

# تشخیص ویژگی‌های ضمنی با استفاده از قواعد نحوی زبان فارسی و خوشه‌بندی صفات

عاطفه محمدی<sup>۱</sup>، دانشجوی دکتری؛ مهدی یزدیان دهکردی<sup>۲</sup>، استادیار؛ محمدعلی نعمت‌بخش<sup>۳</sup>، دانشیار

۱- گروه مهندسی کامپیوتر، پردیس فنی و مهندسی، دانشگاه یزد، یزد، ایران، atefehmohammadi@stu.yazd.ac.ir

۲- گروه مهندسی کامپیوتر، پردیس فنی و مهندسی، دانشگاه یزد، یزد، ایران، yazdian@yazd.ac.ir

۳- دانشکده مهندسی کامپیوتر، دانشگاه اصفهان، اصفهان، ایران، nematbakhsh@eng.ui.ac.ir

**چکیده:** به‌طور معمول وقتی فردی قصد خرید یک محصول برخط را دارد، نظرات و یادداشت‌های نوشته‌شده توسط سایر افراد در مورد محصول را بررسی می‌کند و این مسئله تاثیر بسزایی در تصمیم‌گیری فرد برای خرید دارد. تجزیه و تحلیل نظرات کاربران که با عنوان تجزیه و تحلیل احساس یا نظرکاوی شناخته می‌شود، یکی از داغ‌ترین موضوعات تحقیقاتی در علم کامپیوتر است. هدف اصلی نظرکاوی، استخراج نظرات افراد درباره‌ی ویژگی‌های یک موجودیت یا کالا است. در این تحقیق، یک راه‌کار نظرکاوی بدون نظارت و مبتنی بر استخراج ویژگی‌های ضمنی برای محصولات در زبان فارسی ارائه شده است. استخراج ویژگی‌های ضمنی یکی از مراحل دشوار در تحلیل احساسات مبتنی بر ویژگی می‌باشد. در بیشتر پژوهش‌های پیشین از اطلاعات آماری برای ایجاد ماتریس هم‌رخداد و سپس تشخیص ویژگی‌های ضمنی استفاده شده است. در این پژوهش در کنار اطلاعات آماری، از قواعد نحوی زبان و خوشه‌بندی صفات جهت بهبود ماتریس هم‌رخداد بین کلمات (ویژگی‌ها و احساسات) بهره گرفته شده است. ارزیابی‌های انجام‌شده بر روی داده‌های واقعی و استخراج‌شده از نظرات کاربران در سایت دیجی‌کالا، نشان می‌دهند که روش ارائه‌شده در این مقاله به نرخ بازخوانی و دقت بهتری نسبت به کارهای قبلی دست یافته است.

**واژه‌های کلیدی:** نظرکاوی، ویژگی ضمنی، ماتریس هم‌رخداد، خوشه‌بندی.

## Identification of Implicit Features using Persian Language Rules and Sentiments Clustering

A. Mohammadi<sup>1</sup>, PhD Student; M. Yazdian-Dehkordi<sup>2</sup>, Assistant Professor;  
M. Nematbakhsh<sup>3</sup>, Associated Professor

1- Faculty of Engineering, Yazd University, Yazd, Iran, Email: atefehmohammadi@stu.yazd.ac.ir

2- Faculty of Engineering, Yazd University, Yazd, Iran, Email: yazdian@yazd.ac.ir

3- Faculty of Computer Engineering, University of Isfahan, Isfahan, Iran, Email: nematbakhsh@eng.ui.ac.ir

**Abstract:** Typically, when someone wants to buy an online product, he/she reviews the comments and notes written by others about the product. Clearly, it has a profound impact on the person's decision to buy or not to buy the product. Sentiment analysis or opinion mining is one of the hot research topics in computer science. The main purpose of sentiment analysis is to extract opinions of individuals about the characteristics of an entity such as a product. In this research, an unsupervised opinion mining is proposed for Persian products based on implicit feature extraction which is a critical step in sentiment analysis. In most previous studies, statistical information is utilized to create a co-occurrence matrix and determine the implicit features. In this paper, we benefit from syntactic rules and sentiment clustering in conjunction with statistical information to construct an efficient co-occurrence matrix between features and sentiment words. The evaluation results provided on a real-world dataset, extracted from Digikala website, indicates that the proposed method achieves higher recall and precision compared to the previous studies.

**Keywords:** Opinion mining, implicit feature, co-occurrence matrix, clustering.

تاریخ ارسال مقاله: ۱۳۹۷/۰۵/۰۶

تاریخ اصلاح مقاله: ۱۳۹۷/۰۹/۲۵ و ۱۳۹۷/۱۱/۱۳

تاریخ پذیرش مقاله: ۱۳۹۷/۱۱/۲۰

نام نویسنده مسئول: مهدی یزدیان دهکردی

نشانی نویسنده مسئول: ایران، یزد، صفائیه، بلوار دانشگاه، دانشگاه یزد، پردیس فنی مهندسی، گروه مهندسی کامپیوتر.

## ۱- مقدمه

با رشد چشم‌گیر محبوبیت وب، تعداد نظرات برخط در مورد محصولات مختلف با سرعت بالایی در حال افزایش است. تعداد زیادی از وبلاگ‌ها، تالارهای پرس‌مان و وب سایت‌هایی مانند دیجی‌کالا به کاربران اجازه می‌دهند تا نظراتشان را در مورد محصولات و سرویس‌های مختلف بیان کنند. استخراج دانش از نظرات ارائه شده در مورد محصول توسط یک سیستم نظرکاوی [۱، ۲] انجام می‌شود و پیشنهاد یک محصول مناسب به کاربر توسط یک سیستم توصیه‌گر [۳، ۴] انجام می‌شود.

نظرکاوی، تولیدکنندگان را قادر می‌سازد تا بتوانند ارزیابی مناسبی از رضایت کاربران داشته و با توجه به آن محصولات خود را ارتقا دهند [۱]. از آن جایی که بررسی تمامی نظرات برای کاربران امکان‌پذیر نیست، ارائه روش‌هایی جهت استخراج دانش از نظرات و ارائه آن به صورت خلاصه برای کاربران اهمیت زیادی پیدا کرده است.

از این‌رو، بسیاری از محققان، الگوریتم‌های موثری را برای کشف دانش از نظرات، تحت عنوان نظرکاوی<sup>۱</sup>، ارائه کرده‌اند. امروزه تجزیه و تحلیل احساسات<sup>۲</sup> (که نظرکاوی هم نامیده می‌شود)، مطالعه‌ی نظرات، احساسات و نگرش‌های افراد درباره‌ی موجودیت‌هایی مانند محصولات و سرویس‌ها است.

تجزیه و تحلیل احساسات در سه سطح سند، جمله و ویژگی مورد بررسی قرار می‌گیرد [۲، ۵، ۶]. در تجزیه و تحلیل احساسات در سطح سند، گرایش کلی سند با فرض این‌که هر سند حاوی نظرانی روی یک موجودیت منفرد است، به صورت مثبت یا منفی تعیین می‌شود. تجزیه و تحلیل احساسات در سطح جمله<sup>۳</sup>، با فرض این‌که هر جمله، دارای یک ویژگی یا موجودیت است گرایش هر جمله، اعم از مثبت، منفی یا خنثی را مشخص می‌کند. در صورتی‌که، یک جمله ممکن است حاوی نظرانی درباره‌ی چندین ویژگی از یک موجودیت باشد. تجزیه و تحلیل احساسات در سطح سند و جمله به‌طور دقیق مشخص نمی‌کنند که افراد به چه مشخصه‌ای از محصول تمایل و یا عدم تمایل بیشتری دارند. در تجزیه و تحلیل احساسات در سطح ویژگی یا موجودیت، هر نظر شامل احساسات مثبت یا منفی و مجموعه‌ای از ویژگی‌ها است. درک اهمیت هدف‌های حاوی نظر به درک بهتر مسئله‌ی تجربه و تحلیل احساسات کمک می‌کند. به‌عنوان مثال، جمله‌ی نظری «گرچه ظاهرش جالب نبود، اما من هنوز عاشق این گوشی هستم»، دارای لحنی مثبت است؛ اما نمی‌توانیم بگوییم که این جمله به‌طور کامل دارای گرایش مثبت است. در واقع جمله دارای گرایش مثبت برای «گوشی» و دارای گرایش منفی برای «ظاهر گوشی» است. هدف اصلی تجزیه و تحلیل احساسات در سطح ویژگی، کشف احساسات مربوط به موجودیت‌ها و یا ویژگی‌های آن‌ها است [۷].

