

تحلیل ادراکات کاربران درباره خرید تلفن همراه در سایت دیجی کالا

فاطمه عباسی
آمنه خدیور
محسن یزدی نژاد

چکیده

امروزه افراد برای خرید محصولات و خدمات آنلاین از نظرات دیگران در شبکه‌های اجتماعی جهت تصمیم‌گیری استفاده می‌نمایند. همچنین شرکت‌های ارائه دهنده محصولات از تحلیل ادراکات و نظرات کاربران و مشتریان برای اتخاذ تصمیمات آگاهانه و ارائه محصولات جدید استفاده می‌نمایند. تحلیل ادراکات از جمله رویکردهای نوین در استخراج نظرات می‌باشد. تحلیل ادراکات، استفاده از روش‌های متن کاوی و پردازش زبان طبیعی برای شناسایی، استخراج و بررسی اطلاعات ذهنی می‌باشد. اطلاعات حاصل از تحلیل ادراکات می‌تواند بر انتخاب موثر مشتریان تاثیر بسزایی داشته باشد. در این پژوهش مدلی جهت تحلیل ادراکات کاربران در ارتباط با خرید تلفن همراه از سایت دیجی کالا ارائه شده است. تحقیق حاضر از لحاظ هدف کاربردی است و جامعه مورد بررسی شامل نظرات کاربران در سایت دیجی کالا می‌باشد و نمونه آماری نظرات کاربران تلفن همراه سایت دیجی کالا است. جهت تحلیل و پیاده سازی، رویکرد یادگیری نظارت شده و از پکیج‌های متن کاوی پایتون استفاده شده است. نتایج نشان می‌دهند مدل پیشنهادی با دقت ۰,۸۹۲ می‌تواند نظرات کاربران را دسته بندی نماید. همچنین در مجموع نظرات کاربران در مورد سهولت استفاده و امکانات و قابلیت‌های تلفن همراه مثبت و در مورد ارزش خرید نسبت به قیمت، نوآوری، طراحی و ظاهر و کیفیت ساخت گوشی‌ها نظر کاربران منفی می‌باشند. مدل پیشنهادی می‌تواند در سایت‌های تجارت الکترونیک مانند دیجی کالا پیاده سازی شود و خروجی آن به صورت سیستماتیک توسط کاربران قابل مشاهده باشد که در نهایت می‌تواند منجر به تصمیم‌گیری آگاهانه برای خریداران و شرکت‌های ارائه دهنده محصولات باشد.

کلیدواژه‌گان: تحلیل ادراکات، شبکه‌های اجتماعی، یادگیری نظارت شده

محقق پسادکتری، دانشکده علوم اجتماعی و اقتصاد، دانشگاه الزهرا (س)، تهران.

عضو هیات علمی، گروه مدیریت، دانشکده علوم اجتماعی و اقتصاد دانشگاه الزهرا (س)، تهران، ایران. (نویسنده مسئول):

a.khadivar@alzahra.ac.ir

دانشجوی دکتری، هوش مصنوعی، دانشکده کامپیوتر، دانشگاه اصفهان، اصفهان.

تاریخ پذیرش: ۱۳۹۸/۱۰/۳۰

تاریخ دریافت: ۱۳۹۸/۰۶/۲۸

مقدمه

در سال‌های اخیر استفاده از اینترنت رشد فزاینده‌ای داشته است. اطلاعات آماری نشان می‌دهند در سال ۲۰۱۹ بیش از چهار میلیارد کاربر اینترنت در دنیا وجود دارد که نسبت به سال ۲۰۰۰، ۱۱۵۷ درصد رشد داشته است. در این میان ایران با ۶۲ میلیون کاربر اینترنت در رتبه شانزدهم در میان بیست کشور با بیشترین تعداد کاربران اینترنت در دنیا قرار دارد (آمار جهانی اینترنت^۱، ۲۰۱۹). سرعت زیاد گسترش سهم اینترنت در زندگی انسان‌ها و نیز تحت تأثیر قرار دادن شبکه زندگی آنها با تبدیل روش‌های سنتی به مدرن، فصل جدیدی را در ابعاد رفتاری نمایان کرده است. در این میان تجارت الکترونیک یکی از کلیدی‌ترین مفاهیم اضافه شده به زندگی روزمره انسان‌هاست (بخشی زاده و همکاران، ۱۳۹۷). وب‌سایت‌های تجارت الکترونیک جهت خرید و فروش کالا و خدمات استفاده می‌شوند. از طریق این وبسایت‌ها افراد محصولات مد نظر خود را مشاهده می‌نمایند، نظرات سایر خریداران را بررسی می‌نمایند و در نهایت برای خرید محصول مد نظر خود تصمیم‌گیری می‌نمایند (موروگاوالی و همکاران^۲، ۲۰۱۷). افزایش چشمگیر برنامه‌های کاربردی^۳ گوشی‌های هوشمند و تبلت‌ها به کاربران امکان می‌دهد تا از طریق اینترنت نظرات خود را در شبکه‌های اجتماعی در مورد موضوعات مختلف با دیگران به اشتراک بگذارند (ژانگ و همکاران^۴، ۲۰۱۴). کاربرانی که از طریق وب‌سایت‌های تجارت الکترونیک محصول یا خدمتی را خریداری می‌نمایند می‌توانند نظرات خود را با دیگران به اشتراک بگذارند. در حوزه تجارت الکترونیک اگر کسی بخواهد محصولی را خریداری نماید دیگر محدود به نظرات خانواده و دوستان نیست زیرا نظرات و بحث‌های کاربران در مورد محصولات مختلف در وب در دسترس می‌باشد. برای یک سازمان نیز دیگر نیازی به انجام نظرسنجی و گروه کانونی نیست زیرا چنین اطلاعاتی به وفور در اینترنت در دسترس می‌باشد. با اینحال پیدا کردن و ارزیابی نظرات به دلیل گستردگی و تنوع نظرات کار ساده‌ای نمی‌باشد. هر سایت

1. Internet World Stats
2. Murugavalli et al.
3. Applications
4. Zhang et al.

معمولا دارای حجم زیادی نظرات متنی است که رمزگشایی از آنها کار ساده‌ای نمی‌باشد (لیو، ۲۰۱۲). در حال حاضر رشد محتوای تولید شده توسط کاربران در وبسایت‌ها و شبکه‌های اجتماعی منجر به افزایش قدرت شبکه‌های اجتماعی برای بیان نظرات در مورد خدمات، محصولات و رویدادها شده است (لیو، ۲۰۱۵).

در سال‌های اخیر تجارت الکترونیک، کسب و کارها را قادر ساخته است تا محصولات و خدمات متنوع و شخصی سازی شده ای را به مشتریان ارائه دهند. علی‌رغم توسعه و گسترش سطح شخصی سازی محصولات و خدمات، این روند باعث شده است که حجم اطلاعاتی که مشتریان قبل از خرید با آن مواجه می‌شوند به میزان قابل توجهی افزایش پیدا کند. از سوی دیگر رشد شبکه‌های اجتماعی و به تبع آن افزایش فزاینده محتوای این شبکه‌ها باعث شده است تا افراد برای خرید و استفاده از محصولات، خدمات و یا حتی انتخاب‌های سیاسی خود از نظرات سایر افراد برای تصمیم‌گیری استفاده نمایند. بیشتر اطلاعاتی که از طریق شبکه‌های اجتماعی مبادله و ذخیره می‌شوند در فرمت متنی هستند که با توجه به حجم عظیم داده‌های متنی تحلیل و کشف دانش از این داده‌ها بسیار ارزشمند می‌باشد (لیو، ۲۰۱۲). با توجه به استفاده روزافزون کاربران از وبسایت‌های تجارت الکترونیک جهت خرید و فروش محصولات و خدمات، تحلیل نظرات مشتریان می‌تواند به شناخت ترجیحات و افزایش رضایت مشتریان کمک نماید. قطعا تحلیل و پردازش چنین داده‌های متنی و ساختارنیافته ای از اهداف کلیدی جهت هوشمندسازی فرآیند انتخاب و تصمیم‌گیری وبسایت‌های تجارت الکترونیک می‌باشد. در سال‌های اخیر تکنیک‌های داده کاوی و متن کاوی امکان تحلیل حجم عظیم داده‌ها، پیش بینی نیازهای آتی مشتریان و تحلیل ادراکات و نظرات مشتریان را فراهم نموده اند.

