

راهکاری به منظور طبقه‌بندی حالات عاطفی در شبکه اجتماعی توییتر

احمد آقا کاردان^۱، یوسف آقاجانلو^۲

^۱ استادیار دانشگاه صنعتی امیرکبیر، دانشکده مهندسی کامپیوتر و فناوری اطلاعات
aakardan@aut.ac.ir

^۲ دانشجوی کارشناسی ارشد دانشگاه صنعتی امیرکبیر، دانشکده مهندسی کامپیوتر و فناوری اطلاعات
yousef9190@gmail.com

چکیده

شبکه‌های اجتماعی به عنوان یک رسانه محبوب به منظور به اشتراک گذاری نظرات، افکار، اطلاعات و تجارب مطرح شده‌اند. از سویی افراد معمولاً به طور غیرمستقیم رفتارهای احساسی خود را در متون و گفتارشان بازتاب می‌کنند. در این راستا استخراج حالت‌های عاطفی افراد از منابع برخط همانند شبکه‌های اجتماعی و سپس طبقه‌بندی این حالت‌ها در دسته‌های از پیش تعیین شده می‌تواند اطلاعات غنی درباره موضوعات مطرح و اعضای شبکه اجتماعی مورد نظر به دست دهد و در موارد گوناگون همانند تطبیق محتوا با عواطف هر فرد، کمپین‌های بازاریابی، نظارت بر پاسخ‌ها در اتفاقات محلی و عمومی و کشف روندهای حالت‌های افراد مورد استفاده قرار گیرد. در مقاله پیش‌رو از شبکه اجتماعی توییتر به دلیل در دسترس بودن داده‌های آن و کاربرد بسیار آن در حوزه طبقه‌بندی عواطف، استفاده شده است. در این مقاله قدم‌های مورد نیاز برای به اجرا رساندن پژوهش طی شده و نتایج مربوطه مطرح گردیده است. در پایان نیز راهکاری در قالب یک سیستم بدین منظور ارائه شده است.

کلمات کلیدی

تشخیص عواطف، تحلیل احساسات، طبقه‌بندی متون، شبکه اجتماعی، انتخاب ویژگی، دسته‌های عواطف

فضای مسئله می‌باشد. این واژه برای اولین بار توسط ناسوکاوا و ایی^۲ در سال ۲۰۰۳ مطرح گردید [1].

با توسعه شبکه‌های اجتماعی، مردم شروع به بیان نظرات و عواطفشان بر روی شبکه‌ای همانند توییتر نمودند. در همین راستا کارهایی به منظور تجزیه و تحلیل توییت‌های کاربران و تشخیص عواطف کاربران بر روی توییتر انجام شد و تشخیص عواطف روی توییتر تبدیل به کانونی برای محققان به منظور تمرکز بر روی آن گردید [2].

دلایل متعددی برای توجه به این حوزه وجود دارد. اول اینکه این حوزه تقریباً در هر دامنه‌ای، کاربردهای بسیاری دارد. به‌طور مثال تحلیل احساسات در صنعت با توجه به گسترش برنامه‌های تجاری در حال رونق گرفتن است. دوم اینکه راه‌حلی را برای بسیاری از مشکلات پژوهشی که قبلاً مورد

۱- مقدمه

انسان همواره به دنبال تعاملات اجتماعی بوده و این تعاملات بر افکار و اعمال وی تأثیر می‌گذارند. نمونه‌ای از این تعاملات، شبکه‌های اجتماعی می‌باشد که می‌تواند این تأثیر را پویاتر نماید. تأثیر تعاملات بر عواطف و احساسات از جمله مهم‌ترین تأثیراتی است که در این شبکه‌ها وجود دارد.

تحلیل احساسات^۱ عنوانی است که شامل تحلیل نظرات، احساسات، نگرش‌ها و عواطف افراد نسبت به محصولات، سرویس‌ها، سازمان‌ها، افراد، موضوعات، وقایع و ویژگی‌های آن‌ها است و این موارد نشان‌دهنده بزرگی

از طرفی در طول وقایعی همانند بحران‌های طبیعی، تعداد زیادی از کاربران محتوایی همانند توییت‌ها، پست‌های بلاگ و پیام‌های فروم تولید می‌نمایند. برنیلسن و همکاران^۷ توییت‌های افراد در طول بحران طوفان شن در سال ۲۰۱۲ را دریافت نموده و ۴ دسته مثبت، خشم، ترس و دیگر را به منظور طبقه‌بندی تعیین کردند. ایده کلی برای سازمان‌های مدیریت بحران بدین صورت می‌باشد که آن‌ها می‌توانند با پیگیری مطالبی که افراد منتشر می‌نمایند، استراتژی‌هایشان را در راه برآورده کردن انتظارات و نیازهای مردم تنظیم نمایند. در اینجا نیز پس از تعریف بردار ویژگی‌ها از روش‌های ماشین بردار پشتیبان و بیز ساده به منظور طبقه‌بندی احساسات استفاده شد [7].

