzing and comparing opinions on the web," in *Proceedings of the 14th international conference on World Wide Web*, 2005.
- 4.E. Cambria, B. Schuller, Y. Xia and

2010.

13.D. M. Blei, A. Y. Ng and M. I. Jordan, "Latent dirichlet allocation," *Journal of machine Learning research*, vol. 3, pp. 993-1022, 2003.

14.Z. Qiao, X. Zhang, M. Zhou, G. A. Wang and W. Fan, "A Domain Oriented LDA Model for Mining Product Defects from Online Customer Reviews," 2017.

15.C. Brun and C. Hagege, "Suggestion Mining: Detecting Suggestions for Improvement in Users' Comments.," *Research in Computing Science*, vol. 70, pp. 199-209, 2013.

16.L. Zhang and B. Liu, "Aspect and entity extraction for opinion mining," in *Data mining and knowledge discovery for big data*, Springer, 2014, pp. 1-40.

17.X. Zhang, Z. Qiao, L. Tang, W. Fan, E. Fox and G. Wang, "Identifying Product Defects from User Complaints: A Probabilistic Defect Model," 2016.

18.A. Liaw, M. Wiener and others, "Classification and regression by randomForest," *R news*, vol. 2, pp. 18-22, 2002.

learning system for web opinion mining and extraction," in *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, 2009.

9.F. Li, C. Han, M. Huang, X. Zhu, Y.-J. Xia, S. Zhang and H. Yu, "Structure-aware review mining and summarization," in *Proceedings of the 23rd international conference on computational linguistics*, 2010.

10.S. Moghaddam and M. Ester, "The FLDA model for aspect-based opinion mining: addressing the cold start problem," in *Proceedings of the 22nd international conference on World Wide Web*, 2013.

11.W. X. Zhao, J. Jiang, H. Yan and X. Li, "Jointly modeling aspects and opinions with a MaxEnt-LDA hybrid," in *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing*, 2010.

12.S. Brody and N. Elhadad, "An unsupervised aspect-sentiment model for online reviews," in *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*,