دیجی کالا نیز به عنوان یکی از وبسایت‌های فروش آنلاین در ایران این امکان را فراهم آورده تا کاربران بتوانند نظرات خود را در مورد محصولات مختلف این سایت از جمله تلفن همراه، بدون هیچ محدودیتی ارائه نمایند. این نظرات به سایر کاربران در تصمیم‌گیری در مورد خرید یا عدم خرید کتاب کمک می‌نماید. مشکلی که امروزه اغلب کاربران با آن مواجه هستند،

کمبود وقت است که باعث می شود اغلب نظرات را مطالعه نکنند و تنها بر امتیازات سایر کاربران برای تصمیم گیری تکیه نمایند که این خود می تواند گمراه کننده باشد. امتیاز کاربران برای یک محصول خاص تجربه کلی کاربر از استفاده از آن محصول را منتقل می کند و بستر و زمینه ای را که منجر به آن تجربه شده است را انتقال نمی دهد. به عنوان مثال امتیازی که کاربران به یک تلفن همراه در سایت دیجی کالا می توانند بدهند به ترتیب ۱،۲،۳،۴،۵ است که بیان کننده علاقه یا عدم علاقه به آن گوشی است و نمی تواند به خوبی نظرات و ادراکات کاربر را از جهت کیفیت و یا قیمت بیان نماید. از این جهت نظرات کاربران با اطمینان بیشتری بیان کننده عقیده و ادراک کاربران در مورد محصول یا خدمتی خاص می باشند. با افزایش نظرات کاربران لازم است روش هایی برای مشخص نمودن هدف نویسنده نظرات ایجاد گردد. درک ادراک و نظرات کاربران و مشتریان محصولات به بازاریابی بهتر و کسب اعتبار در محیط آنلاین کمک می نماید. تحلیل و بررسی این نظرات از راهکارهای نوین برای سایت های فروش آنلاین جهت افزایش کیفیت و بهبود خدمات به مشتریان و در نتیجه تصمیم گیری آگاهانه تر خریداران این سایت ها می باشد. با توجه اهمیت توسعه و استفاده از رویکردهای متن کاوی و تحلیل ادراکات برای جامعه علمی و نیز نقش موثر این موضوع از جهت کاربرد برای کسب و کارهای مختلف و خلاء شناسایی شده در این زمینه، مساله اصلی این پژوهش ارائه مدلی برای تحلیل ادراکات خریداران تلفن همراه می باشد. منظور از ادراک در این مقاله تفسیر پژوهشگر از کامنت یا اظهار نظر کاربر است. با استفاده از رویکرد پیشنهادی شرکت ها نیازی به ارائه پرسشنامه جهت استخراج نظرات مشتریان ندارند و بر اساس تحلیل نظرات متنی کاربران، ترجیحات و نظرات آنها نسبت به محصول مشخص می گردد و مجموعه ای از اقلام پیشنهادی به مشتریان ارائه می گردد که می تواند گامی ارزشمند در جهت افزایش رضایت و بهبود ارتباط با مشتریان که مهمترین سرمایه سازمانی هستند باشد. به علاوه جهت استخراج و تحلیل نظرات مشتریان با استفاده از روش تحلیل ادراکات، مشتریان محدود به یکسری پرسش مشخص نمی باشند و معایب پرسشنامه برطرف می گردد. در این روش کاربر با توجه به علاقه، نظرات واقعی خود را در ارتباط با یک موضوع خاص بیان می کند. در نتیجه صحت خروجی نتایج و

اعتبار پیشنهادات بالا می باشد. در ادامه ابتدا پیشینه ای از پژوهش های صورت گرفته در این حوزه ارائه می شود. در بخش های بعدی به ترتیب روش شناسی پژوهش، یافته های پژوهش و در نهایت نتیجه گیری و پیشنهادها جهت تحقیقات آتی ارائه می گردد.

پیشینه پژوهش

پیشینه نظری پژوهش

عقاید دیگران در موقع تصمیم گیری یا انتخاب یک گزینه از میان چندین گزینه می تواند بسیار حیاتی باشد. وقتی این گزینه ها حاوی منابع ارزشمند باشند (برای مثال صرف زمان و هزینه برای خرید محصولات یا سرویس ها)، مردم اغلب به تجربیات قبلی خود تکیه می کنند. در گذشته منابع مهم اطلاعاتی، دوستان و مجلات و وبسایت های تخصصی بودند. امروزه شبکه های اجتماعی ابزار جدیدی برای به اشتراک گذاری ایده ها با افراد متصل به شبکه ی جهانی وب فراهم می کنند. دریافت عقاید عمومی درباره رویدادهای اجتماعی، فعالیت های بازاریابی و اولویت های محصول، علاقه جوامع علمی و جهان تجارت را به خود جلب کرده است. زمینه های ادغام شده حاصل ادراک سنجی و تحلیل ادراکات است (کامبریا و همکاران^۱، ۲۰۱۲). داده های متنی یکی از مهمترین اشکال داده های بدون ساختار^۲ می باشند. این نوع از داده ها، مخزن بزرگی از اطلاعات هستند که می توانیم محتوای مناسب را از آن استخراج نماییم. یکی از راه های استخراج اطلاعات از داده های متنی، تحلیل ادراکات می باشد (ژان و فانگ^۳، ۲۰۱۵). تحلیل ادراکات اشاره به فرآیند متن کاوی برای استخراج لحن نوشته شده توسط کاربر دارد که این لحن مثبت یا منفی می باشد (موروگاوالی و همکاران، ۲۰۱۷). تحلیل ادراکات از زبان طبیعی، چالش برانگیز است زیرا نیاز به درک عمیقی از قوانین کار و ضمنی، منظم و نامنظم و نحوی و معنایی زبان دارد. تحقیقات تحلیل ادراکات با مشکلات حل نشده پردازش زبان طبیعی^۴ یعنی وضوح هم ارجاعی، وضوح تکرار و وضوح حساس به کلمه دست و پنجه نرم می کند. تحلیل

1. Cambria et al.

2. Unstructured Data

3. Zhan & Fang

4. Natural Language Processing

ادراکات یک مسئله‌ی بسیار محدود پردازش زبان طبیعی است زیرا سیستم فقط نیاز به درک ادراکات مثبت و منفی از هر جمله و موجودیت‌ها و موضوعات مقصد دارد (کامبریا و همکاران، ۲۰۱۲). تاکنون روش‌های مختلفی برای تحلیل ادراکات معرفی شده است که این روش‌ها به سه رویکرد اصلی تقسیم می‌شوند: رویکرد مبتنی بر یادگیری ماشینی^۱، مبتنی بر شبکه‌ی واژگانی^۲ و رویکرد ترکیبی (پاندی و همکاران^۳، ۲۰۱۷). رویکرد یادگیری ماشینی نظارت شده تلاش دارد تا ادراکات متن یا سند را بر اساس اطلاعات جمع‌آوری شده یا آموخته شده، طبقه‌بندی و پیش‌بینی نماید (وانگ^۴، ۲۰۱۷). در یادگیری بدون نظارت برخلاف یادگیری نظارت شده، داده‌های مشخصی از قبل وجود ندارد و هدف، ارتباط ورودی و خروجی نیست، بلکه تنها دسته‌بندی آنها مهم است و این یادگیرنده است که بایستی در داده‌ها به دنبال ساختاری خاص بگردد. از آنجا که کلمات ادراکی فاکتور مهمی برای طبقه‌بندی برداشت‌ها می‌باشند، دور از تصور نیست که این کلمات و عبارات برای طبقه‌بندی عواطف در یک روش بدون نظارت استفاده شوند (لیو، ۲۰۱۵). در رویکرد ترکیبی داده‌هایی که در ابتدا با تکیه بر روش مبتنی بر لغت‌نامه بار ادراکی آنها تعیین می‌شوند، به عنوان داده‌های ورودی برای آموزش مدل در الگوریتم‌های یادگیری ماشینی در روش تحت نظارت مورد استفاده قرار می‌گیرند (فیلهو و پاردو^۵، ۲۰۱۳).

