

مجمع آموزش عالی فنی و مهندسی اسفراین
دانشکده مهندسی برق ، کامپیوتر ، صنایع

متن کاوی در متون ادبیات فارسی

پایان نامه کارشناسی رشته مهندسی کامپیوتر (گرایش نرم افزار)

نگارنده

علی زاهدی سهرورزانی

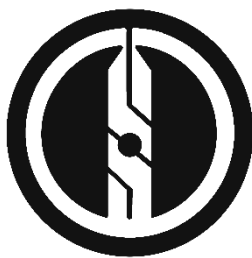
استاد راهنما

دکتر حجت اله موید

بهمن ۱۴۰۲

به نام خداوند جان و خرد

کزین برتر اندیشه برنگذرد



مجمع آموزش عالی فنی و مهندسی اسفراین
دانشکده مهندسی برق ، کامپیوتر ، صنایع

پایان نامه ی کارشناسی رشته ی مهندسی کامپیوتر آقای علی زاهدی سهر و فروزانی
تحت عنوان:

متن کاوی در متون ادبیات فارسی

در تاریخ توسط کمیته ی تخصصی زیر مورد بررسی و تایید نهایی قرار گرفت.

۱- استاد راهنما: دکتر حجت اله موید امضاء.....

۲- مدیریت گروه مهندسی کامپیوتر: دکتر سمیه گرایلو امضاء.....



قدردانی

با تشکر از استاد راهنمای گرامی، جناب دکتر حجت اله موید، بابت توصیه‌های علمی و تکنیکی خود در طول انجام پروژه به من کمک فراوان کردند و بدون حمایت ایشان، این پروژه به سختی به اتمام می‌رسید.

همچنین، من از دانشکده مهندسی برق کامپیوتر و صنایع، مجتمع آموزش عالی فنی و مهندسی اسفرا این قدردانی می‌کنم که با فراهم کردن امکانات لازم و حمایت‌های علمی، به من در انجام این پروژه کمک کردند.

با سپاس فراوان



کلیه ی حقوق مادی و معنوی مربوط به نتایج مطالعات، ابتکارات و نوآوری های ناشی از

تحقیق موضوع این پایان نامه متعلق به "مجمع آموزش عالی فنی و مهندسی اسفراين" است



تقدیم می کنم به خانواده عزیزم که همواره از حمایت من دست نکشیدند



فهرست مطالب

فهرست مطالب	ث
فهرست جداول	خ
فهرست اشکال	د
چکیده	ذ
فصل ۱ مقدمه	۱
۱-۱ پیشنهاد تحقیق و اهمیت موضوع	۱
۲-۱ هدف از انجام پایان نامه	۲
۳-۱ روش تحقیق	۳
۴-۱ ساختار پایان نامه	۴
فصل ۲ مروری بر تحقیقات	۵
۱-۲ مقدمه	۵
۲-۲ چالش ها	۶
۳-۲ اهمیت موضوع	۷
۴-۲ روش شناسی	۷
۵-۲ جوانب پروژه	۸
فصل ۳ مفاهیم پایه	۱۲
۱-۳ متن کاوی	۱۲
۲-۳ پیش پردازش	۱۳
۳-۳ کلمات توقف	۱۴
۴-۳ پاک سازی دیتاست	۱۵
۵-۳ بازیابی اطلاعات	۱۶
۶-۳ مدل های بازیابی اطلاعات	۱۷
۷-۳ مفاهیم پایه بازیابی اطلاعات	۱۸
۸-۳ روش ها و مراحل بازیابی اطلاعات	۱۹
۹-۳ کاربردها و دامنه های استفاده	۱۹



۳-۱۰ برخی از تکنیک‌های بازیابی اطلاعات ۲۰

۳-۱۱ معیارهای امتیازدهی ۲۲

۳-۱۱-۱ تی اف ۲۲

۳-۱۱-۱-۱ رابطه ریاضی ۲۳

۳-۱۱-۲ دی اف ۲۳

۳-۱۱-۲-۱ رابطه ریاضی ۲۴

۳-۱۱-۳ مقایسه DF و TF ۲۴

۳-۱۱-۴ ای دی اف ۲۵

۳-۱۱-۴-۱ رابطه ریاضی ۲۶

۳-۱۱-۴-۲ مقایسه با معیارهای قبلی ۲۶

۳-۱۱-۵ تی اف ای دی اف ۲۶

۳-۱۱-۵-۱ رابطه ریاضی ۲۸

۳-۱۱-۵-۲ مقایسه با معیارهای قبلی ۲۹

فصل ۴ مطالعات و تحقیقات ۳۰

۴-۱ روند کلی ۳۰

۴-۲ پیش‌پردازش ۳۱

۴-۳ رویکرد ما در عملیات بایس کردن استاپ وردها ۳۳

۴-۴ محاسبه TF-IDF ۳۳

۴-۵ اعلان متغیر ۳۴

۴-۶ مدل ذخیره‌سازی TF-IDF ۳۴

۴-۷ تست داده‌های جدید ۳۶

۴-۸ روش‌های تست ورودی ۳۶

۴-۹ نحوه تشخیص احساس از TF-IDF ۳۷

۴-۱۰ دقت مدل ۳۷

فصل ۵ نتیجه‌گیری و پیشنهادات ۳۹



۳۹ ۱-۵ نتیجه گیری

۴۳ ۲-۵ پیشنهادات

۴۶ مراجع



فهرست جداول

جدول ۴-۱ داده‌های اصلی.....	۳۱
جدول ۴-۲ فایل‌های پیش پردازش.....	۳۲
جدول ۴-۳ کلمه بر حسب TF-IDF برای هر شش احساس.....	۳۵
جدول ۴-۴ نمونه از داده های تست.....	۳۷
جدول ۴-۵ کد نمونه‌ها.....	۳۸



فهرست اشکال

شکل ۱-۲ فراوانی سند احساسی..... ۸



چکیده

با توجه به عنوان مسئله و ارتباط مستقیم موضوع با بحث بازیابی اطلاعات صورت مسئله یافتن الگوهای زبانی (ترکیبی از کلمات) در یک سند متنی به زبان فارسی است بدین شکل که با دریافت چندین سند متنی از ورودی کاربر پیرامون یک موضوع خاص به دنبال الگوها می گردد، این تشخیص الگو به وسیله ابزارهای آماری در بازیابی اطلاعات انجام می شود نظیر فراوانی کلمه در سند و مجموعه و دیگر ابزار که در ادامه به آنها اشاره خواهد شد.

پس از تشخیص الگو با در نظر گرفتن روش بهینه، سند دیگری دریافت و اقدام به واکاوی متون از آن می کند که انتظار می رود همان الگویی باشد که در اسناد اولیه پردازش شده یافته بودیم.

اهمیت موضوع از آنجا نمایان می شود که تا کنون پژوهش مشابهی پیرامون آنچه توضیح داده شد در زبان فارسی صورت نگرفته است. از همین رو اهمیت پژوهش دو چندان خواهد بود و منابعی که می توان به آن مراجعه کرد همگی جهت استفاده در زبان انگلیسی و دیگر زبان های روز دنیا خواهند بود که با به پایان رسیدن این پروژه متن کاوی در متون ادبیات فارسی به شکل کاملاً عملی و علمی صورت پذیر خواهد بود.

رویکرد ما مبتنی بر قرار دادن فرایندها و ابزارها جهت رسیدن به نتایج، جمع آوری و تحلیل و تفسیر می باشد. در این مسیر از روش های کمی و کیفی، ترکیبی، تجربی، توصیفی، تحلیل استفاده می شود.