۳- نظریه‌های پایه

در این بخش مجموعه مفاهیم و نظریه‌های پایه‌ای استفاده شده در سیستم پیشنهادی معرفی شده است. این مجموعه مفاهیم شامل توصیف انواع دسته‌های عواطف، توصیف مجموعه دادگانی که می‌توانند به عنوان نمونه مورد استفاده قرار گیرد، روش‌های برچسب‌گذاری مجموعه دادگان، موضوع انتخاب ویژگی‌ها و توصیف انواع روش‌های طبقه‌بندی رایج می‌باشد که در زیر بخش‌های بعدی به اجمال مورد بررسی قرار گرفته‌اند.

۳-۱- انواع دسته‌های عواطف

عواطف انسانی خود را در قالب‌های حالات چهره، بیان، گفتار، نوشته‌ها و در حرکات بدن و فعالیت‌ها آشکار می‌کنند. روانشناسان از حالات صورت انسان که از آن می‌توان عواطف را برداشت نمود، استفاده کرده و آن‌ها را به عواطف انسانی پایه‌ای نگاشت نموده‌اند. اکمان^۸ در سال ۱۹۹۲ عواطف پایه‌ای را تعریف نموده است. این عواطف از طریق حالات چهره تعیین شده و مورد قبول جهانی قرار گرفته است. شش نوع عواطف پایه‌ای شامل شاد، ناراحت، ترس، عصبانی، تنفر و تعجب می‌باشد. در طول دهه‌های اخیر نیز، محققان انواع مختلفی از دسته‌های عاطفی را شناسایی و معرفی نموده‌اند [2].

مدل عمده برای نمایش عواطف وجود دارد؛ مدل دسته‌ای و مدل ابعادی. هر یک از این مدل‌ها به پوشش جنبه‌های مختلف عواطف انسانی کمک می‌نمایند [8].

مدل دسته‌ای فرض می‌کند دسته‌های عاطفی، همانند شش دسته‌بندی پایه‌ای اکمان می‌باشند. هر عاطفه به وسیله یک مجموعه ویژگی‌های خاص و شرایط بروز آن مشخص می‌شود. بیشتر کارها بر روی شش دسته عواطف پایه‌ای متمرکز شده است، با این حال بسیاری از محققان مجموعه عواطف مختلفی را برای حوزه‌های مختلف استفاده کرده‌اند. این نوع مدل دسته‌بندی، مدلی غالب و قطعی می‌باشد و با توجه به سادگی و شهرت آن تغییرات زیادی در آن به وجود آمده و به‌وفور مورد استفاده قرار می‌گیرد [8].

رویکرد دوم برای نمایش عواطف، استفاده از مدل ابعادی است. حالات عاطفی در این مدل به وسیله یک مجموعه از ابعاد رایج، به هم مرتبط بوده و به صورت عمومی در یک فضای دو یا سه‌بعدی تعریف می‌شوند. هر یک از عواطف، مکانی را در این فضا اشغال می‌کنند. مثالی از این مدل، مدل راسل^۹ می‌باشد که به وسیله یک شکل و یک مجموعه نقاط که نشان‌دهنده عواطف هستند معرفی شده است. شکل (۱) این مدل را نشان می‌دهد. واژه‌های مرتبط با عواطف در یک فضای دوبعدی سازمان‌دهی شده‌اند که قادر می‌سازد هر

مطالعه قرار نگرفته است ارائه می‌دهد. سوم اینکه برای اولین بار در تاریخ بشر ما اکنون با حجم عظیمی از داده‌ها در رسانه‌های اجتماعی بر روی وب مواجه هستیم و بدون این داده‌ها بسیاری از تحقیقات امکان‌پذیر نمی‌باشند. از این رو تحقیقات در تحلیل احساسات نه تنها تأثیر مهمی بر روی پردازش زبان طبیعی دارا می‌باشد بلکه تأثیرات عمیقی نیز بر روی علوم مدیریتی، علوم سیاسی، اقتصاد و علوم اجتماعی داشته و همه این‌ها تحت تأثیر نظرات مردم می‌باشند [1].

بنابراین استخراج احساسات پنهان شده در متون موجود در شبکه‌های اجتماعی و طبقه‌بندی آن‌ها مسئله مورد نظر در این مقاله و پژوهش‌های مرتبط است. در بخش دوم مروری بر کارهای انجام گرفته در این حوزه شده و ویژگی‌های هر یک از این کارها بیان گردیده است. در بخش سوم نظریه‌های پایه و مفاهیم مرتبط در این زمینه مطرح شده است. در بخش چهارم سیستم پیشنهادی به منظور طبقه‌بندی حالات عاطفی در شبکه اجتماعی تویتر معرفی شده است. این سیستم ترکیبی از دو روش معروف طبقه‌بندی یعنی روش‌های نظارت‌شده و بدون نظارت می‌باشد و ویژگی‌های خاصی را دارا است که می‌تواند دقت کارهای پیشین را بهبود بخشد. در پایان نیز جمع‌بندی از مطالب آورده شده است.