روش‌های تشخیص ویژگی به‌طور کلی به دو دسته‌ی بانظارت<sup>۴</sup> و بدون نظارت<sup>۵</sup> تقسیم می‌شوند [۷، ۸]. روش‌های بانظارت، روش‌هایی هستند که در آن‌ها، مجموعه داده‌ی آموزشی دارای برچسب صحیح و از قبل مشخص هستند. از جمله این تکنیک‌ها به بیز ساده، ماشین برداری پشتیبان<sup>۶</sup> و حداکثر انتروپی<sup>۷</sup> می‌توان اشاره کرد. از آن جایی که مسئله‌ی استخراج بانظارت ویژگی از متون، یک مسئله‌ی طبقه‌بندی است؛ این مدل‌ها برای استخراج ویژگی‌های صریح<sup>۸</sup> به‌کار می‌رود. در ابتدا، مجموعه‌ای از برچسب‌ها با داشتن مجموعه‌ای از کلمات در یک جمله تهیه می‌شوند که ویژگی یا ویژگی نبودن یک کلمه را مشخص می‌کند. در گام بعدی، فرایند برچسب‌گذاری واژگان بر روی جملات انجام می‌شود. پس از ایجاد این مجموعه داده، مدل حداکثر انتروپی به‌عنوان مثال آموزش داده می‌شود و بر روی داده‌های بدون برچسب آزمایش می‌شود. هدف از به‌کارگیری مجموعه‌ی آموزشی، تشخیص و استخراج ویژگی‌های مناسب توسط سیستم است [۹، ۱۰]. هرچند که تهیه‌ی این مجموعه داده وظیفه‌ی سخت و طاقت‌فرسا و نیازمند صرف زمان زیادی است اما در عوض نتایج تولیدشده توسط این روش از دقت بیشتری برخوردار هستند. مدل‌های بانظارت وابسته به دامنه بوده و در صورت انتقال به دامنه‌ی دیگر بایستی مجدداً آموزش ببینند. مدل‌های بدون نظارت، نیاز به مجموعه داده‌ی آموزشی برچسب‌خورده ندارند و در مقایسه با روش‌های بانظارت از دقت کمتری برخوردار هستند [۷، ۸].

ویژگی‌ها می‌توانند به دو روش صریح و ضمنی ظاهر شوند. اسم یا عبارات اسمی که به‌طور صریح در جملات رخ می‌دهند، ویژگی صریح نامیده می‌شوند. به‌عنوان مثال در جمله‌ی «کیفیت ساخت خوبی دارد»، «کیفیت ساخت» به عنوان ویژگی صریح استخراج می‌شود. ویژگی که اسم یا عبارت اسمی نباشد را ویژگی ضمنی<sup>۹</sup> می‌نامند که به‌طور معمول به شکل صفت یا قید ظاهر می‌شود. به‌عنوان مثال در جمله‌ی «این گوشی خیلی گران است»، به‌طور ضمنی به ویژگی «قیمت» اشاره دارد.

در زمینه تشخیص ویژگی‌های صریح کارهای زیادی نسبت به ویژگی‌های ضمنی انجام شده است. متأسفانه طبق بررسی‌های ما، در زبان فارسی تعداد کارهای انجام شده بر روی ویژگی‌های ضمنی انگشت‌شمار بوده است. در داده‌های واقعی وب، وجود ویژگی‌های ضمنی در نظرات کاربران غیرقابل اجتناب است. بنابراین انجام پژوهش در حوزه ویژگی‌های ضمنی نیز از اهمیت زیادی برخوردار است. در روش‌های استخراج ویژگی ضمنی با توجه به عدم وجود ویژگی در جملات نسبت به روش استخراج ویژگی صریح دارای چالش بیشتری هستند. در این پژوهش، برای تشخیص ویژگی‌های ضمنی از ماتریس هم‌رخداد<sup>۱۰</sup> استفاده می‌شود. در پژوهش‌های پیشین از اطلاعات آماری برای ایجاد ماتریس هم‌رخداد استفاده شده است. در این پژوهش برای استفاده از اطلاعات آماری، علاوه بر معیار LRT، از معیار فاصله‌ی ویژگی صریح از کلمه‌ی احساسات، قواعد نحوی زبان، گراف وابستگی جملات و

متنی استفاده کرده‌اند. تفاوت یادگیری عمیق نسبت به روش‌های مرسوم در این حوزه این است که یادگیری عمیق برای یادگیری ویژگی‌ها نیازمند داده‌های زیاد است.

ع سگریان و همکاران [۱۴]، یک wordnet فارسی کامل، با دقت و بازخوانی مناسب ارائه داده‌اند. آن‌ها با به کارگیری روش‌های انتخاب ویژگی شناخته شده مانند  $CHI^{14}$ ،  $IG^{15}$ ،  $MI^{16}$ ، ۲۰۰۰ ویژگی با بالاترین امتیاز را استخراج کرده‌اند و سپس با کمک روش‌های یادگیری ماشین، یک طبقه‌بندی احساس برای نظرات متنی فارسی ارائه داده‌اند. در برخی از روش‌ها از روش یادگیری نیمه‌نظارتی نیز استفاده شده است.

روش‌های بانظارت نیازمند داده‌های آموزشی برجسب‌دار هستند که تهیه آن‌ها در داده‌های حجیمی مانند وب بسیار پرهزینه است. به علاوه این روش‌ها محدود به دامنه نوع محصولات هستند. در مقابل روش‌های بانظارت، روش‌های بدون نظارت برای استخراج ویژگی نیز ارائه شده‌اند. در سال ۲۰۱۸، دراگونی [۱۵] برای استخراج ویژگی از یک پیکره‌ی نظرات بزرگ استفاده کرده‌اند که شامل ۳۵ میلیون نظر است. در این روش برای استخراج ویژگی‌ها از یک تکنیک پردازش زبان طبیعی بی‌نظارت استفاده شده که شامل دو گام است. در گام اول، برجسب‌زن کلمات و وابستگی‌های بین کلمات از متن استخراج شده و در گام بعدی، واژه‌هایی که به‌عنوان اسم هستند به کلمات احساس مرتبط می‌شوند. سپس، ویژگی‌های صریح خوشه‌بندی می‌گردند.

در روش ارائه شده توسط سو و همکاران [۱۶]، صفات به‌عنوان احساس و اسم، عبارات اسمی و عبارات‌های فعلی به‌عنوان ویژگی‌های کاندید محصول در نظر گرفته می‌شوند. در گام بعدی، ویژگی‌های کاندید، هرس و ویژگی‌های محصول و کلمات احساس<sup>۱۷</sup> دسته‌بندی می‌شوند. در گام آخر، ارتباطات احساسی پنهان<sup>۱۸</sup> بین گروه‌های ویژگی محصول و کلمات نظر کاوش شده و یک مجموعه‌ی انجمنی<sup>۱۹</sup>، با شناسایی n تا از قوی‌ترین پیوند احساسی بین دو گروه از اشیاء داده‌ای، ایجاد و با استفاده از مجموعه‌ای از پیش ساخته‌شده، ویژگی‌های ضمنی شناسایی می‌گردند.

یک روش کاوش قوانین انجمن هم رخداد<sup>۲۰</sup> جهت تشخیص ویژگی‌های ضمنی توسط‌های و همکاران [۱۷] ارائه شد. کلمات احساس و ویژگی‌ها از جملات صریح موجود در پیکره استخراج و یک ماتریس هم‌رخداد ایجاد می‌شود. سپس برای هر کلمه‌ی احساس، قواعد انجمنی از ماتریس هم‌رخداد ساخته می‌شوند و با استفاده از اندازه‌گیری‌های پشتیبان و اطمینان، قوانین ضعیف به‌طور مناسبی هرس می‌گردند. برای کلمه‌ی احساس داده‌شده‌ای که دارای ویژگی صریح نیست، یک لیست تطبیق یافته از قوانین انجمنی مورد جست‌وجو قرار گرفته تا ویژگی ضمنی مرتبط از آن استخراج شوند.

