

فهرست مطالب

۱	مقدمه
۳	مروری بر کارهای پیشین
۳	۲.۱ معرفی
۶	۲.۲ تجزیه و تحلیل احساسات
۷	۲.۲.۱ سطوح تحلیل احساسات
۸	۲.۳ تشخیص احساسات
۸	۲.۳.۱ مدل های هیجانی / نظریه های هیجانی
۱۴	۲.۴ مروری بر روش های پیشین
۱۴	۲.۴.۱ شناختن داده ها
۱۶	۲.۴.۲ پیش پردازش داده ها
۱۷	۲.۴.۳ استخراج ویژگی
۱۹	۲.۴.۴ تکنیک های تجزیه و تحلیل احساسات و تشخیص احساسات
۲۰	۲.۴.۵ تکنیک های تحلیل احساسات
۲۵	۲.۴.۶ تکنیک های تشخیص احساسات
۲۹	۳ روش پیشنهادی
۳۰	۳.۱ شناختن داده ها
۳۲	۳.۲ پیش پردازش داده ها
۳۲	۳.۲.۱ حذف کردن url
۳۴	۳.۲.۲ حذف mentions
۳۶	۳.۲.۳ تمیز کردن متن
۳۹	۳.۳ تقسیم داده ها
۴۲	۳.۴ ساختن مدل
۴۴	۳.۵.۱ SVC
۴۷	۳.۵.۲ Random Forest
۵۰	۳.۵.۳ XGBoost
۵۳	۳.۵.۴ Naive Bayes
۵۶	۴ ارزیابی روش
۵۶	۴.۱ accuracy
۵۶	۴.۲ Recall
۵۷	۴.۳ F1-Score
۵۷	۴.۴ Specificity
۵۷	۴.۵ False Positive Rate (FPR)
۵۷	۴.۶ False Negative Rate (FNR)
۵۸	۴.۷ Area Under the Receiver Operating Characteristic (ROC AUC)
۶۰	۵ نتایج
۶۱	۶ جمع بندی و کارهای آینده
۶۲	۷ مراجع

فهرست شکلها

- شکل 1: مدل بعدی احساسات.....11
- شکل 2: تصویرسازی انواع مدل های عاطفی با برخی حالات روانی.....13
- شکل 3: تکنیک هایی برای تجزیه و تحلیل احساسات و تشخیص احساسات.. 20.

فهرست جدولها

- جدول 1: مدل های هیجانی که توسط روانشناسان مختلف تعریف شده است..

۱ مقدمه

در عصر دیجیتال امروزی، حجم وسیعی از داده‌های متنی موجود به صورت آنلاین، پنجره‌ای به احساسات، نظرات و احساسات جمعی افراد و جوامع می‌دهد. تجزیه و تحلیل احساسات، همچنین به عنوان نظر کاوی شناخته می‌شود، یک تکنیک پردازش زبان طبیعی (NLP) است که به دنبال استخراج و درک احساسات بیان شده در متن است. نقش مهمی در رمزگشایی اینکه آیا یک قطعه متن معین احساسات مثبت، منفی یا خنثی را منتقل می‌کند، ایفا می‌کند.

تجزیه و تحلیل احساسات به ابزاری قدرتمند با کاربرد در زمینه‌های مختلف، از تجارت و بازاریابی گرفته تا نظارت بر رسانه‌های اجتماعی و تحلیل سیاسی تبدیل شده است. با خودکار کردن فرآیند تشخیص احساسات، سازمان‌ها می‌توانند بینش‌های ارزشمندی به دست آورند، تصمیمات مبتنی بر داده‌ها را اتخاذ کنند و به بازخورد مشتریان پاسخ مؤثری دهند و در نهایت محصولات، خدمات و استراتژی‌های خود را بهبود بخشند.

در این عصر دیجیتال، که داده‌های متنی فراوان است و ارتباطات عمده‌تاً مبتنی بر متن است، تجزیه و تحلیل احساسات به عنوان ابزاری حیاتی برای درک چشم انداز احساسی دنیای آنلاین است. این نمای کلی مقدماتی زمینه را برای کاوش در پیچیدگی‌ها، روش‌شناسی‌ها و کاربردهای دنیای واقعی تحلیل احساسات فراهم می‌کند.

با رشد پلتفرم های مختلف در زمینه شبکه های اجتماعی این پلتفرم های به ابزاری قدرتمند برای افراد تبدیل شده اند تا افکار، نظرات و احساسات خود را در مورد مسائل مختلف بیان کنند. ما در این پروژه قصد داریم تکنیک های مختلف تحلیل احساسات را به کار گرفته و احساسات را در ۱۲ گروه مختلف طبقه بندی کنیم

این تحقیق با یک مرحله پیش پردازش کامل شروع می شود، شامل حذف URL ها و اشاره ها، پاک کردن متن، تبدیل آن به حروف کوچک، حذف فضاها، خالی اضافی، نشانه گذاری متن، حذف کلمات توقف، علائم نقطه گذاری، و اعداد، و در نهایت واژه سازی یا ریشه کردن کلمات. هدف این مراحل پیش پردازش آماده سازی داده های متنی برای تجزیه و تحلیل بعدی است.

مرحله بعدی شامل رمزگذاری برچسب های احساسات برای تسهیل آموزش و ارزیابی مدل های تحلیل احساسات است. سپس داده ها به مجموعه های آموزشی و آزمایشی برای توسعه و ارزیابی مدل تقسیم می شوند. برای نمایش موثر داده های متنی، از تکنیک بردارسازی TF-IDF استفاده می شود که اهمیت اصطلاحات را در اسناد نشان می دهد.

چندین مدل یادگیری ماشین در این تحقیق استفاده شده است، از جمله طبقه بندی بردار پشتیبانی (SVC)، جنگل تصادفی، XGBoost و Naive Bayes. این مدل ها برای طبقه بندی احساسات بیان شده در پیام های توییتر بر روی مجموعه داده ها آموزش داده شده و آزمایش می شوند.

ارزیابی عملکرد مدل ها با محاسبه دقت با استفاده از طبقه بندی گزارش انجام می شود، که بینش هایی درباره دقت، یادآوری، امتیاز $F1$ و سایر معیارهای ارزیابی برای هر کلاس احساسی ارائه می دهد. این

تجزیه و تحلیل به ارزیابی اثربخشی و محدودیت‌های مدل‌ها در گرفتن تفاوت‌های ظریف احساسات در پیام‌های تویتر کمک می‌کند.

یافته‌های این تحقیق به حوزه تحلیل احساسات، به ویژه در زمینه پیام‌های شبکه اجتماعی تویتر کمک می‌کند. این مطالعه بینش‌هایی را در مورد عملکرد مدل‌های مختلف یادگیری ماشین و توانایی آنها در طبقه‌بندی احساسات در توییت‌هایی که با ۱۲ احساس مختلف حاشیه‌نویسی شده‌اند، ارائه می‌کند. نتایج در زمینه‌هایی مانند نظارت بر رسانه‌های اجتماعی، مدیریت شهرت برند و تحلیل افکار عمومی پیامدهای عملی دارند.

۲ مروری بر کارهای پیشین

۲.۱ معرفی

درک زبان انسانی و تولید زبان انسانی دو جنبه پردازش زبان طبیعی (NLP) هستند. با این حال، اولی به دلیل ابهامات در زبان طبیعی دشوارتر است. تشخیص گفتار، خلاصه‌سازی اسناد، پاسخ به سؤال، سنتز گفتار، ترجمه ماشینی و سایر کاربردها همگی از NLP استفاده می‌کنند (ایتانی و همکاران، ۲۰۱۷). دو حوزه مهم پردازش زبان طبیعی، تحلیل احساسات و تشخیص احساسات است. اگرچه گاهی اوقات این دو نام به جای یکدیگر استفاده می‌شوند، اما از چند جنبه با هم تفاوت دارند. تحلیل احساسات وسیله‌ای برای ارزیابی مثبت، منفی یا خنثی بودن داده‌ها است.

در مقابل، تشخیص احساسات ابزاری برای شناسایی انواع احساسات انسانی متمایز مانند خشمگین، شاد و یا افسرده. «تشخیص احساسات»، «محاسبات عاطفی»، «تحلیل احساسات» و «شناسایی احساسات» همگی عباراتی هستند که گاهی به جای یکدیگر استفاده می‌شوند (Munezero et al. ۲۰۱۴). از زمانی که خدمات اینترنتی بهبود یافته است، مردم از رسانه‌های اجتماعی برای انتقال احساسات خود استفاده می‌کنند. در رسانه‌های اجتماعی، مردم آزادانه احساسات، استدلال‌ها، نظرات خود را در مورد طیف گسترده‌ای از موضوعات بیان می‌کنند. علاوه بر این، بسیاری از کاربران در سایت‌های مختلف تجارت الکترونیک، محصولات و خدمات مختلف را بازخورد و بررسی می‌کنند. رتبه‌بندی‌ها و بررسی‌های کاربران در چندین پلتفرم، فروشندگان و ارائه‌دهندگان خدمات را تشویق می‌کند تا سیستم‌ها،

کالاها یا خدمات فعلی خود را بهبود بخشند. امروزه تقریباً هر صنعت یا شرکتی در حال گذراندن دوره‌ای از انتقال دیجیتالی است که منجر به افزایش حجم عظیمی از داده‌های ساختاریافته و بدون ساختار می‌شود. وظیفه بزرگ شرکت‌ها تبدیل داده‌های بدون ساختار به بینش‌های معنادار است که می‌تواند به آنها در تصمیم‌گیری کمک کند (احمد و همکاران، ۲۰۲۰).

به عنوان مثال، در دنیای تجارت، فروشندگان از پلتفرم‌های رسانه‌های اجتماعی مانند اینستاگرام، یوتیوب، توییتر و فیس بوک برای پخش اطلاعات در مورد محصول خود و جمع‌آوری کارآمد بازخورد مشتریان استفاده می‌کنند (Agbehadji and Ijabad-eniyi، ۲۰۲۱). بازخورد فعاله مردم نه تنها برای بازاریابان کسب‌وکار برای سنجش رضایت مشتری و پیگیری رقابت، بلکه برای مصرف‌کنندگانی که می‌خواهند قبل از خرید درباره یک محصول یا خدمات بیشتر بدانند ارزشمند است. آی تی. تجزیه و تحلیل احساسات به بازاریابان کمک می‌کند تا دیدگاه‌های مشتریان خود را بهتر درک کنند تا تغییرات لازم را در محصولات یا خدمات خود ایجاد کنند (Jang et al، ۲۰۱۳). اجراوی و همکاران (۲۰۲۱). در هر دو حالت پیشرفته و در حال ظهور کشورها، تأثیر احساسات تجاری و مشتریان بر عملکرد بازار سهام را می‌توان مشاهده کرد. علاوه بر این، ظهور رسانه‌های اجتماعی تعامل سرمایه‌گذاران در بازار سهام را آسان‌تر و سریع‌تر کرده است. در نتیجه، احساسات سرمایه‌گذاران بر تصمیم‌های سرمایه‌گذاری آنها تأثیر می‌گذارد که می‌تواند به سرعت در شبکه گسترش پیدا کند و بزرگ‌تر شود و بازار سهام تا حدی تغییر کند (احمد، ۲۰۲۰). در نتیجه، تحلیل احساسات و عواطف، روشی را که ما برای کسب و کار انجام می‌دهیم تغییر داده است (بهاردواج و همکاران، ۲۰۱۵).

در بخش مراقبت‌های بهداشتی، رسانه‌های اجتماعی آنلاین مانند توییتر به منابع ضروری اطلاعات مرتبط با سلامتی تبدیل شده‌اند که توسط متخصصان مراقبت‌های بهداشتی و شهروندان ارائه می‌شود. برای مثال، مردم افکار، نظرات و احساسات خود را در مورد همه‌گیری کووید-۱۹ به اشتراک گذاشته‌اند (گارسیا و برتون، ۲۰۲۱). به بیماران دستور داده شد که از عزیزان خود جدا بمانند که به سلامت روان آنها آسیب می‌رساند. برای نجات بیماران از مسائل مربوط به سلامت روانی مانند افسردگی، پزشکان بهداشت باید از تجزیه و تحلیل خودکار احساسات و عواطف استفاده کنند (سینگ و همکاران، ۲۰۲۱). مردم معمولاً احساسات یا باورهای خود را در سایت‌ها از طریق پست‌هایشان به اشتراک می‌گذارند، و

اگر فردی افسرده به نظر می‌رسد، مردم می‌توانند برای کمک به او مراجعه کنند، بنابراین از بدتر شدن شرایط سلامت روان جلوگیری می‌کنند.

تجزیه و تحلیل احساسات و عواطف نقش مهمی در بخش آموزش، هم برای معلمان و هم برای دانش آموزان ایفا می‌کند. کارآمدی یک معلم نه تنها با مدرک تحصیلی او تعیین می‌شود، بلکه با اشتیاق، استعداد و فداکاری او نیز تعیین می‌شود. گرفتن بازخورد به موقع از دانش آموزان موثرترین تکنیک برای معلم برای بهبود رویکردهای تدریس است (سان-گیتا و پرابها). (۲۰۲۰). مشاهده بازخورد متنی با پایان باز دشوار است و نتیجه گیری دستی نیز چالش برانگیز است. یافته های تجزیه و تحلیل احساسات و تحلیل احساسات به معلمان و سازمان ها در انجام اقدامات اصلاحی کمک می کند. از زمان تأسیس سایت اجتماعی، مؤسسات آموزشی به طور فزاینده ای به رسانه های اجتماعی مانند فیس بوک و توییتر برای اهداف بازاریابی و تبلیغات متکی هستند. دانشجویان و سرپرستانه تحقیقات آنلاین قابل توجهی را انجام می دهند و در مورد مؤسسه، دوره ها و اساتید بالقوه اطلاعات بیشتری کسب می کنند. آنها از وبلاگ ها و دیگر انجمن های گفتگو برای تعامل با دانشجویانی که علایق مشترک دارند و ارزیابی کیفیت کالج ها و دانشگاه های احتمالی استفاده می کنند. بنابراین، استفاده از تحلیل احساسات و عواطف می تواند به دانش آموز کمک کند تا بهترین مؤسسه یا معلم را در فرآیند ثبت نام خود انتخاب کند (Archana Rao و Baglodi). (۲۰۱۷).

تحلیل احساسات و عواطف کاربردهای گسترده ای دارد و با استفاده از روش شناسی های مختلف قابل انجام است. سه نوع تکنیک تحلیل احساسات و عواطف وجود دارد: مبتنی بر واژگان، مبتنی بر یادگیری ماشینی و مبتنی بر یادگیری عمیق. هر کدام مزایا و معایب خاص خود را دارند. علیرغم تشخیص احساسات و عواطف متفاوت محققان با چالش های مهمی از جمله برخورد با زمینه، تمسخر، جملاتی که چندین احساسات را منتقل می کنند، گسترش عامیانه وب و ابهامات واژگانی و نحوی مواجه هستند. علاوه بر این، از آنجایی که هیچ قانون استاندارد برای برقراری ارتباط احساسات در چندین پلتفرم وجود ندارد، برخی آن ها را با تأثیری باورنکردنی بیان می کنند، برخی احساسات خود را خفه می کنند و برخی پیام خود را به صورت منطقی ساختار می دهند. بنابراین، توسعه تکنیکی که بتواند به طور موثر در همه حوزه ها کار کند، چالش بزرگی برای محققان است.

در این مقاله مروری، بخش ۲، تجزیه و تحلیل احساسات و سطوح مختلف آن، تشخیص احساسات و مدل های روانی-منطقی را معرفی می کند. بخش ۳ مراحل متعددی را که در تجزیه و تحلیل احساسات و عواطف دخیل هستند، از جمله مجموعه داده ها، پیش پردازش متن، تکنیک های استخراج ویژگی، و رویکردهای مختلف تحلیل احساسات و عواطف را مورد بحث قرار می دهد. بخش ۴ چالش های متعددی را که محققان در طول تجزیه و تحلیل احساسات و عواطف با آن مواجه هستند، بررسی می کند. در نهایت، بخش ۵ کار را به پایان می رساند.

۲.۲ تجزیه و تحلیل احساسات

در حال حاضر بسیاری از مردم در سراسر جهان از وبلاگ ها، انجمن ها و سایت های رسانه های اجتماعی مانند توییتر و فیس بوک برای به اشتراک گذاشتن نظرات خود با سایر نقاط جهان استفاده می کنند. رسانه های اجتماعی به یکی از موثرترین رسانه های ارتباطی موجود تبدیل شده اند. در نتیجه، مقدار زیادی داده تولید می شود که داده های بزرگ نامیده می شود و تجزیه و تحلیل احساسات برای تجزیه و تحلیل موثر و کارآمد این کلان داده معرفی شد (ناگامانجولا و پتالاکشمی ۲۰۲۰). درک احساسات کاربر برای صنعت یا سازمان بسیار مهم است. تجزیه و تحلیل احساسات، که اغلب به عنوان عقیده کاوی شناخته می شود، روشی برای تشخیص مثبت یا منفی بودن دیدگاه نویسنده یا کاربر در مورد یک موضوع است. تحلیل احساسات به عنوان فرآیند به دست آوردن اطلاعات و معناشناسی معنادار از متن با استفاده از تکنیک های پردازش طبیعی و تعیین نگرش نویسنده که ممکن است مثبت، منفی یا خنثی باشد، تعریف می شود (Onyenwe et al. ۲۰۲۰). از آنجایی که هدف از تجزیه و تحلیل احساسات، تعیین قطبیت و طبقه بندی متون دارای نظر مثبت یا منفی است، محدوده کلاس مجموعه داده درگیر در تحلیل احساسات فقط به مثبت یا منفی محدود نمی شود. می تواند موافق یا مخالف باشد، خوب یا بد. همچنین می توان آن را در مقیاس ۵ درجه ای تعیین کرد: کاملاً مخالفم، مخالفم، خنثی، موافقم، یا کاملاً موافقم (پرابوو و تلوال ۲۰۰۹). به عنوان مثال، (Ye et al. ۲۰۰۹) تجزیه و تحلیل احساسات را روی نظرات در مقاصد اروپایی و ایالات متحده با برچسب در مقیاس ۱ تا ۵ اعمال کردند. آنها نظرات ۱ یا ۲ ستاره را با قطبیت منفی و نظرات بیش از ۲ ستاره را با قطبیت مثبت مرتبط کردند. گرابنر و همکاران (۲۰۱۲) یک واژگان مخصوص دامنه را ساخت که از نشانه هایی با ارزش احساسی

آنها تشکیل شده است. این توکن ها بودناز نظرات مشتریان در حوزه گردشگری برای طبقه بندی احساسات به رتبه بندی های ۵ ستاره از وحشتناک تا عالی در حوزه گردشگری جمع آوری شده است. علاوه بر این، تحلیل احساسات از متن می تواند در سه سطح مورد بحث در بخش زیر انجام شود. سالینکا (۲۰۱۵) الگوریتم های یادگیری ماشینی را روی مجموعه داده Yelp اعمال کرد که شامل بررسی ارائه دهندگان خدمات از ۱ تا ۵ است. تجزیه و تحلیل احساسات را می توان در سه سطح طبقه بندی کرد که در بخش بعدی اشاره شد.