۲- مروری بر کارهای مرتبط

همان‌طور که در بخش مقدمه اشاره شد در سال‌های اخیر توجهات زیادی به حوزه تحلیل احساسات شده و پژوهش‌هایی با روش‌ها و ویژگی‌های مختلف صورت گرفته است. در اینجا به توصیف برخی از این پژوهش‌ها پرداخته‌ایم. برخی محققان عواطف را در یک چهارچوب وسیع‌تری از حالات خصوصی مطالعه نموده‌اند. ویب و همکاران^۳ بر روی برچسب‌گذاری دستی حالات خصوصی شامل عواطف، نظرات و احساسات در یک مجموعه‌ای شامل هزار جمله خبری پرداختند [3]. در کار [4] تمرکز بر روی یادگیری عواطف خاص از متن بوده و آلم و همکارانش^۴ طبقه‌بندی خودکار جملات در متون داستان‌های مربوط به کودکان را بر اساس عواطف پایه‌ای اکمان ایجاد کردند. در مقاله [5] یک روش طبقه‌بندی احساسات دومارحله‌ای برای تویتر پیشنهاد شده است که ابتدا پیام‌ها تحت عنوان ذهنی و عینی طبقه‌بندی می‌شوند و نشان داده می‌شود که توییت‌های ذهنی مثبت هستند یا منفی. از برچسب‌های دارای نویز به عنوان داده‌های آموزشی استفاده می‌شود. در این کار از روش یادگیری ماشین استفاده شده و ویژگی‌های نحوی تشکیل‌دهنده ماتریس ویژگی‌ها می‌باشند. در مقاله [6] پارور و بترسبای^۵ از روشی موسوم به نظارت از راه دور^۶ برای تجزیه و تحلیل احساسات و برچسب‌گذاری در شبکه‌های اجتماعی استفاده نمودند. ویژگی این روش، استفاده از نشانگرهای مرسوم عاطفی (شکلک‌ها و هشتک‌ها) در متون بود که می‌توانست به عنوان جانشینی برای برچسب‌های صریح در نظر گرفته شود. مجموعه آزمایش‌ها برای طبقه‌بندی در ۶ دسته پایه‌ای اکمان انجام شد. از داده‌های شبکه اجتماعی تویتر که به صورت تصادفی انتخاب شده‌اند و روش‌های یادگیری ماشین و از ویژگی‌های unigram که شامل تمامی کلمات و هشتک‌ها بعد از حذف urlها و نام‌های کاربری تویتر می‌باشد، به منظور پیاده‌سازی استفاده گردید. نشانگرها نیز به‌عنوان یک برچسب در این آزمایش استفاده شد. نهایتاً با انجام آزمایش‌ها مشخص گردید که این روش برای برخی عواطف (شاد، ناراحت و عصبانی) مناسب‌تر بوده و کمتر قادر به تشخیص عواطف دیگر می‌باشد [6].

پس از دریافت به منظور پردازش، توییت‌های غیر انگلیسی، ریتوییت‌ها، توییت‌های تکرار شده، توییت‌های دارای نویز، توییت‌های هدف‌دار و خشی و توییت‌های عینی که در آن‌ها واقعیت‌ها بیان شده است حذف خواهند شد و توییت‌های باقیمانده در فرمت CSV در پایگاه داده‌ای ذخیره می‌گردند [7].

۳-۳- برچسب‌گذاری مجموعه دادگان

به منظور استفاده از داده جمع‌آوری شده به عنوان داده‌های آموزشی در مدل‌های طبقه‌بندی، نیاز به برچسب‌گذاری دادگان بر مبنای یک احساس مورد اعتماد زمینی می‌باشد. یک احساس مورد اعتماد زمینی احساسی است که توسط مردم درک می‌شود و مورد توافق بیشتر افراد گیرنده در یک زمان است. یک سیستم تعاملی کامپیوتر-انسان باید قادر به تشخیص، تفسیر و پردازش عواطف انسانی باشد و اولین گام شناخت عواطف انسانی و برچسب‌گذاری داده‌ها بر پایه این شناخت می‌باشد [9].

سه رویکرد را به منظور برچسب‌گذاری مجموعه دادگان می‌توان در پیش گرفت: (۱) مجموعه دادگان را بین تعدادی حاشیه‌نویس تقسیم نمود و هر یک به صورت مستقل عملیات برچسب‌گذاری را روی دادگان خود انجام دهند. هر حاشیه‌نویس، عاطفه متناسب با هر متن را بر مبنای درصد حضور کلمات و عبارات مرتبط با هر دسته عاطفی مشخص می‌کند. (۲) تعدادی حاشیه‌نویس تعیین می‌شود و هر توییت توسط دو حاشیه‌نویس برچسب‌گذاری می‌شود. اگر هر دو بر روی عاطفه موجود در توییت توافق داشته باشند، توییت برچسب عاطفی موردنظر را می‌خورد در غیر این صورت برچسب خشی به آن توییت اختصاص می‌یابد. (۳) تمام توییت‌ها در اختیار تمام حاشیه‌نویسان قرار می‌گیرند و آن‌ها به صورت مستقل عملیات برچسب‌گذاری را انجام می‌دهند و پرتکرارترین برچسب روی یک توییت، به عنوان حالت عاطفی آن توییت انتخاب می‌گردد [9].