جامعه هدف ما ترکیبی از ادبیات معاصر با در نظر داشتن حالات و احساسات و رفتارها می باشد، نحوه انتخاب نمونه بصورت تصادفی و هدفمند از مجموعه داده ای تحت عنوان "مجموعه ۳۰ هزار جمله ای" که هر جمله از آن را با در نظر گرفتن ۶ احساس اصلی امتیاز بندی شده است استخراج شده از پلتفرم توییتر می باشد.

واژگان کلیدی: بازیابی اطلاعات، فراوانی کلمات، متن کاوی



فصل ۱

مقدمه

۱-۱ پیشینه تحقیق و اهمیت موضوع

متن کاوی در ادبیات فارسی به عنوان یک زمینه تحقیقی مهم در حوزه مطالعات ادبیات و پردازش زبان طی دهه‌های اخیر به شدت توجه پژوهشگران جلب کرده است. ادبیات فارسی با بی‌شماری آثار و متونی از دوران کلاسیک تا ادب معاصر، از غنای فراوانی از احساسات و نگرانی‌ها نمایانگر است. با این حال، تشخیص و تحلیل احساسات موجود در این متون برای محققان همچون ادب‌شناسان، متن‌شناسان، و همچنین برنامه‌نویسان کار سختی است.

تحلیل احساسات و نوع جمله در متون ادبیات فارسی اهمیت زیادی در درک عمیق‌تر این متون دارد. ادبیات به عنوان یک نمایه از فرهنگ و تاریخ یک جامعه، شامل انواع مختلفی از احساسات و جریانات انسانی است. توانایی تشخیص کمیته‌های احساسی و نوع جمله در متن‌ها، به ادب‌شناسان امکان می‌دهد تا درک بهتری از نوع و مضمون اثرات ادبی داشته باشند.



برای مطالعه تطبیقی ادبیات، تحلیل احساسات و نوع جمله در متون می‌تواند به مقایسه متن‌ها از منظر احساساتی و ادبی کمک کند. این تحلیل می‌تواند به پژوهشگران کمک کند تا تأثیر تاریخ، فرهنگ و محیط اجتماعی بر متون ادبی را بهتر درک کنند.

در کاربردهای کنونی متن‌کاوی و پردازش زبان طیف گسترده‌ای از امکانات فراهم کرده است. از مدل‌های یادگیری عمیق برای تشخیص احساسات در متن تا تحلیل خودکار محتوای ادبی، این پروژه می‌تواند از مزایای پیشرفت‌های فناوری در زمینه پردازش زبان برای تحلیل متون ادبی استفاده کند.

با تأکید بر اهمیت تحلیل احساسات و نوع جمله در متون ادبیات فارسی و توانایی‌های پیشرفته پردازش زبان، این پروژه می‌تواند به افزایش فهم عمیق‌تر از متون ادبی و کاربردهای علمی و فرهنگی آنها در جامعه کمک منجر شود.

۱-۲ هدف از انجام پایان نامه

هدف اصلی این پایان‌نامه ارتقای فهم عمیق‌تر و دقیق‌تر متون ادبیات فارسی با بهره‌گیری از تحلیل احساسات و نوع جمله است. این پروژه با توجه به غنای فراوان از متون ادبی و انواع مختلفی از احساسات و جریانات انسانی در این متون، تلاش می‌کند تا توانایی تشخیص و تحلیل احساسات مختلف موجود در این متون را بهبود بخشد.

از اهمیت‌ها و دلایل می‌توان به موارد زیر اشاره کرد: فهم عمیق‌تر: تحلیل احساسات متون ادبیات فارسی به ما امکان می‌دهد تا فهم عمیق‌تری از طبیعت و مضمون این متون پیدا کنیم. متون ادبیاتی به‌تنهایی تنها کلمات نیستند؛ آنها حاوی احساسات، افکار و تجربیات انسانی عمیقی هستند که با تحلیل دقیق‌تر احساسات متن، می‌توانیم به فهم عمیق‌تری از آثار ادبی دست یابیم.

تحلیل متون مقایسه‌ای: با تحلیل احساسات و نوع جمله در متون ادبیات فارسی، می‌توانیم به تحلیل متون مقایسه‌ای پردازیم و اثرات تاریخ، فرهنگ و محیط اجتماعی را در متون مقایسه کنیم. این می‌تواند به ما کمک کند تا به مفاهیم عمیق‌تری از تطبیق متون ادبی دست یابیم.



کاربردهای فناوری: در دنیای امروزی، فناوری‌های پردازش زبان و متن کاوی پیشرفت‌های چشمگیری داشته‌اند.

این پروژه از این تکنولوژی‌ها بهره‌مند شده و به تحلیل خودکار متون ادبی به راحتی و دقت کمک می‌کند.

این پایان‌نامه امیدوار است که به پیشرفت‌های درک متون ادبی فارسی و کاربردهای علمی و فرهنگی آنها کمک

کند و در راستای بهترین شناخت و تجزیه و تحلیل متون ادبی به کمک محققان و علاقه‌مندان در این حوزه باشد.

۳-۱ روش تحقیق

تشخیص احساسات از متون یک مسئله مهم در حوزه پردازش زبان طبیعی و متن کاوی است. با توجه به افزایش

روزافزون متون موجود در شبکه‌های اجتماعی و پلتفرم‌های آنلاین، توانایی تشخیص احساسات موجود در متون، به ما

امکان می‌دهد تا علاوه بر درک مفاهیم متن، نگاهی به جریانات احساسی مرتبط با آن داشته باشیم.

ما از یک مجموعه داده‌ای استفاده می‌کنیم که شامل متون از شبکه‌های اجتماعی به زبان فارسی با برچسب‌های

مربوط به شش احساس مختلف (خشم، ترس، خوشحالی، نفرت، غم و تعجب) است. ما از این داده‌ها بهره‌مند

می‌شویم تا با استفاده از روش TF-IDF^۱ وزن دار وزن‌دهی به کلمات مربوط به هر یک از این احساسات را انجام

دهیم [1].

روش TF-IDF یک روش معمول برای بررسی وزن کلمات در یک متن است. با اعمال این روش به داده‌های

ما، ما می‌توانیم وزن‌دهی به کلمات مربوط به هر احساس از آن‌ها انجام دهیم. این کار امکان تفکیک کلماتی که

بیشترین ارتباط با هر احساس دارند را فراهم می‌کند.

^۱ Term Frequency-Inverse Document Frequency



هدف اصلی این است که بتوانیم با استفاده از این معیار، کلمات مربوط به هر احساس را به دقت تشخیص داده و از آنها برای تحلیل احساسات در متون مختلف استفاده کنیم. این اطلاعات می تواند به ما کمک کند تا درک بهتری از تغییرات احساساتی در متون داشته باشیم و نهایتاً به تفهیم بهتری از مضمون و مفاهیم متون ادبی دست یابیم.

۴-۱ ساختار پایان نامه

ابتدا با معرفی موضوع و اهمیت تشخیص احساسات از متون به عنوان مسئله اصلی، به معرفی داده های مورد استفاده پرداختیم. سپس به مرور ادبیات در زمینه تشخیص احساسات و تفسیر مفاهیم مشابه پرداختیم. در مرحله بعد، متدولوژی تحقیق شرح داده شد و مراحل پیش پردازش متن، تعیین وزن به کلمات با استفاده از TF-IDF وزن دار و اجرای آزمایش ها تشریح گردید. پس از اجرای آزمایش ها، نتایج به دست آمده مورد تجزیه و تحلیل قرار گرفتند و با تحلیل مفصل به نتیجه گیری رسیدیم. در نهایت، توصیه ها و پیشنهادات برای تحقیقات آینده در این زمینه ارائه گردید.