پیشینه تجربی پژوهش

پیکری و همکاران در سال ۱۳۹۴ با استفاده از تکنیک‌های متن کاوی و تحلیل ادراکات فوت مرتضی پاشایی را در شبکه اجتماعی توییتر مورد بررسی قرار داده‌اند. در این پژوهش محققین این نوشتار محتوای بارگذاری شده را در پنج مقوله رده‌بندی نموده‌اند (پیکری و همکاران، ۱۳۹۴). عباسی و همکاران مدلی را جهت دسته‌بندی ادراکات خریداران کتاب سایت آمازون ارائه نموده است. در این پژوهش از روش وزن‌دهی با دقت ۰/۸۰۰۳ استخراج ادراکات استفاده

1. Machine Learning
2. Lexicon-based Approach
3. Pandey et al.
4. Wang
5. Filho & Pardo

شده است (عباسی و همکاران، ۱۳۹۶). نجف‌زاده و همکاران چارچوب نیمه نظارتی مبتنی بر لغت‌نامه جهت تحلیل نظرات فارسی ارائه نموده‌اند. در این پژوهش برچسب گذاری ادراکات با استفاده از یک لغت‌نامه و بدون نیاز به خبره انسانی انجام گرفته است. رضایی و همکاران نظرات کاربران سایت دیجی کالا را در مورد برندهای مختلف تلفن همراه با استفاده از متن کاوی مورد بررسی قرار می‌دهند. نتایج این پژوهش نشان می‌دهد از دیدگاه مشتریان هیچ برند تلفن همراهی برتر نیست و هرکدام ویژگی‌های مثبت و منفی دارند (نجف‌زاده و همکاران، ۱۳۹۷). حاج سیدجوادی در سال ۱۳۹۵ با ترکیب روش‌های یادگیری ماشین و شباهت معنایی رویکرد جدیدی را جهت تحلیل ادراکات ارائه نموده‌اند. در این پژوهش الگوریتم ژنتیکی جهت کاهش ابعاد فضای ویژگی‌ها طراحی شده است (حاج سید جوادی و جلالی، ۱۳۹۴).

بات^۱ در سال ۲۰۱۵ سیستمی را برای دسته بندی نظریات مشتریان بر اساس تحلیل ادراکات نظرات پیشنهاد می‌نماید. تحلیل بر روی نظرات کاربران در سایت آمازون در ارتباط با تلفن همراه آیفون ۵ انجام می‌شود (بهاث و همکاران^۲، ۲۰۱۵). لیانگ و همکاران در پژوهشی تاثیر نظرات مشتریان را بر فروش برنامه‌های موبایلی ارزیابی نموده‌اند. در این پژوهش تحلیل ادراکات چند وجهی بر روی نظرات مشتریان انجام می‌شود. یافته‌های این پژوهش نشان می‌دهد نظرات مشتریان در مورد کیفیت محصولات و خدمات بر رتبه‌بندی فروش تاثیر بسزایی دارد (لیانگ و همکاران^۳، ۲۰۱۵). ویدیا و همکاران در سال ۲۰۱۵ از طریق تحلیل ادراکات بر روی نظرات کاربران اعتبار برندهای تلفن همراه را مورد ارزیابی قرار داده است. جهت مدل‌سازی از رویکرد یادگیری نظارت شده استفاده شده است (ویدیا و همکاران^۴، ۲۰۱۵). مورگوالی و همکاران^۵ در سال ۲۰۱۷ با استفاده از تحلیل ادراکات، بازخورد کاربران در زمینه تجارت الکترونیک را تجزیه و تحلیل نموده‌اند. در این پژوهش الگوریتمی برای تشخیص قطبیت ادراکات بازخورد کاربران در مورد محصولات ارائه شده است (مورگوالی و

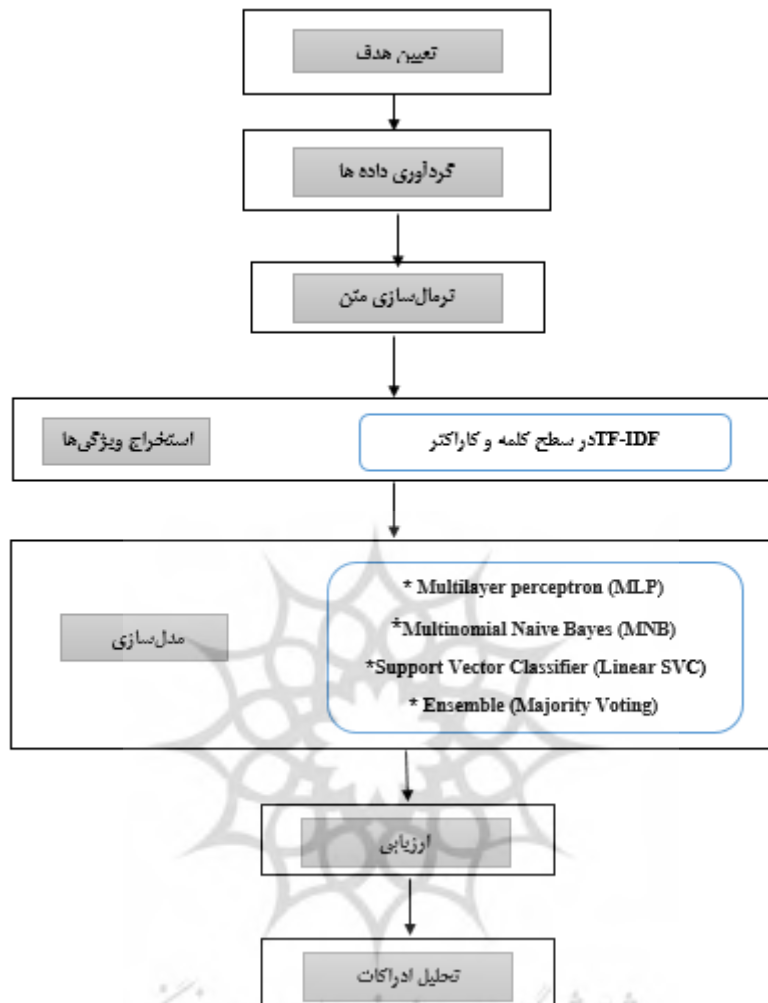
1. Bhatt
2. Bhatt et al.
3. Liang et al.
4. Vidya et al.
5. Murugavalli et al.

همکاران، ۲۰۱۷). ژانگ و همکاران در سال ۲۰۱۹ از ترکیب چندین لغت‌نامه جهت تحلیل ادراکات نظرات کاربران در مورد برنامه‌های موبایلی استفاده نموده است. همچنین با استفاده از روش‌های آماری و ویژگی‌های نظرات کاربران و ارزش ادراکات آنها مورد تجزیه و تحلیل قرار گرفته است. نتایج این پژوهش نشان می‌دهند که انحراف استاندارد ارزش ادراکات، تغییرات ادراکات میان مردم را منعکس می‌نماید (ژانگ و همکاران، ۲۰۱۹).

روش‌شناسی پژوهش

پژوهش حاضر از نظر هدف کاربردی می‌باشد. زیرا نتایج این تحقیق به مشتریان و کاربران آنلاین جهت تصمیم‌گیری برای انتخاب کالا و خدمات بر اساس تجربیات سایر افراد و نیز به مدیران و تصمیم‌گیران این قبیل مشاغل جهت بهبود عملکرد و اخذ تصمیمات آگاهانه‌تر کمک می‌نماید. در این پژوهش پس از گردآوری و پیش‌پردازش داده‌ها با استفاده از رویکرد یادگیری نظارت شده و سه الگوریتم پرسپترون چند لایه^۲، بیزین ساده چند جمله‌ای^۳ و طبقه‌بندی کننده بردار پشتیبان خطی^۴ و رویکرد ترکیبی (رای اکثریت)^۵ مدلسازی انجام شده است. در گام‌های بعد پس از ارزیابی مدل، ادراکات مثبت و منفی نظرات کاربران استخراج گردید. جامعه مورد بررسی شامل نظرات کاربران در سایت دیجی کالا می‌باشد و نمونه آماری نظرات کاربران تلفن همراه سایت دیجی کالا است. در شکل ۱ مراحل و رویکردهای هر مرحله جهت استخراج ادراکات و نظرات کاربران نشان داده شده است.

1. Zhang et al
2. Multilayer perceptron (MLP)
3. Multinomial Naive Bayes (MNB)
4. Linear Support Vector Classification (Linear SVC)
5. Ensemble (Majority Voting)



شکل ۱: چارچوب پژوهش

تعیین هدف

جهت مشخص نمودن هدف لازم است مساله پژوهش به دقت مورد بررسی قرار گیرد. هدف پژوهش حاضر بررسی ادراکات و دیدگاه‌های کاربران در ارتباط با کالاهای خریداری شده از سایت دیجی کالا می‌باشد و بر این اساس تحلیل ادراکات و ادراک سنجی به عنوان رویکرد مطلوب برای تجزیه و تحلیلی نظرات متنی کاربران انتخاب گردیده است.