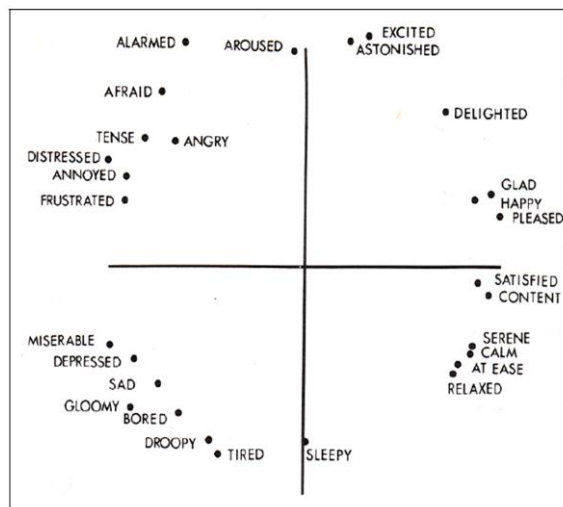
برای انجام حاشیه‌گذاری می‌توانیم یک واسط حاشیه‌گذاری مبتنی بر وب توسعه دهیم که می‌تواند شامل یک میله کشویی برای شش حالت عاطفی و یک میله کشویی برای نشان دادن مقدار آن حالت عاطفی باشد که بین صفر تا صد است. همچنین تمام حاشیه‌نویسان، از قبل می‌توانند قواعد و دستورالعمل‌های مربوطه را دریافت نمایند که این قواعد باعث می‌شود اطمینان پیدا کنیم هر کس توییت‌ها را به شیوه‌ای یکسان تگ نموده است. به طور مثال این قواعد می‌تواند شامل تعریف هر یک از دسته‌های عاطفی و اطلاعاتی درباره آن‌ها باشد.

۳-۴- انتخاب ویژگی‌ها

تبدیل یک تکه متن به یک بردار ویژگی، قدم پایه‌ای در هر رویکردی برای تجزیه و تحلیل عواطف می‌باشد. بدین منظور از یک مجموعه ویژگی‌های قوی که نمایش چکیده‌ای از یک توییت می‌باشد، استفاده می‌شود که در آن از اطلاعات مرتبط با کلمات و ویژگی‌های خاص موجود در یک توییت استفاده شده است. بردار ویژگی یک بردار n بعدی از ویژگی‌های عددی می‌باشد که نشان‌دهنده اشیاء یا ویژگی‌های مختلف است. در اینجا می‌توان به ۴ نوع از ویژگی‌هایی که در کارهای قبلی استفاده شده است، اشاره نمود: شامل نحوی^۱، معنایی^۲، مبتنی بر پیوند^۳ و سبکی^۴ [10].

ویژگی‌های نحوی عمدتاً شامل n -gram، برچسب‌های POS، الگوهای POS و علائم نگارشی (نقطه‌گذاری‌ها) می‌باشند. اگر یک نمایش uni-

مطلب مکانی را بین هر دو واژه مرتبط با عواطف گسسته انتخاب نماید. داده‌های عددی از موقعیت نسبی نقاط در یک فضای دوقطبی و دوبعدی (والنس-اروسال)^۱ به دست می‌آیند. بعد والنس، عواطف مثبت و منفی در دو سمت آن و بعد اروسال حالات آرام در مقابل هیجانی را مشخص می‌کنند [8].



شکل (۱): مدل عواطف چندبعدی راسل [8]

۳-۲- مجموعه دادگان

به منظور جمع‌آوری داده، از شبکه اجتماعی توییت و توییت‌های انگلیسی مربوط به این شبکه اجتماعی استفاده خواهد شد. داده‌های توییت یک منبع غنی به منظور دریافت اطلاعات در هر موضوع قابل تصویری می‌باشد. این داده‌ها می‌تواند در موارد مختلفی مورد استفاده قرار بگیرد از جمله: یافتن روندهای مرتبط با یک کلمه کلیدی خاص، جمع‌آوری بازخورد درباره سرویس‌ها و محصولات جدید و اندازه‌گیری احساسات و عواطف.

توییت‌ها می‌تواند به صورت تصادفی، در موضوعات مختلف و فرمت‌های متفاوتی استخراج شوند. موضوعاتی که قرار است عواطف موجود در آن‌ها استخراج شود، می‌تواند در حوزه‌های مختلف سیاسی، اقتصادی، اجتماعی و غیره باشد. همچنین مجموعه دادگان می‌تواند در دو فرم JSON و CSV تهیه شوند. فرمت JSON شامل تاریخ تولید توییت، نام کاربر، محتوای توییت، محل زندگی کاربر، تعداد دنبال کننده کاربر، ریتوییت و اطلاعات دیگر است. این فرمت بیشتر به درد تحلیل شبکه اجتماعی می‌خورد. در فرمت CSV اطلاعات فقط شامل محتوای توییت است. این مجموعه داده، محتوای توییت‌های جمع‌آوری شده بدون توجه به اطلاعات دیگر مانند نام کاربر و غیره می‌باشد. این فرمت بیشتر به درد متن‌کاوی و نظرکاوی می‌خورد.