فصل ۲

مروری بر تحقیقات

۱-۲ مقدمه

مروری بر اکثر کارهای گذشته نشان می‌دهد که روش‌های مورد استفاده برای سیستم‌های تحلیل احساسات را

می‌توان به گروه‌های زیر تقسیم کرد:

- رویکردهای کلیدواژه
- رویکردهای یادگیری ماشینی
- رویکردهای ترکیبی

اولین رویکرد محققان روی شناسایی و استفاده از کلمات کلیدی که مستقیماً احساسات بیان شده در متن را

منعکس می‌کنند، کار کرده‌اند. برای شناسایی این کلمات کلیدی از چندین روش استفاده شده است. دقت این

روش‌ها به وجود کلمات کلیدی بستگی دارد. یکی از معایب این روش وابستگی آن به مناطق خاص است [2].

عموماً در این رویکرد فهرستی از کلمات با برچسب‌های معنایی عواطف مثبت و منفی یا احساس غم، شادی و...

تهیه می‌شود.

گاهی اوقات فهرست با دادن نکاتی ارائه می شود که نشان دهنده شدت یک احساس است. سپس متن به کلمات جداگانه تقسیم می شود. در مرحله بعد، کلماتی که احساسات را پوشش می دهند مشخص می شود. سپس شدت احساسات اندازه گیری می شود.

در مرحله بعد بحث می شود که آیا فعل منفی است یا خیر. فعل منفی بودن احساس تعیین شده را تغییر می دهد. برخلاف رویکردهای مبتنی بر کلیدواژه، رویکردهای مبتنی بر یادگیری ماشینی، با استفاده از سیستم های طبقه بندی از پیش آموزش دیده، سعی در شناسایی احساسات بیان شده در متن دارند. عدم وابستگی به حوزه های خاص منجر به پیشرفت سریع روش های یادگیری ماشینی تحت نظارت و بدون نظارت شده است.

۲-۲ چالش ها

به دلیل کمبود منابع احساسی برچسب گذاری شده در دسترس، صفحات تویتر و تگ ها مورد استفاده کاربران به عنوان منابع پردازشی موجود در بسیاری از مطالعات مورد استفاده قرار گرفته اند.

این مجموعه داده ها می توانند به عنوان مجموعه داده های فعال برای پردازش استفاده شوند. آنچه در بیشتر تحقیقات رایج است استفاده از مجموعه داده ها، ویژگی ها و روش های طبقه بندی مختلف به منظور دستیابی به دقت بهتر است.

به عنوان مثال، در یک کار مشابه، دقت ۷۳ درصد با استفاده از الگوریتم ماشین بردار پشتیبان به دست آمده است که از ویژگی هایی مانند کلیدواژه ها معنایی برای روابط معنایی کلیدواژه ها استفاده شده است.



۲-۳ اهمیت موضوع

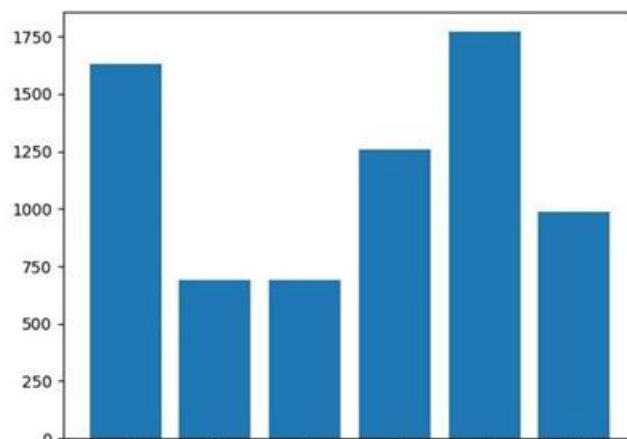
تشخیص احساسات، متن را به یک حوزه مهم برای مطالعه تبدیل می کند. با توجه به مشکلات موجود، این کار بسیار چالش برانگیز به نظر می رسد. به عنوان مثال، تشخیص احساسات واقعی انسان گاهی اوقات دشوار است به علت اینکه گاهی اوقات افراد از کلماتی استفاده می کنند که احساس واقعی شان را نشان نمی دهد.

نکته دیگر اینکه سبک نوشتن برای هر فرد یکسان نیست و تشخیص مستقیم احساسات واقعی آن ها در سبک های رسمی و غیررسمی دشوار است. این موارد نمونه هایی هستند که تشخیص احساسات را آن طور که نگاه اول به نظر می رسد آسان نمی کنند.

۲-۴ روش شناسی

تعیین نوع احساس بیان شده در متن روش پیشنهادی است. با استفاده از مجموعه توییت های رتبه بندی شده بر اساس عاطفه، ترس، اندوه، شادی، خشم و حیرت است. در ابتدا، یکی از مشکلاتی که به وجود می آید تحلیل احساسات متون تولید شده توسط کاربران در رسانه هایی مانند توییتر است. یکی از وظایف اصلی تحلیل احساسات طبقه بندی کلمات در دسته بندی احساسات می باشد.

ارزیابی جملات در منظرهای مختلف برای دریافت احساسات متعددی که یک متن واحد به وجود می آیند، تازگی این پروژه است. در تصویر زیر فراوانی جمع کل امتیاز احساسات را به تفکیک مشاهده می کنید.



شکل ۱-۲ فراوانی سند احساسی

۲-۵ جوانب پروژه

استفاده از مجموعه داده‌های مناسب یکی از عواملی است که بر دقت اندازه‌گیری سیستم‌های پردازش زبان طبیعی تأثیر زیادی دارد. برای این منظور یکی از اولین گام‌های برداشته شده، انتخاب مجموعه مناسب از داده‌های مرتبط با هدف انتخاب شده است. به عنوان مثال در متون علمی، همان‌طور که از نام آن‌ها پیداست، استفاده از جملات با ساختار عاطفی عجیب به نظر می‌رسد. همچنین استفاده از جملات احساسی در متون روایی کاملاً طبیعی به نظر می‌رسد. از سوی دیگر اگر چه منابع مکتوب زیادی وجود دارد، اما منابع احساسی برای پردازش زبان طبیعی به‌ویژه برای زبان فارسی بسیار کم است. در تحقیقات قبلی که در مورد زبان فارسی انجام شده است معمولاً توسط نظرات کاربران در مورد یک محصول خاصی اخبار روزانه انجام شده است. در واقع هیچ مجموعه داده کافی و جامعی برای ارزیابی نتایج سیستم‌ها پیشنهادی وجود ندارد. این امر منجر به اولین گام در ارائه داده‌های احساسی از مجموعه توییت‌های فارسی شد که یکی از محبوب‌ترین مجموعه‌ها در زبان فارسی است. از میان منابع فارسی موجود، مجموعه توییت‌های فارسی به دلیل تنوع در مفاد اسناد، برای برچسب‌گذاری احساسی انتخاب شد. مجموعه توییت‌های فارسی مجموعه‌ای از متون فارسی است که به صورت حرفه‌ای با برچسب‌های دسته‌بندی نحوی برچسب‌گذاری شده‌اند. این



مجموعه شامل شش برچسب احساسی است. این مجموعه توسط زبان‌شناسان به صورت دستی ساخته شده و تهیه آن زمان بسیار طول کشیده.

باتوجه به کلماتی که افراد برای بیان احساسات خود به طور مستقیم و غیرمستقیم استفاده می‌کنند، مفاهیم عاطفی به دو اصطلاح تحت‌اللفظی و مجازی تقسیم می‌شوند. اگر معنای آن در نگاه اول آشکار نباشد و معنای آن ضمنی باشد این بنا را اصطلاح مجازی می‌نامند.