گردآوری داده‌ها

در این پژوهش جامعه مورد بررسی نظرات کاربران در سایت دیجی کالا می‌باشد. نمونه مورد بررسی نظرات کاربران تلفن همراه در سایت دیجی کالا است که از طریق سایت دیتاهارت^۱ این داده‌ها تهیه شده است. این مجموعه داده مشتمل بر ۱۵۰۰ نظر به زبان فارسی می‌باشد که از طریق سایت دیجی کالا در بازه زمانی سال ۲۰۱۵ تا سال ۲۰۱۶ توسط کاربران به اشتراک گذاشته شده‌اند. نمونه‌ای از داده‌ها در جدول ۱ ارائه شده است.

جدول ۱: نمونه‌ای از داده‌های گردآوری شده جهت تحلیل

ردیف	نظر کاربران
۱	اپل گوشی خوش ساختی است.
۲	روان بودن سیستم عامل عالی است
۳	اسپیکر استریو و صدا قوی تر شده است.
۴	کیفیت صدای طرف مقابل در زمان مکالمه خوب نیست.
۵	کسانی که قیمت کمتر می‌خواهند سامسونگ بهتر است.
۶	از لحاظ کیفیت خیلی از آن راضی هستم.
۷	اگر با قیمت آن مشکل ندارید، حتی یک لحظه در خریدش شک نکنید.
۸	کیفیت صفحه رویایی است.
۹	لوازم جانبی به خوبی نمی‌توانند از صفحه کاملا محافظت کنند.
۱۰	گران بودن آن عیب محسوب می‌شود.

نرمال سازی متن

یکی از مهمترین مراحل متن کاوی، پیش پردازش است که در آن کاراگرها، کلمات و جملات شناسایی می‌شوند (کانان و گوروسامی^۲، ۲۰۱۴). پیش پردازش متن شامل تکنیک‌هایی چون حذف ایست واژه‌ها، ریشه‌یابی کلمات، برچسب گذاری اجزای واژگانی کلام است

1. <http://dataheart.ir>
2. Kannan & Gurusamy
3. Part of Speech Tagging (POS tagging)

(ویجایارانی و همکاران^۱، ۲۰۱۵). پردازش متون فارسی از جهاتی با پردازش متون انگلیسی متفاوت می‌باشد. در زبان انگلیسی تمامی حروف و کلمات جدا از هم و با قانونی مشخص نوشته می‌شوند. در حالیکه در زبان فارسی برخی از حروف به هم چسبیده‌اند، برخی جدا نوشته می‌شوند، بعضی از کلمات یکپارچه‌اند و بعضی از کلمات با فاصله یا نیم فاصله به دو یا چند بخش تقسیم می‌شوند (زمانی و سرخپور، ۱۳۹۳). تمامی حوزه‌های مرتبط با پردازش زبان طبیعی به‌نحوی با متون واقعی سروکار دارند. صورتهای غیراستاندارد نویسه‌ها و کلمات به وفور در این نوع متون دیده می‌شوند. قبل از اینکه بتوان از این متون به منظور استفاده در سیستم‌های تبدیل متن به گفتار، ترجمه ماشینی، بازشناسی حروف فارسی، خلاصه ساز فارسی، جستجو در متون فارسی استفاده کرد و یا در پایگاه داده ذخیره‌شان کرد، باید ابتدا پیش پردازشی روی آنها انجام گیرد تا صورت‌های غیراستاندارد به شکل استاندارد تبدیل گردند. اگر حروف، نشانه‌های نگارشی و کلمات فارسی به شکل یکسانی نوشته نشوند، متون مورد استفاده قابل تحلیل توسط سامانه‌های رایانه‌ای نخواهند بود. طی فرایند نرمال سازی علائم نگارشی، حروف، فاصله‌های بین کلمات، اختصارات و غیره بدون ایجاد تغییرات معنایی در متن به شکل استانداردشان تبدیل می‌گردند (عبدی قویدل و همکاران، ۱۳۹۱). جهت پیش پردازش و نرمال سازی متن در این پژوهش مراحل ذیل انجام گرفته است:

- حذف ایست واژه‌ها: در متون فارسی کلمات پرتکراری چون حروف تعریف و حروف ربط ارزش معنایی ندارند که در این مرحله این کلمات حذف می‌شوند.
- پاکسازی متن: در این مرحله کاراکترهایی چون برچسب‌ها، html، "@"، ویرگول، علامت سوال، نقطه و یگول و کلیه کاراکترهای نامطلوب حذف می‌شوند.
- حذف نویسه «-» که برای کشش نویسه‌های چسبان مورد استفاده قرار می‌گیرد.
- حذف فاصله و نیم‌فاصله‌های اضافه بکار رفته در متن

بخش می‌باشد (وانگ، ۲۰۱۷). این مدل بسته کلمات^۱ نیز نامیده می‌شود زیرا وقوع کلمات را می‌شمارد و دستور زبان و ترتیب کلمات را نادیده می‌گیرد (شونلاو و گوئنتر^۲، ۲۰۱۶). در این پژوهش n مربوط به n-grams در بازه ۱ تا ۸ و بر اساس کد زیر مورد انتخاب و بررسی قرار گرفته است. بر اساس این کد ۲۰۰۰ ویژگی با TF-IDF بالاتر را مورد تحلیل قرار می‌دهد.

```
from sklearn.feature_extraction.text import TfidfVectorizer
TFIDF = TfidfVectorizer(ngram_range=(1,8), max_features=2000)
X = TFIDF.fit_transform(df.Text)
```

در جدول ۲ نمونه ای از استخراج ویژگی‌ها در بازه ۱ تا ۸ نشان داده شده است.

فراوانی کلمه در مقابل فراوانی سند^۳

تکنیک استخراج ویژگی‌ها در طبقه‌بندی متن، بازیابی اطلاعات، تشخیص موضوع و خلاصه‌سازی سند مورد استفاده قرار می‌گیرد. روش‌های اصلی در این تکنیک، فراوانی کلمه در مقابل فراوانی سند (TF-IDF) اطلاعات به دست آمده^۴ (IG)، آماره مربع کای^۵ (CHI) می‌باشند (چاکرابورتی^۶، ۲۰۱۳). هدف از این کار کم کردن تأثیر لغاتی است که در تحلیل اهمیت کمتری دارند. این روش کلمات را بر اساس اهمیتشان وزن دهی می‌کند. برای محاسبه این ضریب، فراوانی کلمه (اهمیت کلمه در سند) با فراوانی معکوس سند (اهمیت کلمه در کل اسناد) ترکیب می‌شود تا وزن هر کلمه در هر سند مشخص شود. فراوانی کلمه، نسبت تعداد تکرار یک کلمه در سند به تعداد کلمات موجود در سند است. در مقابل فراوانی معکوس سند، لگاریتم نسبت تعداد سندها به تعداد اسنادی است که کلمه مورد نظر را دارا است. این روش باعث می‌شود تا کلماتی که در تعداد کمتری از اسناد هستند وزن بیشتری پیدا کنند و کلماتی که در اغلب اسناد موجودند وزن کمتری پیدا کنند. بدین ترتیب کلمات با وزن کمتر در تحلیل وارد نمی‌شوند (خدیور و عباسی، ۱۳۹۸؛ لطفی آذر داریان و جاویدان، ۱۳۹۵).

1. Bag of Words
2. Schonlau & Guenther
- 3 Term Frequency Inverse Document Frequency (TFIDF)
4. Information Gain
5. Chi-Square Statistics
6. Chakraborty

رابطه (۱)

$$W_t = TF(d, t) IDF(t)$$

براساس این فرمول، بیشترین مقدار وزن زمانی است که کلمه t به تعداد دفعات زیاد در یک سند یا در تعداد محدودی سند رخ می دهد و کمترین مقدار زمانی است که کلمه t تعداد دفعات کمتر در سند یا در بسیاری از اسناد رخ دهد. در ضمن مقدار وزن زمانی که کلمه در همه اسناد رخ داده باشد حداقل می باشد. شباهت میان اسناد از فرمول زیر محاسبه می شود که معادل روش شباهت کسینوسی می باشد (خدییور و عباسی، ۱۳۹۸).