در سیستم پیشنهادی به منظور جمع‌آوری داده از جریان توییت API استفاده می‌شود. API مخفف Application Programming Interface یا واسط برنامه‌نویسی کاربردی می‌باشد و ابزاری است که تعاملی را بین برنامه‌های کامپیوتر و سرویس‌های وب ایجاد می‌کند. بسیاری از سرویس‌های وب، API را برای توسعه‌دهندگان به منظور تعامل با سرویس‌هایشان و دسترسی به داده‌ها فراهم می‌کنند. در نتیجه Twitter API نیز برای دانلود توییت‌های مرتبط با کلمات کلیدی خاص استفاده می‌شود.



۳-۵- روش های طبقه بندی عواطف

رویکردهای طبقه بندی عواطف را می توان به طور کلی در دو بخش تقسیم بندی نمود: (۱) طبقه بندی مبتنی بر یادگیری نظارت شده (۲) طبقه بندی مبتنی بر یادگیری بدون نظارت. رویکرد نظارت شده شامل روش های یادگیری ماشین همانند ماشین بردار پشتیبان و بیز ساده بوده و رویکرد بدون نظارت نیز شامل روش لغوی ساده^{۱۶} و روش های NLP^{۱۷} می باشد که در ادامه به صورت مفصل به این روش ها اشاره شده است.

۳-۵-۱- طبقه بندی مبتنی بر یادگیری نظارت شده

روش های یادگیری ماشین که برای تجزیه و تحلیل عواطف قابل اجرا هستند عمدتاً متعلق به دسته بندی نظارت شده و به خصوص روش های طبقه بندی متون می باشند از این رو یادگیری نظارت شده نامیده می شوند. یادگیری نظارت شده نوعی از یادگیری است که در آن ورودی و خروجی مشخص است و ناظر 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 و مجموعه آزمایشی. مجموعه آموزشی به وسیله یک دسته بند خودکار برای یادگیری ویژگی های متمایز از متون استفاده می شود و مجموعه آزمایشی برای ارزیابی دسته بند خودکار استفاده می گردد [10].

سابقه استفاده از روش های یادگیری نظارت شده در تحلیل احساسات به اولین مقاله نوشته شده در سال ۲۰۰۲ برمی گردد. اصولاً روش های یادگیری نظارت شده کاربرد وسیع و گسترده ای داشته و استفاده از آن بسیار مورد توجه می باشد. پرکاربردترین روش یادگیری ماشین نیز ماشین بردار پشتیبان است که بهترین عملکرد را در بین این الگوریتم ها ارائه می دهد.

۳-۵-۲- طبقه بندی مبتنی بر یادگیری بدون نظارت

در یادگیری بدون نظارت برخلاف یادگیری نظارت شده، داده های مشخصی از قبل وجود ندارد و هدف، ارتباط ورودی و خروجی نیست، بلکه تنها دسته بندی آن ها مهم است و این یادگیرنده است که بایستی در داده ها به دنبال ساختاری خاص بگردد. از آنجا که کلمات احساسی فاکتور مهمی برای طبقه بندی عواطف می باشند، دور از تصور نیست که این کلمات و عبارات برای طبقه بندی عواطف در یک روش بدون نظارت استفاده شوند. رویکرد بدون نظارت شامل دو روش اصلی لغوی ساده و روش های NLP می باشد که در ادامه به این روش ها اشاره شده است [11].

در روش لغوی ساده با استفاده از یک لغت نامه که تمام لغات مربوط به دامنه کاری ما در آن طبق مثبت و منفی بودن، امتیاز داده شده اند، متن را تجزیه و تحلیل می کنیم. این روش بسیار ساده بوده و در متن، لغاتی مثل خوب، عالی، بد، زشت و غیره را در نظر گرفته و امتیاز آن ها را جمع می نماید. نتیجه نهایی، امتیاز جمله به دست آمده است که اگر مثبت باشد جمله مثبت و اگر منفی باشد جمله را منفی در نظر می گیریم. این روش با وجود سادگی، به دلیل پیچیدگی ساختار زبانی به ندرت مورد استفاده قرار می گیرد؛ زیرا به

gram استفاده شود، کلمات به تنهایی به عنوان ویژگی استفاده می شوند و اگر یک نمایش bia-gram استفاده شود، جفت هایی از کلمات به عنوان ویژگی استفاده می شوند. برچسب گذاری اجزای کلام یا POS Tag، عمل انتساب برچسب به کلمات و نشانه های تشکیل دهنده یک متن بر اساس نقش آن ها در جمله می باشد. الگوهای POS همانند الگوی n+aj (اسم و به دنبال آن یک صفت مثبت) نشان دهنده یک احساس مثبت و الگوی n+dj (اسم و به دنبال آن یک صفت منفی) نشان دهنده یک احساس منفی است.