باقری و همکاران [۷] یک مدل بدون ناظر تشخیص ویژگی و احساس را ارائه داده‌اند که قادر به استخراج ویژگی‌های صریح و ضمنی بر روی زبان انگلیسی است. در این مدل، اسم‌ها بیان‌گر ویژگی و

خوشه‌بندی صفات، برای ایجاد کاراتر ماتریس هم رخداد استفاده شده است؛ تا با به کارگیری این ماتریس بتوان ویژگی‌های ضمنی مناسب‌تری را استخراج نمود. برای ارزیابی روش ارائه شده، از داده‌های واقعی که از سایت معروف و پربازدید فرو شگاه اینترنتی دیجی کالا استخراج شده است، استفاده گردید. نتایج به دست آمده بر روی این داده‌ها نشان می‌دهند که روش ارائه شده توانسته است با استخراج ویژگی‌های ضمنی مناسب‌تر، کارایی بهتری نسبت به روش‌های پیشین به همراه داشته باشد.

در ادامه‌ی پژوهش، در بخش ۲ مروری بر تاریخچه موضوع و کارهای پیشین صورت گرفته است. در بخش ۳، به ارائه روش پیشنهادی و مراحل انجام کار پرداخته شده و در بخش ۴، نتایج ارزیابی روش پیشنهادی بیان خواهند شد و در نهایت در بخش ۵، جمع بندی و کارهای آینده شرح داده می‌شوند.

## ۲- مروری بر کارهای پیشین

همان‌گونه که در بخش مقدمه عنوان شد، روش‌های استخراج ویژگی در سیستم‌های نظرکاوی را می‌توان به دو دسته عمده شامل روش‌های بانظارت و بدون نظارت تقسیم‌بندی کرد.

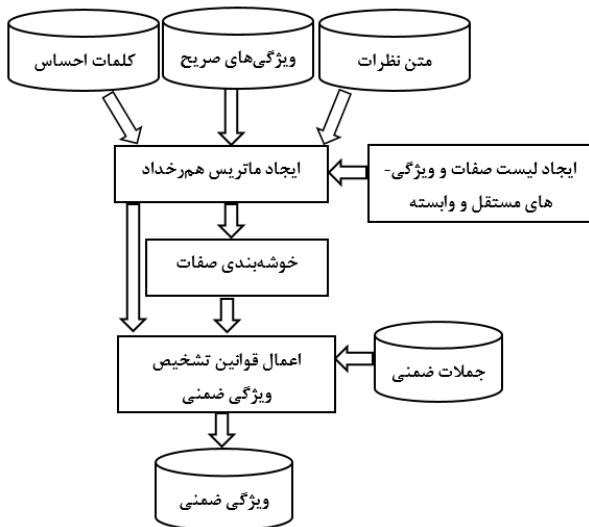
در روش بانظارت داده‌های آموزش دارای برجسب هستند و به‌طور معمول از طبقه‌بندی‌کننده<sup>۱۱</sup> برای ساخت سیستم استفاده می‌شود.

شیخ حسنی و منصوری‌زاده [۱۱] با استفاده از میدان تصادفی شرطی (CRF)<sup>۱۲</sup>، جنبه‌های نهفته در متن نظرات را استخراج کرده‌اند. با توجه به انتخاب مدل میدان تصادف شرطی، آن‌ها هر جمله از نظر را به کلمات تشکیل‌دهنده‌اش بخش‌بندی کرده و برای کلماتی غیر از کلمات لیست ایست‌واژه<sup>۱۳</sup> و لیست غیر الفبایی، یک بردار ویژگی ساختند. به این صورت که بر اساس مجموعه دادگان برجسب‌خورده، کلاس هر کلمه را مشخص کرده و به‌دنبال آن، ویژگی‌های استخراج‌شده را قرار داده‌اند. سپس، کلمه‌های هر جمله را با یک خط خالی از هم جدا نموده و مجموعه داده‌های آموزشی را به‌عنوان بردار ورودی به طبقه‌بندی‌کننده CRF داده‌اند.

در سال ۲۰۱۷ اصغر و همکاران [۱۲] چارچوبی پیشنهاد داده‌اند که شامل چهار ماژول اصلی استخراج ویژگی- احساس، گروه‌بندی ویژگی، دسته‌بندی ویژگی- احساس و خلاصه‌سازی ویژگی است. آن‌ها تعدادی الگوی جدید ارائه داده‌اند که در آن اسم، عبارات اسمی و فعل‌ها، بیان‌گر ویژگی‌های کاندید هستند. پس از استخراج ویژگی‌ها، گروه‌بندی ویژگی‌ها را بر اساس اندازه‌گیری شباهت معنایی پیشنهاد داده‌اند. به‌منظور انتساب امتیاز به نظرات داده‌شده در بای ویژگی‌های استخراج‌شده، یک تکنیک ترکیبی برای امتیازدهی احساسات پیشنهاد داده شد که مفاهیم مبتنی بر پیکره و مبتنی بر لغت را ترکیب و در نهایت نقاط مثبت و منفی هر ویژگی را بیان می‌کند.

لیانگ و همکاران [۱۳] در سال ۲۰۱۷ از یادگیری عمیق که یک روش جدید و موفق در حوزه متن‌کاوی بوده، برای استخراج ویژگی‌های

فرایند کلی فلوجارت الگوریتم پیشنهادی که به استخراج ویژگی ضمنی می‌پردازد در شکل ۱، آورده شده است و در ادامه نیز به تفصیل هر یک از مراحل خواهیم پرداخت.



شکل ۱: فلوجارت روش پیشنهادی استخراج ویژگی ضمنی

### ۳-۱- پیش پردازش جملات

در ابتدا، لازم است که متون نظرات پیش‌پردازش شوند. نظرات مربوط به محصول توسط ابزار هضم  $h_{[20]}$  نرمال سازی شده، ریشه‌یابی و سپس برچسب‌زده می‌شوند. به‌عنوان مثال خروجی این عملیات برای جمله‌ی «دوربین مناسبی دارد ولی کیفیت ساخت خوبی ندارد» به صورت «دوربین/Ne/ مناسب/ AJ/ داشت/ V/ ولی/ CONJ/ کیفیت/ Ne/ ساخت/ Ne/ خوب/ AJ/ نداشت/ V» خواهد بود. سپس از ابزار هضم، برای تعیین گراف وابستگی جملات استفاده می‌شود که این گراف، وابستگی بین کلمات یک جمله را از نظر نحوی مشخص می‌کند. جهت بهبود دقت گراف وابستگی، نظرات طولانی با استفاده از کلماتی که برچسب «V» دارند به یک یا چند جمله شکسته می‌شوند که تنها حاوی یک فعل باشند. به‌عنوان مثال، نظر بالا به دو جمله‌ی «دوربین/Ne/ مناسب/ AJ/ داشت/ V» و «ولی/ CONJ/ کیفیت/ Ne/ ساخت/ Ne/ خوب/ AJ/ نداشت/ V» شکسته می‌شود. این گراف وابستگی و نقش نحوی استخراج‌شده از ابزار هضم، در پایگاه داده ذخیره می‌شوند. به‌عنوان مثال نقش نحوی جمله‌ی اول به‌صورت «دوربین/Ne/ SBJ/ مناسب/ NPOSTMOD/ داشت/ ROOT» مشخص می‌شود. گراف وابستگی همان جمله، به‌صورت شکل ۲ آورده شده است.



شکل ۲: نمونه گراف وابستگی جمله

صفت‌ها، قیدها و فعل‌ها بیان‌گر کلمات احساس هستند. ویژگی‌های چندکلمه‌ای توسط روشی به نام FMLR<sup>۲۱</sup> شناسایی می‌شوند. پس از نهایی شدن لیست ویژگی‌ها از دو روش هرس پشتیبان بالا مجموعه و زیر مجموعه برای حذف ویژگی‌های غیر مفید و افزونه استفاده می‌شود. در نهایت، از یک گراف برای تشخیص ویژگی ضمنی استفاده می‌شود.