رابطه (۲)

$$\text{sim}(d_1, d_2) = \frac{d_1 \cdot d_2}{|d_1| |d_2|}$$

در این پژوهش از n-grams و فراوانی کلمه در مقابل فراوانی سند جهت استخراج ویژگی ها استفاده شده است که در جدول ۲ نمونه ای از استخراج ویژگی ها با استفاده از دو روش نام برده ارائه شده است.

جدول ۲: نمونه ای از استخراج ویژگی ها

ردیف	استخراج ویژگی ها
۱	'با این'
۲	'با این سیستم'
۳	'با این سیستم عامل'
۴	'با این سیستم عامل واقعا'
۵	'با این سیستم عامل واقعا لذت'
۶	'با این سیستم عامل واقعا لذت بخش'
۷	'با این سیستم عامل واقعا لذت بخش است'

یافته‌های پژوهش

مدل‌سازی

پس از تبدیل داده‌های متنی ساختارنیافته به داده‌های ساختاریافته، کلمات و ویژگی‌های استخراج شده جهت مدل‌سازی مورد استفاده قرار گرفته‌اند. در این پژوهش از رویکرد یادگیری نظارت شده جهت مدل‌سازی استفاده شده است. جهت مدل‌سازی از سه الگوریتم پرسپترون چند لایه^۱، بیزین ساده چندجمله‌ای^۲ و طبقه بندی کننده بردار پشتیبان خطی^۳ و رویکرد ترکیبی (رای اکثریت)^۴ از کتابخانه NLTK^۵ پایتون جهت مدل‌سازی استفاده شده است.

طبقه بندی کننده بیزین ساده چندجمله‌ای: این دسته بندی کننده به میزان وسیعی برای دسته بندی متون به کار گرفته می‌شود. این دسته بندی کننده بر مبنای اصول یادگیری بیزین می‌باشد و فرض می‌کند که توزیع کلمات در اسناد توسط مدل‌های پارامتری خاصی ایجاد می‌شود (ژائو و همکاران^۶، ۲۰۱۶).

طبقه بندی کننده پرسپترون چندلایه: این دسته بندی کننده در دسته شبکه‌های عصبی مصنوعی^۷ قرار می‌گیرد. یک دسته بندی کننده پرسپترون چندلایه حداقل شامل سه لایه گره می‌باشد. به جز گره‌های ورودی، هر گره یک نورون است که از یک تابع فعال‌سازی غیرخطی استفاده می‌نماید. پرسپترون چند لایه جزء شبکه‌های عصبی پیش‌خور^۸ می‌باشد که با انتخاب تعداد لایه‌ها و نورون‌های بهینه، عملکرد با دقت دلخواه را ارائه نماید (پاپسکیو و همکاران^۹، ۲۰۰۹).

1. Multilayer Perceptron (MLP)
2. Multinomial Naive Bayes (MNB)
3. Linear Support Vector Classification (Linear SVC)
4. Ensemble (Majority Voting)
5. Natural Language Toolkit
6. Zhao et al.
7. Artificial Neural Network (ANN)
8. Feedforward Artificial Neural Network
9. Popescu et al.

طبقه بندی کننده بردار پشتیبان خطی: این روش در زمره روش های یادگیری نظارت شده قرار می گیرد که از روش های نسبتا جدید است که از کارایی بیشتری نسبت به روش های قدیمی تر برخوردار است. مبنای کار این روش دسته بندی خطی داده ها می باشد و در تقسیم خطی داده ها تلاش می شود تا خطی انتخاب شود تا حاشیه اطمینان بیشتری داشته باشد (میتچل^۱، ۲۰۱۵).

رویکرد ترکیبی^۲:

رویکرد یادگیری ماشین نظارت شده نوعی از یادگیری است که در آن ورودی و خروجی مشخص است و ناظری وجود دارد که اطلاعاتی را در اختیار یادگیرنده قرار می دهد و به این ترتیب سیستم سعی می کند تا تابعی را از ورودی به خروجی فراگیرد. در یک طبقه بندی مبتنی بر یادگیری ماشین دو مجموعه از متون موجود می باشد: مجموعه آموزشی^۳ و مجموعه آزمایشی^۴. مجموعه آموزشی به وسیله یک دسته بندی کننده خودکار برای یادگیری ویژگی های متمایز از متون استفاده می شود و مجموعه آزمایشی برای ارزیابی دسته بند خودکار استفاده می گردد (مدهات و همکاران^۵، ۲۰۱۴). یکی از رویکردهای یادگیری نظارت شده روش ترکیبی است. یادگیری ترکیبی که چندین دسته بندی کننده را برای حل مساله آموزش می دهد. به دو دلیل از رویکرد ترکیبی برای یادگیری استفاده می شود:

۱. آماری^۶: زمانی که نتیجه بر ترکیب دسته بندی کننده ها متکی باشد، احتمال انتخاب طبقه بندی اشتبه کاهش می یابد.

۲. محاسباتی^۷: برخی از الگوریتم های یادگیری مبتنی بر جستجوی محلی می باشند که احتمال گیر کردن در نقاط بهینه محلی وجود دارد. با استفاده از رویکرد ترکیبی می توان به نتایج بهتری

1. Mitchell
2. Ensemble Method
3. Training Set
4. Test Set
5. Medhat et al.
6. Statistical
7. Computational

دست یافت (دیتریچ^۱، ۲۰۰۰).

رویکرد ترکیبی مبتنی بر ترکیب مجموعه ای از مدل‌های دسته بندی است که ترکیب آن‌ها صحت عملکرد بیشتری را نسبت به استفاده منفرد و مستقل از آن‌ها در اختیار کاربر خواهد گذاشت. مدل جمعی دسته بندی کننده، مجموعه ای از مدل‌های دسته بندی است که تصمیمات منفرد آن‌ها را با یک روش شناسی خاص ترکیب می‌شود (الهی و همکاران، ۱۳۹۳). رویکرد ترکیبی با پوشش نواحی که هر یک از دسته بندی کننده‌ها نتوانستند پیش بینی نمایند، دقت مدل پیش بینی کننده تحلیل ادراکات افزایش باید. روش‌های رای اکثریت^۲ و رتبه بندی از روش‌های متداول برای ترکیب نتایج در رویکردهای ترکیبی است. روش‌های رای گیری زمانی که هر دسته بندی کننده یک برچسب منحصر بفرد برای هر کلاس را نمایش می‌دهد، استفاده می‌شود. روش رای اکثریت بر این اساس می‌باشد که کلاس برنده بیش از نیمی از آراء را به خود اختصاص می‌دهد و در این روش همه دسته بندی کننده‌ها دارای وزن یکسان هستند (سوئن و لام^۳، ۲۰۰۰). به طور کلی دسته بندی کننده رای اکثریت به صورت مجموعه ای از آراء تعریف می‌شود که به شرح زیر می‌باشد:

رابطه ۳)

$$C(X) = \arg \max \sum_{i=1}^n v_{i,j}$$

در رویکرد رای اکثریت نمونه‌های بدون برچسب بر اساس دسته ای که بیشترین تعداد آراء را به دست آورده است، برچسب گذاری می‌شوند (عبدماناف و همکاران^۴، ۲۰۱۷).

ارزیابی

برای ارزیابی مدل‌های دسته بندی کننده، یک مجموعه تصادفی از اسناد را که مستقل از

1. Dietterich
2. Majority Voting
3. Suen & Lam
4. Abdmanaf et al.

مجموعه آموزش^۱ می باشد را به عنوان مجموعه تست^۲ در نظر می گیریم. پس از آموزش مدل با مجموعه داده های آموزش، مجموعه تست را دسته بندی می کنیم و برچسب های برآورد شده را با برچسب های واقعی مقایسه می نماییم و عملکرد مدل دسته بندی کننده را ارزیابی می نماییم (صنعی آباده و محمودی، ۱۳۹۴). پس از پیاده سازی مدل با استفاد از رویکرد ترکیبی لازم است دقت، بازخوانی و معیار ارزیابی F تخمین زده شود. دقت^۳ و بازخوانی^۴ معیارهای کاربردی در حوزه ارزیابی اطلاعات هستند که میزان تناسب اسناد ارزیابی شده توسط سیستم را با نیاز کاربر تعیین می کنند (راقاوان و جی وانگ^۵، ۱۹۸۹).