بحث انتزاع بودن احساس با گفتار مجازی بیان می‌شود. یکی از ویژگی‌های انتزاعی گفتار مجازی، استفاده مکرر از استعاره است. مثلاً احساس خشم در جمله از عصبانیت منفجر شد منفجر شدن یک انتزاع از عصبانیت در نظر گرفته شده است. استعاره‌های مفهومی این دو حوزه را از نظر برخی ویژگی‌های مشترکی که دارند به هم مرتبط می‌کنند. یک ناحیه معمولاً عینی‌تر و قابل لمس‌تر از دیگری است. علاوه بر استعاره‌های عاطفی، از مفاهیم کنایه نیز برای بیان احساسات استفاده می‌شود. در این حالت بر خلاف استعاره‌ها به دو حوزه مربوط نمی‌شود و فقط یک حوزه را نشان می‌دهد. هدف آن فراهم کردن دسترسی ذهنی به یک دامنه خاص از طریق بخشی از همان دامنه است یا سعی می‌کند با استفاده از بخشی از یک دامنه خاص به بخش دیگری از همان دامنه ارتباط برقرار کند. به عنوان مثال، عصبانیت یکی از ویژگی‌های افزایش دمای بدن است، بنابراین زمانی که فرد قرمز می‌شود، نشان‌دهنده ساختار چهره وضعیت فیزیکی فرد است.

باتوجه به آنچه در علوم شناختی در مورد عاطفه و چگونگی بازنمایی احساس توسط زبان شناخته شده است، با فرض وجود رابطه بین عاطفه و زبان نفرت، می‌توان دیدگاه‌های متعددی را مورد بررسی قرارداد.

دیدگاه اول آن است که کلمات می‌توانند احساسات درونی افراد را نشان دهند.

دیدگاه دوم نحوه بیان اصطلاحات عاطفی در سطوح آوایی، واژگانی، نحوی، معنایی و بافتی زبان است.



دیدگاه سوم آن است که سازه‌های احساسات شخصی را نشان می‌دهند؛ بنابراین ویژگی‌های مورد استفاده در

سیستم پیشنهادی به شرح زیر است:

- کلمات کلیدی احساسی
- مقولات نحوی عاطفی
- ترکیب‌های احساسی

هنگام استفاده از زبان، انسان‌ها سعی می‌کنند محتوای ذهنی پیچیده خود را در قالب عینی تر و قابل دسترس بیان

کنند. اگر بین عاطفه و زبان رابطه مستقیم یا غیرمستقیم وجود دارد، این رابطه چگونه نمایش داده می‌شود؟ این رابطه

را به خوبی می‌توان به این صورت خلاصه کرد:

کلماتی مانند "او"، "هورا" و "وای" در متن به ترتیب نشان‌دهنده غم، شادی و تعجب هستند.

این واژه‌ها را واژه‌های بیانی می‌گویند. این کلمات بیانگر معنای احساسات در درون خود هستند.

کلمات کلیدی و احساسی افراد بسته به احساسی که دارند، در موقعیت‌های مختلف از کلمات متفاوتی استفاده

می‌کنند. به عنوان مثال، کلماتی که افراد هنگام عصبانیت استفاده می‌کنند با کلماتی که در هنگام خوشحالی استفاده

می‌کنند متفاوت است.

در این قسمت واژه‌های بیانی و تحت‌اللفظی مورد تجزیه و تحلیل قرار می‌گیرد؛ بنابراین فهرستی از کلمات

کلیدی معنایی مرتبط با هر احساس می‌توان یافت.

بر اساس منطق انسان، واژه‌هایی که بیشترین ارتباط را با متون احساسی داشته و فراوانی بالایی در متون دارند را

مرتبط می‌دانیم. سپس کلیدواژه‌های تهیه شده مورد بررسی قرار گرفت که به کدام دسته از احساسات تعلق دارند.

برخی از برترین کلمات کلیدی مرتبط با احساسات در زبان فارسی هستند.



این مجموعه کلماتی که احساسات را توصیف می کنند، کلمات توصیفی نامیده می شوند. کلماتی مانند "خشم"، "شادی"، "شگفتی" کلمات توصیفی هستند. در بین کلمات توصیفی، همه کلمات به یک اندازه اهمیت ندارند و برخی از آنها نمونه بهتری هستند و نمونه اولیه احساسات به حساب می آیند. این کلمات به عنوان کلمات اصلی شناخته می شوند و گروهی دیگر واژه های غیر پایه نامیده می شوند.

در این کار پژوهشی سعی شده است از ویژگی های متن کاوی برای دستیابی به دقت بالا استفاده شود. با توجه به ساختارهای احساسی و کلیدواژه ها در متون فارسی، جنبه بدیع این پروژه است. امید است شناخت کلیدواژه های عاطفی کارایی مناسبی کسب کند.

با توجه موارد فوق، هدف این کار ارائه سیستمی برای شناخت عاطفه بر اساس ویژگی های متن کاوی بود تا بتواند عملکردی مشابه عملکرد ذهن انسان داشته باشد. همچنین نشان می دهد که در عصر مدرن، استفاده از علم آمار در مطالعه قسمت های زیرین زبان بسیار کمک کننده است.

ساختار بقیه مقاله به شرح زیر است:

در فصل بعدی به بررسی معیارهای متن کاوی و بازیابی اطلاعات می پردازیم.

در فصل آخر، سیستم پیشنهادی با توجه به ویژگی های آن شرح داده شده است. در این بخش به بررسی ویژگی های تعریف شده، نحوه تهیه و استفاده از آنها می پردازد و به میزان تأثیر و ویژگی های تعریف شده بر شناسایی احساسات مختلف می پردازد.

توجه به کلیدواژه ها از جمله مواردی بود که در این کار بررسی شد.



فصل ۳

مفاهیم پایه

۳-۱ متن کاوی

متن کاوی یک فرایند تحلیلی است که با استفاده از الگوریتم‌ها و تکنیک‌های مختلف، اطلاعات مخفی یا الگوهای پنهان در داده‌ها را شناسایی می‌کند. این فرایند، معمولاً به منظور استخراج اطلاعات مفید و نهایی از داده‌های حجیم و پیچیده انجام می‌شود.

در مرحله ابتدایی متن کاوی، داده‌های موردنظر از منابع مختلف جمع‌آوری می‌شوند. سپس با استفاده از الگوریتم‌های مختلف، داده‌ها تحلیل و بررسی می‌شوند تا الگوها، روابط، و اطلاعات مهم شناسایی شوند. این اطلاعات ممکن است از نظرات مشتریان، رفتارهای بازار، یا دیگر جنبه‌های مهم کسب‌وکار باشند.

یکی از روش‌های متداول در متن کاوی، تحلیل متنوعی از داده‌های متنی است. با بهره‌گیری از تکنیک‌های پردازش زبان طبیعی، متن‌های بلند مانند نظرات مشتریان، مقالات، یا ارسال‌های رسانه‌های اجتماعی مورد تجزیه و تحلیل قرار می‌گیرند. این تجزیه و تحلیل متن می‌تواند به شناخت احساسات، موضوعات محبوب، یا حتی شناسایی احتمالی مشکلات و نقاط قوت در محتواها منجر شود [3].



در نهایت، نتایج به دست آمده از متن کاوی، به تصمیم گیری های مؤثر در کسب و کار کمک می کند. این اطلاعات می توانند در بهبود تجربه مشتریان، افزایش کارایی عملیات، و حتی شناسایی فرصت های جدید در بازار کمک کنند.

در نتیجه، متن کاوی به عنوان یک ابزار قدرتمند در دنیای اطلاعاتی امروزی، برای سازمان ها و تصمیم گیرندگان ارزش زیادی ایجاد می کند. این فرایند تحلیلی نه تنها امکان شناسایی الگوها و اطلاعات مهم را فراهم می کند؛ بلکه به بهبود استراتژی ها و عملکرد کلان کسب و کار نیز کمک می کند.

۳-۲ پیش پردازش

در پیش پردازش متن چالش های متعددی برای زبان فارسی وجود دارد. یکی از چالش های سیستم های پردازش زبان فارسی، پیچیدگی سیستم نوشتاری آن است. به عنوان مثال، استفاده نامناسب از فاصله به جای نیم فاصله قابل ذکر است. توجه به استفاده از علائم نگارشی فارسی و استفاده از حروف غیرعربی می تواند از پیچیدگی سیستم نوشتاری بکاهد.