۲.۲.۱ سطوح تحلیل احساسات

تحلیل احساسات در سه سطح امکان پذیر است: سطح جمله، سطح سند و سطح جنبه. در تحلیل احساسات در سطح جمله یا عبارت، اسناد یا پاراگراف ها به جملات تقسیم می شوند و قطبیت هر جمله مشخص می شود (Prabhakar و Meena، ۲۰۰۷؛ آرولموروگان و همکاران ۲۰۱۹؛ شیرست و همکاران ۲۰۱۹). در سطح سند، احساس از کل سند یا رکورد شناسایی می شود (Pu et al. ۲۰۱۹). ضرورت تجزیه و تحلیل احساسات در سطح سند، استخراج احساسات جهانی از متون طولانی است که حاوی الگوهای محلی اضافی و نویز زیادی هستند. چالش برانگیزترین جنبه طبقه بندی احساسات در سطح سند، در نظر گرفتن پیوند بین کلمات و عبارات و زمینه کامل اطلاعات معنایی برای انعکاس ترکیب سند است (رائو و همکاران ۲۰۱۸؛ لیو و همکاران ۲۰۲۰ a). این امر مستلزم درک عمیق تر ساختار درونی درونی احساسات و کلمات وابسته است (لیو و همکاران ۲۰۲۰ b). در سطح جنبه، تحلیل احساسات، نظر در مورد یک جنبه یا ویژگی خاص تعیین می شود. به عنوان مثال، سرعت پردازنده بالا است، اما این محصول قیمت بالایی دارد. در اینجا سرعت و هزینه دو جنبه یا دیدگاه هستند. سرعت در جمله اشاره شده است، از این رو جنبه صریح نامیده می شود، در حالی که هزینه یک جنبه ضمنی است. تجزیه و تحلیل احساسات سطح جنبه کمی سخت تر از دو مورد دیگر است زیرا شناسایی ویژگی های ضمنی دشوار است. دوی سری ناندینی و پرادیپ (۲۰۲۰) الگوریتمی را برای استخراج جنبه های ضمنی از اسناد بر اساس فراوانی همزمانی جنبه با شاخص ویژگی و با بهره برداری از رابطه بین کلمات نظری و جنبه های صریح پیشنهاد کرد. ما و همکاران (۲۰۱۹) به دو موضوع مربوط به تحلیل سطح جنبه رسیدگی کرد: جنبه های مختلف در یک جمله با قطبیت های متفاوت و موقعیت صریح بافت در یک جمله نظری. نویسندگان یک مدل دو مرحله ای بر اساس LSTM با مکانیزم توجه

برای حل این مسائل ایجاد کردند. آنها این مدل را بر اساس این فرض پیشنهاد کردند که کلمات بافت نزدیک به جنبه مرتبط تر هستند و نیاز به توجه بیشتری نسبت به کلمات بافت دورتر دارند. در مرحله یک، مدل از جنبه های متعدد در یک جمله یک به یک با مکانیسم توجه موقعیت بهره برداری می کند. سپس در حالت دوم، جفت ها (جنبه، جمله) را با توجه به موقعیت جنبه و زمینه اطراف آن شناسایی می کند و قطبیت هر تیم را به طور همزمان محاسبه می کند. همانطور که قبلاً گفته شد، تحلیل احساسات و تحلیل احساسات اغلب به جای یکدیگر توسط محققان مورد استفاده قرار می گیرند. با این حال، آنها از چند جهت متفاوت هستند. در تجزیه و تحلیل احساسات، قطبیت نگرانی اصلی است، در حالی که، در تشخیص احساسات، حالت یا خلق عاطفی یا روانی تشخیص داده می شود. تجزیه و تحلیل احساسات به طور استثنایی ذهنی است، در حالی که تشخیص احساسات عینی تر و دقیق تر است. بخش ۲.۲ همه چیز در مورد تشخیص احساسات را با جزئیات شرح می دهد.

۲.۳ تشخیص احساسات

احساسات جزء جدایی ناپذیر زندگی انسان هستند. این احساسات بر تصمیم گیری انسان تأثیر می گذارد و به ما کمک می کند تا به روشی بهتر با جهان ارتباط برقرار کنیم. تشخیص عواطف، که به عنوان تشخیص احساسات نیز شناخته می شود، فرآیند شناسایی احساسات یا عواطف مختلف یک فرد (به عنوان مثال، شادی، غم یا خشم) است. محققان در چند سال گذشته سخت کار کرده اند تا تشخیص احساسات را خودکار کنند. با این حال، برخی از فعالیت های فیزیکی مانند ضربان قلب، لرز دست ها، تعریق و زیر و بمی صدا نیز وضعیت عاطفی فرد را منتقل می کنند (Kratzwald et al. ۲۰۱۸)، اما تشخیص احساسات از روی متن بسیار سخت است. علاوه بر این، هر روز که می گذرد، ابهامات مختلف و اصطلاحات یا اصطلاحات عامیانه جدید معرفی می شود، تشخیص احساسات از متن را چالش برانگیزتر می کند. علاوه بر این، تشخیص احساسات فقط به شناسایی شرایط روانی اولیه (شادی، غمگین، خشم) محدود نمی شود. در عوض، بسته به مدل احساسی، به مقیاس ۶ یا ۸ می رسد.

۲.۳.۱ مدل های هیجانی / نظریه های هیجانی

در زبان انگلیسی، کلمه "احساس" در قرن هفدهم به وجود آمد که از کلمه فرانسوی "emo-tion" به معنای اختلال جسمانی گرفته شده است. قبل از قرن نوزدهم، اشتیاق، اشتها و محبت ها به عنوان حالت های روانی طبقه بندی می شدند. در قرن نوزدهم، واژه «هیجان» یک اصطلاح روانشناختی در نظر گرفته شد (دیکسون ۲۰۱۲). در روانشناسی، حالات پیچیده احساس منجر به تغییر در افکار، اعمال، رفتار و شخصیت می شود که به آن احساسات می گویند. به طور کلی، مدل های روان شناختی یا هیجانی به دو دسته طبقه بندی می شوند: بعدی و مقوله ای.

مدل عاطفه بعدی این مدل احساسات را بر اساس سه پارامتر نشان می دهند ظرفیت، برانگیختگی و قدرت (باکر و همکاران ۲۰۱۴). والانس به معنای قطبیت است و برانگیختگی به معنای هیجان انگیز بودن یک احساس است. به عنوان مثال، خوشحال بودن بیشتر هیجان انگیز است تا خوشحال. قدرت یا تسلط به معنای محدودیت بر احساسات است. این پارامترها موقعیت حالات روانی را در فضای دو بعدی تعیین می کنند، همانطور که در شکل ۱ نشان داده شده است. ۱.

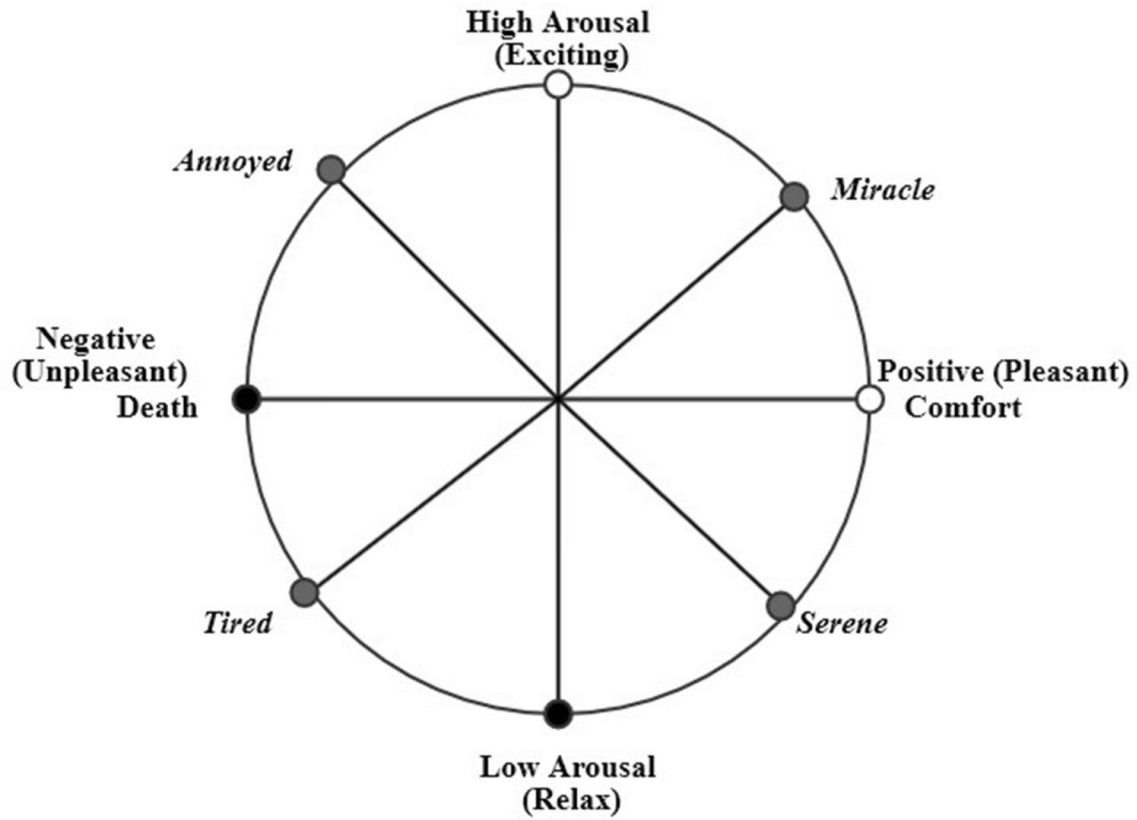
مدل احساسات طبقه بندی شده

در مدل مقوله ای، احساسات به صورت گسسته تعریف می شوند، مانند خشم، شادی، غم و ترس. بسته به در مدل طبقه بندی خاص، احساسات به چهار، شش یا هشت دسته طبقه بندی می شوند.

جدول ۱ مدل های هیجانی متعددی را نشان می دهد که ابعادی و مقوله ای هستند. در حوزه تشخیص احساسات، اکثر محققان مدل احساسات اکمن و پلاچیک را اتخاذ کردند. حالت های احساسی تعریف شده توسط مدل ها مجموعه ای از برجسب ها را تشکیل می دهد که برای حاشیه نویسی جملات یا اسناد استفاده می شود. Batbaatar و همکاران. (۲۰۱۹) بکر و همکاران. (۲۰۱۷) جین و همکاران. (۲۰۱۷) شش احساس اساسی اکمن را پذیرفت. سای لونا و الحاج (۲۰۱۹) از مدل های اکمن برای حاشیه نویسی توییت ها استفاده کرد. برخی از محققان با گسترش مدل با یک یا دو حالت اضافی، از مدل های احساسی سفارشی استفاده کردند. رابرتز و همکاران (۲۰۱۲) از مدل اکمن برای حاشیه نویسی توییت ها با حالت "عشق" استفاده کرد. احمد و همکاران (۲۰۲۰) چرخ احساسات مدل سازی شده

توسط پلوچیک را برای برچسب‌گذاری جملات هندی با ۹ حالت مختلف مدل پلوچیک، کاهش سردرگمی معنایی، در میان کلمات دیگر، اتخاذ کرد. حالت های مدل پلاچیک و اکمن نیز در واژگان دست ساز مختلف مانند WordNet-Affect استفاده می شود (Strapparava et al. ۲۰۰۴) و NRC (محمد و تورنی ۲۰۱۳) واژگان کلمه-احساس. لوبرت و پارلامیس (۲۰۱۹) به مدل Shaver به دلیل ساختار سلسله مراتبی سه سطحی احساسات اشاره کرد. ظرفیت یا قطبیت در سطح اول ارائه می شود، سپس سطح دوم شامل پنج احساس است و سطح سوم ۲۴ حالت احساسی گسسته را نشان می دهد. برخی از محققان به هیچ مدلی اشاره نکردند و مجموعه داده را به سه احساس اصلی طبقه بندی کردند: خوشحال، غمگین یا عصبانی.

شکل ۲ حالت های عاطفی متعددی را که در مدل های مختلف یافت می شود به تصویر می کشد. این حالت ها با در نظر گرفتن مدل پلاچیک به عنوان مدل پایه بر روی چهار محور ترسیم می شوند. رارایج ترین حالت های هیجانی مورد استفاده در مدل های مختلف شامل خشم، ترس، شادی، تعجب و انزجار است که در شکل بالا نشان داده شده است. از شکل می توان دریافت که احساسات در دو طرف محور همیشه مخالف یکدیگر نخواهند بود. مثلاً غم و شادی متضاد یکدیگرند، اما خشم مخالف ترس نیست.

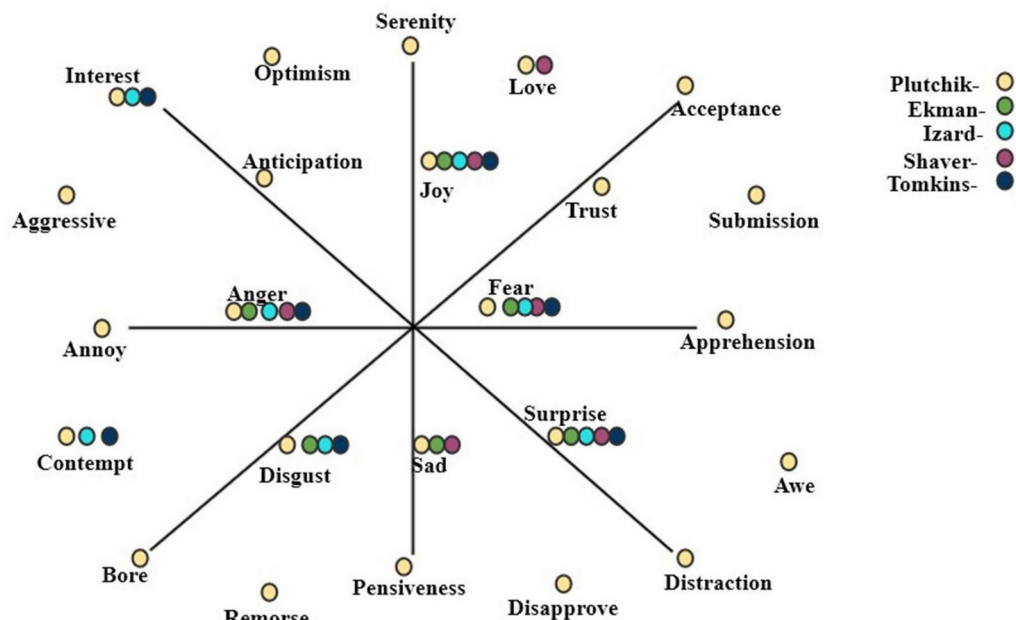


شکل ۱: مدل بعدی احساسات

جدول ۱: مدل های هیجانی که توسط روانشناسان مختلف تعریف شده است

Emotions	Type of model	No, of states	Psychological states	Representations	Discussion
Ekman model (Ekman 1992)	Categorical	6	Anger, disgust, fear, joy, sadness, surprise	-	Ekman's model consisted of six emotions, which act as a base for other emotion models like Plutchik model
Plutchik Wheel of Emotions (Plutchik 1982)	Dimensional	-	Joy, pensiveness, ecstasy, acceptance, sadness, fear, interest, rage, admiration, amazement, anger, vigilance boredom, annoyance, submission, serenity, apprehension, contempt, surprise, disapproval, distraction, grief, loathing, love, optimism, aggressiveness, remorse, anticipation, awe, terror, trust, disgust	wheel	Plutchik considered two types of emotions: basic (Ekman model + Trust +Anticipation) and mixed emotions (made from the combination of basic emotions). Plutchik represented emotions on a colored wheel
Izard model (Izard 1992)	-	10	Anger, contempt, disgust, anxiety, fear, guilt, interest, joy, shame, surprise	-	-
Shaver model (Shaver et al. 1987)	Categorical	6	Sadness, joy, anger, fear, love, surprise	tree	Shaver represented the primary, secondary and tertiary emotions in a hierarchal manner. The top-level of the tree presents these six emotions
Russell's circumplex model (Russell 1980)	Dimensional	-	Sad, satisfied, Afraid, alarmed, frustrated, angry, happy, gloomy, annoyed, tired, relaxed, glad, aroused, astonished, at ease, tense, miserable, content, bored, calm, delighted, excited, depressed, distressed, serene, droopy, pleased, sleepy	-	Emotions are presented over the circumplex model
Tomkins model (Tomkins and		9		-	

McCarter 1964)	Categorical		Disgust, surprise-Startle, anger-rage, anxiety, fear-terror, contempt, joy, shame, interest-Excitement	Tomkins identified nine different emotions out of which six emotions are negative. Most of the emotions are defined as a pair
Lövheim Model (Lövheim 2012)	Dimensional	-	Anger, contempt, distress, enjoyment, terror, excitement, humiliation, startle	cube Lövheim arranged the emotions according to the amount of three substances (Noradrenaline, dopamine and Serotonin) on a ۳-D cube



شکل ۲: تصویرسازی انواع مدل های عاطفی با برخی حالات روانی

۲.۴ مروری بر روش‌های پیشین

در پردازش احساسات به صورت کلی مراحل زیر طی می‌شوند.

برای پیاده‌سازی هرکدام از این مراحل روش‌های بسیار متنوعی وجود دارند که در ادامه شرح داده میشوند.

۲.۴.۱ شناختن داده‌ها

قبل از انجام هر مرحله پیش پردازش در پردازش زبان طبیعی (NLP) درک ویژگی‌ها و ویژگی‌های مجموعه داده از اهمیت بالایی برخوردار است. مجموعه داده، که به عنوان پایه‌ای برای تجزیه و تحلیل عمل می‌کند، بینش‌های ارزشمندی را ارائه می‌دهد که تصمیمات پیش پردازش را هدایت می‌کند و در نهایت بر کیفیت نتایج تأثیر می‌گذارد.