دسته دوم ویژگی ها یعنی ویژگی های معنایی روش های حاشیه گذاری خودکار، نیمه خودکار یا دستی را برای اضافه نمودن قطبیت یا امتیاز مرتبط به کلمه ها یا عبارات ترکیب می کند. این ویژگی ها تکیه بر معنای کلمات دارند و شامل دو رویکرد مبتنی بر امتیاز و مبتنی بر لغت می باشند. در روش مبتنی بر امتیاز می توان از فرمول PMI^{۱۵} برای محاسبه امتیاز جهت گیری معنایی یک کلمه یا عبارت استفاده کرد. در روش مبتنی بر لغات نیز از گروه های ارزیابی برای تعیین معنای کلمات یا عبارات استفاده می شود و هر یک از کلمات یا عبارات در دسته های ارزیابی مخصوص به خود قرار می گیرند.

دسته سوم ویژگی ها، ویژگی های مبتنی بر پیوند می باشند. ویژگی های مبتنی بر پیوند از تجزیه و تحلیل های پیوندها و ارجاعات به منظور تعیین عواطف برای متون وب استفاده می کند. از این ویژگی ها به ندرت استفاده می شود و دلیل آن نیز نبود ساختار پیوند مشخص در بسیاری از محتوای وب می باشد.

ویژگی های سبکی نیز شامل استفاده از ویژگی های ساختاری و لغاتی می باشد که این ویژگی ها نیز به دلیل عدم کارایی لازم به ندرت مورد استفاده قرار گرفته است [10].

در جدول (۱) به صورت خلاصه ویژگی های ذکر شده در این بخش آورده شده است.

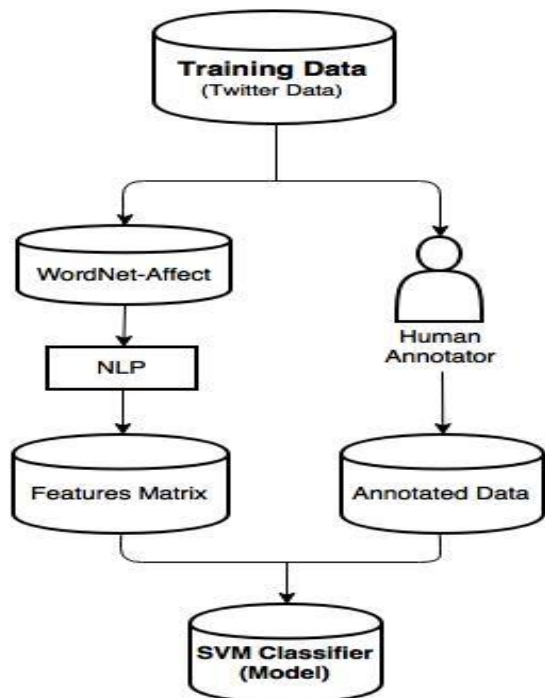
جدول (۱): انواع ویژگی ها

نحوی	معنایی	مبتنی بر پیوند	سبکی
* برچسب های POS (part-of-speech)	* استفاده از فرمول PMI برای محاسبه امتیاز جهت گیری معنایی یک کلمه یا عبارت	استفاده از از تجزیه و تحلیل های پیوندها و ارجاعات لغاتی	استفاده از ویژگی های ساختاری و لغاتی
* الگوهای n-gram			
* الگوهای POS	* استفاده از گروه های ارزیابی برای تعیین معنای کلمات یا عبارات		
* علائم نگارشی (نقطه گذاری ها)			

بنابراین ویژگی ها بسیار بارز بوده و در فرایند طبقه بندی مهم و مفید می باشند. در میان این ویژگی ها، ویژگی های نحوی استفاده بیشتری در تجزیه و تحلیل احساسات داشته و ویژگی ها هر چه قدر وابسته به دامنه باشند، مفیدتر بوده و دقت نتایج به دست آمده را افزایش می دهند.

۴- سیستمی برای راهبرد پیشنهادی

سیستمی که در اینجا به منظور طبقه‌بندی حالات عاطفی در شبکه اجتماعی تویتر پیشنهاد می‌شود، ترکیبی از روش‌های نظارتی و بدون نظارت و همچنین ترکیبی از ویژگی‌های استفاده شده در پژوهش‌های قبلی و ویژگی‌های جدید می‌باشد. برای اینکه بتوانیم دیدی از روش ترکیبی به دست آوریم، در شکل (۳) چهارچوب پیشنهادی برای طبقه‌بندی عواطف با استفاده از الگوریتم ماشین بردار پشتیبان (SVM) و لغت‌نامه WordNet-Affect نمایش داده شده است.



شکل (۳): سیستم پیشنهادی به منظور طبقه‌بندی عواطف

داده‌هایی که برای تحلیل و طبقه‌بندی عواطف در نظر گرفته می‌شود شامل توییت‌های مربوط به هر یک از موضوعات مختلف مطرح در شبکه اجتماعی تویتر می‌باشد که در فرمت CSV و پس از پیش‌پردازش‌های لازم در سیستم استفاده می‌گردد.