در این پژوهش جهت تخمین معیارهای فوق از ماتریس درهم ریختگی^۶ استفاده شده است. ماتریس درهم ریختگی عملکرد الگوریتم های مربوطه را نشان می دهند. در جدول ۳ ماتریس درهم ریختگی به تصویر کشیده شده است که بر اساس اطلاعات این ماتریس ارزیابی صورت می گیرد.

جدول ۳: ماتریس درهم ریختگی الگوریتم های طبقه بندی

مقدار پیش بینی شده ^۷		مقدار واقعی ^۸	
دسته مثبت	دسته منفی		
مثبت نادرست ^{۱۰}	مثبت درست ^۹	دسته مثبت	مقدار واقعی ^۸
منفی درست ^{۱۲}	منفی نادرست ^{۱۱}	دسته منفی	

1. Train set
2. Test set
3. Precision
4. Recall
5. Raghavan & Gwang
6. Confusion matrix
7. Predicted
8. Actual
9. True Positive(TP)
- 10 False Positive(FP)
11. False Negative(FN)
12. True Negative(TN)

مثبت درست (TP): تعدادی از داده‌ها که به درست به عنوان دسته مثبت شناسایی شده‌اند.

مثبت نادرست (FP): تعدادی از داده‌ها که به غلط به عنوان دسته مثبت شناسایی شده‌اند.

منفی نادرست (FN): تعدادی از داده‌ها که به غلط به عنوان دسته منفی شناسایی شده‌اند.

منفی درست (TN): تعدادی از داده‌ها که به درست به عنوان دسته منفی شناسایی شده‌اند.

مقادیر ماتریس همواره عددی بین ۰ و ۱ هستند که هر چه به ۱ نزدیکتر باشند عملکرد روش پیشنهادی بهتر خواهد بود (سوکولاوا و لاپالمه، ۲۰۰۹). معیار بازخوانی دقت دسته بندی کننده را با توجه به کل رکوردها نشان می‌دهد. معیار دقت، دقت دسته بندی کننده را با توجه به کل مواردی نشان می‌دهد که توسط دسته بندی کننده پیشنهاد شده است. معیار بازخوانی کارایی دسته بنده را نشان می‌دهد و معیار دقت اساساً مبتنی بر دقت پیش بینی دسته بندی کننده می‌باشد و مبین آن است که به چه میزان می‌توان به خروجی دسته بندی کننده اعتماد نمود. معیار F که ترکیب دو معیار بازخوانی و دقت است در مواردی مورد استفاده قرار می‌گیرد که نتوان اهمیت ویژه ای را برای هر یک از دو معیار بازخوانی و دقت نسبت به یکدیگر قائل شد. معیار نرخ دقت که دقت یک دسته بندی کننده را مورد ارزیابی قرار می‌دهد و نشان می‌دهد دسته بند طراحی شده چند درصد از کل مجموعه رکوردهای آزمایشی را به درستی طبقه بندی می‌کند. معیار نرخ خطا که برعکس نرخ دقت است که کمترین مقدار آن برابر صفر (بهترین کارایی) و بیشترین مقدار آن یک (ضعیف ترین کارایی) می‌باشد (صنعی آباده و محمودی، ۱۳۹۴).

معیارهای ارزیابی:

$$\frac{TN}{FP+TN} = \text{بازخوانی}$$

رابطه ۴)

$$\frac{TN}{TN+FN} = \text{دقت} \quad \text{رابطه (۵)}$$

$$\frac{2*Recall*Precision}{Recall+Precision} = F \text{ معیار} \quad \text{رابطه (۶)}$$

معیارهای ارزیابی الگوریتم دسته بندی کننده:

$$\frac{TN+TP}{TN+FN+TP+FP} = \text{نرخ دقت}^1 \quad \text{رابطه (۷)}$$

$$\frac{FN+FP}{TN+FN+TP+FP} = 1 - \text{نرخ دقت} = \text{نرخ خطا}^2 \quad \text{رابطه (۸)}$$

همچنین در این پژوهش از سطح زیر نمودار^۳ (AUC) جهت ارزیابی مدل استفاده شده است. سطح زیر نمودار یکی دیگر از معیارهای مهمی است که جهت مشخص نمودن میزان کارایی الگوریتم های دسته بندی کننده استفاده می شود. این معیار نشان دهنده سطح زیر منحنی مشخصه عملکرد گیرنده^۴ (ROC) می باشد. هرچه معیار AUC برای یک دسته بندی کننده بزرگ تر باشد کارایی الگوریتم بهتر است. جهت رسم نمودار ROC از دو معیار DR و FAR استفاده می شود که توجه به دسته مثبت را نشان می دهند.

$$\frac{TP}{FN+TP} = DR \quad \text{رابطه (۹)}$$

$$\frac{TP}{TP+FP} = FAR \quad \text{رابطه (۱۰)}$$

در این نمودار محور X نشان دهنده FPR یا همان نرخ تشخیص غلط دسته مثبت^۵ و محور Y نشان دهنده TPR یا نرخ تشخیص صحیح دسته مثبت^۶ می باشد. منحنی های ROC

1. Accuracy-Rate
2. Error Rate
3. Area Under Curve
4. Receiver Operating Characteristic (ROC) Curve
5. False Positive Rate (FPR)
6. True Positive Rate (TPR)

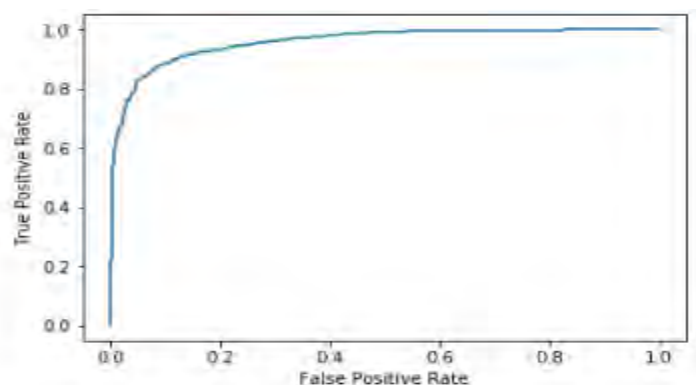
رفتار یک دسته را بدون توجه به توزیع دسته یا هزینه خطا نشان می دهند و کارایی دسته بندی کننده را از این معیارها جدا می کنند. در صورتیکه الگوریتم دسته بندی کننده ایده آل بوده و کلیه نمونه های مثبت را درست دسته بندی نماید مقدار معیار AUC برای آن ۱ خواهد بود (فاوکت^۱، ۲۰۰۳). جهت انتخاب الگوریتم بهینه جهت مدل سازی برای تحلیل ادراکات ابتدا هر یک از الگوریتم ها پیاده سازی گردیدند و سپس با استفاده از روش از روش رای اکثریت دسته بندی کننده ها با یکدیگر ترکیب گردیدند. همانطور که نتایج جدول زیر نشان می دهند دسته بندی کننده پرسپترون چند لایه با دقت ۰٫۸۹۲ دارای دقت بالاتری نسبت به سایر دسته بندی کننده ها و رویکرد ترکیبی می باشد و جهت مدلسازی برای تحلیل ادراکات نظرات استفاده می شود.

در جدول ۴ نتایج ارزیابی هر یک از الگوریتم های دسته بندی کننده و رویکرد ترکیبی ارائه شده است.

جدول ۴: ارزیابی دسته بندی کننده ها

نرخ خطا	نرخ دقت	معیار F	بازخوانی	دقت	دسته بندی کننده
۰/۱۱۲	۰/۸۸۸	۰/۸۸۷	۰/۸۸۵	۰/۸۸۸	Linear Support Vector Classification (Linear SVC)
۰/۱۲۲	۰/۸۷۸	۰/۸۷۷	۰/۸۷۷	۰/۸۷۹	Multinomial Naive Bayes (MNB)
۰/۱۰۸	۰/۸۹۲	۰/۸۹۲	۰/۸۹۲	۰/۸۹۲	Multilayer perceptron (MLP)
۰/۱۱۵	۰/۸۸۵	۰/۸۸۵	۰/۸۸۴	۰/۸۸۶	Ensemble (Majority Voting)

همچنین در شکل ۲ نمودار ROC دسته بندی کننده منتخب جهت تحلیل ادراکات نشان داده شده است.