قبل از هر چیز، دانستن مجموعه داده‌ها به محققان اجازه می‌دهد تا درک جامعی از داده‌های متنی که با آنها کار می‌کنند به دست آورند. با بررسی مجموعه داده‌ها، می‌توان زبان مورد استفاده، دامنه یا موضوع متون، وجود اصطلاحات خاص یا اصطلاحات عامیانه و ساختار و قالب کلی داده‌ها را شناسایی کرد. این دانش به تطبیق مراحل پیش پردازش با زمینه خاص کمک می‌کند و اطمینان حاصل می‌کند که تحلیل‌ها و مدل‌های بعدی به خوبی با مجموعه داده سازگار هستند.

علاوه بر این، درک مجموعه داده‌ها به شناسایی چالش‌ها و مسائل بالقوه‌ای که ممکن است در طول پیش پردازش ایجاد شود کمک می‌کند. به عنوان مثال، وجود داده‌های پر سر و صدا، مانند اشتباهات تایپی، جملات ناقص، یا ناهماهنگی در قالب‌بندی، ممکن است به تکنیک‌های ویژه‌ای نیاز داشته باشد. به همین ترتیب، مجموعه داده ممکن است شامل کاراکترهای خاص، ایموجی‌ها یا نمادهای غیر

استاندارد باشد که باید به درستی مورد توجه قرار گیرند. با آگاهی از این چالش ها از قبل، محققان می توانند خطوط لوله پیش پردازش قوی ایجاد کنند که این مسائل را به طور موثر کاهش داده یا به آنها رسیدگی کند.

علاوه بر این، درک مجموعه داده ها امکان تصمیم گیری آگاهانه در مورد انتخاب و کاربرد تکنیک های پیش پردازش را فراهم می کند. تکنیک های مختلف، مانند نشانه سازی، حذف کلمه توقف، ریشه یابی، واژه سازی یا عادی سازی، بسته به ویژگی های آن می توانند تأثیرات متفاوتی بر داده ها داشته باشند. به عنوان مثال، اگر مجموعه داده شامل محتوای تولید شده توسط کاربر از پلتفرم های رسانه های اجتماعی باشد، ممکن است به مراحل پیش پردازش خاصی نیاز داشته باشد تا هشتگ ها، اشاره ها یا URL ها را به درستی مدیریت کند. با دانستن مجموعه داده، محققان می توانند تصمیمات آگاهانه ای در مورد اینکه کدام تکنیک ها را به کار گیرند و چگونه آنها را برای مجموعه داده خاص بهینه کنند، اتخاذ کنند که منجر به نتایج دقیق تر و معنی داری شود.

علاوه بر این، درک مجموعه داده، محققان را قادر می سازد تا سوگیری ها یا محدودیت های بالقوه ای را که ممکن است بر تحلیل تأثیر بگذارد، شناسایی کنند. سوگیری ها می توانند از فرآیند جمع آوری مجموعه داده ها ناشی شوند، مانند سوگیری های نمونه گیری یا سوگیری های جمعیتی در محتوای تولید شده توسط کاربر. با آگاهی از این سوگیری ها، محققان می توانند اقدامات مناسبی را برای رسیدگی یا کاهش آنها در طول پیش پردازش و تجزیه و تحلیل انجام دهند و از یک نتیجه عینی تر و نماینده تر اطمینان حاصل کنند.

در نتیجه، دانستن مجموعه داده قبل از پیش پردازش برای اطمینان از دقت، ارتباط و قابلیت اطمینان تجزیه و تحلیل بعدی حیاتی است. با به دست آوردن درک کامل از مجموعه داده، محققان می توانند

مراحل پیش پردازش را با زمینه خاص تنظیم کنند، به چالش ها و سوگیری ها بپردازند، و تصمیمات آگاهانه ای در مورد تکنیک های مورد استفاده اتخاذ کنند. در نهایت، این دانش کیفیت پیش پردازش را افزایش می دهد و پایه محکمی برای تحلیل و مدل سازی موفق NLP ایجاد می کند.

۲.۴.۲ پیش پردازش داده ها

در رسانه های اجتماعی، مردم معمولاً احساسات و عواطف خود را به روش هایی بدون زحمت بیان می کنند. در نتیجه، داده های به دست آمده از پست ها، ممیزی ها، نظرات، اظهارات و انتقادات این پلت فرم رسانه های اجتماعی بسیار ساختارمند نیستند و تجزیه و تحلیل احساسات و احساسات را برای ماشین ها دشوار می کنند. در نتیجه، پیش پردازش یک مرحله حیاتی در پاکسازی داده ها است، زیرا کیفیت داده ها به طور قابل توجهی روی بسیاری از رویکردهای پیش پردازش تأثیر می گذارد.

سازماندهی یک مجموعه داده نیاز به پیش پردازش دارد، از جمله توکن سازی، حذف کلمه توقف، برچسب گذاری POS و غیره (عبدی و همکاران، ۲۰۱۹؛ باسکار و همکاران، ۲۰۱۵). برخی از این تکنیک های پیش پردازش می توانند منجر به از دست رفتن اطلاعات حیاتی برای تجزیه و تحلیل احساسات و عواطف شوند که باید مورد توجه قرار گیرند.

توکن سازی فرآیند تجزیه کل سند یا پاراگراف یا فقط یک جمله به تکه هایی از کلمات به نام نشانه است (ناگارا جان و گاندی، ۲۰۱۹). به عنوان مثال، جمله "این مکان بسیار زیبا است" را در نظر بگیرید و پس از توکنیزاسیون، تبدیل به "این"، "مکان"، "بسیار زیبا" است. عادی سازی متن برای دستیابی به یکنواختی در داده ها با تبدیل متن به فرم استاندارد، تصحیح املای کلمات و غیره ضروری است (آهو جلا و همکاران، ۲۰۱۹).

کلمات غیر ضروری مانند مقاله ها و برخی حرف های اضافه که به شناخت احساسات و تحلیل احساسات کمک نمی کنند باید حذف شوند. به عنوان مثال، کلمات توقفی مانند "، "an"، "at"، "is"، "the" هیچ ربطی به احساسات ندارند، بنابراین برای جلوگیری از محاسبات غیر ضروری، باید حذف شوند (Bhaskar et al. ۲۰۱۹؛ عبدی و همکاران، ۲۰۱۹). برچسب گذاری POS راهی برای شناسایی بخش های مختلف گفتار در یک جمله است. این مرحله برای یافتن جنبه های مختلف یک

جمله مفید است که عموماً با اسم‌ها یا عبارات اسمی توصیف می‌شوند در حالی که احساسات و عواطف با صفت‌ها منتقل می‌شوند (سان و همکاران، ۲۰۱۷).

ریشه زایی و ریشه یابی دو مرحله مهم پیش پردازش هستند. در **stemming**، کلمات با کوتاه کردن پسوندها به شکل ریشه خود تبدیل می‌شوند. به عنوان مثال، اصطلاحات «استدلال» و «استدلال» به «استدلال» تبدیل می‌شود. این فرآیند محاسبه ناخواسته جملات را کاهش می‌دهد (کراتزوالد و همکاران، ۲۰۱۸؛ آکیلاندرسواری و جوتی، ۲۰۱۸). **Lemmatization** شامل تجزیه و تحلیل مورفولوژیکی برای حذف پایان‌های عطفی از یک نشانه و تبدیل آن به کلمه پایه لم است (قنبری-ادیوی و مصلح). (۲۰۱۹). به عنوان مثال، اصطلاح "گرفتار" به "گرفتن" تبدیل می‌شود (آهوجلا و همکاران، ۲۰۱۹). **Symeonidis** و همکاران. (۲۰۱۸) عملکرد چهار مدل یادگیری ماشین را با مطالعه ترکیبی و فرسایشی تکنیک‌های مختلف پیش پردازش روی دو مجموعه داده، یعنی **SS-Tweet** و **SemEval** مورد بررسی قرار داد. نویسندگان به این نتیجه رسیدند که حذف اعداد و واژه‌سازی دقت را افزایش می‌دهد، در حالی که حذف علائم نگارشی بر دقت تأثیری ندارد.

۲.۴.۳ استخراج ویژگی

دستگاه متن را بر حسب اعداد درک می‌کند. فرآیند تبدیل یا نگاشت متن یا کلمات به بردارهای با ارزش واقعی را واژه برداری یا جاسازی کلمه می‌نامند. این یک تکنیک استخراج ویژگی است که در آن یک سند به جملاتی تقسیم می‌شود که بیشتر به کلمات تقسیم می‌شوند. پس از آن، نقشه یا ماتریس ویژگی ساخته می‌شود. در ماتریس به دست آمده، هر ردیف نشان دهنده یک جمله یا سند است در حالی که هر ستون ویژگی یک کلمه را در فرهنگ لغت، و مقادیر موجود در سلول‌های نقشه ویژگی به طور کلی تعداد کلمه در جمله یا سند را نشان می‌دهد. برای انجام استخراج ویژگی، یکی از ساده‌ترین روش‌های مورد استفاده، "**Bag of Words** (BOW)" است که در آن یک بردار با طول ثابت از شمارش تعریف می‌شود که در آن هر ورودی با یک کلمه در فرهنگ لغت از پیش تعریف‌شده کلمات مطابقت دارد. اگر کلمه در یک جمله در فرهنگ لغت از پیش تعریف شده وجود نداشته باشد، تعداد ۰ اختصاص داده می‌شود، در غیر این صورت بسته به تعداد دفعات ظاهر شدن آن در جمله، تعداد آن بیشتر یا مساوی ۱ است. به همین دلیل است که طول بردار همیشه با کلمات موجود در فرهنگ لغت برابر است. مزیت این تکنیک اجرای آسان آن است اما دارای اشکالات قابل توجهی است زیرا منجر به

یک ماتریس پراکنده می شود، ترتیب کلمات را در جمله از دست می دهد. ۲۰۱۷؛ عبدی و همکاران (۲۰۱۹). به عنوان مثال، برای نشان دادن متن "آیا از خواندن لذت می برید" از فرهنگ لغت از پیش تعریف شده I, Hope, you, are, enjoying, reading خواهد بود (۰,۰,۱,۱,۱,۱). با این حال، این نمایش ها را می توان با پیش پردازش متن و با استفاده از TF-IDF، n-gram، بهبود بخشید.

روش N-gram یک گزینه عالی برای حل ترتیب کلمات در نمایش برداری جمله است. در یک نمایش برداری n-gram، متن به عنوان یک همکاری از n-gram منحصر به فرد به معنای گروه هایی از n واژه یا کلمه مجاور نشان داده می شود. مقدار n می تواند هر عدد طبیعی باشد. به عنوان مثال، جمله "آموزش برای همیشه لمس کردن است" را در نظر بگیرید و $n = 3$ به نام تریگرام "آموزش دادن است"، "تدریس کردن است"، "لمس کردن است"، "لمس کردن است" لمس کردن است" ایجاد می کند. یک زندگی، "یک زندگی برای همیشه". به این ترتیب می توان نظم جمله را حفظ کرد (آهوجا و همکاران ۲۰۱۹). ویژگی های N-gram بهتر از رویکرد BOW عمل می کنند، زیرا الگوهای نحوی، از جمله اطلاعات مهم (Chaffar و Inkpen) را پوشش می دهند. ۲۰۱۱). با این حال، اگرچه n-gram ترتیب کلمات را حفظ می کند، اما ابعاد و پراکندگی داده بالایی دارد (Le and Mikolov ۲۰۱۴).

اصطلاح فرکانس معکوس سند فرکانس، که معمولاً به اختصار TFIDF نامیده می شود، روش دیگری است که معمولاً برای استخراج ویژگی استفاده می شود. این روش متن را به صورت ماتریسی نشان می دهد، که در آن هر عدد مقدار اطلاعاتی را که این اصطلاحات در یک سند معین حمل می کنند را کمیت می دهد. بر این فرض ساخته شده است که اصطلاحات نادر دارای اطلاعات زیادی در سند متنی هستند (Liu et al. ۲۰۱۹). فراوانی عبارت تعداد دفعاتی است که یک کلمه w در یک سند تقسیم می شود بر تعداد کل کلمات W در سند و IDF عبارت است از \log (تعداد کل اسناد (N) تقسیم بر تعداد کل اسنادی که کلمه w در آنها ظاهر می شود. (ن)) (سونگبو و جین ۲۰۰۸). آهوجا و همکاران (۲۰۱۹) شش تکنیک پیش پردازش را اجرا کرد و دو تکنیک استخراج ویژگی را برای شناسایی بهترین رویکرد مقایسه کرد. آنها از شش الگوریتم یادگیری ماشین استفاده کردند و از n-gram با $n = 2$ و TF-

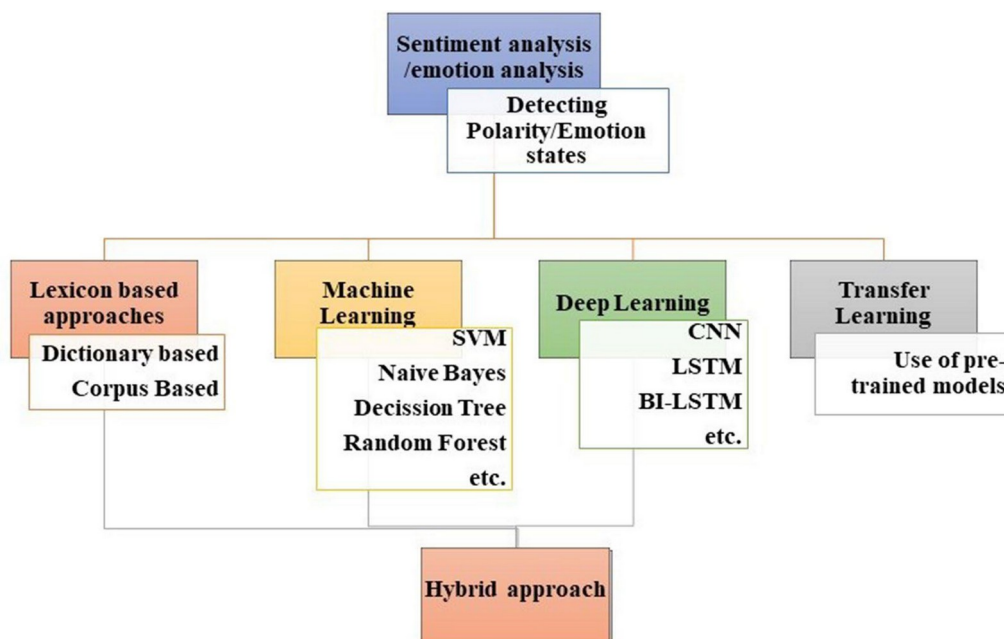
IDF برای استخراج ویژگی روی مجموعه داده SS-tweet استفاده کردند و به این نتیجه رسیدند که TF-IDF عملکرد بهتری نسبت به n-gram دارد.

در دسترس بودن حجم وسیعی از داده ها به یک شبکه یادگیری عمیق اجازه می دهد تا بازنمایی های برداری خوبی را کشف کند. استخراج ویژگی با جاسازی کلمه بر اساس شبکه های عصبی آموزنده تر است. در جاسازی کلمه مبتنی بر شبکه عصبی، کلمات با معنایی یکسان یا مرتبط با یکدیگر با بردارهای مشابه نمایش داده می شوند. این در پیش بینی کلمات محبوب تر است زیرا معنای کلمات را حفظ می کند. تیم تحقیقاتی گوگل به سرپرستی توماس میکولوف، مدلی به نام Word2Vec را برای جاسازی کلمه توسعه دادند. با Word2Vec، می توان برای ماشینی فهمید که نمایش برداری «ملکه» + «مونث» + «مرد» مانند بازنمایی بردار «شاه» است (سوما و همکاران، ۲۰۱۹).

نمونه های دیگر از مدل های جاسازی کلمه مبتنی بر یادگیری عمیق عبارتند از GloVe که توسط محققان دانشگاه استنفورد توسعه یافته و FastText که توسط Face-book معرفی شده اند. آموزش بردارهای GloVe سریعتر از Word2vec است. بردارهای FastText در مقایسه با بردارهای Word2Vec با چندین معیار متفاوت دقت بهتری دارند. یانگ و همکاران (۲۰۱۸) ثابت کرد که انتخاب جاسازی کلمه مناسب بر اساس شبکه های عصبی می تواند منجر به پیشرفت های قابل توجهی حتی در مورد کلمات خارج از واژگان (OOV) شود. نویسندگان انواع جاسازی کلمات را که با استفاده از توئیتر و ویکی پدیا به عنوان مجموعه آموزش دیده بودند، با جاسازی کلمه TF-IDF مقایسه کردند.

۲.۴.۴ تکنیک های تجزیه و تحلیل احساسات و تشخیص احساسات

شکل ۴ تکنیک های مختلفی را برای تجزیه و تحلیل احساسات و تشخیص احساسات ارائه می دهد که به طور گسترده به رویکرد مبتنی بر فرهنگ لغت، رویکرد مبتنی بر یادگیری ماشین، رویکرد مبتنی بر یادگیری عمیق طبقه بندی می شود. رویکرد ترکیبی ترکیبی از رویکردهای آماری و یادگیری ماشینی برای غلبه بر اشکالات هر دو رویکرد است. یادگیری انتقالی نیز زیرمجموعه ای از یادگیری ماشین است که امکان استفاده از مدل از پیش آموزش دیده را در سایر حوزه های مشابه فراهم می کند.