مدل ارائه شده توسط شوتن و همکاران [۱۸] از ماتریس هم‌رخداد برای تشخیص ویژگی‌های ضمنی استفاده می‌کند. ابتدا داده‌ی آموزشی، بررسی شده و سه لیست ایجاد می‌شوند. (۱) لیستی به نام F از تمامی ویژگی‌های ضمنی (۲) یک لیست به نام O از تمامی کلمات (به‌جای به‌کارگیری کلمات احساس، از همه‌ی کلمات برای تشخیص ویژگی ضمنی استفاده می‌کند) و (۳) یک ماتریس به نام V از هم‌رخدادهای بین ویژگی‌های ضمنی و کلماتی که با هم در یک جمله ظاهر می‌شوند. در داده‌ی آزمایشی، برای هر ویژگی ضمنی  $f_i$  بر اساس مقادیر لیست‌های F، O، و ماتریس C یک نمره محاسبه‌شده که از مجموع هم‌رخداد هر کلمه تقسیم بر تعداد تکرار آن کلمه در جملات به‌دست می‌آید.

باباعلی و همکارش [۱۹] یک روش نظرکاوی بدون نظارت مبتنی بر ویژگی برای محصولات در زبان فارسی ارائه داده‌اند که شامل مراحل شناسایی ویژگی‌های صریح، ویژگی‌های ضمنی و تعیین جهت معنایی نظرات است. در مرحله‌ی استخراج ویژگی صریح، ابتدا ویژگی‌های جذاب و پرتکرار و سپس ویژگی‌های نادر شناسایی می‌شوند. در مرحله‌ی استخراج ویژگی‌های پرتکرار، ویژگی‌های تک و چندکلمه‌ای با حداقل پشتیبان یک درصد به‌دست می‌آیند. سپس، ویژگی‌های ضمنی با استفاده از کاوش قوانین انجمنی هم‌رخداد ارائه‌شده توسط‌های [۱۷] استخراج می‌شوند.

با وجود تحقیقات بسیاری که در زمینه نظر کاوی مبتنی بر ویژگی‌های صریح انجام شده، اما هنوز در حوزه نظرکاوی مبتنی بر ویژگی‌های ضمنی با چالش‌های بسیاری روبه‌رو بوده و نیازمند پژوهش‌های بیشتری در این زمینه است.

### ۳-۲ روش پیشنهادی

از آنجایی که تمرکز این پژوهش بر استخراج ویژگی‌های ضمنی است، در ابتدا به اختصار به توصیف ویژگی‌های ضمنی پرداخته و سپس فرایند پیشنهادی ارائه می‌شود.

ویژگی ضمنی، ویژگی است که به‌وضوح در جملات ظاهر نمی‌شود و با استفاده از کلمات احساس در جمله‌ی دارای احساس، قابل‌شناسایی است. ویژگی‌های ضمنی را می‌توان با استفاده از ویژگی‌های صریح شناخته‌شده از نظرات استخراج نمود. در جدول ۱، نمونه‌هایی از ویژگی‌های ضمنی آورده شده است.

جدول ۱: سه مثال برای ویژگی ضمنی

ویژگی ضمنی	نظر
قیمت	خیلی گران است
قیمت	اصلا به صرفه نیست
ظاهر	خیلی خوشم اومد، خیلی شیک است

و تحلیل متن شده است و دلایل استفاده از آن به اختصار در ادامه آورده شده است:

- برای متن‌های کوتاه نیز نتایج خوبی را تولید می‌کند.
  - در محاسبه‌ی وابستگی بین دو واژه، به وابستگی بین کلمات نادر و همچنین کلمات پرتکرار اهمیت می‌دهد [۲۱].
- برخلاف سایر روش‌ها در ایجاد ماتریس هم‌رخداد، که ویژگی‌های صریح و کلمات احساس تنها از طریق اطلاعات آماری به ماتریس اضافه می‌شوند، در این جا، همه‌ی ویژگی‌های صریح و احساسات موجود در جملات بر اساس روابط وابستگی بین کلمات، قواعد نحوی زبان فارسی و اطلاعات آماری به این ماتریس افزوده می‌شوند. در روش پیشنهادی برای استفاده از اطلاعات آماری، علاوه بر معیار LRT، از معیار فاصله‌ی ویژگی صریح از کلمه‌ی احساس، برخلاف سایر روش‌ها استفاده شد که بهبود خوبی در نتایج ارزیابی داشت. به عبارتی، برای هر زوج «ویژگی صریح-کلمه‌ی احساس»، در صورتی که معیار LRT این دو زوج از حد آستانه بیشتر باشد و فاصله‌ی این دو زوج در محدوده [۱-۴] باشد به ماتریس هم‌رخداد اضافه می‌شوند.

شبه کد ایجاد ماتریس هم‌رخداد در شکل ۳ آورده شده که شامل سه مرحله می‌باشد که در ادامه آورده شده است:

- ۱- ایجاد گراف وابستگی
- ۲- تعیین قواعد نحوی زبان
- ۳- استخراج اطلاعات آماری بین ویژگی‌های صریح و کلمات احساس

```

Input: opinion Set  $\{O\}$ , feature set  $\{F\}$ 
Output: co-occurrence matrix  $M_{OF}$ 
for each opinion word  $o_i$  in  $\{O\}$  do
  for each feature  $f_j$  in  $\{F\}$  do
    if ( $o_i$  and  $f_j$  are identified dependent by graph) then
      add  $o_i$  and  $f_j$  to  $M_{OF}$ ;
    end
    else if ( $o_i$  and  $f_j$  are identified dependent by syntactic rules) then
      add  $o_i$  and  $f_j$  to  $M_{OF}$ ;
    end
    sort LRT( $o_i$  and  $f_j$ ) to descending order ;
    else if ( $o_i$  and  $f_j$  are identified dependent by LRT criterion and are in range of (-1,4)) then
      add  $o_i$  and  $f_j$  to  $M_{OF}$  ;
    end;
  end;
end;
return  $M_{OF}$  co-occurrence matrix;
  
```

شکل ۳: شبه کد پیشنهادی تشکیل ماتریس هم‌رخداد

در مرحله‌ی اول، گراف وابستگی جملات در نظر گرفته می‌شود. در این گراف، در صورتی که یک کلمه‌ی احساس به یک ویژگی صریح وابسته باشد، در قالب زوج «ویژگی صریح-کلمه‌ی احساس» به ماتریس هم‌رخداد اضافه خواهد شد. برای تشکیل زوج‌های مذکور، از لیست‌های مستقل و وابسته که قبلاً ایجاد شده‌اند، استفاده شده است. به‌عنوان مثال در جمله‌ی «ظاهر گوشه‌ی زیبا و جذاب است.» در

پس از تعیین نقش نحوی و تشکیل گراف وابستگی، لیست صفات و ویژگی‌های صریح مستقل و وابسته ایجاد می‌شوند. در ابتدا، برای هر یک از جملات پیکره که حاوی ویژگی یا صفت هستند، دو لیست از صفات به نام لیست صفات مستقل و وابسته و دو لیست از ویژگی‌ها به نام لیست ویژگی‌های مستقل و وابسته ایجاد می‌شوند. از این لیست‌ها در فرایند تشکیل ماتریس هم‌رخداد استفاده شده است. در ادامه هر یک به اختصار بیان شده‌اند.

- صفات مستقل و وابسته: تمامی صفت‌هایی که در یک جمله رخ می‌دهند، صفت مستقل خواهند بود مگر اینکه صفتی (صفت‌هایی) بعد از حرف ربط «و» و یا علامت «ویرگول» آمده باشند و قبل از آن حروف نیز، یک صفت باشد. در این صورت، اولین صفت قبل از حرف ربط «و» و یا علامت «ویرگول» صفت مستقل و صفت‌های بعد از آن‌ها، صفت‌های وابسته خواهند بود. به‌عنوان مثال در جمله‌ی «دوربین و باتری خوب است»، «خوب» صفت مستقل و یا در جمله‌ی «ظاهر گوشه‌ی زیبا و جذاب است»، «جذاب» صفت مستقل و «زیبا» صفت وابسته است.
- ویژگی‌های صریح مستقل و وابسته: دسته‌بندی ویژگی‌ها دقیقاً مشابه فرایند ارائه‌شده برای دسته‌بندی صفات است، یعنی همه‌ی ویژگی‌ها مستقل خواهند بود مگر این که ویژگی (ویژگی‌هایی) بعد از حرف ربط «و» و یا علامت «ویرگول» بیاید و قبل از آن حروف نیز یک ویژگی آمده باشد. به‌عنوان مثال در جمله‌ی «دوربین و باتری خوب است»، «دوربین» ویژگی مستقل و «باتری» ویژگی وابسته است.