شکل ۲: نمودار ROC مربوط به دسته بندی کننده بردار پشتیبان تصمیم

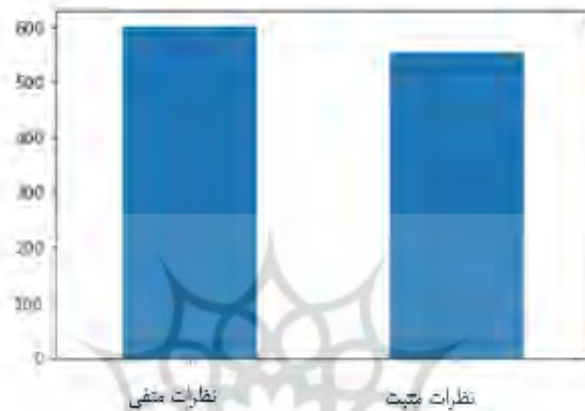
تحلیل ادراکات

در مرحله پایانی، مدل منتخب بر روی نظرات مستخرج از سایت دیجی کالا پیاده سازی می شود و نظرات مثبت و منفی مشخص می گردند. در جدول ۵ نمونه از امتیاز حاصل از تحلیل ادراکات نظرات کاربران ارائه شده است. ستون "دسته بندی موضوعی نظر" در جدول ۵ بر اساس گروه بندی نظرات از سایت دیجی کالا ارائه شده است.

جدول ۵: نمونه ای از نظرات مثبت و منفی

دسته بندی موضوعی نظر	امتیاز نظر (مثبت / منفی)	نظر
کیفیت ساخت	۰ (مثبت)	اپل گوشی خوش ساختی است.
امکانات و قابلیت ها	۱ (منفی)	کیفیت صدای طرف مقابل در زمان مکالمه خوب نیست.
امکانات و قابلیت ها	۱ (منفی)	اصلا گوشی خوبی نیست سیستم عامل تکراری دارد.
طراحی و ظاهر	۱ (منفی)	بیرون زده گی لنز دوربین در طراحی خوب نیست.
امکانات و قابلیت ها	۰ (مثبت)	گوشی سامسونگ VS قابلیت هایش بیشتر است.
امکانات و قابلیت ها	۰ (مثبت)	این گوشی را با گوشی هایی مانند سامسونگ مقایسه کنند از نظر کار آمدی به طور فاحشی شکست می خورند.

در شکل ۳ تعداد نظرات مثبت و منفی ارائه شده است که نشان می‌دهد در مجموع نظرات مشتریان در تلفن‌های همراه ارائه شده در سایت دیجی کالا مثبت نمی‌باشد. همانطور که اطلاعات جدول ۶ نشان می‌دهد در مجموع نظر کاربران در مورد سهولت استفاده و امکانات و قابلیت‌های تلفن‌های همراه مثبت است و در مورد سایر دسته‌بندی‌ها یعنی ارزش خرید نسبت به قیمت، نوآوری، طراحی و ظاهر و کیفیت ساخت گوشی‌ها نظر کاربران مثبت نمی‌باشد.



شکل ۳: تعداد نظرات مثبت و منفی

در جدول ۶ درصد بازخورد مثبت و منفی هر یک از دسته‌بندی‌های موضوعی ارائه شده است.

جدول ۶: نمونه‌ای از نظرات مثبت و منفی

درصد منفی	درصد مثبت	دسته‌بندی موضوعی نظر
۹۳	۷	ارزش خرید نسبت به قیمت
۴۵	۵۵	امکانات و قابلیت‌ها
۴۲	۵۸	سهولت استفاده
۵۲	۴۸	طراحی و ظاهر
۹۰	۱۰	نوآوری
۵۱	۴۹	کیفیت ساخت

شکل ۴ و ۵ نمودار ابر کلمات^۱ در مورد نظرات مثبت و منفی نشان داده شده است. همانطور که از شکل مشخص است کلماتی چون "دوربین"، "کیفیت"، "طراحی" و "سیستم عامل" در نظرات مثبت و در کلماتی مانند "باتری"، "قیمت"، "صفحه نمایش" در نظرات منفی دارای فراوانی بیشتری می‌باشند.



شکل ۵: نمودار ابر کلمات نظرات منفی

شکل ۴: نمودار ابر کلمات نظرات مثبت

نتیجه‌گیری و پیشنهادها

طی چند سال گذشته با توسعه سریع رسانه‌های اجتماعی مردم بیشتر ادراکات، عقاید و نگرش‌های خود را در محیط وب بیان می‌کنند که این موضوع باعث ایجاد داده‌های متنی حجیم در شبکه‌های اجتماعی در ارتباط با موضوعات مختلف شده است. دسترسی به این داده‌ها مسیر فرآیند تصمیم‌گیری افراد را تغییر داده است. این نظرات و داده‌ها بصورت ساختار نیافته می‌باشند که حاوی اطلاعات ارزشمند در ارتباط با ادراکات و ترجیحات کاربران می‌باشند. امروزه تحلیل داده‌های متنی و غیرساخت یافته در محیط وب و شبکه‌های

1. Word Cloud Diagram

اجتماعی از چالش‌های اصلی در تجارت الکترونیک و فروش محصولات دیجیتال می‌باشد. خودکارسازی استخراج ادراکات و نظرات از راهکارهای پیشنهادی جهت استخراج دانش نهفته در این داده‌ها و استفاده بهینه از آنها می‌باشد.

روش‌های مورد استفاده در این تحقیق برای تحلیل ادراکات و نشان دادن قطبیت ادراکات مفید می‌باشند. در این پژوهش تلاش گردید تا با استفاده از رویکرد یادگیری نظارت شده و استفاده از سه الگوریتم پرسپترون چند لایه، بیزین ساده چند جمله‌ای و طبقه بندی کننده بردار پشتیبان و رویکرد ترکیبی (رای اکثریت) مدلی جهت دسته بندی نظرات مثبت و منفی خریداران تلفن همراه سایت دیجی کالا ارائه شود. نتایج پژوهش نشان می‌دهند که دسته بندی کننده پرسپترون چند لایه با دقت ۰.۸۹۲، مدل بهتری را جهت تحلیل ادراکات و نظرات کاربران ارائه می‌دهد. همچنین در مجموع نظرات کاربران در مورد تلفن همراه ارائه شده در سایت دیجی کالا منفی می‌باشد. آیت‌هایی چون ارزش خرید نسبت به قیمت، نوآوری، طراحی و ظاهر و کیفیت ساخت دارای بیشترین میزان نظر منفی می‌باشند که لازم است برای انتخاب محصولات جدید بر این موارد تمرکز بیشتری داشته باشد. مدل پیشنهادی می‌تواند در سایت‌های تجارت الکترونیک مانند دیجی کالا پیاده سازی شود و خروجی آن به صورت سیستماتیک توسط کاربران قابل مشاهده باشد که در نهایت می‌تواند منجر به تصمیم‌گیری آگاهانه برای خریداران و شرکت‌های ارائه دهنده محصولات باشد. این روش می‌تواند برای شناسایی الگوهای روانشناختی مشتریان برای افزایش رضایت و کسب بازارهای جدید مورد استفاده قرار گیرد.

در این پژوهش از رویکرد کلاسیک استخراج ویژگی‌ها جهت دسته بندی نظرات استفاده شده است که در تحقیقات آتی می‌توان از سایر روش‌های ترکیبی چون روش وزن‌دهی و یادگیری عمیق جهت دست‌یابی به نتایج با دقت بالاتر استفاده نمود. به علاوه نتایج تحلیل نظرات و ادراکات کاربران منبع ارزشمندی برای تصمیم‌گیری شرکت‌ها جهت

تدوین استراتژی برای ارائه محصولات و خدمات جدید به مشتریان می‌باشد. همچنین نتایج تحلیل ادراکات نظرات کاربران می‌تواند به عنوان ورودی برای طراحی سیستم‌های هوشمند و تصمیم‌یار باشد. برای تحلیل قطبیت کلمات، پیشنهاد می‌گردد از روش‌های جدید استفاده گردد و بجای در نظر گرفتن +۱ و -۱ بازه اعداد بین این دو عدد در نظر گرفته شود. به علاوه استفاده از واژگان پویا^۱ بر بهبود عملکرد تحلیل موثر می‌باشد. همچنین پیشنهاد می‌گردد این موضوع که آیا ایموجی‌ها^۲ با ادراکاتی که در متون ارائه می‌شود، مطابقت دارد یا خیر در پژوهش‌های آتی مورد بررسی قرار گیرد. شناسایی نظرات ساختگی^۳ نیز می‌تواند بر بهبود تحلیل‌ها اثرگذار باشد.