شکل ۳: تکنیک‌هایی برای تجزیه و تحلیل احساسات و تشخیص احساسات

۲.۴.۵ تکنیک های تحلیل احساسات

رویکرد مبتنی بر فرهنگ لغت‌تین روش یک فرهنگ لغت دارد که در آن به هر کلمه مثبت و منفی یک ارزش احساسی اختصاص داده می‌شود. سپس از مجموع یا میانگین ارزش‌های احساسی برای محاسبه احساس کل جمله یا سند استفاده می‌شود. با این حال، جورک و همکاران. (۲۰۱۵) رویکرد متفاوتی به نام تابع نرمال سازی را برای محاسبه دقیق ارزش احساسات نسبت به این جمع پایه و تابع میانگین امتحان کرد. رویکرد مبتنی بر فرهنگ لغت و رویکرد مبتنی بر پیکره دو نوع رویکرد مبتنی بر واژگان مبتنی بر واژگان احساسی هستند. به طور کلی، یک فرهنگ لغت کلمات برخی از زبان‌ها را به صورت سیستمی حفظ می‌کند، در حالی که یک مجموعه نمونه تصادفی متن در برخی از زبان‌ها است. معنای دقیق در اینجا در رویکرد مبتنی بر فرهنگ لغت و رویکرد مبتنی بر پیکره اعمال می‌شود. در رویکرد مبتنی بر فرهنگ لغت، یک فرهنگ لغت از کلمات بذر نگهداری می‌شود (Schouten و ۲۰۱۵Frasincar). برای ایجاد این فرهنگ لغت، اولین مجموعه کوچک از کلمات احساسی، احتمالاً با زمینه‌های بسیار کوتاه مانند نفی، همراه با برجسب‌های قطبیت آن جمع‌آوری شده است (Bern- ۲۰۲۰ abé-Moreno et al.). سپس فرهنگ لغت با جستجوی مترادف آنها (کلمات با قطبیت یکسان) و متضاد (کلمات با قطبیت مخالف) به روز می‌شود. دقت تحلیل احساسات از طریق این

رویکرد به الگوریتم بستگی دارد. با این حال، این تکنیک دارای ویژگی دامنه نیست. رویکرد مبتنی بر پیکره محدودیت‌های رویکرد مبتنی بر فرهنگ لغت را با گنجاندن کلمات احساسی خاص دامنه که در آن برچسب قطبیت به کلمه احساس با توجه به زمینه یا دامنه آن اختصاص داده می‌شود، حل می‌کند. این یک رویکرد داده محور است که در آن می‌توان به کلمات احساسی همراه با زمینه دسترسی داشت. این رویکرد مطمئن‌تر می‌تواند یک رویکرد مبتنی بر قانون با برخی از تکنیک‌های تجزیه NLP باشد. بنابراین رویکرد مبتنی بر پیکره تمایل به تعمیم ضعیف دارد، اما می‌تواند عملکرد عالی را در یک حوزه خاص به دست آورد. از آنجایی که رویکرد مبتنی بر فرهنگ لغت، زمینه پیرامون کلمه احساسی را در نظر نمی‌گیرد، منجر به کارایی کمتری می‌شود. بنابراین، چو و همکاران. (۲۰۱۴) به صراحت قطبیت زمینه ای را به کار گرفت تا فرهنگ لغت را در حوزه‌های متعدد با رویکرد داده محور سازگار کند. آنها گرفتند

استراتژی سه مرحله‌ای: لغت‌نامه‌های مختلف را ادغام کنید، کلماتی را که به طبقه‌بندی کمک نمی‌کنند حذف کنید، و قطبیت را بر اساس یک دامنه خاص تغییر دهید.

SentiWordNet (اصولی و سباستینی-۲۰۰۶) و فرهنگ لغت آگاه و منطق کننده احساسات (VADER) (هوتو و گیلبرت) (۲۰۱۴) واژگان محبوب در احساسات هستند. جها و همکاران (۲۰۱۸) سعی کرد با ایجاد فرهنگ لغت احساسات به نام فرهنگ لغت آگاه چند دامنه ای هندی (HMDSAD) برای تجزیه و تحلیل احساسات در سطح سند، برنامه واژگان را در حوزه‌های چندگانه گسترش دهد. از این فرهنگ لغت می‌توان برای حاشیه نویسی نظرات به مثبت و منفی استفاده کرد. روش پیشنهادی ۲۴ درصد کلمات بیشتری نسبت به واژگان عمومی سستی هندی (HSWN) Sentiwordnet، یک واژگان خاص دامنه، برچسب‌گذاری کرد. روابط معنایی بین کلمات در واژگان سستی مورد بررسی قرار نگرفته است، که عملکرد طبقه بندی احساسات را بهبود می‌بخشد. بر اساس این فرض، ویگاس و همکاران. (۲۰۲۰) واژگان را با گنجاندن اصطلاحات اضافی پس از استفاده از جاسازی کلمات برای کشف مقادیر احساسات برای این کلمات به طور خودکار به روز کرد. این ارزش‌های احساسی از تعبیه‌های واژه‌ای «نزدیک» از واژه‌های موجود در واژگان مشتق شده‌اند.

رویکرد مبتنی بر یادگیری ماشین رویکرد دیگری برای تحلیل احساسات وجود دارد که به آن رویکرد یادگیری ماشینی می‌گویند. کل مجموعه داده برای اهداف آموزشی و آزمایشی به دو بخش تقسیم می‌شود: مجموعه داده آموزشی و مجموعه داده آزمایشی. مجموعه داده آموزشی اطلاعاتی است که برای آموزش مدل با ارائه ویژگی‌های نمونه‌های مختلف یک آیت استفاده می‌شود. سپس از مجموعه داده آزمایشی استفاده می‌شود تا ببینیم مدل مجموعه داده آموزشی با چه موفقیتی آموزش داده شده است. به طور کلی، الگوریتم‌های یادگیری ماشین مورد استفاده برای تجزیه و تحلیل احساسات تحت طبقه‌بندی نظارت شده قرار می‌گیرد. انواع مختلفی از الگوریتم‌های مورد نیاز برای طبقه‌بندی احساسات ممکن است عبارتند از Naïve Bayes، ماشین بردار پشتیبان (SVM)، درخت‌های تصمیم‌گیری، و غیره که هر کدام مزایا و معایب خود را دارند. گامون (۲۰۰۴) یک ماشین بردار پشتیبانی را روی ۴۰۸۸۴ بازخورد مشتری جمع‌آوری شده از نظرسنجی‌ها اعمال کرد. نویسندگان ترکیب‌های مختلفی از مجموعه ویژگی‌ها را پیاده‌سازی کردند و به دقت ۸۵.۴۷٪ رسیدند. یه و همکاران (۲۰۰۹) با SVM، مدل N-gram و Naïve Bayes روی احساسات و بررسی هفت مقصد محبوب اروپا و ایالات متحده آمریکا کار کرد که از yahoo.com جمع‌آوری شد. نویسندگان با مدل n-gram به دقت ۸۷.۱۷ درصد دست یافتند. تورفتگی بوچار و همکاران. (۲۰۱۸) واژگانی به نام ۱۰۰ JOB ایجاد کرد و مجموعه‌های خبری به نام ۱۰۰ SentiNews را برای تجزیه و تحلیل احساسات در متون اسلوونیایی نامگذاری کرد. ۱۰۰ JOB شامل ۲۵۵۲۴ کلمه اصلی است که با مقیاس بندی احساسات از

- ۵ تا ۵ بر اساس مدل AFINN. برای ساخت مجموعه‌ها، داده‌ها از رسانه‌های خبری مختلف وب حذف شدند. سپس، پس از پاکسازی و پیش پردازش داده‌ها، از حاشیه نویسان خواسته شد تا ۱۰۴۲۷ سند را در مقیاس ۱ تا ۵ حاشیه نویسی کنند که یک به معنای منفی و ۵ به معنای بسیار مثبت است. سپس این اسناد با برچسب‌های مثبت، منفی و خنثی مطابق با میانگین امتیاز مقیاس خاص برچسب‌گذاری شدند. نویسندگان مشاهده کردند که Naïve Bayes در مقایسه با ماشین بردار پشتیبانی (SVM) بهتر عمل کرد. Naive Bayes به امتیاز $F1$ بالای ۹۰٪ در طبقه‌بندی باینری و امتیاز $F1$ بالای ۶۰٪ برای طبقه‌بندی احساسات سه کلاسه دست یافت. تیواری و همکاران (۲۰۲۰) سه الگوریتم یادگیری ماشین به نام‌های SVM، Naive Bayes و حداکثر آنتروپی را با روش استخراج ویژگی n گرم روی مجموعه داده‌های گوجه‌فرنگی پوسیده پیاده‌سازی کرد. مجموعه داده آموزش و

آزمایش شامل ۱۶۰۰ مرور در هر یک بود. نویسندگان کاهش دقت را با مقادیر بالاتر n در n گرم مشاهده کردند مانند $n =$ چهار، پنج و شش. سومیا و پرامود (۲۰۲۰) ۳۱۸۴ توییت مالایالام را با استفاده از بردارهای ویژگی های مختلف مانند BOW، Unigram و Sentiwordnet و غیره به نظرات مثبت و منفی طبقه بندی کردند. نویسندگان الگوریتم های یادگیری ماشینی مانند جنگل تصادفی و Naïve Bayes را پیاده سازی کردند و مشاهده کردند که جنگل تصادفی با دقت ۹۵.۶ درصد عملکرد بهتری دارد. Unigram Sentiwordnet با در نظر گرفتن کلمات نفی.

رویکرد مبتنی بر یادگیری عمیق در سال های اخیر، الگوریتم های یادگیری عمیق بر سایر رویکردهای سنتی برای تحلیل احساسات تسلط دارند. این الگوریتم ها احساسات یا نظرات متن را بدون انجام مهندسی ویژگی تشخیص می دهند. الگوریتم های یادگیری عمیق متعددی وجود دارد، یعنی شبکه های عصبی مکرر و شبکه های عصبی کانولوشنال، که می توانند برای تحلیل احساسات اعمال شوند و نتایج دقیق تری نسبت به مدل های یادگیری ماشینی ارائه دهند. این رویکرد باعث می شود انسان ها از ساختن ویژگی های متن به صورت دستی رها شوند، زیرا مدل های یادگیری عمیق آن ویژگی ها یا الگوها را خودشان استخراج می کنند. ژیان و همکاران (۲۰۱۰) از یک مدل مبتنی بر عصبی استفاده کرد

فناوری شبکه برای دسته بندی احساسات که شامل ویژگی های احساسی، بردارهای وزن و ویژگی و پایگاه دانش قبلی است. نویسندگان این مدل را برای بررسی داده های فیلم کرنل به کار بردند. نتایج تجربی این مقاله نشان داد که سطح دقت مدل I در مقایسه با HMM و SVM فوق العاده است. پاسوپا و آیوتایا (۲۰۱۹) اعتبارسنجی متقاطع پنج برابری را روی مجموعه داده های داستان کودکان (تایلندی) اجرا کرد و سه مدل یادگیری عمیق به نام های LSTM، CNN و Bi-LSTM را مقایسه کرد. این مدل ها با یا بدون ویژگی ها اعمال می شوند: برچسب گذاری POS (تکنیک پیش پردازش برای شناسایی بخش های مختلف گفتار). Thai2Vec (جاسازی کلمه آموزش دیده از ویکی پدیا تایلندی)؛ sentic (برای درک احساس کلمه). نویسندگان بهترین عملکرد را در مدل CNN با هر سه ویژگی که قبلاً اشاره شد مشاهده کردند. همانطور که قبلاً گفته شد، پلتفرم های رسانه های اجتماعی به عنوان منبع مهمی از داده ها در زمینه تحلیل احساسات عمل می کنند. داده های جمع آوری شده از این سایت های اجتماعی به دلیل شیوه نوشتن رایگان کاربران، سر و صدای زیادی دارد. بنابراین، آرورا و کانسال (۲۰۱۹) مدلی به نام Conv-char-Emb را پیشنهاد کرد که می تواند مشکل داده های نویزدار را حل کند و از فضای

کوچک حافظه برای جاسازی استفاده کند. برای جاسازی، از شبکه عصبی کانولوشن (CNN) استفاده شده است که از پارامترهای کمتری در نمایش ویژگی استفاده می کند. دشتی پور و همکاران (۲۰۲۰) چارچوب یادگیری عمیق را برای انجام تحلیل احساسات در زبان فارسی پیشنهاد کرد. محققان به این نتیجه رسیدند که شبکه های عصبی عمیق مانند LSTM و CNN از الگوریتم های یادگیری ماشین موجود در مجموعه داده های بررسی هتل و محصول بهتر عمل می کنند.

رویکرد یادگیری انتقالی و رویکرد ترکیبی یادگیری انتقالی نیز بخشی از یادگیری ماشینی است. یک مدل آموزش داده شده بر روی مجموعه داده های بزرگ برای حل یک مشکل می تواند برای سایر مسائل مرتبط اعمال شود. استفاده مجدد از یک مدل از پیش آموزش دیده در حوزه های مرتبط به عنوان نقطه شروع می تواند باعث صرفه جویی در زمان و ایجاد نتایج کارآمدتر شود. ژانگ و همکاران (۲۰۱۲) با مدل سازی مستقیم توزیع بین حوزه های مختلف، یک روش یادگیری نمونه جدید را پیشنهاد کرد. نویسندگان مجموعه داده را طبقه بندی کردند: بررسی محصول آمازون و مجموعه داده های توییتر به احساسات مثبت و منفی. تائو و نیش (۲۰۲۰) گسترش روش های طبقه بندی اخیر در تحلیل احساسات مبتنی بر جنبه را به طبقه بندی چند برچسبی پیشنهاد کرد. نویسندگان همچنین مدل های یادگیری انتقالی به نام XLNet و Bert را توسعه دادند و رویکرد پیشنهادی را در مجموعه داده های مختلف Yelp ارزیابی کردند. رویکردهای یادگیری عمیق و یادگیری ماشین نتایج خوبی دارند، اما رویکرد ترکیبی می تواند نتایج بهتری به همراه داشته باشد زیرا بر محدودیت های هر مدل سنتی غلبه می کند. ملادنووویچ و همکاران (۲۰۱۶) یک تکنیک کاهش ویژگی، یک چارچوب ترکیبی ساخته شده از واژگان احساسات و wordnet صربی را پیشنهاد کرد. نویسندگان هر دو فرهنگ لغت را با افزودن برخی واژه های احساسی صرفی گسترش دادند تا از دست دادن اطلاعات مهم در حین ریشه گیری جلوگیری کنند. آل امرانی و همکاران (۲۰۱۸) مدل هیبریدی ساخته شده خود را مقایسه کردند SVM و مدل جنگل تصادفی، یعنی RFSVM، در بررسی محصول ama-zon. نویسندگان به این نتیجه رسیدند که RFSVM، با سطح دقت ۸۳.۴ درصد، عملکرد بهتری نسبت به SVM با دقت ۸۲.۴ درصد و جنگل تصادفی با دقت ۸۱ درصد به صورت جداگانه در مجموعه داده ۱۰۰۰ بررسی دارد. القریوتی و همکاران (۲۰۲۰) ترکیبی از رویکرد مبتنی بر قانون و واژگان دامنه را برای تشخیص احساسات در سطح جنبه برای درک نظرات مردم در مورد برنامه های هوشمند دولتی پیشنهاد کرد. نویسندگان به این نتیجه رسیدند

که تکنیک پیشنهادی ۵ درصد از سایر مدل‌های پایه مبتنی بر واژگان برتری دارد. ری و چاکرابارتی (۲۰۲۰) رویکرد مبتنی بر قانون را برای استخراج جنبه‌ها با یک مدل یادگیری عمیق ۷ لایه CNN برای برچسب گذاری هر جنبه ترکیب کرد. مدل ترکیبی دقت ۸۷ درصدی را به دست آورد، در حالی که مدل‌های منفرد ۷۵ درصد دقت مبتنی بر قانون و ۸۰ درصد دقت با مدل CNN داشتند.

۲.۴.۶ تکنیک های تشخیص احساسات

رویکرد مبتنی بر واژگان رویکرد مبتنی بر واژگان یک رویکرد جستجوی مبتنی بر کلمه کلیدی است که به جستجوی کلمات کلیدی احساسات اختصاص داده شده به برخی از حالات روانی می‌پردازد (رایبا و همکاران ۲۰۱۷). واژگان محبوب برای تشخیص احساسات عبارتند از Word-Net-Affect (Strapparava et al. ۲۰۰۴) و NRC واژگان کلمه-احساس (محمد و تورنی ۲۰۱۳). WordNet -Affect یک فرم توسعه یافته از WordNet است که از کلمات عاطفی تشکیل شده است که با برچسب های احساسات حاشیه نویسی شده اند. واژگان NRC شامل ۱۴۱۸۲ کلمه است که هر کدام به یک احساس خاص و دو احساس اختصاص دارد. این واژگان، واژگان دسته بندی هستند که هر کلمه را با یک حالت احساسی برای طبقه بندی احساسات برچسب گذاری می کنند. با این حال، با نادیده گرفتن شدت احساسات، این واژگان سنتی کمتر آموزنده و کمتر سازگار می شوند. بنابراین، لی و همکاران. (۲۰۲۱) یک استراتژی موثر برای به دست آوردن توزیع هیجان در سطح کلمه پیشنهاد کرد تا با ادغام یک فرهنگ لغت بعدی به نام NRC-Valence arosal dominance، احساسات با شدت به کلمات احساسی اختصاص یابد. EmoSenticNet (پوریا و همکاران ۲۰۱۴) همچنین شامل تعداد زیادی است که به برچسب های کیفی و کمی اختصاص داده شده است. به طور کلی، محققان واژگان خود را تولید می کنند و مستقیماً آنها را برای تجزیه و تحلیل احساسات به کار می برند، اما واژگان می توانند برای اهداف استخراج ویژگی نیز استفاده شوند. عبادوی و همکاران (۲۰۱۷) از استفاده از ابزارهای ترجمه آنلاین برای ایجاد واژگان فرانسوی به نام FEEL (فرانسوی واژگان عاطفی گسترش یافته) استفاده کرد که شامل بیش از ۱۴۰۰۰ کلمه با دو قطبیت و برچسب احساسات است. این واژگان توسط افزایش تعداد کلمات در واژگان احساسات NRC و ترجمه نیمه خودکار با استفاده از شش مترجم آنلاین. آن مدخل‌هایی که از حداقل سه مترجم

به دست می آمدند، از پیش تأیید شده در نظر گرفته شدند و سپس توسط مترجم دستی تأیید شدند. بندکاوی و همکاران (۲۰۱۷) یک واژگان خاص دامنه را برای فرآیند استخراج ویژگی در تحلیل احساسات به کار برد. نویسندگان به این نتیجه رسیدند که ویژگی های مشتق شده از واژگان پیشنهادی آنها از سایر ویژگی های پایه بهتر عمل می کند. براون و همکاران (۲۰۲۱) یک مجموعه چند زبانه به نام MEMoFC که مخفف عبارت Multilingual Emotional Football Corpus است، ساخته است که شامل گزارش های فوتبال از وبسایت های انگلیسی، هلندی و آلمانی و آمار مسابقات خزیده شده از Goal.com است. مجموعه با ایجاد دو جدول فراداده ایجاد شد: یکی توضیح جزئیات یک مسابقه مانند تاریخ، مکان، تیم های شرکت و غیره، و جدول دوم شامل اختصارات باشگاه های فوتبال است. نویسندگان مجموعه را با رویکردهای مختلف برای دانستن تأثیر گزارش ها بر نتایج بازی نشان دادند.