دیگر چالش ها شامل کلماتی است که می توان به روش های مختلفی نوشت. شناسایی واژه های مرکب و تک واژه ها در نظر گرفتن آنها از دیگر چالش هایی است که برای نگارش زبان فارسی وجود دارد. به عنوان مثال، پیدا کردن ریشه کلمات برای در نظر گرفتن یک واحد برای کلمات مربوطه است؛ بنابراین به راحتی بر عملکرد سیستم تأثیر می گذارد.

کلمات توقف شامل کلمات دستوری مانند حروف ربط و حروف اضافه است. این کلمات فاقد اطلاعات معنایی و متنی هستند. در مورد مطالعاتی که بر فرکانس کلمه تمرکز دارند، فراوانی بالای کلمات توقف که در اسناد وجود دارد باعث خطا در نتایج پردازش می شود.

۳-۳ کلمات توقف

در پردازش متن، کلمات توقف به کلماتی اطلاق می‌شود که در تحلیل و استخراج اطلاعات از متن، ارزش زیادی ندارند و بیشتر برای ارتقای کارایی و کاهش حجم داده‌ها مورد استفاده قرار می‌گیرند. این کلمات به صورت متداول در زبان فارسی نیز وجود دارند و در فرایند پردازش متن به منظور بهبود دقت و سرعت تحلیل، استفاده می‌شوند.

کلمات توقف معمولاً از دسته‌های کلمات متداول و پرتکرار در زبان فارسی هستند. این شامل حروف اضافه، حروف ربط، ضمائر و کلمات پرتکراری مانند "و"، "در"، "به" و "از" می‌شود. این کلمات توقف برخلاف کلمات مهم مانند اسم‌ها یا صفات، به تنهایی اطلاعات زیادی ارائه نمی‌دهند.

از اهمیت کلمات توقف می‌توان به موارد زیر اشاره کرد:

۱. کاهش حجم داده‌ها: حذف کلمات توقف باعث کاهش حجم داده‌ها می‌شود و این امکان را فراهم می‌کند تا در تحلیل داده‌ها سریع‌تر عمل شود.

۲. بهبود دقت تحلیل: با حذف کلمات توقف، تمرکز بر روی کلمات اصلی و معنادارتر متن افزایش می‌یابد که منجر به بهبود دقت تحلیل متن می‌شود.

برای درک بهتر، مثالی از حذف کلمات توقف در یک جمله فارسی را مورد بررسی قرار دهیم: جمله اصلی:

"در یک روز زیبا، آسمان آبی بود و افرازه‌های کوه‌ها با برف پوشیده بودند."

جمله پس از حذف کلمات توقف: "روز زیبا، آسمان آبی، افرازه‌های کوه‌ها برف پوشیده."

در این مثال، کلمات توقف مانند "در"، "و" و "حذف شده‌اند تا متن کلی‌تر و معنادارتری به دست آید.

استفاده از کلمات توقف در پردازش متن فارسی بهبود سرعت و دقت تحلیل را تضمین می‌کند. این اقدام مؤثر در استخراج اطلاعات معنی‌دار از متن‌های حجیم است و در تحلیل دقیق‌تر و کارآمدتر داده‌ها نقش بسزایی ایفا می‌کند.

۳-۴ پاک‌سازی دیتاست

تمیز کردن دیتاست متنی یک گام اساسی در پردازش داده است و می‌تواند تأثیر زیادی در کیفیت و قابلیت استفاده از داده‌ها داشته باشد. در ادامه، روش‌ها و مراحل برای تمیز کردن دیتاست متنی را بررسی می‌کنیم:

۱. حذف سطرهای خالی: حذف ردیف‌هایی که هیچ اطلاعاتی ندارند از جمله اولین گام‌هاست. این سطرها ممکن است اطلاعات ناقص یا بدون ارزشی داشته باشند.

۲. پر کردن مقادیر ناقص: اگر بخشی از داده‌ها ناقص باشد، می‌توانید از روش‌های پر کردن مقادیر ناقص مانند میانگین، مد، یا مقدار پیش‌فرض استفاده کنید. یا می‌توانید در صورت برخورد با چنین نوع رکوردی کل آن رکورد را حذف کنید.

۳. حذف داده‌های تکراری: اگر در دیتاست داده‌های تکراری وجود داشته باشد، حذف آنها می‌تواند به کاهش انحراف و بهبود دقت در تحلیل کمک کند.

۴. تشخیص و حذف داده‌های نویزی: داده‌های نویزی ممکن است تحلیل دقیق را مختل کنند. با شناسایی و حذف داده‌های نویزی می‌توانید کیفیت تحلیل خود را افزایش دهید.

۵. حذف اطلاعات غیرضروری: اگر بخشی از داده‌ها برای هدف خاصی که در دست استفاده می‌شود ارزشی ندارد، آن بخش را حذف کنید.

۶. تبدیل متن به حالت نرمال: برخی از پیش پردازش های متنی مانند تبدیل به حالت نرمال، تبدیل کلمات به حروف کوچک (در نوشتار انگلیسی)، و حذف علائم نگارشی اضافی مانند ضمه، فتحه، کسره و ... می توانند بهبودی در تحلیل متن داشته باشند.

۷. تشخیص و حذف داده های پرتکرار: داده های پرتکرار ممکن است موثریت تحلیل را کاهش دهند. با تشخیص و حذف این داده ها، دقت تحلیل افزایش می یابد.

۸. استفاده از کلمات توقف: حذف کلمات توقف که در توضیح قبلی آمده، به بهبود دقت تحلیل کلمات کمک می کند.

۹. تصحیح املائی: اگر در داده ها املائی اشتباه وجود دارد، از ابزارهای تصحیح املائی استفاده کنید تا دقت اطلاعات افزایش یابد.

۱۰. تبدیل داده های متنی به بردارهای عددی: برای استفاده از داده های متنی در مدل های عددی، می توانید از روش های تبدیل به بردارهای عددی استفاده کنید.

تمیز کردن دیتاست متنی مرحله مهمی در استفاده مؤثر از داده ها است. هر گاه که با دیتاست متنی روبرو می شوید، می توانید از این روش ها و تکنیک ها برای بهبود کیفیت و قابلیت تحلیل داده ها بهره مند شوید.

۳-۵-۵ بازایی اطلاعات

بازایی اطلاعات یک حوزه مهم در علوم کامپیوتر و علم اطلاعات است که به بررسی و تحقیق در مورد روش ها و تکنیک های بازایی اطلاعات اختصاص دارد. در این حوزه، هدف اصلی یافتن اسناد یا اطلاعات مناسب در مجموعه های بزرگ داده ها با استفاده از یک سیستم کامپیوتری است



۳-۶ مدل‌های بازیابی اطلاعات

۱. مدل محتوا^۱: در این روش، بازیابی اطلاعات بر اساس محتوای اسناد صورت می‌گیرد. ویژگی‌های مختلف مانند کلمات کلیدی، عناصر مشخص، یا ویژگی‌های محتوایی دیگر برای تطابق و جستجوی اطلاعات استفاده می‌شوند.

۲. مدل مبتنی بر کاربر^۲: این روش به اطلاعات مربوط به کاربران توجه دارد. سابقه جستجوها، نظرات یا فعالیت‌های کاربر در بازیابی اطلاعات تأثیرگذار هستند. سیستم‌های توصیه نیز معمولاً از این نوع مدل استفاده می‌کنند.

۳. مدل ترکیبی^۳: در این روش، از ترکیب مدل‌های محتوایی و مبتنی بر کاربر برای بهبود دقت و کارایی در بازیابی اطلاعات استفاده می‌شود.