### ۳-۲- ایجاد ماتریس هم‌رخداد

اولین گام از فرایند کلی استخراج ویژگی ضمنی، ایجاد ماتریس هم‌رخداد است. برای ایجاد ماتریس هم‌رخداد ابتدا ویژگی‌های صریح و احساس به‌عنوان ورودی به سیستم داده می‌شوند. این اطلاعات از نتیجه پژوهش‌های قبلی استخراج شده است.

در ماتریس هم‌رخداد برای هر ویژگی صریح یک سطر و برای هر کلمه‌ی احساس یک ستون در نظر گرفته می‌شود. برای هر زوج «ویژگی صریح-کلمه‌ی احساس» موجود در این ماتریس، معیار آزمون نسبت درست‌نمایی<sup>۲۳</sup> (LRT) [۲۱] و تعداد رخداد ویژگی صریح و کلمه‌ی احساس با هم محاسبه می‌شوند.

نسبت درست‌نمایی یا LRT از نسبت احتمال رخداد دو واژه به احتمال کل نظرات محاسبه می‌شود که میزان وابستگی بین دو واژه را نشان می‌دهد. هر چه مقدار این معیار بزرگتر باشد، وابستگی دو واژه به یکدیگر بیشتر خواهد بود [۲۱]. عنصر  $(i,j)$  در ماتریس هم‌رخداد میزان LRT و همچنین تعداد رخداد ویژگی صریح  $\lambda_m$  با احساس  $\lambda_n$  را نشان می‌دهد. استفاده از معیار LRT منجر به نتایج آماری بهتری در تجزیه

«باتری-عالی» را محاسبه و به ترتیب نزولی مرتب می‌کند و زوج با بالاترین مقدار LRT را در صورتی که در بازه‌ی [۱-۴] باشد به ماتریس هم‌خدای اضافه می‌کند. بر اساس جمله فوق دو زوج ویژگی «دوربین-خوب» و «باتری-عالی» به ماتریس اضافه می‌شوند. دلیل انتخاب مقدار ۱-، این است که گاهی اوقات، صفتی که برای بیان احساس یک ویژگی به کار می‌رود، قبل از آن ویژگی ظاهر می‌شود. به عنوان نمونه در جمله‌ی نظری «در کل اگر یک گوشی ارزان قیمت و با امکانات عالی می‌خواهید توی خرید این گوشی شک نکنید»، صفت «ارزان» که برای بیان احساس ویژگی «قیمت» به کار می‌رود، قبل از آن ظاهر می‌شود. در شکل ۴، مثالی برای تشکیل ماتریس هم‌رخداد آورده شده است.

<p>(۱) تاچ، سرعت و قدرت گوشی واقعا رضایت بخش و عالی است! (۲) هزینه‌اش یکم گران است ولی ظاهر زیبا و قشنگی دارد.</p> <p>(۳) دوربین چندان هم خوب نبود ولی باتری آن به نسبت عالی است. (۴) قیمتش خیلی ارزان است.</p> <p>(۵) قیمتش کمی گران است ولی دوربینش عالی است. (۶) هزینه‌اش یکم گران است.</p>							
تاج	سرعت	قدرت	هزینه	ظاهر	دوربین	قیمت	باتری
			۲				
				۱			
					۱		
						۱	
							۱

شکل ۴: مثالی برای تشکیل ماتریس هم‌رخداد

### ۳-۳- خوشه‌بندی صفات

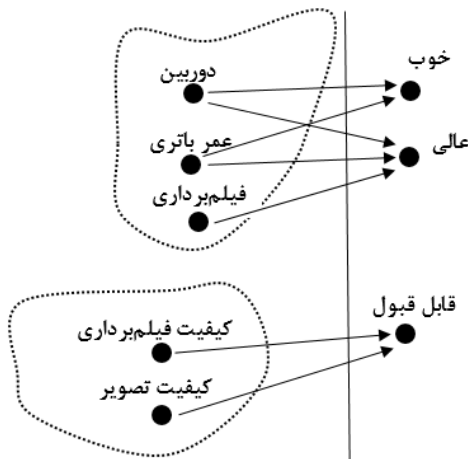
در گام بعدی، صفات‌های موجود در ماتریس هم‌خدای بر اساس شباهت کسینوسی خوشه‌بندی می‌شوند. در خوشه‌بندی صفات، صفتهایی که از نظر مفهومی و معنایی به هم نزدیک هستند، در یک خوشه قرار می‌گیرند. در ابتدا، برای هر صفتی که در ماتریس هم‌خدای موجود است، تعداد تکرار آن صفت با تمامی ویژگی‌های

صورتی که گراف تشخیص دهد که کلمه‌ی احساس «زیبا» به ویژگی «ظاهر» وابسته است، با استفاده از لیست صفات وابسته، علاوه بر زوج «ظاهر-زیبا»، زوج «ظاهر-جذاب» هم به ماتریس اضافه خواهد شد؛ چرا که صفت «جذاب» به عنوان وابسته‌ی صفت «زیبا» دسته‌بندی شده است. برای ویژگی‌های صریح وابسته‌ی موجود در جمله نیز روال مشابهی انجام می‌شود. به عنوان مثال، از جمله‌ی «دوربین و باتری خوب است»، دو زوج «دوربین-خوب» و «باتری-خوب» به ماتریس هم‌خدای اضافه خواهند شد؛ چرا که «باتری» ویژگی وابسته‌ی «دوربین» است.

ممکن است که ابزار هضم نتواند برخی روابط وابستگی بین کلمات احساس و ویژگی‌های صریح را نشان دهد. در مرحله‌ی دوم، قواعد نحوی زبان بررسی می‌شوند. به عنوان نمونه، در زبان فارسی، هر صفتی که بلافاصله پس از یک ویژگی آورده شود یا یک قید مقدار (مانند خیلی، بسیار، کمی و غیره) بین صفت و ویژگی بیاید، آن صفت مربوط به ویژگی است. صفتهای و ویژگی‌هایی که از قاعده‌ی فوق پیروی کنند، به صورت زوج «ویژگی صریح-کلمه‌ی احساس» به ماتریس هم‌خدای اضافه می‌شوند. به عنوان مثال، جمله‌ی «صفحه نمایش و کیفیت تصویر خیلی خوب است.» را در نظر بگیرید. بر اساس این قانون نحوی، صفت «خوب» به ویژگی «کیفیت تصویر» وابسته است، زیرا این صفت بلافاصله بعد از ویژگی و قید مقدار آمده است؛ در نتیجه، زوج «کیفیت تصویر-خوب» به ماتریس هم‌خدای اضافه می‌شود، علاوه بر این، زوج «صفحه نمایش-خوب» به دلیل وابسته بودن ویژگی «کیفیت تصویر» به ویژگی «صفحه نمایش» نیز به ماتریس هم‌خدای اضافه خواهد شد.

مرحله‌ی سوم، تمامی زوج‌های ویژگی صریح و صفتهایی که در دو مرحله قبلی استخراج نشده‌اند، بررسی می‌شوند. در ابتدا معیار LRT برای هر زوج ویژگی صریح-کلمه‌ی احساس محاسبه می‌شود که نشان‌دهنده‌ی میزان هم‌رخداد واژه‌ها در جملات پیکره است. پس از محاسبه‌ی LRT هر زوج ویژگی صریح-کلمه‌ی احساس، لیست زوج ویژگی صریح-کلمه‌ی احساس به ترتیب نزولی با توجه به مقدار LRT مرتب شده و لیست مرتب به دست آمده جهت افزودن زوج ویژگی صریح-کلمه‌ی احساس به ماتریس هم‌خدای، از بالا پیمایش می‌شود. سپس آن عنصر در صورتی که ماتریس هم‌خدای اضافه می‌شود که در بازه‌ی [۱-۴] باشد. در محاسبه‌ی بازه کلمات مانع<sup>۲۴</sup> نادیده گرفته می‌شوند. دلیل استفاده از بازه این است که صفتی که برای توصیف احساس یک ویژگی به کار می‌رود، در نزدیکی آن ویژگی ظاهر می‌شود. نتایج به دست آمده نشان می‌دهد که استفاده از بازه باعث افزایش دقت و بازخوانی در ارزیابی نتایج می‌شود. هر بار که عنصری از این لیست به ماتریس افزوده می‌شود سایر زوج ویژگی صریح-کلمات احساس موجود در لیست که شامل ویژگی یا کلمه‌ی احساس اضافه شده به ماتریس باشند، از لیست حذف می‌شوند. به عنوان مثال، در جمله‌ی «دوربین چندان هم خوب نبود ولی باتری آن نسبتا عالی است»، مقدار LRT تمامی زوج‌های «دوربین-خوب»، «دوربین-عالی»، «باتری-خوب» و

ضمنی نمی‌کند. از این رو از دسته‌بندی احساسات بر اساس بردار ویژگی پیشنهاد شده است.