تکنیک های مبتنی بر یادگیری ماشین تشخیص یا طبقه بندی احساسات ممکن است به انواع مختلفی از مدل های یادگیری ماشین نیاز داشته باشد، مانند Naïve Bayes، ماشین بردار پشتیبان، درخت های تصمیم، و غیره. (۲۰۱۷) احساسات را از متون چند زبانه جمع آوری شده از سه حوزه مختلف استخراج کرد. نویسندگان از یک رویکرد جدید به نام خلاصه سایت غنی برای جمع آوری داده ها استفاده کردند و از الگوریتم های یادگیری ماشینی SVM و Naïve Bayes برای طبقه بندی احساسات متن توییت استفاده کردند. نتایج نشان داد که سطح دقت ۷۱.۴٪ با الگوریتم Naïve Bayes به دست آمد. حسن و همکاران (۲۰۱۹) الگوریتم های یادگیری ماشینی مانند Naïve Bayes، SVM و درخت های تصمیم را برای شناسایی احساسات در پیام های متنی ارزیابی کرد. این کار به دو زیرکار تقسیم می شود: وظیفه ۱ شامل مجموعه ای از مجموعه داده از توییت و برچسب زدن خودکار مجموعه داده با استفاده از هشتگ ها و آموزش مدل است. ۲ Task در حال توسعه EmotexStream دو مرحله ای است که توییت های بدون احساس را در مرحله اول جدا می کند و با استفاده از مدل های آموزش دیده در task ۱، احساسات را در متن شناسایی می کند. نویسندگان دقت ۹۰ درصدی را در طبقه بندی احساسات مشاهده کردند. اصغر و همکاران (۲۰۱۹) با هدف اعمال چندین مدل یادگیری ماشین در مجموعه داده ISEAR برای یافتن بهترین طبقه بندی کننده. آنها دریافتند که مدل رگرسیون لجستیک بهتر از سایر طبقه بندی کننده ها با ارزش یادآوری ۸۳ درصد عمل می کند.

یادگیری عمیق و تکنیک ترکیبیحوزه یادگیری عمیق بخشی از یادگیری ماشینی است که اطلاعات یا سیگنال ها را به همان روشی که مغز انسان انجام می دهد پردازش می کند. مدل های یادگیری عمیق حاوی چندین لایه نورون هستند. هزاران نورون به یکدیگر متصل هستند، که سرعت پردازش را به صورت موازی افزایش می دهد. چت ترجی و همکاران (۲۰۱۹) مدلی به نام تشخیص احساسات و احساسات معنایی (SSBED) با تغذیه احساسات و بازنمایی های معنایی به ترتیب به دو لایه LSTM ایجاد کرد. این نمایش ها سپس به هم متصل می شوند و سپس برای طبقه بندی به شبکه مش منتقل می شوند. رمان رویکرد مبتنی بر احتمال وجود احساسات متعدد در جمله است و از بازنمایی معنایی و احساسی برای طبقه بندی بهتر احساسات استفاده می کند. نتایج بر روی مجموعه داده های ساخته شده خودشان با جفت های مکالمه توییت ارزیابی می شوند و مدل آنها با سایر مدل های پایه مقایسه می شود. خو و همکاران (۲۰۲۰) با استفاده از مدل های دو ترکیبی به ترتیب با نام های حافظه کوتاه مدت متحرک- طولانی سه بعدی (۳DCLS) و CNN-RNN احساسات را از ویدیو و متن استخراج کرد. در همان زمان، نویسندگان SVM را برای طبقه بندی احساسات مبتنی بر صوتی پیاده سازی کردند. نویسندگان نتایج را با ادغام ویژگی های صوتی و تصویری در سطح ویژگی با تکنیک MKL fusion و ترکیب بیشتر نتایج آن با نتایج طبقه بندی احساسات مبتنی بر متن به نتیجه رسیدند. این روش دقت بهتری نسبت به سایر تکنیک های همجوشی چندوجهی ارائه می کند و قصد دارد احساسات بررسی های دارویی نوشته شده توسط بیماران در پلتفرم های رسانه های اجتماعی را تحلیل کند. بصیری و همکاران (۲۰۲۰) دو مدل را با استفاده از تئوری تصمیم گیری سه طرفه پیشنهاد کرد. مدل اول ادغام سه طرفه یک مدل یادگیری عمیق با روش یادگیری سستی (۳W۱DT) است، در حالی که مدل دیگر تلفیقی سه طرفه از سه مدل یادگیری عمیق با روش یادگیری معمولی (۳W۳DT) است. نتایج به دست آمده با استفاده از مجموعه داده Drugs.com نشان داد که هر دو چارچوب بهتر از تکنیک های یادگیری عمیق سستی عمل می کنند. علاوه بر این، عملکرد مدل فیوژن اول در مقایسه با مدل دوم از نظر دقت و متریک F۱ بسیار بهتر بود. در روزهای اخیر، پلتفرم های رسانه های اجتماعی مملو از پست های مرتبط با کووید-۱۹ هستند. سینگ و همکاران (۲۰۲۱) تجزیه و تحلیل تشخیص احساسات را روی توییت های کووید-۱۹ جمع آوری شده از کل جهان و هند تنها با مدل های رمزگذار دوطرفه از ترانسفورماتورها (BERT) روی مجموعه داده های توییت اعمال کرد و تقریباً ۹۴ درصد دقت را به دست آورد.

رویکرد یادگیری انتقالی در رویکردهای سنتی، فرض رایج این است که مجموعه داده از یک حوزه است. با این حال، زمانی که دامنه تغییر می کند، نیاز به یک مدل جدید وجود دارد. رویکرد یادگیری انتقال به شما امکان می دهد از مدل های از پیش آموزش دیده موجود در حوزه هدف مجدداً استفاده کنید. برای مثال احمد و همکاران. (۲۰۲۰) از تکنیک یادگیری انتقالی به دلیل کمبود منابع برای تشخیص احساسات در زبان هندی استفاده کرد. محققان از قبل آموزش دیدند

مدل بر روی دو مجموعه داده انگلیسی مختلف: SemEval-۲۰۱۸، تجزیه و تحلیل احساسات، و یک مجموعه داده هندی با برچسب های مثبت، خنثی، درگیری و منفی. آنها امتیاز f_1 ۰.۵۳ را با استفاده از یادگیری انتقال و ۰.۴۷ با استفاده از مدل های پایه CNN و Bi-LSTM با جاسازی کلمه متقابل زبانی به دست آوردند. هزاریکا و همکاران (۲۰۲۰) یک مدل TL-ERC ایجاد کرد که در آن مدل از قبل بر روی مکالمات چند نوبتی منبع آموزش داده شد و سپس بر روی وظیفه طبقه بندی احساسات در پیام های رد و بدل شده منتقل شد. نویسندگان بر موضوعاتی مانند کمبود داده های برچسب گذاری شده در مکالمات چندگانه با چارچوب مبتنی بر یادگیری انتقال استقرایی تأکید کردند.

اکثر محققان مدل ها را با ترکیب تکنیک های یادگیری ماشین و یادگیری عمیق با تکنیک های مختلف استخراج ویژگی پیاده سازی کردند. بیشتر مجموعه داده ها به زبان انگلیسی در دسترس هستند. با این حال، برخی از محققان مجموعه داده های زبان منطقه ای خود را ساختند. به عنوان مثال، ساسیذر و همکاران. (۲۰۲۰) مجموعه داده کد هندی-انگلیسی را با سه احساس اصلی ایجاد کرد: خوشحال، غمگین و عصبانی، و مشاهده شد که CNN-BILSTM عملکرد بهتری نسبت به دیگران دارد.

۳ روش پیشنهادی

در این روش ما تمام مراحل را که برای پردازش داده‌های دیتاست استفاده شده است پیاده‌سازی میکنیم. مراحل که ما استفاده کردیم عبارتند از:

برای پیاده‌سازی یک مدل روی یک مجموعه داده ما قدم‌های زیر را طی خواهیم کرد:

۱. شناختن داده‌ها : در این بخش ما تعدادی بررسی آماری روی مجموعه داده انجام خواهیم داد که داده‌هایی که قرار است با آنها کار کنیم را بهتر شناخته و به یک دید کلی نسبت به داده‌ها برسیم.
۲. پیش‌پردازش داده‌ها : در بخش پیش‌پردازش ما از تکنیک‌هایی استفاده میکنیم که بتوانیم به ویژگی‌های مناسبی که مد نظر ما هستند و تاثیر گذاری مستقیم در خروجی دارند برسیم و داده‌های اضافی که تاثیری در خروجی ندارند و یا ما نمیخواهیم تاثیر آنها را مورد بحث و بررسی قرار دهیم حذف کنیم. همچنین داده‌ها را به فرم مناسبی در آورده که برای مدل‌هایی که میخواهیم پیاده‌سازی کنیم قابل پردازش باشند.

۳. تقسیم کردن داده‌ها : از آنجایی که بعد از ساختن مدل ما نیاز به ارزیابی مدل و بحث و بررسی در مورد دقت مدل خواهیم داشت میبایست داده‌ها را به دو بخش داده یادگیری و تست تقسیم کنیم که

بتوان بعد از خوراندن داده های یادگیری به مدل و ساخته شدن مدل با استفاده از داده های تست دقت مدل را مورد ارزیابی قرار داد

۴. ساختن مدل : در نهایت میتوان مدل خود را با الگوریتم های متفاوت ساخت و مدل های ساخته شده را مورد بحث و بررسی قرار داد و دقت هرکدام را سنجید. انتخاب مدل به عوامل متعددی وابسته است که مهم ترین آن ها پراکندگی آماری داده ها و تعداد کلاس های موجود هستند.

۳.۱ شناختن داده ها

شناختن داده ها برای کار کردن با آن ها از اهمیت بالایی است در ادامه به چندی از آن ها اشاره میکنیم:

۱. **تضمین کیفیت داده:** هنگام کار با وظایف NLP، باید از کیفیت داده هایی که استفاده می کنید مطمئن باشید. این شامل درک منبع داده ها، فرآیند جمع آوری و سوگیری های بالقوه است. به عنوان مثال، اگر در حال ساخت یک مدل تجزیه و تحلیل احساسات هستید، دانستن اینکه آیا منبع داده تویتر است یا مقاله های خبری رسمی می تواند تا حد زیادی بر زبان و عبارات احساسات در مجموعه داده تأثیر بگذارد.
۲. **تشخیص سوگیری و کاهش:** درک مجموعه داده ها به تشخیص سوگیری ها کمک می کند. مدل های NLP می توانند ناخواسته تعصب های موجود در داده های آموزشی را یاد بگیرند و منتشر کنند. دانستن منبع و زمینه مجموعه داده شما را قادر می سازد تا چنین سوگیری ها را شناسایی و کاهش دهید و از عدالت و برابری در برنامه های NLP خود اطمینان حاصل کنید.
۳. **حجم و تنوع داده ها:** وظایف مختلف NLP به مقادیر متفاوتی از داده نیاز دارند. دانستن اندازه و تنوع مجموعه داده ها برای انتخاب الگوریتم ها و مدل های مناسب بسیار مهم است. برای مثال، مدل های یادگیری عمیق اغلب به مجموعه داده های بزرگ برای آموزش مؤثر نیاز دارند، در حالی که مدل های ساده تر ممکن است برای مجموعه داده های کوچک تر کافی باشند.
۴. **برچسب گذاری و حاشیه نویسی داده ها:** در بسیاری از وظایف NLP، مانند طبقه بندی متن یا شناسایی موجودیت نام گذاری شده، برچسب گذاری دستی یا حاشیه نویسی داده ها ضروری

است. درک مجموعه داده‌ها می‌تواند به تعریف دستورالعمل‌های برچسب‌گذاری واضح و اطمینان از ثبات در حاشیه‌نویسی کمک کند، که برای آموزش مدل ضروری است.

۵. **استراتژی‌های پیش پردازش داده‌ها:** بسته به مجموعه داده، ممکن است به استراتژی‌های پیش پردازش متفاوتی نیاز داشته باشید. به عنوان مثال، متن از رسانه‌های اجتماعی ممکن است در مقایسه با مقالات دانشگاهی نیاز به تمیز کردن و کاهش نویز گسترده تری داشته باشد. دانستن ماهیت داده‌ها، تصمیم‌گیری در مورد تکنیک‌هایی مانند ریشه‌یابی، واژه‌سازی یا مدیریت کاراکترهای خاص را تعیین می‌کند.

۶. **دانش اختصاصی دامنه:** وظایف NLP اغلب شامل اصطلاحات و زمینه خاص دامنه است. دانستن دامنه مجموعه داده (به عنوان مثال، پزشکی، حقوقی، مالی) برای استخراج موثر ویژگی، سفارشی سازی مدل، و درک ظرافت های متن ضروری است.

۷. **تنوع زبان:** مناطق و جوامع مختلف ممکن است زبان را به طور متفاوتی استفاده کنند. تشخیص منشأ مجموعه داده و زبان های استفاده شده در آن، به ویژه برای کارهایی مانند شناسایی زبان یا تشخیص گویش، مهم است.

۸. **تشخیص ناهنجاری:** در برخی از برنامه های NLP، مانند تشخیص تقلب یا فیلتر کردن هرزنامه، دانستن ویژگی های داده های عادی و غیرعادی بسیار مهم است. درک مجموعه داده به شناسایی آنچه رفتار یا زبان عادی را تشکیل می دهد و در نتیجه، تشخیص ناهنجاری ها کمک می کند.

۹. **افزایش داده ها:** تکنیک های افزایش داده ها، مانند جایگزینی مترادف یا ترجمه متن، زمانی که درک خوبی از محتوا و زمینه مجموعه داده داشته باشید، می توانند به طور موثرتری اعمال شوند.

در اصل، دانستن مجموعه داده مانند درک مواد خامی است که در NLP با آن کار می کنید. این شما را در تصمیم گیری آگاهانه در هر مرحله از خط لوله NLP، از جمع آوری داده ها و پیش پردازش گرفته تا آموزش و ارزیابی مدل، راهنمایی می کند. این درک به شما اجازه می دهد تا رویکرد خود را با ویژگی ها و چالش های خاص ناشی از داده ها تطبیق دهید و در نهایت منجر به نتایج و برنامه های کاربردی بهتر NLP می شود.

مجموعه داده‌ای که ما با آن کار خواهیم کرد مشخصات زیر را دارد:

- ستون اول **user ID** کاربر
- ستون دوم متن خام توییت (۴۰.۰۰۰)
- ستون سوم احساس (۱۳ مدل)

این مجموعه داده دامنه عمومی از پلتفرم **data.world** جمع آوری شده است. با تشکر، **data.world** برای انتشار آن.

۳.۲ پیش پردازش داده‌ها

برای آماده سازی داده‌ها ما کارهای زیر را انجام دادیم:

۱. **cleaning**
۲. **tokenization**
۳. **removing stopwords**
۴. **lemmatization**
۵. **stemming**

کدی که برای این مرحله زدیم به شرح زیر است:

۳.۲.۱ حذف کردن **url**

حذف **URL** ها از داده های متنی، به ویژه در زمینه **NLP** برای توییت ها یا رسانه های اجتماعی، به چند دلیل یک مرحله پیش پردازش رایج است:

- **کاهش نویز:** **URL** های موجود در توییت ها اغلب حاوی نویز زیادی از جمله کاراکترها، نمادها و اعداد تصادفی هستند. حذف آنها به کاهش نویز در داده های متن کمک می کند و آن را تمیزتر و مناسب تر برای تجزیه و تحلیل می کند.

- **یکنواختی:** حذف URL ها تضمین می کند که داده های متنی سازگارتر هستند. از آنجایی که URL ها می توانند از نظر طول و محتوا به طور قابل توجهی متفاوت باشند، حضور آنها می تواند تنوعی را در مجموعه داده ایجاد کند که به وظیفه NLP مرتبط نیست.
- **تمرکز بر محتوای متنی:** وظایف NLP معمولاً بر تجزیه و تحلیل و درک محتوای متنی داده ها متمرکز است. URL ها که آدرس های وب هستند، اغلب حاوی اطلاعات معناداری برای تجزیه و تحلیل خود متن نیستند. حذف URL ها به تمرکز روی کلمات و عباراتی که پیام یا احساس موجود در توییت را منتقل می کنند کمک می کند.
- **کاهش ابعاد:** URL ها می توانند بسیار طولانی باشند، و ممکن است اغلب در توییت ها، به خصوص در ریتوییت ها یا هنگام اشتراک گذاری لینک ها، ظاهر شوند. گنجاندن URL ها در داده های متنی می تواند ابعاد مجموعه داده را به میزان قابل توجهی افزایش دهد. حذف آنها ابعاد را کاهش می دهد و آن را برای کارهای NLP قابل مدیریت و کارآمدتر می کند.
- **ثبات در توکن سازی:** زمانی که متن را نشانه گذاری می کنید (آن را به کلمات یا نشانه ها تقسیم می کنید)، URL ها اغلب به چندین توکن تقسیم می شوند که هر کدام شامل بخش هایی از URL است. حذف نشانی های اینترنتی تضمین می کند که توکن سازی به جای تکه هایی از URL ها، به کلمات معنادار منجر می شود.
- **حریم خصوصی و امنیت:** URL ها گاهی اوقات می توانند به وب سایت ها یا منابع خارجی منجر شوند که ممکن است پیامدهای حفظ حریم خصوصی داشته باشند. با حذف URL ها، می توانید از حریم خصوصی کاربران محافظت کنید و اطمینان حاصل کنید که اطلاعات حساس به طور ناخواسته در معرض دید قرار نمی گیرند یا به آنها دسترسی پیدا نمی شود.
- **عملکرد مدل بهبود یافته:** مدل های NLP اغلب زمانی که بر روی داده های متنی تمیز و از پیش پردازش شده آموزش می بینند، عملکرد بهتری دارند. حذف URL ها یک مرحله در خط لوله پیش پردازش است که می تواند با کاهش نویز و اطلاعات نامربوط به بهبود عملکرد مدل کمک کند.
- **افزایش خوانایی:** برای خوانندگان انسانی یا تحلیلگرانی که متن را بررسی می کنند، URL ها چندان آموزنده نیستند و می توانند جریان فهم را مختل کنند. حذف آنها خوانایی متن را افزایش می دهد.

به طور خلاصه، حذف URL ها از داده های متنی در NLP برای توییت ها یک روش معموله برای تمیز کردن و آماده سازی داده ها برای تجزیه و تحلیل، بهبود عملکرد مدل و تمرکز بر محتوای متنی مربوطه توییت ها است. این بخشی از مجموعه گسترده تری از مراحل پیش پردازش متن است که با هدف استخراج اطلاعات معنی دار از متن و در عین حال کاهش نویز و اطلاعات نامربوط انجام می شود.