-
1. Dynamic Lexicon
 2. Emojis
 3. Fake

منابع

- بخشی زاده برج، ک.، حاجی جعفر، ع.، و نصیری، ح. (۱۳۹۷). ترسیم نقشه ذهنی مشتریان فروشگاه اینترنتی دیجی کالا با استفاده از تکنیک استخراج استعاره‌ای زالتمن (زیمت). مدیریت بازرگانی، ۴۹۷۲.
- پیکری، ن.، یعقوبی، س.، و طاهری، ح. (۱۳۹۴). تحلیل احساسات در شبکه اجتماعی توییتر با تکنیک متن کاوی. کنفرانس بین‌المللی وب پژوهی. تهران: دانشگاه علم و فرهنگ.
- حاج سید جوادی، ش.، و جلالی، م. (۱۳۹۴). ارائه روشی کارا برای تجزیه و تحلیل احساسات توییتر براساس ترکیب روش‌های یادگیری ماشین و شباهت معنایی. دومین کنفرانس بین‌المللی فن آوری.
- خدیور، آ.، و عباسی، ف. (۱۳۹۸). متن کاوی با تمرکز بر تحلیل احساسات. تهران: نگاه دانش. زمانی، م.، و سرخپور، ب. (۱۳۹۳). داده کاوی متون فارسی با نگرش مدیریت دانش. هفتمین کنفرانس ملی و اولین کنفرانس بین‌المللی مدیریت دانش. تهران.
- صنّعی آباده، م.، و محمودی، س. (۱۳۹۴). داده کاوی کاربردی. تهران: نیاز دانش.
- عباسی، ف.، سهرابی، ب.، مانیان، ا.، و خدیور، آ. (۱۳۹۶). ارائه مدلی جهت دسته بندی احساسات خریداران کتاب با استفاده از رویکرد ترکیبی. مطالعات مدیریت کسب و کار هوشمند، ۹۴-۶۵.
- عبدی قویدل، ه.، وزیر نژاد، ب.، و بحرانی، م. (۱۳۹۱). برجسب زنی موضوعی متون فارسی. چهارمین کنفرانس فناوری اطلاعات و دانش. بابل.
- لطفی آذری داریان، س.، و جاویدان، ر. (۱۳۹۵). استفاده از روشهای داده‌کاوی به منظور تسهیل جستجو در موتورهای جستجوگر متنی. بیست و چهارمین کنفرانس برق ایران، (ص. ۲۸۱۷-۲۸۰۹). شیراز.
- نجف زاده، م.، راحتی قوچانی، س.، و قائمی، ر. (۱۳۹۷). یک چارچوب نظارتی مبتنی بر لغت نامه وفقی خودساخت جهت تحلیل نظرات فارسی. پردازش علائم و داده‌ها، ۱۰۱-۸۹.

الهی، ش.، قدس الهی، ا.، و ناجی، ح. (۱۳۹۳). ارائه مدل ترکیبی شبکه های عصبی با بهره گیری از یادگیری جمعی به منظور ارزیابی ریسک اعتباری. *انجمن فناوری اطلاعات و ارتباطات ایران*، ۲۸-۱۱.

Abdmanaf, S., Mustapha, N., Sulaiman, M., Azura Husin, N., Zainuddin, M., & Shafri, H. (2017). Majority Voting of Ensemble Classifiers to Improve Shoreline Extraction of Medium Resolution Satellite Images. *Journal of Theoretical and Applied Information Technology*, pp: 4394-4405.

Bhatt, A., Patel, A., Chheda, H., & Gawande, K. (2015). Amazon Review Classification and Sentiment Analysis. *International Journal of Computer Science and Information Technologies*, pp: 5107-5110.

Cambria, E., Havasi, C., & Hussain, A. (2012). SenticNet 2: A semantic and affective resource for opinion mining and sentiment analysis. *Association for the Advancement of Artificial*, pp:202-207.

Chakraborty, R. (2013). Domain Keyword Extraction Technique: A New Weighting Method Based on Frequency Analysis. *Computer Science & Information Technology*, pp: 109-118.

Dietterich, T. (2000). *Ensemble methods in machine learning*. Berlin, Heidelberg: Springer.

Fang, X., & Zhan, u. (2015). Sentiment analysis using product review data. *Journal of Big Data*, 2(5), pp:1-14.

Fawcett, T. (2003). *ROC Graphs: Notes and Practical Considerations for Data Mining Researchers*. Intelligent Enterprise Technologies Laboratory.

Filho, P., & Pardo, T. (2013). NILC USP: A Hybrid System for Sentiment Analysis in Twitter Messages. *Second Joint Conference on Lexical and Computational Semantics (*SEM)* (pp: 568-572). Association for Computational Linguistics.

Internet World Stats. (2019). *Internet World Stats Usage and Population Statics*. <https://www.internetworldstats.com/top20.htm>

Kannan, S., & Gurusamy, V. (2014). *Preprocessing Techniques for Text Mining*.

- Liang, T.-P., Li, X., Yang, C.-T., & Wang, M. (2015). What in Consumer Reviews Affects the Sales of Mobile Apps: A Multifacet Sentiment Analysis Approach. *International Journal of Electronic Commerce*, pp:226-260.
- Liu, B. (2012). *Sentiment Analysis and Opinion Mining (Synthesis Lectures on Human Language Technologies)*. Williston: Morgan & Claypool Publishers.
- Liu, B. (2015). *Sentiment analysis: Mining opinions, sentiments, and emotions*. Cambridge University Press.
- Mitchell, T. (2015). *Generative and Discriminative Clasifiers: Naive bayes and logestic regression*. T. Mitchell, Machine Learning .
- Murugavalli, S., Bagirathan, U., Saiprassanth, R., & Arvindkumar, S. (2017). Feedback analysis using Sentiment Analysis for E-commerce. *International Journal of Latest Engineering Research and Applications (IJLERA)*, pp: 84-90.
- Pandey, A., Rajpoor, D., & Saraswat, M. (2017). Twitter sentiment analysis using hybrid cuckoo search method. *Information Processing and Management*, pp:764-769.
- Popescu, M.-C., Balas, V., Perescu-Popescu, L., & Mastorakis, N. (2009). Multilayer Perceptron and Neural Networks. *WSEAS Transactions on Circuits and Systems*, 579-588.
- Raghavan, V., & Gwang, J. (1989). A Critical Investigation of Recall and Precision as Measures of Retrieval System Performance. *ACM Transactions on Information Systems*, pp: 206-229.
- Schonlau, M., & Guenther, N. (2016). Text Mining Using N-Grams. *SSRN Electronic Journal*.
- Sokolova, M., & Lapalme, G. (2009). A systematic analysis of performance measures for classification tasks. *Information Processing and Management*, pp:427-437.
- Suen, C., & Lam, L. (2000). *Multiple Classifier Combination Methodologies*. MCS 2000, LNCS 1857, 52-66.
- Vidya, N., Fanany, M., & Budi, I. (2015). Twitter Sentiment to Analyze Net Brand Reputation of Mobile Phone Providers. *Procedia Computer Science*, pp:519-526.

- Vijayarani, S., Ilamathi, M., & Nithya, M. (2015). Preprocessing Techniques for Text Mining - An Overview. *International Journal of Computer Science & Communication Networks*, pp: 7-16.
- Wang, Z. (2017). *The Evaluation of Ensemble Sentiment Classification Approach on Airline Services Using Twitter*. Dublin: Dublin Institute of Technology.
- Zhang, L., Hua, K., Wang, H., Qian, G., & Zhang, L. (2014). Sentiment Analysis on Reviews of Mobile Users. *Procedia Computer Science*, pp:458-465.
- Zhang, Y., Ren, W., Zhu, T., & Faith, E. (2019). *MoSa: A Modeling and Sentiment Analysis System for Mobile Application Big Data*. Symmetry.
- Zhao, L., Huang, M., Yao, Z., Su, R., Jiang, Y., & Zhu, X. (2016). Semi-Supervised Multinomial Naive Bayes for Text Classification by Leveraging Word-Level Statistical Constraint. *Proceeding of the Thirtieth AAAI Conference on Artificial Intelligence*, pp: 2877-2883.