توصیه می‌شود در این سیستم از الگوریتم SVM که الگوریتم نظارتی است، استفاده شود. SVM در برابر نویز داده‌ها مقاوم است، می‌تواند با تعداد زیادی ویژگی کار کند و در کارهای مشابه مانند دسته‌بندی متن عملکرد خوبی دارد؛ از این رو ابزار مناسبی به شمار می‌رود. برای استفاده از الگوریتم SVM باید داده‌ها برچسب دار شوند. برچسب‌گذاری این داده‌ها توسط عامل‌های خبره انسانی و مطابق دستورالعمل‌های تهیه‌شده برای آن‌ها و به وسیله روش‌های بیان‌شده در بخش ۳-۳ انجام می‌گردد.

دسته‌های عاطفی می‌توانند هر یک از دسته‌های عواطف معرفی شده در این مقاله باشند ولی پیشنهاد ما تطبیق دادن دسته‌ها با موضوع توییت‌های جمع‌آوری شده و ایجاد دسته‌های وابسته به دامنه است.

به منظور به دست آوردن ویژگی‌ها از لغت‌نامه WordNet-Affect استفاده می‌کنیم و به منظور بالا بردن دقت نتایج، ترکیبی از ویژگی‌های

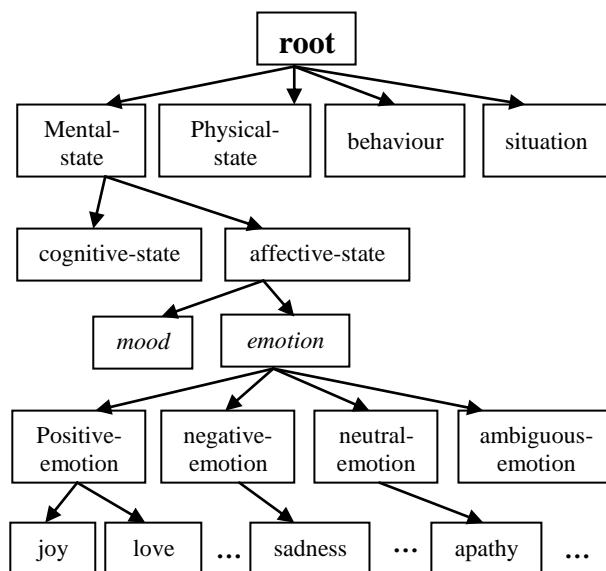
طور مثال یک جمله می‌تواند دارای شماری زیادی از کلمات منفی باشد ولی معنی مثبتی داشته باشد.

در روش‌های NLP به کلمات در یک مقیاس، ارزشی معادل با مثبت و یا منفی بودن آن‌ها داده شده و متن با استفاده از NLP تجزیه و تحلیل می‌شود. روش NLP را می‌توان تکامل‌یافته روش لغوی ساده دانست. این روش چهار مرحله کلی دارد:

۱. شناسایی بخش‌های مهم متن به منظور طبقه‌بندی عواطف
۲. تشخیص ساختار جمله و تشکیل درخت گرامری
۳. امتیازدهی به کلمات مهم تشخیص داده‌شده در متن (features)
۴. امتیازدهی نهایی به متن

در نهایت با محاسبه امتیاز همه ویژگی‌ها برای متن، دو امتیاز کلی مثبت و منفی محاسبه می‌شود که می‌تواند از جمع امتیازهای مثبت و منفی ویژگی‌ها و یا میانگین امتیازهای مثبت و منفی ویژگی‌ها به دست آید [11].

در هر دوی این روش‌ها از یک لغت‌نامه به منظور نگاشت کلمات به یک دسته‌بندی (مثبت، منفی، خنثی و ...) و یا به یک امتیاز احساسی استفاده شده است که این منبع به وسیله الگوریتمی برای تعیین احساس کلی موجود در متن استفاده می‌شود. لغت‌نامه‌های زیادی برای تعیین قطبیت کلمه‌ها وجود دارد. یکی از کاربردی‌ترین و مهم‌ترین این لغت‌نامه‌ها، لغت‌نامه WordNet-Affect می‌باشد که در سیستم پیشنهادی نیز از آن استفاده شده است. این منبع توسعه‌یافته WordNet بوده و در آن، مجموعه‌های مترادف با یک مفهوم عاطفی به‌عنوان A-LABELS تعریف می‌شوند. به‌عنوان مثال واژه‌های joy و elation با مفهوم JOY برچسب‌گذاری می‌شوند و یا در سطوح بالاتر euphoria با مفهوم positive-emotion برچسب‌گذاری می‌شود. در واقع نگاشت بر مبنای یک سلسله‌مراتب برچسب‌گذاری عواطف خودکار و مستقل از دامنه با تأکید بر روابط موجود در WordNet انجام می‌شود که بخشی از آن را در شکل (۲) مشاهده می‌کنیم [12].



شکل (۲): گراف جزئی از WordNet-Affect

Computational Linguistics, China, p.36-44, August 23-27, 2010.