ما با استفاده از این قطعه کد لینک ها را از مجموعه داده خود حذف کردیم:

```
# Function to remove URLs from text
def remove_urls(text):
    url_pattern = re.compile(r'https?://\S+|www\.\S+')
    return url_pattern.sub('', text)
```

۳.۲.۲ حذف mentions

اشاره در توییت ها که به نام های کاربری یا دسته های توییت نیز شناخته می شوند، اشاره ای به سایر کاربران توییت هستند. قبل از آنها علامت "@" و پس از آن نام کاربری وجود دارد. در اینجا توضیح مختصری در مورد اینکه چرا اشاره ها معمولاً در طول پیش پردازش متن برای وظایف NLP شامل توییت ها به روشی خاص برخورد می شوند آورده شده است:

۱. **شناسایی تعاملات کاربر:** اشاره ها برای درک نحوه تعامل کاربران با یکدیگر در توییت بسیار مهم هستند. آنها نشان می دهند که یک کاربر مستقیماً به کاربر دیگری آدرس می دهد یا به او ارجاع می دهد. تجزیه و تحلیل اشاره ها می تواند به مطالعه تعاملات اجتماعی، احساسات و نفوذ در جامعه توییت کمک کند.

۲. **حذف اطلاعات خاص کاربر:** در برخی از کارهای NLP، ممکن است بخواهید برای محافظت از حریم خصوصی کاربر یا تمرکز صرفاً بر محتوای متنی توییت ها، اشاره ها را حذف کنید. از آنجایی که اشاره ها شامل نام های کاربری خاصی هستند، حذف آنها می تواند به حفظ ناشناس بودن و تعمیم تحلیل کمک کند.

۳. **کاهش ابعاد:** مشابه URL ها، اشاره می تواند اغلب در توییت ها، به خصوص در مکالمات، پاسخ ها، یا ریتوییت ها اتفاق بیفتد. گنجاندن نام های موجود در داده های متنی می تواند ابعاد مجموعه داده را به میزان قابل توجهی افزایش دهد. حذف نام ها ابعاد را کاهش می دهد و می تواند منجر به مدل سازی کارآمدتر شود.

۴. **ثبات توکن سازی:** هنگام توکن کردن متن، اشاره ها اغلب به نشانه های جداگانه تقسیم می شوند، یکی برای نماد "@" و دیگری برای نام کاربری. حذف نام ها تضمین می کند که توکن سازی به جای تکه هایی از نام های کاربری منجر به کلمات معنادار می شود.

۵. **ساده کردن تجزیه و تحلیل متن:** در برخی از وظایف NLP، ممکن است علاقه ای به تجزیه و تحلیل تعاملات بین کاربران خاص یا مطالعه اشاره ها به عنوان تمرکز اصلی نداشته باشید. حذف اشاره ها با تمرکز بر محتوا و احساسات توییت بدون در نظر گرفتن تعاملات کاربر، تجزیه و تحلیل متن را ساده می کند.

۶. **استخراج محتوا:** برای تجزیه و تحلیل احساسات یا مدل سازی موضوع، اشاره ها ممکن است نویز در نظر گرفته شوند زیرا معمولاً حاوی احساسات یا اطلاعات موضوعی نیستند. حذف نام ها به شما امکان می دهد تا محتوای اصلی توییت را برای تجزیه و تحلیل استخراج کنید.

۷. **احترام به داده های کاربر:** حذف نام ها می تواند راهی برای احترام به حریم خصوصی و مالکیت داده های کاربران توییت باشد. در تحقیقات یا برنامه های کاربردی، مهم است که داده های کاربر را مسئولانه و مطابق با دستورالعمل های حفظ حریم خصوصی مدیریت کنید.

قطعه کدی که برای حذف mentions بود به شرح زیر است:

```
# Function to remove mentions from text
def remove_mentions(text):
    mention_pattern = re.compile(r'@\w+')
    return mention_pattern.sub('', text)
```

۳.۲.۳ تمیز کردن متن

پاکسازی متن یک مرحله پیش پردازش ضروری در پردازش زبان طبیعی (NLP) است که شامل فرآیند تهیه داده های متن خام برای تجزیه و تحلیل است. در اینجا توضیحی در مورد اینکه چرا پاکسازی متن بسیار مهم است و مراحل مربوطه آورده شده است:

اهمیت پاکسازی متن در NLP:

۱. **کاهش نویز:** داده های متنی اغلب حاوی نویز به شکل کاراکترهای خاص، علائم نگارشی، تگ های HTML و نمادهای نامربوط هستند. پاک کردن متن این عناصر را حذف می کند و داده ها را برای تجزیه و تحلیل مناسب تر می کند.
۲. **سازگاری:** تمیز کردن تضمین می کند که داده های متنی سازگار و استاندارد هستند. تغییرات در حروف بزرگ، املا و قالب بندی را می توان کاهش داد، که منجر به یک مجموعه داده یکنواخت تر می شود.
۳. **توکن سازی بهبودیافته:** توکن سازی، فرآیند تقسیم متن به کلمات یا نشانه ها، بر روی متن تمیز مؤثرتر است. متن پاک شده منجر به توکن سازی دقیق تری می شود که برای کارهای پایین دستی NLP حیاتی است.
۴. **حذف کلمات توقف:** تمیز کردن اغلب شامل حذف کلمات توقف (کلمات رایج مانند "و"، "the"، "است") است که معنای زیادی ندارند. این مرحله ابعاد داده ها را کاهش می دهد و روی کلمات معنی دارتر تمرکز می کند.
۵. **خوانایی پیشرفته:** خواندن و درک متن پاک شده آسان تر است، هم برای انسان و هم برای مدل های NLP. درک بهتر محتوای متن را تسهیل می کند.
۶. **کاهش ابعاد:** با حذف کاراکترها و نمادهای غیر ضروری، تمیز کردن می تواند ابعاد داده ها را به میزان قابل توجهی کاهش دهد. این امر به ویژه هنگام کار با مجموعه داده های بزرگ مهم است.
۷. **عملکرد مدل بهبود یافته:** مدل های NLP اغلب زمانی که بر روی متن تمیز و از قبل پردازش شده آموزش می بینند بهتر عمل می کنند. داده های پاک ابهام و نویز را کاهش می دهد که می تواند منجر به پیش بینی دقیق تر مدل شود.

مراحل پاکسازی متن:

تمیز کردن متن معمولاً شامل یک سری مراحل است که ممکن است شامل موارد زیر باشد:

۱. **حروف کوچک:** تبدیل تمام متن به حروف کوچک برای اطمینان از سازگاری و کاهش تأثیر

حروف بزرگ.

۲. **حذف کاراکترهای خاص:** حذف کاراکترهای خاص، علائم نقطه گذاری و نمادهایی که حاوی

اطلاعات معنی دار نیستند.

۳. **مدیریت اعداد:** بسته به کار، می توانید اعداد را در متن نگه دارید، حذف کنید یا جایگزین کنید.

۴. **حذف برچسب های HTML:** هنگام کار با داده های وب، حذف برچسب های HTML

برای استخراج فقط محتوای متنی بسیار مهم است.

۵. **ریشه یا Lemmatization:** کاهش کلمات به شکل پایه یا ریشه آنها. ساقه زنی

تهاجمی تر است، در حالی که lemmatization دقیق تر است.

۶. **Stopword Removal:** حذف کلمات توقف متداول که کمک زیادی به معنای متن

ندارند.

۷. **بررسی املا:** تصحیح اشتباهات املائی، در صورت وجود.

۸. **Tokenization:** تقسیم متن به کلمات یا نشانه های جداگانه برای تجزیه و تحلیل

بیشتر.

۹. **تمیز کردن سفارشی:** بسته به مجموعه داده و کار خاص، مراحل تمیز کردن اضافی ممکن است

لازم باشد، مانند مدیریت ایموجی ها، URL ها، یا ذکرها.

ما در این پروژه از قطعه کد زیر برای تمیز کردن متن استفاده کردیم:

```
# Function to cleaning the text
```

```
def clean_text(text):
```

```
    # Remove special characters, hashtags, and mentions
```

```
    cleaned_text = re.sub(r"[^a-zA-Z0-9\s]", "", text)
```

```
    #cleaned_text = re.sub(r"#\w+", "", cleaned_text)
```

```
    cleaned_text = re.sub(r"@w+", "", cleaned_text)
```



```
# Convert text to lowercase
cleaned_text = cleaned_text.lower()

# Remove extra whitespace
cleaned_text = re.sub(r"\s+", " ", cleaned_text).strip()

return cleaned_text

# Function to tokenize
def tokenize_text(text):
    # Tokenize the text into individual words/tokens
    tokens = word_tokenize(text)
    return tokens

# Function to remove stop words from text
def remove_stopwords(tokens):
    stop_words = set(stopwords.words('english'))
    filtered_tokens = [token for token in tokens if token.lower() not in
stop_words]
    return filtered_tokens

# Function to remove punctuations and numbers
def remove_punctuations_and_numbers(text):
    # Remove punctuation marks
    text = text.translate(str.maketrans('', '', string.punctuation))

    # Remove numbers
    text = re.sub(r'\d+', '', text)

    return text

# Lemmatizer Function
```

```
def perform_lemmatization(tokens):
    lemmatizer = WordNetLemmatizer()
    lemmatized_tokens = [lemmatizer.lemmatize(token) for token in tokens]
    return lemmatized_tokens

# Stemmer Function
def perform_stemming(tokens):
    stemmer = PorterStemmer()
    stemmed_tokens = [stemmer.stem(token) for token in tokens]
    return stemmed_tokens
```

در مرحله بعدی ما داده‌ها را در فایل جداگانه ذخیره می‌کنیم که در مرحله بعدی از فایل جدید استفاده کنیم.

```
# Save to pre_processed.csv
df.to_csv('pre_processed.csv', index=False)
```

۳.۳ تقسیم داده‌ها

تقسیم یک مجموعه داده به زیرمجموعه‌های آموزشی و آزمایشی، به چند دلیل مهم، یک عمل اساسی در یادگیری ماشین، از جمله پردازش زبان طبیعی (NLP) است:

۱. **ارزیابی مدل:** هدف اصلی از تقسیم مجموعه داده، ارزیابی عملکرد یک مدل یادگیری ماشین است. با رزرو بخشی از داده‌ها برای آزمایش (مجموعه تست)، می‌توانید ارزیابی کنید که مدل شما چقدر به نمونه‌های جدید و دیده نشده تعمیم می‌یابد. بدون یک مجموعه آزمایشی جداگانه، راه قابل اعتمادی برای اندازه‌گیری عملکرد مدل خود در داده‌های دنیای واقعی نخواهید داشت.

۲. **تشخیص بیش‌برازش:** تقسیم داده‌ها به تشخیص اضافه‌برازش کمک می‌کند، که زمانی اتفاق می‌افتد که یک مدل داده‌های آموزشی را خیلی خوب یاد می‌گیرد اما نمی‌تواند به داده‌های جدید تعمیم یابد. اگر یک مدل در داده‌های آموزشی عملکرد فوق‌العاده‌ای داشته باشد اما در داده‌های آزمایشی ضعیف عمل کند، این نشانه بیش از حد برازش است. این به شما اطلاع می‌دهد که

مدل ممکن است خیلی پیچیده باشد یا داده های آموزشی را به جای یادگیری از آن به خاطر بسپارد.

۳. **تنظیم فرایارامتر:** در طول توسعه مدل های یادگیری ماشین، اغلب نیاز به تنظیم دقیق فرایارامترها (به عنوان مثال، نرخ یادگیری، نقاط قوت منظم سازی) دارید تا به بهترین عملکرد برسید. شما از مجموعه تست برای ارزیابی تاثیر این تغییرات فرایارامتر بر عملکرد مدل استفاده می کنید، بدون اینکه بایاس از داده های آموزشی وارد کنید.

۴. **انتخاب مدل:** در برخی موارد، ممکن است بخواهید چندین مدل را با هم مقایسه کنید تا مشخص کنید کدام یک در داده های دیده نشده بهترین عملکرد را دارد. مجموعه تست به شما امکان می دهد این مقایسه را به صورت عینی انجام دهید و مناسب ترین مدل را برای کار NLP خود انتخاب کنید.

۵. **جلوگیری از نشت داده ها:** نگه داشتن یک مجموعه آزمایشی مجزا به شما کمک می کند تا اطمینان حاصل کنید که در طول آموزش مدل، سهواً از هیچ اطلاعاتی از داده های آزمایش استفاده نکنید. نشت داده ها می تواند منجر به تخمین های عملکرد بیش از حد خوش بینانه شود، و باعث می شود مدل شما بهتر از آنچه هست ظاهر شود.

۶. **شبیه سازی دنیای واقعی:** مجموعه تست نحوه عملکرد مدل شما را در یک سناریوی واقعی که در آن با داده های جدید و دیده نشده مواجه می شود، شبیه سازی می کند. این برای ارزیابی اینکه آیا مدل NLP شما برای استقرار آماده است یا خیر حیاتی است.

۷. **اعتبارسنجی متقاطع:** علاوه بر یک تقسیم آزمایش ساده قطار، می توانید تکنیک های پیشرفته تری مانند اعتبارسنجی متقاطع k-fold را انجام دهید. اعتبارسنجی متقابل داده ها را به زیر مجموعه های متعدد تقسیم می کند و به شما این امکان را می دهد که مدل خود را چندین بار روی ترکیب های مختلف مجموعه های آموزشی و اعتبارسنجی آموزش و آزمایش کنید. این یک تخمین قوی تر از عملکرد مدل ارائه می دهد.

۸. **ارزیابی تعمیم:** تقسیم داده ها به شما امکان می دهد توانایی مدل خود را برای تعمیم اندازه گیری کنید. تعمیم به ظرفیت مدل برای به کارگیری آنچه از داده های آموزشی آموخته است برای پیش بینی های دقیق بر روی نمونه های جدید و نادیده اشاره دارد - یک جنبه حیاتی از یادگیری ماشین.

به طور خلاصه، تقسیم یک مجموعه داده به زیرمجموعه های آموزشی و آزمایشی یک گام مهم در گردش کار یادگیری ماشین، از جمله NLP است. این به شما امکان می دهد عملکرد مدل را ارزیابی کنید، بیش از حد برازش را تشخیص دهید، فرایامترها را تنظیم کنید، بهترین مدل را انتخاب کنید، و اطمینان حاصل کنید که مدل شما می تواند به داده های جدید و نادیده تعمیم یابد - جنبه های اساسی ساخت مدل های یادگیری ماشین موثر و قابل اعتماد.

در این پروژه ما با استفاده کد های زیر ابتدا ۴ لایه از احساسات موجود در دیتاست را جدا میکنیم.

```
df = pd.read_csv('pre_processed.csv')
desiered_labels = ['neutral', 'sadness', 'worry', 'happiness']
filtered_df = df[df['emotions'].isin(desiered_labels)]
filtered_df.to_csv('4_emotions_data.csv')
```

در ادامه label ها را encode کرده

```
# Load data frame
df = pd.read_csv('4_emotions_data.csv')
df.dropna(inplace=True)
X = df['joined_lemmatized']
y = df['emotions']

# Encode the labels into numerical values
label_encoder = LabelEncoder()
encoded_labels = label_encoder.fit_transform(y)
```

و سپس داده ها را به دو بخش **test, train** تقسیم میکنیم

```
# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, encoded_labels,
test_size=0.2, random_state=42)
```

مرحله بعدی **vectorize** کردن داده ها است که قطعه کد زیر آین امر را میسر میکند

```
# Create a TF-IDF vectorizer
vectorizer = TfidfVectorizer()
```

```
# Vectorize the text data
```

```
X_train_vectorized = vectorizer.fit_transform(X_train)
```

```
X_test_vectorized = vectorizer.transform(X_test)
```

بعد از طی کردن این مراحل تمام کار ما با پیش پردازش داده‌ها به اتمام رسیده و دو مجموعه x_train, y_train آماده تزریق به الگوریتم‌هایی هستند که کار یادگیری را انجام می‌دهند.

۳.۴ ساختن مدل

در زمینه یادگیری ماشینی، مدل یک نمایش ریاضی یا محاسباتی است که الگوها و روابط را از داده‌ها یاد می‌گیرد. این به عنوان یک سیستم الگوریتمی عمل می‌کند که می‌تواند پیش‌بینی‌ها، طبقه‌بندی‌ها یا دیگر استنباط‌ها را بر روی داده‌های جدید و نادیده بر اساس آنچه در طول آموزش آموخته است انجام دهد.

یک مدل اساساً تابعی است که داده‌های ورودی را به پیش‌بینی‌های خروجی نگاشت می‌کند. به عنوان مثال، در NLP، یک مدل ممکن است یک قطعه از متن را بگیرد و احساسات، دسته‌بندی یا هر اطلاعات مرتبط دیگری را پیش‌بینی کند.

ایجاد یک مدل:

۱. **انتخاب یک معماری مدل:** اولین قدم در ایجاد یک مدل، انتخاب یک معماری مناسب است. در

NLP، این می‌تواند مدل‌های مختلفی باشد، مانند:

۲. **مدل‌های کلاسیک** (مانند Naive Bayes، SVM): این مدل‌ها بر اساس اصول آماری و

ریاضی سستی هستند.

۳. **شبکه‌های عصبی** (به عنوان مثال، LSTM، RNN، ترانسفورماتور): این مدل‌ها بر اساس

شبکه‌های عصبی مصنوعی هستند که از ساختار مغز انسان الهام گرفته شده‌اند.

۴. **آماده‌سازی داده‌ها:** قبل از آموزش یک مدل، داده‌ها باید آماده شوند. این شامل کارهایی مانند

تمیز کردن متن، توکن‌سازی، برداری (تبدیل متن به فرمت عددی) و تقسیم داده‌ها به مجموعه

های آموزشی و آزمایشی است.

۵. **مهندسی ویژگی (اختیاری):** در NLP، مهندسی ویژگی شامل تبدیل داده‌های متن خام به قالبی است که برای مدل انتخابی مناسب است. این ممکن است شامل تکنیک‌هایی مانند جاسازی کلمه، TF-IDF یا سایر نمایش‌ها باشد.

۶. **کامپایل و آموزش مدل:** برای مدل‌های کلاسیک، این مرحله شامل تنظیم مدل، تعیین فرآپارامترها و سپس برآزش آن با داده‌های آموزشی است. برای شبکه‌های عصبی، این شامل تعریف معماری، کامپایل مدل و آموزش آن با استفاده از یک الگوریتم بهینه‌سازی است.