شکل ۶: نمونه خوشه‌بندی ویژگی‌ها بر اساس بردار احساس

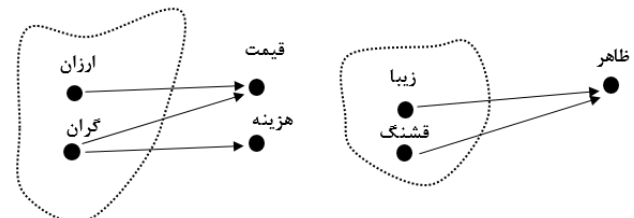
### ۳-۴- اعمال قوانین تشخیص ویژگی ضمنی

پس از انجام مراحل فوق، از صفات خوشه‌بندی شده و ماتریس هم‌رخداد برای تشخیص ویژگی‌های ضمنی استفاده می‌شود. تشخیص ویژگی ضمنی برای جملاتی انجام می‌شود که فاقد ویژگی صریح و دارای کلمه‌ی احساس است. برای هر صفت مانند  $s$  در این گونه جملات، از چهار قانون برای استخراج ویژگی ضمنی مناسب استفاده می‌شود که در ادامه توضیح داده شده‌اند و شبه کد آن نیز در شکل ۷ آورده شده است.

- در صورتی که صفت  $s$  فقط با یک ویژگی صریح مانند  $a$  در ماتریس هم‌رخدادی رخ داده باشد، ویژگی  $a$  به‌عنوان ویژگی ضمنی جمله‌ای که دارای احساس و فاقد ویژگی است انتخاب می‌شود. به‌عنوان مثال در جمله‌ی نظری «خیلی گران است» در صورتی که تنها زوج «قیمت-گران» در ماتریس هم‌رخدادی وجود داشته باشد، ویژگی ضمنی «قیمت» برای صفت «گران» انتخاب می‌شود.
- در صورتی که صفت  $s$  در مجموعه صفات خوشه‌بندی شده موجود باشد، خوشه‌ی دربردارنده‌ی صفت  $s$  انتخاب می‌شود. سپس ویژگی صریح مشترک بین تمامی صفات داخل خوشه از ماتریس هم‌رخدادی انتخاب و به‌عنوان ویژگی ضمنی صفت  $s$  انتساب داده می‌شود. به‌عنوان مثال اگر صفات «ارزان»، «گران» و «عادلان» در یک خوشه باشند و صفت «ارزان» و «گران» با ویژگی‌های «قیمت» و «هزینه» و صفت «عادلان» با ویژگی «قیمت» در ماتریس هم‌رخدادی موجود باشند. در این صورت اگر جمله‌ی ضمنی «خیلی گران است» را داشته باشیم، ویژگی مشترک بین این سه صفت یعنی «قیمت» به‌عنوان ویژگی ضمنی صفت مربوطه تعیین می‌شود.

اگر برای صفت  $s$  موجود در جمله، دو ویژگی  $a_1$  و  $a_2$  در ماتریس هم‌رخدادی وجود داشته باشند، در این حالت، در صورتی که تعداد هم‌رخدادی صفت  $s$  با ویژگی  $a_1$ ، بیش‌تر از دو بار باشد و علاوه بر

صریح به‌دست می‌آید. یعنی برای هر صفت، یک بردار از ویژگی‌ها به‌دست می‌آید که هر عنصر این بردار بیان‌گر تعداد هم‌رخدادی آن صفت با ویژگی‌های صریح موجود در جملات نظرات است. پس از نرمال کردن بردار هر صفت، شباهت کسینوسی را برای هر دو صفت کاندید به‌دست آورده و در صورتی که شباهت کسینوسی برای دو صفت کاندید بزرگتر و یا مساوی مقدار  $0.95$  باشد، در یک خوشه قرار می‌گیرند. خوشه‌های به‌دست‌آمده در پایگاه داده ذخیره می‌شوند تا برای تشخیص ویژگی ضمنی مورد استفاده قرار گیرند. در شکل ۵، خوشه‌بندی معادل برای ماتریس هم‌رخدادی موجود در شکل ۴ نشان داده شده است. در این شکل، صفات «ارزان» و «گران» در یک خوشه و صفات «زیبا» و «قشنگ» در خوشه‌ی دیگر قرار می‌گیرند زیرا بردار ویژگی خوشه‌ی اول شامل «قیمت» و «هزینه» و بردار ویژگی خوشه‌ی دوم شامل «ظاهر» است.



شکل ۵: نمونه خوشه‌بندی احساسات بر اساس بردار ویژگی‌ها

خوشه‌بندی ویژگی‌ها بر اساس بردار احساسات کاربردی نیست زیرا صفات زیادی هستند که برای بیان احساسات ویژگی‌های مختلفی به‌کار می‌روند. به‌عنوان مثال همان‌طور که در جدول ۲، نشان داده شده است، احساس‌هایی مانند «قابل قبول»، «خوب»، «عالی» و غیره برای توصیف ویژگی‌های مختلفی به‌کار می‌روند.

### جدول ۲: مثال‌هایی برای احساسات مشترک مرتبط با ویژگی‌ها

احساس مرتبط با ویژگی	جمله
خوب	دوربین خیلی خوبی دارم من که خیلی دوستش دارم
خوب	عمر باتریش خیلی خوبه، واقعا سامسونگ گل کاشته
عالی	عمر باتریش عالی است
قابل قبول	کیفیت تصویرش قابل قبول است
عالی	دوربین عالی
عالی	فیلم برداریش عالی
قابل قبول	کیفیت فیلم برداریش قابل قبول است

با توجه به جدول ۲، خوشه‌بندی ویژگی‌ها بر اساس بردار احساسات به‌صورت شکل ۶ است. همان‌طور که مشاهده می‌شود، ویژگی‌های «دوربین»، «عمر باتری» و «فیلم برداری» در یک خوشه و ویژگی‌های «کیفیت فیلم برداری» و «کیفیت تصویر» در خوشه‌ی دیگر قرار می‌گیرند. پس خوشه‌بندی ویژگی‌ها بر اساس بردار احساسات، نتایج درستی ایجاد نکرده و کمکی به حل مسئله‌ی تشخیص ویژگی‌های

تلفن همراه و سه نوع لپ‌تاپ [۲۲] است. اطلاعات مربوط به این داده ها در جدول ۳ ارائه شده است.

جدول ۳: اطلاعات آماری مجموعه داده‌ها

تعداد جملات	تعداد نظرات	شرح کالا	گروه کالا
۱۸۴۷	۳۹۴	-Fujitsu LifeBook AH-532-B -Huawei Ascend Y51 -Samsung Galaxy Young S6310	تلفن همراه
۶۶۵	۱۶۸	- ASUS K46CB-C - ASUS N550JV - Nokia Lumia 520	لپ‌تاپ
۲۵۱۲	۵۶۲	۶	تعداد کل

تمامی جملات مرتبط با محصولات توسط ابزار هضم نرمال شده و فعل‌های محاوره (مانند همیشه) به فعل رسمی (مانند می‌شود) تبدیل شده و همچنین تا حدودی غلط‌های املائی اصلاح می‌شوند. سپس نقش لغوی، نحوی و گراف وابستگی تمامی جملات توسط این ابزار استخراج می‌گردد. این مجموعه داده‌ها به‌عنوان ورودی الگوریتم پیشنهادی استفاده می‌شود.

روش پیشنهادی ما بر روی ویژگی‌های ضمنی عمل می‌کند. برای ارزیابی صحت خروجی، از دو کاربر خبره خواسته شد تا با خواندن تمامی نظرات، ویژگی‌های ضمنی هر جمله را (در صورت وجود) استخراج کنند. این ویژگی‌ها با ویژگی‌های ضمنی استخراج شده توسط روش پیشنهادی ما مقایسه شده است. به‌علاوه، برای استخراج قواعد رایج نحوی زبان فارسی (مرحله ۲ از روش پیشنهادی) از چند مرجع معتبر در این حوزه استفاده شده است.