۴. مدل مبتنی بر ساختار^۴: این روش در بازیابی اطلاعات از اطلاعات ساختاری و ارتباطات میان اسناد استفاده می‌کند. به عنوان مثال، اسناد متنی با ساختارهای مختلف مانند وب‌سایت‌ها یا مقالات علمی تفاوت دارند.

۵. مدل مبتنی بر احتمال^۵: این روش از مدل‌های احتمالاتی برای تعیین احتمال اطلاعات مورد جستجو استفاده می‌کند. مدل‌های مانند مدل فضای برداری احتمالاتی^۶ در این دسته قرار می‌گیرند.

۶. مدل مبتنی بر هوش مصنوعی^۷: با پیشرفت هوش مصنوعی، مدل‌های مبتنی بر الگوریتم‌ها و شبکه‌های عصبی عمیق نیز در بازیابی اطلاعات مورد استفاده قرار می‌گیرند.

¹ Content-Based Model

² User-Based Model

³ Combination Model

⁴ Structure-Based Model

⁵ Probabilistic Model

⁶ Probabilistic Vector Space Model

⁷ AI-Based Model

۷. مدل مبتنی بر ویژگی‌های مکانی^۱: در مواردی که اطلاعات مرتبط با مکان مهم هستند، مدل‌های مبتنی بر موقعیت جغرافیایی و ویژگی‌های مکانی به کار می‌روند.

هر یک از این روش‌ها و مدل‌ها بر اساس نیازها و موارد خاص مسائل بازیابی اطلاعات مورد استفاده قرار می‌گیرند. انتخاب مناسب‌ترین روش بستگی به خصوصیات داده‌ها و هدف نهایی بازیابی اطلاعات دارد.

۳-۷ مفاهیم پایه بازیابی اطلاعات

همان‌طور که در توضیحات قبلی گفته شد بازیابی اطلاعات یک حوزه گسترده در علوم کامپیوتر و علم اطلاعات است که به توسعه و بهینه‌سازی فرایند جستجو و بازیابی داده‌ها از مجموعه‌های حجیم اطلاعاتی می‌پردازد. هدف اصلی این حوزه، یافتن اسناد یا اطلاعات مرتبط با یک پرسش خاص یا یک نیاز مشخص است که در ادامه به بررسی مفاهیم اساسی و پایه این حوزه می‌پردازیم.

۱. پرس و جو^۲: در بازیابی اطلاعات، کلمات یا عباراتی که توسط کاربر به منظور جستجو از سیستم وارد می‌شوند، به عنوان کلمه کلیدی شناخته می‌شوند.

۲. سند^۳: هرگونه اطلاعات یا محتوا که ممکن است به عنوان پاسخ به یک جستجو بازیابی شود، به عنوان سند معرفی می‌شود.

۳. اطلاعات مرتبط^۴: در زمینه بازیابی اطلاعات، اطلاعاتی که با توجه به کلمه کلیدی وارد شده توسط کاربر، بیشترین تطابق را دارند و به میزان مرتبطی با جستجو، به عنوان اطلاعات مرتبط شناخته می‌شوند.

¹ Location-Based Model

² Query

³ Document

⁴ Relevance

۳-۸ روش‌ها و مراحل بازیابی اطلاعات

۱. پرسمان‌سازی^۱: در این مرحله، کاربر یا سیستم جستجوی متنی و یا تصویری خود را با استفاده از کلمات کلیدی مناسب مشخص می‌کند.
۲. جستجو و اندازه‌گیری تطابق^۲: سیستم با استفاده از الگوریتم‌های بازیابی، اسنادی که بیشترین تطابق با پرسمان دارند را انتخاب می‌کند.
۳. مرتب‌سازی نتایج^۳: اطلاعات انتخاب شده بر اساس معیارهایی نظیر تطابق، اهمیت و احتمال مرتب می‌شوند تا نتایج به کاربر به ترتیب مرتب سپرده شوند.
۴. نمایش نتایج^۴: نتایج به کاربر نمایش داده می‌شود. این ممکن است شامل متن، تصاویر، یا سایر فرمت‌های متناسب با نوع اطلاعات باشد.

۳-۹ کاربردها و دامنه‌های استفاده

۱. موتورهای جستجوی وب: موتورهای جستجو مانند گوگل از تکنیک‌های بازیابی اطلاعات برای ارائه نتایج دقیق و مرتبط به کاربران استفاده می‌کنند.
۲. سیستم‌های مدیریت اسناد: در سازمان‌ها و شرکت‌ها برای مدیریت و بازیابی اطلاعات از داکيومنت‌ها و اسناد به وسیله سیستم‌های مدیریت اسناد^۵ از تکنیک‌های بازیابی استفاده می‌شود.

¹ Query Formulation
² Search and Matching
³ Ranking Results
⁴ Presentation of Results
⁵ DMS



۳. پرتال‌های خبری و اطلاعاتی: سیستم‌های بازیابی اطلاعات در پرتال‌های اخباری جهت ارائه اطلاعات مرتبط

با جستجوهای کاربران به کار می‌روند.

۴. پرسش و پاسخ هوشمند: در سامانه‌های پرسش و پاسخ هوشمند، از تکنیک‌های بازیابی اطلاعات برای یافتن

پاسخ‌های مرتبط با سوالات کاربران استفاده می‌شود.

بازیابی اطلاعات به عنوان یک حوزه اساسی در علوم کامپیوتر، در جوانب مختلف زندگی و تکنولوژی

تأثیرگذار بوده و با پیشرفت تکنولوژی، روش‌ها و مدل‌های جدیدی برای بهبود عملکرد در این زمینه توسعه می‌یابد.

۳-۱۰ برخی از تکنیک‌های بازیابی اطلاعات

تکنیک‌های بازیابی اطلاعات متنوع و گسترده هستند و با توجه به نوع اطلاعات مورد نظر، میزان دقت مورد نیاز، و

ساختار داده‌ها، از روش‌های مختلف استفاده می‌شوند. در زیر به برخی از تکنیک‌های معروف بازیابی اطلاعات اشاره

می‌شود:

۱. مدل فضای برداری^۱: این مدل اطلاعات را در یک فضای چند بعدی نمایش می‌دهد و از ابزارهای ریاضی

برای اندازه‌گیری تطابق میان پرسمان و اسناد استفاده می‌کند.

۲. مدل BM₂₅: این مدل مبتنی بر احتمالات است و برای اندازه‌گیری تطابق بین پرسمان و اسناد از توزیع

احتمالاتی BM₂₅ استفاده می‌کند.

¹ Vector Space Model

۳. پردازش زبان طبیعی^۱ و تحلیل متن: استفاده از تکنیک‌های پردازش زبان طبیعی مانند تحلیل موضوع، خوشه‌بندی متن، تحلیل احساسات، و استخراج اجزاء مهم متن برای افزایش دقت در بازیابی اطلاعات متنی.
۴. شبکه‌های عصبی: از شبکه‌های عصبی عمیق برای یادگیری ارتباطات پیچیده بین کلمات و ساختارهای مختلف در متن‌ها جهت بهبود بازیابی اطلاعات استفاده می‌شود.
۵. شباهت کسینوسی^۲: این تکنیک بر اساس زاویه میان بردارهای متنوع در فضای بردارهای عددی استفاده می‌شود و میزان شباهت بین پرسمان و اسناد را اندازه‌گیری می‌کند.
۶. الگوریتم‌های هش^۳: برای شتاب‌زدن فرآیند بازیابی اطلاعات، الگوریتم‌های هش برای نمایش سریع و بهینه‌تر داده‌ها مورد استفاده قرار می‌گیرند.
۷. الگوریتم‌های درخت تصمیم^۴: این الگوریتم‌ها به منظور کلاس‌بندی و جستجوی بهینه در فضای ویژگی‌ها برای افزایش دقت و سرعت در بازیابی اطلاعات مورد استفاده قرار می‌گیرند.
۸. مدل‌های ترکیبی^۵: این مدل‌ها از ترکیب چندین مدل مختلف با هدف بهبود دقت و کارایی در بازیابی اطلاعات استفاده می‌کنند.
- هر کدام از این تکنیک‌ها و روش‌ها به مرور زمان با پیشرفت تکنولوژی و علوم داده توسعه می‌یابند و برای حل چالش‌های مختلف در بازیابی اطلاعات مورد استفاده قرار می‌گیرند.