۷. **ارزیابی مدل:** هنگامی که مدل آموزش داده شد، برای ارزیابی عملکرد آن باید در مجموعه تست ارزیابی شود. معیارهای ارزیابی رایج در NLP شامل دقت، دقت، یادآوری، امتیاز F1 و معیارهای خاص‌تر بسته به کار است (به عنوان مثال، امتیاز BLEU برای ترجمه ماشینی).

۸. **تنظیم و بهینه‌سازی:** بسته به نتایج ارزیابی، ممکن است نیاز به تنظیم دقیق مدل داشته باشید. این می‌تواند شامل تنظیم هایپرپارامترها، تغییر معماری، یا انجام تکنیک‌هایی مانند منظم‌سازی برای بهبود عملکرد باشد.

۹. **استقرار و استنتاج:** هنگامی که مدل به طور رضایت بخشی عمل کرد، می‌توان آن را برای استفاده در دنیای واقعی مستقر کرد. این شامل دریافت داده‌های جدید و پیش‌بینی یا طبقه‌بندی بر اساس آموخته‌های آن است.

۱۰. **نظارت و نگهداری:** پس از استقرار، نظارت بر عملکرد مدل در طول زمان بسیار مهم است. اگر توزیع داده‌های اساسی تغییر کند یا اگر عملکرد مدل بدتر شود، ممکن است نیاز به بازآموزی یا تنظیم دقیق باشد.

ما در این پروژه از مدل‌های زیر استفاده میکنیم:

۱. SVC

۲. Random Forest

۳. XGBoost

۴. Naive Bayes

در ادامه به توضیح این مدل‌ها می‌پردازیم.

SVC ۳.۵.۱

SVC مخفف عبارت Support Vector Classification است. این یک الگوریتم یادگیری ماشینی نظارت شده است که برای کارهای طبقه بندی استفاده می شود. طبقه بندی بردار پشتیبان برای طبقه بندی به عنوان SVM (ماشین بردار پشتیبانی) نیز شناخته می شود.

۱. **وظیفه طبقه بندی:** هدف اصلی مدل طبقه بندی بردار پشتیبانی (SVC) طبقه بندی نقاط داده به یکی از دو یا چند کلاس است. این به ویژه برای مسائل طبقه بندی باینری مفید است، جایی که هدف تخصیص نقاط داده به یکی از دو کلاس است، اما می توان آن را به طبقه بندی چند کلاسه گسترش داد.

۲. **حداکثر کردن حاشیه:** SVC با یافتن یک مرز تصمیم (یا ابرصفحه) کار می کند که نقاط داده کلاس های مختلف را به بهترین شکل جدا می کند و در عین حال حاشیه بین کلاس ها را به حداکثر می رساند. حاشیه فاصله بین ابر صفحه و نزدیکترین نقاط داده هر کلاس است. حداکثر کردن حاشیه به بهبود تعمیم مدل کمک می کند و خطر بیش از حد برازش را کاهش می دهد.

۳. **بردارهای پشتیبان:** بردارهای پشتیبان، نقاط داده ای هستند که به مرز تصمیم نزدیک ترین هستند. اینها مهمترین نقاط در مجموعه داده هستند زیرا آنها هستند که بر موقعیت و جهت گیری مرز تصمیم تأثیر می گذارند. الگوریتم "بردار پشتیبانی" نامگذاری شده است زیرا بر این نقاط داده کلیدی متکی است.

۴. **ترفند هسته:** SVC می تواند با استفاده از تکنیکی به نام "ترفند هسته"، هم مجموعه داده های قابل جداسازی خطی و هم مجموعه داده های غیرخطی را مدیریت کند. این تکنیک فضای ویژگی اصلی را به فضایی با ابعاد بالاتر تبدیل می کند که در آن نقاط داده به صورت خطی قابل تفکیک می شوند. توابع رایج هسته شامل هسته خطی، هسته چند جمله ای و هسته تابع پایه شعاعی (RBF) است.

۵. **پارامتر C:** الگوریتم SVC دارای یک فراپارامتر به نام "C" (پارامتر منظم سازی) است که مبادله بین حداکثر کردن حاشیه و به حداقل رساندن خطای طبقه بندی را کنترل می کند. مقدار کمتر C اجازه می دهد تا حاشیه وسیع تری داشته باشد، اما ممکن است منجر به طبقه بندی

اشتباه برخی از نقاط داده شود، در حالی که مقدار بزرگتر C منجر به حاشیه باریک تر اما طبقه بندی اشتباه کمتر می شود.

۶. **Hyperplane**: در یک مسئله طبقه بندی باینری (دو کلاس)، مدل **SVC** یک ابر صفحه را پیدا می کند که داده ها را به دو کلاس جدا می کند. برای مسائل چند کلاسه، از چند ابر صفحه برای جداسازی جفت های مختلف کلاس استفاده می شود.

۷. **پیش بینی**: زمانی که مدل **SVC** بر روی یک مجموعه داده برچسب دار آموزش داده شد، می توان از آن برای پیش بینی نقاط داده جدید و بدون برچسب استفاده کرد. این مدل نقاط داده جدید را به یکی از کلاس ها بر اساس موقعیت آنها نسبت به مرز تصمیم اختصاص می دهد.

مزایای کلیدی SVC:

- **SVC** در فضاهای با ابعاد بالا موثر است.
- می تواند مشکلات طبقه بندی خطی و غیر خطی را حل کند.
- استفاده از بردارهای پشتیبانی به تمرکز بر آموزنده ترین نقاط داده کمک می کند.
- هنگامی که تنظیم مناسب اعمال می شود در برابر بیش از حد برازش قوی است.

محدودیت های SVC:

- **SVC** می تواند از نظر محاسباتی گران باشد، به خصوص با مجموعه داده های بزرگ.
- انتخاب هسته و فرامپارامترها می تواند بر عملکرد مدل تأثیر بگذارد و نیاز به تنظیم دقیق دارد.
- ممکن است در مجموعه داده هایی با کلاس های نویزدار یا همپوشانی خوب عمل نکند.

به طور خلاصه، دسته بندی بردار پشتیبانی (**SVC**) یک الگوریتم یادگیری ماشینی قدرتمند است که برای کارهای طبقه بندی استفاده می شود، به ویژه زمانی که هدف یافتن مرز تصمیم گیری است که حاشیه بین کلاس ها را به حداکثر می رساند. این می تواند مشکلات طبقه بندی خطی و غیر خطی را از طریق استفاده از توابع هسته حل کند و زمانی موثر است که داده ها ابعاد بالایی داشته باشند. با این حال، به تنظیم دقیق پارامتر نیاز دارد و می تواند محاسباتی فشرده باشد.

در این بخش ما با استفاده از قطعه کد زیر یک مدل **SVC** میسازیم و داده‌ها را به آن تزریق می‌کنیم

```
# SVC model
# Create SVC model
svc_model = SVC()

# Train the model
svc_model.fit(X_train_vectorized, y_train)

# Make predictions on the test set
y_pred_svc = svc_model.predict(X_test_vectorized)

# Decode the predicted labels
decoded_pred = label_encoder.inverse_transform(y_pred_svc)

# Save the model to a file
filename_svc = 'svc_model.sav'
joblib.dump(svc_model, filename_svc)
```

در این قطعه کد ما ابتدا مدل یک شی از جنس مدل **SVC** میسازیم.

داده‌های مربوط به یادگیری را به مدل پاس می‌دهیم.

دقت یادگیری مدل را با استفاده از داده‌های تست اندازه‌گیری می‌کنیم.

لیبل‌های داده‌های تست را دیکد می‌کنیم که بعداً بتوانیم برای بررسی آن‌ها را تشخیص دهیم.

و در نهایت هم مدل را ذخیره می‌کنیم.

۳.۵.۲ Random Forest

مدل جنگل تصادفی یک روش یادگیری گروهی است که هم برای کارهای طبقه‌بندی و هم رگرسیون در یادگیری ماشین استفاده می‌شود. این یک الگوریتم همه کاره و قدرتمند است که به دلیل توانایی خود در ارائه پیش بینی های دقیق و قوی محبوبیت پیدا کرده است.

مفاهیم کلیدی:

۱. **یادگیری گروهی:** جنگل تصادفی یک تکنیک یادگیری گروهی است، به این معنی که پیش‌بینی‌های چند مدل پایه (درخت تصمیم) را برای بهبود عملکرد کلی ترکیب می‌کند. به جای تکیه بر یک درخت تصمیم، Random Forest از خرد جمعیت برای پیش‌بینی دقیق‌تر استفاده می‌کند.
۲. **درختان تصمیم:** در هسته یک جنگل تصادفی، درختان تصمیم فردی قرار دارند. درخت‌های تصمیم سازه‌های ساده و درخت‌مانندی هستند که بر اساس مجموعه‌ای از قوانین آموخته شده از داده‌ها تصمیم می‌گیرند. هر درخت تصمیم در Random Forest به طور مستقل ساخته می‌شود و مجموعه قوانین خاص خود را دارد.
۳. **تصادفی سازی:** "تصادفی" در جنگل تصادفی به دو منبع کلیدی تصادفی مورد استفاده در ساخت آن اشاره دارد:
۴. **نمونه‌گیری تصادفی:** هر درخت تصمیم بر روی زیرمجموعه‌ای تصادفی از داده‌های آموزشی با جایگزینی (نمونه‌گیری بوت استرپ) آموزش داده می‌شود. این باعث ایجاد تنوع در بین درختان می‌شود.
۵. **انتخاب ویژگی تصادفی:** در هر گره درخت تصمیم، یک زیرمجموعه تصادفی از ویژگی‌ها هنگام تصمیم‌گیری تقسیم در نظر گرفته می‌شود. این باعث تنوع بیشتر و کاهش بیش از حد برازش می‌شود.
۶. **رای گیری یا میانگین گیری:** برای کارهای طبقه‌بندی، Random Forest پیش‌بینی‌های درختان منفرد را با رای اکثریت ترکیب می‌کند. به عبارت دیگر، کلاسی را انتخاب می‌کند که بیشترین رای را از درختان جداگانه دریافت کند. برای وظایف رگرسیون، میانگین مقادیر پیش‌بینی شده را از همه درخت‌ها می‌گیرد.

مزایای جنگل تصادفی:

- **دقت بالا:** Random Forest به دلیل دقت پیش‌بینی بالا معروف است. با جمع‌آوری پیش‌بینی‌ها از چندین درخت، خطر بیش از حد برازش را کاهش می‌دهد و نتایج قوی ارائه می‌کند.
- **انتخاب ویژگی ضمنی:** فرآیند انتخاب ویژگی تصادفی می‌تواند به عنوان مکانیزم انتخاب ویژگی ضمنی عمل کند. ویژگی‌هایی که به طور مداوم در بالای درختان تصمیم ظاهر می‌شوند، احتمالاً در پیش‌بینی‌ها مهم هستند.
- **مقاومت در برابر برازش بیش از حد:** ترکیبی از تکنیک‌های تصادفی‌سازی و میانگین‌گیری در میان درختان، جنگل تصادفی را حتی در مجموعه داده‌های پرسر و صدا یا با ابعاد بالا در برابر بیش‌برازش مقاوم می‌کند.
- **مدیریت مقادیر از دست رفته:** Random Forest می‌تواند مقادیر از دست رفته در مجموعه داده را مدیریت کند. این نیازی به نسبت داده‌های از دست رفته ندارد، و همچنان می‌تواند پیش‌بینی‌های معناداری ارائه دهد.
- **اهمیت ویژگی:** می‌تواند اهمیت هر یک از ویژگی‌ها را در مجموعه داده تخمین بزند و به تشخیص اینکه کدام متغیرها بیشترین سهم را در پیش‌بینی‌ها دارند کمک می‌کند.

محدودیت‌های جنگل تصادفی:

- **مدل جعبه سیاه:** جنگل تصادفی را می‌توان یک مدل جعبه سیاه در نظر گرفت، و تفسیر منطق تصمیم خاص مورد استفاده توسط درختان را به چالش می‌کشد.
- **محاسباتی فشرده:** آموزش یک جنگل تصادفی بزرگ با درختان زیاد می‌تواند از نظر محاسباتی گران باشد و ممکن است به حافظه قابل توجهی نیاز داشته باشد.
- **برای همه داده‌ها مناسب نیست:** در حالی که Random Forest یک الگوریتم همه‌کاره است، ممکن است بهترین انتخاب برای همه انواع داده‌ها و مشکلات نباشد.

به طور خلاصه، مدل جنگل تصادفی یک تکنیک یادگیری گروهی قدرتمند است که قدرت درخت‌های تصمیم‌گیری متعدد را برای ارائه پیش‌بینی‌های دقیق و قوی ترکیب می‌کند. این به طور گسترده در برنامه‌های مختلف یادگیری ماشین، از جمله وظایف طبقه‌بندی و رگرسیون، که در آن عملکرد بالا و انعطاف پذیری در برابر بیش از حد مطلوب است، استفاده می‌شود.

قطعه کد استفاده شده برای ساختن این مدل نیز همانند **SVC** است و تفاوتی ندارد که در زیر آمده است.

```
# Random Forest model

# Create a Random Forest model
rf_model = RandomForestClassifier()

# Train the model
rf_model.fit(X_train_vectorized, y_train)

# Make predictions on the test set
y_pred_rf = rf_model.predict(X_test_vectorized)

# Decode the predicted labels
decoded_pred = label_encoder.inverse_transform(y_pred_rf)

# Save the model to a file
filename_rf = 'random_forest_model.sav'
joblib.dump(rf_model, filename_rf)
```

XGBoost ۳.۵.۳

XGBoost، مخفف Extreme Gradient Boosting، یک الگوریتم یادگیری ماشینی است که به خانواده روش‌های تقویت گرادیان تعلق دارد. این به دلیل عملکرد پیش‌بینی خود مشهور است و به طور گسترده در مسابقات مختلف یادگیری ماشین و برنامه‌های کاربردی در دنیای واقعی استفاده می‌شود.

مفاهیم کلیدی:

۱. تقویت گرادیان: XGBoost یک تکنیک یادگیری گروهی است که با ترکیب پیش‌بینی‌های

چند زبان‌آموز ضعیف، معمولاً درخت‌های تصمیم، یک مدل پیش‌بینی قوی ایجاد می‌کند. این یک رویکرد تقویتی را دنبال می‌کند، که در آن هر یادگیرنده جدید بر روی اشتباهات انجام شده توسط افراد قبلی تمرکز می‌کند.

۲. **Decision Trees** به عنوان **Base Learners**: در XGBoost از درخت

های تصمیم به عنوان پایه یا زبان آموزان ضعیف استفاده می‌شود. این درختان کم عمق با تعداد گره و عمق محدود هستند. ایده این است که مجموعه‌ای از چنین درختانی بسازیم که با هم کار کنند تا پیش‌بینی‌های دقیقی انجام دهند.

۳. **بهینه‌سازی گرادیان نزولی**: XGBoost از بهینه‌سازی گرادیان نزول برای به حداقل رساندن

عملکرد تلفات استفاده می‌کند. تابع ضرر تفاوت بین مقادیر پیش‌بینی شده و واقعی را اندازه‌گیری می‌کند. با تنظیم مکرر پارامترهای درختان، XGBoost این تابع از دست دادن را به حداقل می‌رساند و منجر به پیش‌بینی‌های بهتر می‌شود.

۴. **منظم‌سازی**: برای جلوگیری از برازش بیش از حد، XGBoost از تکنیک‌های منظم‌سازی،

از جمله منظم‌سازی L_1 (Lasso) و L_2 (Ridge) استفاده می‌کند. این جریمه‌ها برای پیچیدگی درختان و وزن‌های اختصاص داده شده به ویژگی‌ها اعمال می‌شود.

۵. **اهمیت ویژگی**: XGBoost معیاری از اهمیت ویژگی ارائه می‌دهد که به شناسایی

تأثیرگذارترین ویژگی‌ها در پیش‌بینی‌ها کمک می‌کند. این اطلاعات می‌تواند انتخاب ویژگی و مهندسی را راهنمایی کند.

۶. **Cross-Validation**: اعتبار متقاطع اغلب با XGBoost برای ارزیابی عملکرد آن و تنظیم هایپرپارامترها استفاده می شود. این شامل تقسیم داده ها به زیر مجموعه های متعدد و آموزش مدل بر روی ترکیب های مختلف مجموعه های آموزشی و اعتبار سنجی است.

مزایای XGBoost:

- **دقت پیش‌بینی بالا**: XGBoost برای دقت پیش‌بینی پیشرفته‌اش شناخته شده است. اغلب در طیف وسیعی از کارها از سایر الگوریتم های یادگیری ماشین بهتر عمل می کند.
- **کارایی**: XGBoost برای کارایی طراحی شده است و می تواند مجموعه داده های بزرگ و فضاهای ویژگی با ابعاد بالا را مدیریت کند.
- **استحکام**: به لطف تکنیک های منظم سازی که آن را برای مجموعه داده های پر سر و صدا یا پیچیده مناسب می کند، نسبت به بیش از حد برازش قوی است.
- **انعطاف‌پذیری**: XGBoost را می توان برای کارهای طبقه‌بندی و رگرسیون استفاده کرد و آن را برای مشکلات مختلف یادگیری ماشین همه‌کاره می سازد.
- **اهمیت ویژگی**: این مدل بینش هایی در مورد اهمیت ویژگی ارائه می دهد و به دانشمندان و تحلیلگران داده کمک می کند تا تأثیرگذارترین متغیرها را در داده ها درک کنند.

محدودیت های XGBoost:

- **پیچیدگی**: XGBoost می تواند فرآیندهای زیادی برای تنظیم داشته باشد که بهینه سازی آن را نسبت به الگوریتم های ساده تر پیچیده تر می کند.
- **منابع محاسباتی**: آموزش مجموعه های بزرگ درختان در XGBoost می تواند محاسباتی فشرده باشد و ممکن است به منابع قابل توجهی نیاز داشته باشد.
- **مدل جعبه سیاه**: مانند سایر روش های مجموعه، XGBoost را می توان یک مدل جعبه سیاه در نظر گرفت، که تفسیر منطق تصمیم گیری خاص درختان را به چالش می کشد.

به طور خلاصه، **XGBoost** یک الگوریتم یادگیری ماشینی بسیار موثر و همه کاره است که در دقت پیش بینی برتر است. این به طور گسترده در برنامه های مختلف علوم داده و یادگیری ماشین، از جمله مسابقات **Kaggle** استفاده می شود و برای ارائه عملکرد بالاتر در طیف وسیعی از وظایف، شهرت دارد.