#### ۴-۲- معیارهای ارزیابی

ارزیابی یک مدل نظرکاوی، بر اساس نمونه‌های آزمونی است که توسط فرد خبره برچسب‌گذاری شده باشد. برای ارزیابی مدل، باید برچسبی را که مدل برای متون نظرات (ویژگی‌های ضمنی) در نظر گرفته است، با برچسب فرد خبره مقایسه کرد. بنابراین چهار حالت وجود دارد که جدول ۴، شرایط وقوع این حالت‌ها را نشان می‌دهد.

جدول ۴: وقوع حالت‌های مختلف برای مدل تجزیه و تحلیل متون نظرات

برچسب‌گذاری توسط فرد خبره			برچسب‌گذاری شده توسط مدل
NO	YES	YES	
FP	TP	NO	مدل
TN	FN	NO	

پارامترهای معرفی شده در جدول ۴ به‌صورت زیر برای استخراج ویژگی‌های ضمنی تعریف می‌شوند:

- پارامتر TP: تعداد ویژگی‌هایی است که به‌درستی توسط مدل مورد نظر به‌عنوان ویژگی استخراج شده‌اند.

این مقدار LRT آن دو هم از حد آستانه‌ی  $\gamma$  بیشتر باشد و همچنین تعداد هم‌رخدادی صفت  $s$  با ویژگی  $a_2$ ، کم‌تر از دو بار باشد، ویژگی  $a_1$  به‌عنوان ویژگی ضمنی صفت  $s$  استخراج می‌شود.

• اگر صفت  $s$  با بیش از دو ویژگی در ماتریس هم‌رخدادی وجود داشته باشد، دیگر نمی‌توان ویژگی مشخصی را برای صفت  $s$  در نظر گرفت. به‌عنوان مثال در جمله‌ی «گوشی خیلی عالی است»، احساس «عالی» با ویژگی‌های مختلفی در ماتریس هم‌رخدادی اتفاق می‌افتد. پس نمی‌توان ویژگی مشخصی را برای این جمله تعیین کرد. در این شرایط، این صفت برای کلیت محصول در نظر گرفته می‌شود. به‌عنوان مثال، اگر نظر فوق در مورد گوشی موبایل باشد، صفت «عالی» به «گوشی» نسبت داده می‌شود.

**Input:**  $M_{of}, G_o$

**Output:**  $IM_{of}$

$IS \leftarrow$  implicit sentence set with no feature ;

**For each** opinion words  $o_i$  in  $IS$  **do**

**if** ( $o_i$  just occurred candidate feature ( $cf$ ) in  $M_{of}$ ) **then**

Add  $o_i$  and  $cf$  to  $IM_{of}$ ;

**end;**

**else if** ( $o_i$  is in  $G_o$  cluster) **then**

$cf \leftarrow$  Find common features between all sentiments in  $G_o$

;

add  $o_i$  and  $cf$  to  $IM_{of}$ ;

**end;**

**else if** ( $o_i$  occurred with features  $f_1, f_2$  in  $M_{of}$  matrix)

**then**

**if** (co-occurrence\_number of  $o_i$  with  $f_1 \geq 2$  and  $LRT(o_i, f_1) \geq 5$ )

and (co-occurrence\_number of  $o_i$  with  $f_2 < 2$ ) **then**

add  $o_i$  and  $f_1$  to  $IM_{of}$ ;

**end;**

**else if** (co-occurrence\_number of  $o_i$  with  $f_2 \geq 2$  and  $LRT(o_i, f_2) \geq 5$ )

and (co-occurrence\_number of  $o_i$  with  $f_1 < 2$ ) **then**

add  $o_i$  and  $f_2$  to  $IM_{of}$ ;

**end;**

**end;**

**else if** ( $o_i$  occurred with more than two features) **then**

add  $o_i$  and «کلیت محصول» to  $IM_{of}$ ;

**end;**

**return**  $IM_{of}$ ;

شکل ۷: شبه کد تشخیص ویژگی مناسب برای جمله‌ی ضمنی

#### ۴- نتایج و ارزیابی روش پیشنهادی

در این بخش، ابتدا به معرفی مجموعه دادگان و معیارهای ارزیابی پرداخته و سپس عملکرد روش پیشنهادی، ارزیابی می‌شود.

#### ۴-۱- نحوه آماده‌سازی مجموعه داده‌ها

در این مقاله، از مجموعه نظرات جمع‌آوری شده در مقاله‌ی باباعلی [۱۹] که از سایت دیجی کالا در سال ۱۳۹۴ جمع‌آوری گردیده، جهت ارزیابی سیستم استفاده شده است. این مجموعه نظرات، مربوط به سه نوع



در اطلاعات آماری از ارتباطات نحوی زبان و گراف وابستگی بین کلمات یک جمله برای ایجاد ماتریس هم‌رخداد کارا تر بهره برده است. به‌عنوان مثال جمله‌ی «با این که قیمتش خیلی ارزان است ولی در کل گوشی خوبی است»، را در نظر بگیرید. روش‌های قبلی در صورتی که احساس «خوب» بیش‌تر از احساس «ارزان» با ویژگی «قیمت» رخ داده باشد، احساس «خوب» را به ویژگی «قیمت» نسبت می‌دهند که نادرست است. اما در روش پیشنهادی ما به دلیل استفاده از ارتباطات نحوی زبان، صفت «ارزان» را به ویژگی «قیمت» نسبت خواهد داد.

علیرغم اینکه روش پیشنهادی ما نسبت به روش‌های قبلی بهبود مناسبی داشته است، در برخی موارد به‌درستی ویژگی ضمنی را مشخص نکرده است. در مواردی، این مسئله به‌علت زیادنبودن تعداد قواعد نحوی زبان فارسی رخ داده است. شایان ذکر است که در این پژوهش از قواعد نحوی رایج زبان فارسی استفاده شده است. دلیل دیگر برخی از خطاها نیز نرم‌افزار هضم بوده است که در برخی از موارد مانند تشخیص روابط وابستگی بین کلمات یک جمله، دقت بالایی نداشته است. علاوه بر این، برخلاف سایت‌های انگلیسی مانند آمازون، تعداد نظرات کاربران در سایت دیجی کالا کم است. این مسئله نیز می‌تواند باعث کاهش کارایی در تشکیل ماتریس هم‌رخداد و در نتیجه کاهش دقت در تشخیص ویژگی ضمنی شود.

## ۵- نتیجه‌گیری

در این مقاله، روشی جهت استخراج ویژگی ضمنی مجموعه نظرات کاربران در زبان فارسی ارائه گردید. در این روش یک معیار جهت محاسبه فاصله‌ی صفت مرتبط با ویژگی در اطلاعات آماری پیشنهاد شده است. علاوه بر این گراف وابستگی، قواعد نحوی زبان فارسی و خوشه‌بندی صفات برای تعیین بهتر هم‌رخدادی ویژگی‌ها و کلمات احساس پیشنهاد شده است. ارزیابی انجام‌شده بر روی نظرات زبان فارسی گرفته‌شده از سایت دیجی کالا نشان می‌دهد که این روش بهبود ایجاد ماتریس هم‌رخداد و در نتیجه کارایی تشخیص ویژگی‌های ضمنی را به‌همراه داشته است.

## مراجع

- [1] J. J. Li, H. Yang and H. Tang, "Feature Mining and Sentiment Orientation Analysis on Product Review," in International Conference on Management Information and Optoelectronic Engineering, pp. 79-84, 2016.
- [2] B. Liu, "Sentiment analysis and opinion mining," Synthesis Lectures on Human Language Technologies, vol. 5, pp. 1-167, 2012.
- [3] مصطفی. رجبزاده و رضا. رافع، «ارائه یک سیستم توصیه‌گر ترکیبی برای تجارت الکترونیک»، مجله مهندسی برق دانشگاه تبریز، جلد ۴۵، صفحه ۱۷۷-۱۶۳، ۱۳۹۴.
- [4] سیامک. عبدالعزیز، محمدعلی. بالافرو و لیلی. محمدخانی، «استفاده از خوشه‌بندی و مدل مارکوف جهت پیش‌بینی درخواست آتی کاربر در وب»، مجله مهندسی برق دانشگاه تبریز، جلد ۴۵، صفحه ۱۷۷-۱۶۳، ۱۳۹۴.
- [5] A. Yadollahi, A. G. Shahraki, and O. R. Zaiane, "Current state of text sentiment analysis from opinion to emotion mining," ACM Computing Surveys (CSUR), vol. 50, no. 25, pp. 1-33, 2017.