¹ NLP

² Cosine Similarity

³ Hashing Algorithms

⁴ Decision Trees

⁵ Ensemble Models



۱۱-۳ معیارهای امتیازدهی

معیارهای امتیازدهی در بازیابی اطلاعات جهت ارتباط دقیق تر با نیازهای کاربران مورد استفاده قرار می گیرند. این معیارها عموماً بر اساس تحلیل تکرار و اهمیت واژگان در متن ها شکل می گیرند. امتیازدهی به عنوان یک فرایند انتخابی، با بهره گیری از معیارهای مختلف مانند تعداد تکرار^۱، تعداد اسناد حاوی یک واژه^۲، اهمیت مفهومی یک واژه در متن^۳ و ترکیبی از آنها یعنی TF-IDF، به ما کمک می کند تا اسناد مرتبط را با دقت بالاتری انتخاب کنیم. این معیارها اساسی ترین ابزارها برای بهبود عملکرد سامانه های بازیابی اطلاعات هستند، زیرا امکان تفکیک بین اطلاعات مرتبط و غیر مرتبط را فراهم می کنند.

در ادامه به بررسی دقیق تر معیارها می پردازیم.

۱۱-۳-۱ تی اف

در حوزه بازیابی اطلاعات، معیار TF^4 یکی از مفاهیم کلیدی است که فرکانس تکرار یک کلمه خاص در یک سند را اندازه گیری می کند. این معیار نشان دهنده میزان اهمیت یک کلمه در یک سند خاص است در سیستم های بازیابی اطلاعات، از معیار TF برای رتبه بندی اسناد بر اساس میزان تطابق با کلمات کلیدی پرسمان استفاده می شود، در عین حال این معیار به تنهایی نمی تواند رتبه خوبی ارائه کند. در پردازش زبان طبیعی، معیار TF به عنوان یک ویژگی مهم برای تحلیل و استخراج ویژگی های متنی مورد استفاده قرار می گیرد. در مسائل دسته بندی متون، از معیار TF به عنوان یکی از ویژگی های ورودی برای مدل های یادگیری ماشین استفاده می شود. معیار فرکانس تعداد تکرار کلمه TF یک ابزار در تحلیل محتوای متنی است که با محاسبه تعداد تکرار یک کلمه در یک سند، اطلاعات مهمی در

¹ TF

² DF

³ IDF

⁴ Term Frequency

مورد اهمیت و تأثیر آن کلمه در سند ارائه می‌دهد. این معیار یک روش برای اندازه‌گیری اهمیت یک کلمه در یک سند نسبت به یک مجموعه یا کورپوس از اسناد است. این مقدار می‌تواند به صورت تعداد خام، یا به صورت نرمال شده با تقسیم بر تعداد کل کلمات در سند یا به صورت لگاریتمی یا به صورت نرمال دوگانه، محاسبه شود.

۳-۱۱-۱ رابطه ریاضی

TF بیان می‌کند که i نشان‌دهنده کلمه موردنظر و j نمایانگر سند است. این معیار به صورت زیر محاسبه می‌شود

$$TF = \frac{\text{تعداد تکرار کلمه در سند}}{\text{تعداد کل کلمات در سند}}$$

۳-۱۱-۲ دی اف

معیار DF^1 یکی دیگر از معیارها در حوزه بازیابی اطلاعات است که میزان پراکندگی یا فراوانی یک کلمه در میان اسناد مختلف را اندازه‌گیری می‌کند. این معیار برای تشخیص اهمیت یک کلمه در مجموعه اسناد مورد استفاده قرار می‌گیرد. این معیار نشان‌دهنده تعداد اسنادی است که حاوی یک کلمه خاص هستند و به عبارت دیگر، میزان پراکندگی یا روند فراوانی آن کلمه در کل مجموعه اسناد. از DF به عنوان یکی از معیارهای انتخاب کلمات کلیدی استفاده می‌شود. کلمات کلیدی با DF کمتر ممکن است اطلاعات خاصی را نمایان کنند. در تحلیل و پردازش متن، استفاده از DF به عنوان ویژگی می‌تواند کمک کند تا کلماتی که در اسناد مختلف پراکنده‌تر هستند را شناسایی و

¹ Document Frequency



تحلیل کرد. معیار فراوانی سند ابزاری قدرتمند در تحلیل و بازیابی اطلاعات متنی است. با بهره‌وری استفاده از این معیار، می‌توان نتایج دقیق‌تر و مفیدتری در پروژه‌های مختلف دریافت کرد.

۳-۱۱-۲-۱ رابطه ریاضی

DF برابر با تعداد اسناد حاوی کلمه i است.

تعداد اسناد حاوی کلمه $DF =$

۳-۱۱-۳ مقایسه DF و TF

۱. استفاده در بازیابی اطلاعات

TF: برای اندازه‌گیری اهمیت یک کلمه در یک سند خاص و مرتبط با پرسمان‌ها.

DF: برای انتخاب کلمات کلیدی که در مجموعه اسناد به‌عنوان مهم شناخته شده‌اند.

۲. تاثیر تعداد کلمات

TF: توجه به نسبت تکرار به تعداد کلمات در یک سند.

DF: مستقل از تعداد کلمات در یک سند و فقط به تعداد اسناد حاوی کلمه متمرکز است.

۳. کاربرد در مسائل مختلف

TF: برای مسائل مربوط به تحلیل محتوا و استخراج ویژگی‌ها.

DF: برای مسائل مربوط به انتخاب و تشخیص کلمات کلیدی و اهمیت آنها در میان اسناد.

در کل، هر کدام از معیارهای TF و DF وظایف خاص خود را دارند و در مسائل متفاوت مورد استفاده قرار می گیرند. TF بر روی هر سند تمرکز دارد، در حالی که DF بر روی کلمات موجود در کل مجموعه اسناد تمرکز دارد. در مسائل بازیابی اطلاعات، معمولاً از ترکیب این دو معیار برای بهترین نتایج استفاده می شود.

۳-۱۱-۴ ای دی اف

معیار IDF^۱ یکی از ابزارهای حوزه بازیابی اطلاعات و پردازش متن است. این معیار اهمیت یک کلمه خاص در یک مجموعه اسناد را اندازه گیری می کند. محاسبه این معیار به این اعتقاد بنیان گذاری شده است که کلمات که در تعداد کمی از اسناد مجموعه ظاهر می شوند، اهمیت بیشتری دارند. IDF نشان دهنده عکس میزان تکرار یک کلمه در اسناد مجموعه است. با افزایش تعداد اسنادی که حاوی کلمه مورد نظر هستند، ارزش IDF کم می شود. معیار IDF برای تعیین اهمیت کلمات کلیدی در یک پرونده یا پرسمان استفاده می شود. در مسائل تحلیل متن، معیار IDF برای شناسایی کلمات کلیدی و ویژگی های مهم در متون به کار می رود. این مقدار با تقسیم تعداد کل اسناد، N، بر تعداد اسنادی که کلمه t در آن ها ظاهر شده است، df(t)، به دست می آید. $idf(t) = N / df(t)$. این مقدار نشان می دهد که کلمات کمتر رایج اهمیت بیشتری دارند. برای جلوگیری از تأثیر زیاد اعداد بزرگ، معمولاً از لگاریتم این مقدار استفاده می شود.