کد استفاده شده به شرح زیر است که مشابه با مدل های قبلیست.

```
# XGBoost model
# Create an XGBoost model
xgb_model = xgb.XGBClassifier()

# Train the model
xgb_model.fit(X_train_vectorized, y_train)

# Make predictions on the test set
y_pred_xgb = xgb_model.predict(X_test_vectorized)

# Calculate the accuracy
accuracy_xgb = accuracy_score(y_test, y_pred_xgb)

# Save the model to a file
filename_xgb = 'gradient_boosting_model.sav'
joblib.dump(xgb_model, filename_xgb)
```

۳.۵.۴ Naive Bayes

Naive Bayes یک الگوریتم یادگیری ماشین احتمالی است که معمولاً برای کارهای طبقه بندی استفاده می شود. این بر اساس قضیه بیز است و یک فرض کلیدی برای استقلال ویژگی ایجاد می کند، به همین دلیل است که به آن "ساده لوح" می گویند. با وجود سادگی، Naive Bayes می تواند به طرز شگفت انگیزی در پردازش زبان طبیعی (NLP) و وظایف طبقه بندی متن موثر باشد.

مفاهیم کلیدی:

۱. **قضیه بیز:** در هسته الگوریتم ساده بیز قضیه بیز قرار دارد که یک قضیه اساسی در نظریه احتمالات است. احتمال وقوع یک رویداد را بر اساس دانش قبلی از شرایطی که ممکن است با رویداد مرتبط باشد، توصیف می کند. در طبقه بندی، به محاسبه احتمال یک کلاس با توجه به مجموعه ای از ویژگی ها کمک می کند.

۲. **استقلال مشروط:** فرض "ساده لوح" در Naive Bayes این است که ویژگی ها با توجه به برچسب کلاس مستقل هستند. این بدان معنی است که وجود یک ویژگی مستقل از وجود هر ویژگی دیگر فرض می شود. اگرچه این فرض ممکن است در واقعیت صادق نباشد، اما Naive Bayes هنوز هم می تواند در عمل عملکرد خوبی داشته باشد.

۳. **چند جمله ای و گاوسی ساده بیز:** انواع مختلفی از مدل های نایو بیز وجود دارد.

○ **Multinomial Naive Bayes:** این نوع معمولاً برای کارهای طبقه بندی متن استفاده می شود که در آن ویژگی ها تعداد کلمات یا فرکانس های عبارت را نشان می دهند. این برای داده های گسسته مانند تعداد کلمات در یک سند مناسب است.

○ **Gaussian Naive Bayes:** این نوع فرض می کند که ویژگی ها از توزیع گاوسی (عادی) پیروی می کنند. برای ویژگی های با ارزش پیوسته استفاده می شود.

۴. **تخمین احتمال ویژگی:** Naive Bayes توزیع احتمال ویژگی ها را برای هر کلاس در داده های آموزشی تخمین می زند. به عنوان مثال، در طبقه بندی متن، احتمال وقوع هر کلمه در اسناد یک کلاس خاص را محاسبه می کند.

۵. Class Prior Probability: Naive Bayes همچنین احتمال قبلی هر کلاس را

تخمین می زند. این احتمال تعلق یک سند به یک کلاس خاص بر اساس توزیع کلی کلاس ها در داده های آموزشی است.

مزایای Naive Bayes:

- **کارایی:** Naive Bayes از نظر محاسباتی کارآمد است و با مجموعه داده های بزرگ و فضاهای ویژگی با ابعاد بالا به خوبی مقیاس می شود.
- **آموزش سریع و پیش بینی:** سادگی الگوریتم منجر به آموزش سریع و زمان های پیش بینی می شود و آن را برای برنامه های بلادرنگ یا نزدیک به زمان مناسب می سازد.
- **عملکرد خوب در طبقه بندی متن:** Naive Bayes اغلب در کارهای طبقه بندی متن مانند تشخیص هرزنامه، تجزیه و تحلیل احساسات و دسته بندی اسناد استفاده می شود. این می تواند واژگان بزرگ و داده های متنی با ابعاد بالا را به طور موثر اداره کند.
- **نیاز به داده محدود:** حتی با مقدار کمی از داده های آموزشی می تواند عملکرد مناسبی داشته باشد.

محدودیت های Naive Bayes:

- **فرض مستقل بودن ویژگی:** فرض "ساده لوحانه" استقلال ویژگی ممکن است در همه موارد صادق نباشد. در عمل، اغلب همبستگی بین ویژگی ها وجود دارد.
- **حساسیت به کیفیت داده ها:** Naive Bayes می تواند به کیفیت داده ها حساس باشد. ممکن است با ویژگی های نادر یا دیده نشده در داده های آزمایشی مشکل داشته باشد.
- **فقدان تفسیرپذیری مدل:** مانند سایر مدل های احتمالی، Naive Bayes یک مدل جعبه سیاه در نظر گرفته می شود، به این معنی که ممکن است بینش قابل تفسیری در مورد روابط بین ویژگی ها ارائه نکند.

به طور خلاصه، Naive Bayes یک الگوریتم ساده و در عین حال قدرتمند است که به ویژه برای کارهای طبقه بندی متن مناسب است. کارایی، سرعت و توانایی آن در مدیریت داده های با ابعاد بالا، آن

را به یک انتخاب محبوب تبدیل کرده است، به ویژه در برنامه های **NLP** مانند فیلتر کردن هرزنامه، تجزیه و تحلیل احساسات و دسته بندی اسناد. برای این مدل هم مطابق مدل های قبلی عمل میکنیم.

```
# Naive Bayes model
# Create a Naive Bayes model
nb_model = MultinomialNB()

# Train the model
nb_model.fit(X_train_vectorized, y_train)

# Make predictions on the test set
y_pred_nb = nb_model.predict(X_test_vectorized)

# Calculate the accuracy
report_nb = classification_report(y_test, y_pred_nb)

# Save the model to a file
filename_nb = 'naive_bayes_model.sav'
joblib.dump(nb_model, filename_nb)
```

۴ ارزیابی روش

در بخش‌های قبلی ما در قطعه کد های خود از معیار دقت استفاده کردیم که قطعه کد آن به شرح زیر است:

```
accuracy_xgb = accuracy_score(y_test, y_pred_xgb)
```

اما علاوه بر معیار دقت که تنها به بررسی تعداد حدس های درست بر غلط استوار است معیارهای دیگری نیز برای بررسی وجود دارد

در ادامه به بررسی چندی از روش‌های ارزیابی میپردازیم:

۴.۱ accuracy

دقت پیش بینی های مثبت را اندازه گیری می کند. این نسبت پیش‌بینی‌های مثبت واقعی به تعداد کل پیش‌بینی‌های مثبت (مثبت واقعی + مثبت کاذب) است. دقت زمانی ارزشمند است که هزینه مثبت کاذب بالا باشد.

$$\text{accuracy} = \text{TP} / (\text{TP} + \text{FP})$$

۴.۲ Recall

یادآوری توانایی مدل را در شناسایی صحیح تمام موارد مثبت اندازه گیری می کند. این نسبت پیش بینی های مثبت واقعی به تعداد کل نمونه های مثبت واقعی (مثبت‌های واقعی + منفی‌های کاذب) است. یادآوری زمانی ارزشمند است که از دست دادن موارد مثبت پرهزینه باشد.

$$\text{recall} = \text{TP} / (\text{TP} + \text{FN})$$

F1-Score ۴.۳

F1-Score میانگین هارمونیک دقت و یادآوری است. این یک معیار متعادل از عملکرد یک مدل را ارائه می دهد، به خصوص زمانی که عدم تعادل بین کلاس ها وجود دارد. به صورت زیر محاسبه می شود:

$$F1-Score = 2 * (Precision * Recall) / (Precision + Recall)$$

Specificity ۴.۴

Specificity توانایی مدل را در شناسایی صحیح موارد منفی می سنجد. این نسبت پیش بینی های منفی واقعی به تعداد کل نمونه های منفی واقعی (منفی واقعی + مثبت کاذب) است. زمانی که اجتناب از هشدارهای اشتباه برای موارد منفی بسیار مهم است، ویژگی ارزشمند است.

$$Specificity = TN / (TN + FP)$$

False Positive Rate (FPR) ۴.۵

نسبت پیش بینی های مثبت کاذب به تعداد کل موارد منفی واقعی است. مکمل ویژگی است و به صورت زیر محاسبه می شود:

$$FPR = FP / (TN + FP)$$

False Negative Rate (FNR) ۴.۶

نرخ منفی کاذب (FNR): FNR نسبت پیش بینی های منفی کاذب به تعداد کل نمونه های مثبت واقعی است. نرخی را که مدل در آن موارد مثبت را از دست می دهد اندازه گیری می کند و به صورت زیر محاسبه می شود:

$$FNR = FN / (TP + FN)$$

۴.۷ Area Under the Receiver Operating Characteristic (ROC AUC)

اندازه گیری توانایی مدل در تشخیص نمونه های مثبت و منفی است. این ناحیه زیر منحنی مشخصه عملیاتی گیرنده را کمی می کند، با مقادیر بالاتر که نشان دهنده تبعیض بهتر است.

این معیارها نمای جامعی از عملکرد مدل شما ارائه می دهند، به خصوص در سناریوهایی با مجموعه داده های نامتعادل یا جایی که انواع مختلف خطاهای هزینه های متفاوتی دارند.

برای محاسبه این معیارها در پایتون با استفاده از `scikit-learn`، می توانید از تابع `classification_report` استفاده کنید که خلاصه ای از دقت، یادآوری، امتیاز $F1$ و سایر معیارها را برای هر کلاس در مسئله طبقه بندی شما ارائه می کند.

```
from sklearn.metrics import classification_report

# true_labels are the true class labels from the test set
# predicted_labels are the predicted class labels from your model
report = classification_report(true_labels, predicted_labels)
print(report)
```

به عنوان مثال خروجی به این شکل خواهد بود:

	precision	recall	f1-score	support
Class 0	0.82	0.90	0.86	100
Class 1	0.74	0.65	0.69	50
Class 2	0.91	0.93	0.92	75
accuracy			0.84	225
macro avg	0.82	0.83	0.82	225
weighted avg	0.84	0.84	0.84	225

در این مثال، ما یک مشکل طبقه بندی با سه کلاس داریم: کلاس ۰، کلاس ۱، و کلاس ۲. معنی هر قسمت از گزارش در اینجا آمده است:

- **Precision:** این ستون دقت هر کلاس را نشان می دهد. به عنوان مثال، برای کلاس ۰، دقت ۰.۸۲ است، به این معنی که ۸۲٪ از نمونه های پیش بینی شده به عنوان کلاس ۰ درست بوده اند.
- **Recall:** این ستون فراخوان (یا نرخ مثبت واقعی) را برای هر کلاس نشان می دهد. برای کلاس ۰، فراخوانی ۰.۹۰ است، که نشان می دهد ۹۰٪ از نمونه های کلاس ۰ واقعی به درستی توسط مدل شناسایی شده اند.
- **F1-Score:** میانگین هارمونیک دقت و یادآوری است. تعادلی بین دقت و یادآوری ایجاد می کند. برای کلاس ۰، امتیاز F_1 ۰.۸۶ است.
- **پشتیبانی:** این ستون تعداد نمونه های هر کلاس را در مجموعه داده آزمایشی نشان می دهد.
- **دقت:** در پایین گزارش نشان داده شده است. این دقت کلی مدل است که در این مورد ۰.۸۴ (۸۴٪) است.
- **میانگین ماکرو:** این ردیف میانگین کلان دقت، فراخوانی و امتیاز F_1 را ارائه می کند که در همه کلاس ها میانگین هستند. زمانی مفید است که بخواهید عملکرد مدل را در همه کلاس ها به طور یکسان ارزیابی کنید.
- **میانگین وزنی:** این ردیف با در نظر گرفتن عدم تعادل کلاس، دقت میانگین وزنی، فراخوانی و امتیاز F_1 را ارائه می دهد. زمانی مفید است که کلاس ها اندازه های متفاوتی داشته باشند و بخواهید وزن بیشتری به کلاس های بزرگتر بدهید.

classification_report خلاصه ای مختصر از عملکرد مدل شما برای هر کلاس و به طور کلی است. این یک ابزار ارزشمند برای درک اینکه مدل شما در کجا برتر است و در کجا ممکن است در یک مشکل طبقه بندی چند کلاسه نیاز به بهبود داشته باشد.

۵ نتایج

در این پایان‌نامه و پروژه ما به بررسی مدل‌های مختلفی که برای پردازش زبان طبیعی وجود دارند پرداختیم و انواع الگوریتم‌هایی که میتوان به کمک آن‌ها مدلی بسازیم که نتایج را برای ما پیش‌بینی کنند بررسی کردیم. همچنین تعدادی از آن‌ها همراه با معایب و مزایای هرکدام مورد بحث و بررسی قرار گرفتند.

در انتها میتوان گفت مدل‌های یادگیری ماشین کارایی مناسبی برای پردازش زبان طبیعی دارند و تا زمانی که تعداد دسته بندی‌های ما زیاد نباشد عمل‌کرد خوبی به جا میگذارند اما با افزایش لیبِل‌ها استفاده از الگوریتم‌های عصبی و عمیق خروجی بهتری برای ما خواهند داشت.

از طرفی خروجی ما وابسته به تعداد کلمات موجود در هر توییت است که با بیشتر شدن تعداد کلمات و **vectorize** بهتر میتوان خروجی بهتری گرفت.

در انتها میتوان گفت پارامترهای زیادی در خروجی ما تأثیر گذارند اما در صورتی که ما بتوانیم هر کلمه را به عددی بامعنا مپ کنیم و تبدیل کردن کلمات به اعداد بار احساسی آن‌ها را نیز منتقل کند خروجی ما به مراتب بهتر میشود.

همچنین، تکنیک‌های پیش پردازش و استخراج ویژگی تأثیر قابل توجهی بر عملکرد رویکردهای مختلف تحلیل احساسات و عواطف دارند.

۶ جمع‌بندی و کارهای آینده

کارهایی که در زمینه مپ کردن کلمات به اعداد بامعنا صورت بگیرند به خروجی بهتری منجر می‌شوند پس در آینده به روش‌های مختلف برای پیش پردازش داده‌ها بیشتر پرداخته خواهد شد. همچنین الگوریتم‌های مبتنی بر شبکه‌های عصبی و شبکه‌های عمیق هم مورد بحث و بررسی قرار خواهند گرفت. از سوی دیگر استفاده از سخت‌افزارهای متفاوت برای بهتر شدن زمان یادگیری و ساختن مدل هم می‌تواند بسیاری از مشکلات ما را برای پیاده‌سازی الگوریتم‌های پیچیده از نظر زمانی حل کند.

7 مراجع

- Al Amrani Y, Lazaar M, El Kadiri KE (2018) Random forest and support vector machine based hybrid [1]
approach to sentiment analysis. *Procedia Comput Sci* 127:511–520
- Alqaryouti O, Siyam N, Monem AA, Shaalan K (2020) Aspect-based sentiment analysis using smart government [2]
review data. *Appl Comput Inf.* <https://doi.org/10.1016/j.aci.2019.11.003>
- Alswaidan N, Menai MEB (2020) A survey of state-of-the-art approaches for emotion recognition in text. [3]
Knowl Inf Syst 62(8):1–51
- Archana Rao PN, Baglodi K (2017) Role of sentiment analysis in education sector in the era of big data: a [4]
survey. *Int J Latest Trends Eng Technol* 22–24
- Arora M, Kansal V (2019) Character level embedding with deep convolutional neural network for text [5]
normalization of unstructured data for twitter sentiment analysis. *Soc Netw Anal Min* 9(1):1–14
- Arulmurugan R, Sabarmathi K, Anandakumar H (2019) Classification of sentence level sentiment analysis [6]
using cloud machine learning techniques. *Cluster Comput* 22(1):1199–1209
- Asghar MZ, Subhan F, Imran M, Kundi FM, Shamshirband S, Mosavi A, Csiba P, Várkonyi-Kóczy AR (2019) [7]
Performance evaluation of supervised machine learning techniques for efficient detection of emotions from online
content. *arXiv preprint* [arXiv:1908.01587](https://arxiv.org/abs/1908.01587)
- Bakker I, Van Der Voordt T, Vink P, De Boon J (2014) Pleasure, arousal, dominance: Mehrabian and [8]
Russell revisited. *Curr Psychol* 33(3):405–421
- Balahur A, Turchi M (2014) Comparative experiments using supervised learning and machine translation [9]
for multilingual sentiment analysis. *Comput Speech Lang* 28(1):56–75
- Bandhakavi A, Wiratunga N, Padmanabhan D, Massie S (2017) Lexi-con based feature extraction for [10]
emotion text classification. *Pattern Recogn Lett* 93:133–142
- Basiri ME, Abdar M, Cifci MA, Nemati S, Acharya UR (2020) A novel method for sentiment classification [11]
of drug reviews using fusion of deep and machine learning techniques. *Knowl Based Syst* 198:105949
- Batbaatar E, Li M, Ryu KH (2019) Semantic-emotion neural network for emotion recognition from text. [12]
IEEE Access 7:111866–111878
- Becker K, Moreira VP, dos Santos AG (2017) Multilingual emotion classification using supervised learning: [13]
comparative experiments. *Inf Process Manag* 53(3):684–704
- Bernabé-Moreno J, Tejeda-Lorente A, Herce-Zelaya J, Porcel C, Herrera-Viedma E (2020) A context-aware [14]
embeddings supported method to extract a fuzzy sentiment polarity dictionary. *Knowl-Based Syst* 190:105236
- Bhardwaj A, Narayan Y, Dutta M et al (2015) Sentiment analysis for Indian stock market prediction using [15]
senex and nifty. *Procedia Comput Sci* 70:85–91
- Bhaskar J, Sruthi K, Nedungadi P (2015) Hybrid approach for emotion classification of audio conversation [16]
based on text and speech mining. *Procedia Comput Sci* 46:635–643
- Braun N, van der Lee C, Gatti L, Goudbeek M, Krahmer E (2021) Memofc: introducing the multilingual [17]
emotional football corpus. *Lang Resour Eval* 55(2):389–430
- Bučar J, Žnidaršič M, Povh J (2018) Annotated news corpora and a lexicon for sentiment analysis in [18]
Slovene. *Lang Resour Eval* 52(3):895–919
- Buechel S, Hahn U (2017) Emobank: Studying the impact of annotation perspective and representation [19]
format on dimensional emotion analysis. In: *Proceedings of the 15th conference of the european chapter of the
association for computational linguistics: volume 2, Short Papers*, pp 578–585

Chaffar S, Inkpen D (2011) Using a heterogeneous dataset for emotion analysis in text. In: Butz C, Lingras P [20]
(eds) Advances in artificial