- پارامتر FP: مربوط به تعداد ویژگی‌هایی است که به‌صورت نادرست توسط مدل مورد نظر به‌عنوان ویژگی استخراج شده‌اند.
- پارامتر FN: تعداد ویژگی‌هایی است که ویژگی ضمنی هستند اما توسط مدل مورد نظر به‌عنوان ویژگی ضمنی تشخیص داده نشده‌اند.
- پارامتر TN: تعداد ویژگی‌هایی است که ضمنی نیستند و به‌درستی توسط مدل مورد نظر به‌عنوان ویژگی استخراج نشده‌اند.

جهت ارزیابی کارایی روش پیشنهادی در استخراج ویژگی ضمنی همانند کارهای قبلی از معیارهای دقت، بازخوانی و معیار F استفاده شده است که معیارهای ارزیابی به‌صورت زیر قابل محاسبه هستند:

$$Precision = \frac{TP}{TP + FP} \quad (1)$$

$$Recall = \frac{TP}{TP + FN} \quad (2)$$

$$F - measure = \frac{2 * Precision * Recall}{Precision + Recall} \quad (3)$$

در ادامه با معیارهای فوق نتایج روش پیشنهادی ارزیابی شده است.

## ۴-۳- ارزیابی نتایج

در آزمایش‌های انجام شده روش پیشنهادی، مقدار حد آستانه‌ی ۷ که در استخراج ویژگی‌های ضمنی در بخش ۴-۳ استفاده گردید، به‌صورت تجربی برابر ۵ در نظر گرفته شده است.

مقادیر معیارهای محاسبه شده برای روش پیشنهادی در تشخیص ویژگی ضمنی در مقایسه با نتایج به‌دست‌آمده در روش ارائه‌شده توسط باباعلی [۱۹]، باقری [۷] و دنینگ [۲۱] در جدول ۵ نشان داده شده است.

جدول ۵: مقایسه روش پیشنهادی با روش‌های پیشین

لپ تاپ	تلفن همراه		
۰/۷۵	۰/۷۵	روش پیشنهادی	دقت
۰/۶۹	۰/۶۷	باقری [۷]	
۰/۶۴	۰/۶۰	باباعلی [۱۹]	
۰/۶۰	۰/۵۹	دنینگ [۲۱]	
۰/۶۹	۰/۶۵	روش پیشنهادی	بازخوانی
۰/۵۰	۰/۴۹	باقری [۷]	
۰/۴۵	۰/۴۲	باباعلی [۱۹]	
۰/۴۰	۰/۳۸	دنینگ [۲۱]	
۰/۷۱	۰/۶۹	روش پیشنهادی	معیار F
۰/۵۸	۰/۵۷	باقری [۷]	
۰/۵۲	۰/۴۹	باباعلی [۱۹]	
۰/۴۸	۰/۴۶	دنینگ [۲۱]	

نتایج نشان می‌دهد که روش پیشنهادی ما کارایی تشخیص ویژگی ضمنی را نسبت به روش‌های پیشین [۷، ۱۹، ۲۱] بهبود داده است. زیرا روش ما برخلاف روش‌های قبلی علاوه بر در نظر گرفتن معیار فاصله

- [14] E. Asgarian, M. Kahani, and S. Sharifi, "The impact of sentiment features on the sentiment polarity classification in Persian reviews," *Cognitive Computation*, vol. 10, pp. 117-135, 2018.
- [15] M. Dragoni, "Computational advertising in social networks: an opinion mining-based approach," in *Proceedings of the 33rd Annual ACM Symposium on Applied Computing*, pp. 1798-1804, 2018.
- [16] Q. Su, X. Xu, H. Guo, Z. Guo, X. Wu, and X. Zhang, "Hidden sentiment association in chinese web opinion mining," in *Proceedings of the 17th international conference on World Wide Web*, pp. 959-968, 2008.
- [17] Z. Hai, K. Chang, and J.-j. Kim, "Implicit feature identification via co-occurrence association rule mining," in *International Conference on Intelligent Text Processing and Computational Linguistics, Lecture Notes in Computer Science*, Springer, pp. 393-404, 2011.
- [18] K. Schouten and F. Frasinca, "Implicit Feature Extraction for Sentiment Analysis in Consumer Reviews," in *International Conference on Applications of Natural Language to Data Bases Information Systems, Lecture Notes in Computer Science*, Springer, pp. 228-231, 2014.
- [۱۹] مرضیه. باباعلی و محمدعلی. نعمت‌بخش. "استخراج ویژگی‌های محصول در زبان فارسی." سومین همایش زبان‌شناسی رایانشی، دانشگاه صنعتی شریف، ۱۳۹۳.
- [۲۰] محسن. ایمانی و مجتبی. خلاش. ابزار پردازش زبان فارسی، (http://www.sobhe.ir/hazm). ۱۳۹۲.
- [21] T. Dunning, "Accurate methods for the statistics of surprise and coincidence," *Computational linguistics*, vol. 19, pp. 61-74, 1993.
- [22] Digikala Dataset (2018, December 13), Retrieved February 2, 2019.
- [23] <https://www.uploader.net/files/701edb674a0ef75695b47de35db298d4/Data.rar.html>.
- [6] E. Breck and C. Cardie, *Opinion Mining and Sentiment Analysis*, in *The Oxford Handbook of Computational Linguistics* 2nd edition, 2017.
- [7] A. Bagheri, M. Saraee, and F. De Jong, "Care more about customers: Unsupervised domain-independent aspect detection for sentiment analysis of customer reviews," *Knowledge-Based Systems*, vol. 52, pp. 201-213, 2013.
- [8] Z. Hai, K. Chang, and G. Cong, "One seed to find them all: mining opinion features via association," in *Proceedings of the 21st ACM international conference on Information and knowledge management*, pp. 255-264, 2012.
- [9] N. Jakob and I. Gurevych, "Extracting opinion targets in a single- and cross-domain setting with conditional random fields," in *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pp. 1035-1045, 2010.
- [10] G. Somprasertsri and P. Lalitrojwong, "Automatic product feature extraction from online product reviews using maximum entropy with lexical and syntactic features," in *IEEE International Conference on Information Reuse and Integration*, pp. 250-255, 2008.
- [۱۱] مانده. شیخ حسینی، محرم. منصوری‌زاده و میرحسین. دزفولیان، «نظر کاوی جنبه‌گرا به کمک استخراج روابط معنایی»، بیست و دومین کنفرانس ملی سالانه انجمن کامپیوتر ایران، دانشگاه صنعتی شریف، ۱۳۹۵.
- [12] M. Z. Asghar, A. Khan, S. R. Zahra, S. Ahmad, and F. M. Kundi, "Aspect-based opinion mining framework using heuristic patterns," *Cluster Computing*, <https://doi.org/10.1007/s10586-017-1096-9>, pp. 1-19, 2017.
- [13] H. Liang, X. Sun, Y. Sun, and Y. Gao, "Text feature extraction based on deep learning: a review," *EURASIP journal on wireless communications and networking*, vol. 1, pp. 1-12, 2017.

## زیرنویس‌ها

<sup>13</sup> Stop list

<sup>14</sup> Chi-square

<sup>15</sup> Information Gain

<sup>16</sup> Mutual Infomation

<sup>17</sup> Opinion word

<sup>18</sup> Hidden sentiment links

<sup>19</sup> Association set

<sup>20</sup> Rule mining co-occurrence association approach

<sup>21</sup> Frequency Modified Left Right

<sup>22</sup> Hazm

<sup>23</sup> Likelihood Ratio Tests (LRT)

<sup>24</sup> Stop words

<sup>1</sup> Opinion mining

<sup>2</sup> Sentiment analysis

<sup>3</sup> Sentence Level

<sup>4</sup> Supervised

<sup>5</sup> Unsupervised

<sup>6</sup> Support Vector Machine (SVM)

<sup>7</sup> Maximum Entropy

<sup>8</sup> Explicit features

<sup>9</sup> Implicit feature

<sup>10</sup> Co-occurrence matrix

<sup>11</sup> Classifier

<sup>12</sup> Conditional Random Field (CRF)