- [6] Purver, Matthew, and Stuart Battersby. "Experimenting with distant supervision for emotion classification." Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics. Association for Computational Linguistics, 2012.
- [7] J. Brynielsson, F. Johansson, C. Jonsson, and A. Westling, "Emotion classification of social media posts for estimating people's reactions to communicated alert messages during crises," Security Informatics, vol. 3, no.7, 2014.
- [8] Mac Kim, Sunghwan. "Recognising Emotions and Sentiments in Text." University of Sydney, 2011.
- [9] Constantine, Layale, and Hazem Hajj. "A survey of ground-truth in emotion data annotation." Pervasive Computing and Communications Workshops (PERCOM Workshops), 2012 IEEE International Conference on. IEEE, 2012.
- [10] Rahate, Rohini S., and M. Emmanuel. "Feature Selection for Sentiment Analysis by using SVM." Int J Comput Appl 84.5, 24-32, 2013.
- [11] Taboada, Maite, et al. "Lexicon-based methods for sentiment analysis." Computational linguistics 37.2, p.267-307, 2011.
- [12] Musto, Cataldo, Giovanni Semeraro, and Marco Polignano. "A comparison of Lexicon-based approaches for Sentiment Analysis of microblog posts." Information Filtering and Retrieval, 2014.

زیر نویس ها

- 1 Sentiment Analysis
- 2 Nasukawa and Yi
- 3 Wiebe et al.
- 4 Alm et al
- 5 Purver and Battersby
- 6 Distant Supervision
- 7 Brynielsson
- 8 Ekman
- 9 Russell
- 10 Valence-Arousal
- 11 Syntactic
- 12 Semantic
- 13 Link-base
- 14 Stylistic
- 15 Pointwise Mutual Information
- 16 Basic lexicon analysis
- 17 Natural Language Process

استفاده شده در کارهای قبلی و ویژگی‌های جدیدی همانند استفاده از شکلک‌های موجود در متون در بردار ویژگی عواطف استفاده می‌شود.

نهایتاً برای هر توییت ما یک برداری از ویژگی‌ها خواهیم داشت. برای استفاده از این ویژگی‌ها در الگوریتم SVM حاصل ضرب تکرار هر کلمه در مقدار قطبیت آن را می‌توان در نظر گرفت. هر توییت مجموعه‌ای از کلمه‌ها در نظر گرفته می‌شود؛ در نتیجه هر توییت برداری از کلمه‌ها است. با جست‌وجوی هر کلمه در لغت‌نامه، قطبیت آن تعیین می‌شود و با مقدار عددی به نمایش درمی‌آید. این مقدار در تعداد تکرار هر کلمه ضرب می‌شود و وزن جدیدی را ایجاد می‌کند. در نتیجه این فرایند، برای هر کلمه وزنی ایجاد می‌شود که حاصل ضرب قطبیت آن کلمه در تعداد تکرار آن در توییت می‌باشد.

همچنین به وسیله داده‌های آزمایشی و با استفاده از پارامترهای محاسبه‌شده از داده‌های آموزشی، ارزیابی SVM صورت گرفته و میزان دقت و درستی روش، مورد سنجش قرار می‌گیرد.

۵- نتیجه

در این مقاله، ضمن آشنایی با مفاهیم و روش‌های تشخیص و طبقه‌بندی حالات عاطفی در دسته‌های از پیش تعیین شده، سیستمی نیز بدین منظور ارائه گردید. روش NLP از لحاظ سرعت و قدرت پردازش سریع‌تر از سایر روش‌ها عمل می‌کند اما میزان خطای بالای آن و همچنین الزام نوشتن قواعد بسیار برای زبان باعث سختی و پیچیدگی آن می‌شود. از طرفی روش‌های یادگیری ماشین نیز به دلیل عدم نیاز به دسترسی به ساختار جمله، برای رسیدن به دقت قابل قبول نیاز به حجم داده‌ای بزرگی به منظور طبقه‌بندی دارد که می‌تواند هزینه‌های پردازشی و محاسباتی بالایی داشته باشد.

از این رو سیستم پیشنهادی می‌تواند مزیت‌های زیر را دارا باشد:

- (۱) کاهش هزینه‌ها
- (۲) نتایج خوب برای هر حجم داده‌ای
- (۳) در نظر گرفتن ساختار جملات و قطبیت کلمات
- (۴) افزایش دقت

همچنین با تغییر در ساختار ویژگی‌ها قادر خواهیم بود به دقت مدنظر خود دست یابیم و سیستم پیشنهادی در هر زمان قابل بهبود خواهد بود.

مراجع

- [1] Liu, Bing. "Sentiment analysis and opinion mining." Synthesis Lectures on Human Language Technologies 5.1, 2012.
- [2] Li, Zhaolong. "Analyzing emotion on Twitter for user modeling." Diss. TU Delft, Delft University of Technology, 2013.
- [3] Wiebe, J., Wilson, T., Cardie, "Annotating expressions of opinions and emotions in language. Language Resources and Evaluation 39(2-3), 165-210, 2005.
- [4] Alm, C.O., Roth, D., Sproat, R.: Emotions from text: machine learning for text-based emotion prediction. In: Proc. of the Joint Conf. on Human Language Technology/Empirical Methods in Natural Language Processing, pp. 579-586, 2005.
- [5] Luciano Barbosa, Junlan Feng, Robust sentiment detection on Twitter from biased and noisy data, Proceedings of the 23rd International Conference on