^۱ Inverse Document Frequency

۳-۱۱-۴ رابطه ریاضی

معیار IDF برای یک کلمه از رابطه زیر به دست می آید:

$$IDF = \log\left(\frac{\text{تعداد کل اسناد}}{\text{تعداد اسناد حاوی کلمه} + 1}\right)$$

۳-۱۱-۴ مقایسه با معیارهای قبلی

تفاوت: DF تعداد اسناد حاوی کلمه را نشان می دهد، TF معیار تکرار کلمه در یک سند است، در حالی که IDF

تأثیر فراوانی آن کلمه در کل مجموعه اسناد را نشان می دهد.

استفاده: DF برای انتخاب کلمات کلیدی در یک سند خاص، TF برای تحلیل محتوا و استخراج ویژگی ها

مورد استفاده قرار می گیرد، در حالی که IDF برای انتخاب کلمات کلیدی در مقیاس گسترده تر مورد استفاده قرار می گیرد.

۳-۱۱-۵ تی اف ای دی اف

معیار TF-IDF^۱ یکی از مهم ترین و مؤثرترین ابزارهای استفاده در بازیابی اطلاعات و پردازش متن است. این

معیار اهمیت یک کلمه خاص در یک سند خاص را مشخص می کند و به عنوان یک ویژگی مؤثر در تحلیل و بازیابی

اطلاعات متنی به کار می رود. معیار TF-IDF نشان دهنده ترکیب تأثیر ترم فراوانی و معکوس فراوانی اسناد بر یک

کلمه خاص است. این معیار نمایانگر اهمیت یک کلمه در یک سند خاص و ارتباط آن با سایر اسناد مجموعه است.

معیار TF-IDF به عنوان یک ابزار کامل و مؤثر در تحلیل متن و بازیابی اطلاعات در پروژه های مختلف مورد استفاده

^۱ Inverse Document Frequency - Term Frequency



قرار می گیرد. استفاده صحیح از این معیار می تواند به دقت و اثربخشی سیستم های پردازش متن و بازیابی اطلاعات کمک کند. TF-IDF به عنوان یکی از ویژگی های کلیدی در مدل های بازیابی اطلاعات مورد استفاده قرار می گیرد. این معیار میزان مطابقت میان پرسمان و اسناد را بهبود می بخشد. در تحلیل متن و استخراج اطلاعات مفید، TF-IDF برای تشخیص کلمات کلیدی و اهمیت آنها در یک متن مورد استفاده قرار می گیرد. برای پیاده سازی، می توان از مقادیر ترم فراوانی و معکوس فراوانی اسناد برای هر کلمه در هر سند استفاده کرد و با ضرب دو مقدار حاصل، مقدار TF-IDF را محاسبه کرد.

در واقع هدف این سیستم وزن دهی، نشان دادن اهمیت کلمه در متن است که اغلب در جستجوهای درون بازیابی اطلاعات، متن کاوی و مدل سازی کاربر استفاده می شود. مقدار TF-IDF به تناسب تعداد تکرار کلمه در سند افزایش می یابد و توسط تعداد اسنادی که در مجموعه هستند و شامل کلمه نیز هستند متعادل می شود. به این معنی که اگر کلمه ای در بسیاری از متون ظاهر شود احتمالاً کلمه ای متداول است و ارزش چندانی در ارزیابی متن ندارد. در حال حاضر TF-IDF یکی از محبوب ترین روش های وزن گذاری اصطلاحات می باشد و امروزه بیش از ۸۳ درصد از سامانه های توصیه گر در کتابخانه های دیجیتال از این روش وزن دهی اصطلاحات استفاده می کنند.

این اختلاف بین وزن ها که توسط روش TF-IDF ایجاد می شود توسط بیشتر موتورهای جستجو به عنوان ابزار اصلی رتبه دهی و امتیازدهی اسناد پرس وجو شده کاربر استفاده می شود؛ و همچنین برای فیلتر کردن ایست وازه ها در زمینه های موضوعی مختلف، از جمله خلاصه سازی و دسته بندی متن با موفقیت استفاده شده است. یکی از ساده ترین تابع های رتبه بندی با جمع کردن وزن به دست آمده توسط TF-IDF برای هر اصطلاح پرس وجو محاسبه می شود. بسیاری از توابع رتبه بندی پیچیده تر بر اساس این مدل ساده به وجود آمده اند.

معیار TF-IDF برای کاربردهای مختلفی در حوزه بازیابی اطلاعات، داده کاوی متن، و مدل سازی کاربر مفید

است. برخی از این کاربردها عبارت اند از:

رتبه‌بندی و ارزیابی اسناد: با استفاده از معیار TF-IDF می‌توان اسناد را بر اساس اهمیت کلماتی که در آنها ظاهر شده‌اند، رتبه‌بندی و ارزیابی کرد. این روش به موتورهای جستجو کمک می‌کند تا اسناد مرتبط با پرسمان کاربر را پیدا و نمایش دهند.

خلاصه‌سازی متن: با استفاده از معیار TF-IDF می‌توان کلمات کلیدی یک متن را شناسایی کرد و با استفاده از آنها یک خلاصه از متن ایجاد کرد. این روش به سیستم‌های خلاصه‌سازی متن کمک می‌کند تا محتوای اصلی متن را با حفظ اهمیت کلمات، بازنمایی کنند.

توصیه‌دهی متن: با استفاده از معیار TF-IDF می‌توان شباهت بین اسناد را محاسبه کرد و بر اساس آن اسناد مشابه به کاربر توصیه کرد. این روش به سیستم‌های توصیه‌دهی متن کمک می‌کند تا اسنادی را که کلمات مشترک و مهم با سند موردعلاقه‌ی کاربر دارند، پیشنهاد دهند.

۳-۱۱-۵ رابطه ریاضی

معیار TF-IDF برای یک کلمه در یک سند خاص به‌ازای مجموعه اسناد به‌صورت زیر محاسبه می‌شود:

$$TF - IDF(t, d, D) = TF(t, d) * IDF(t, D)$$

۳-۱۱-۵-۲ مقایسه با معیارهای قبلی

تفاوت DF: تنها تعداد اسناد حاوی کلمه را نشان می‌دهد، TF نمایانگر تعداد تکرار یک کلمه در یک سند است، درحالی‌که TF-IDF ترکیبی از میزان تکرار کلمه در یک سند و اهمیت آن در کل مجموعه اسناد است.

استفاده DF: برای انتخاب کلمات کلیدی در یک سند خاص، TF بیشتر در تحلیل محتوا و استخراج ویژگی‌ها مورد استفاده قرار می‌گیرد، درحالی‌که TF-IDF برای انتخاب کلمات کلیدی با توجه به ارتباطات بین کلمات در کل مجموعه و تحلیل و بازیابی اطلاعات در مقیاس گسترده‌تر مفید است



فصل ۴

مطالعات و تحقیقات

۱-۴ روند کلی

در این فصل به بررسی مدل انتخاب شده می‌پردازیم با توجه به مقدمات توضیح داده شده در فصل قبلی می‌توانیم با استفاده از معیارهای امتیاز دهی اقدام به محاسبه با استفاده از معیار انتخابی کنیم.

نکته حائز اهمیت در این اقدام این است که معیار انتخاب شده بتواند بخوبی فاصله را میان مفاهیم احساسی کلمات ایجاد کند. در فرم اولیه دیتاها هر سطر حاوی ۷ مقدار است که اولین مقدار آن متن توییت است و ۶ مقدار بعدی تشخیص یافته به آن متن بر اساس ۶ احساس می‌باشد.

به فرمت داده‌ها در جدول ۱-۴ توجه کنید.

جدول ۴-۱ داده‌های اصلی

۱	متن	خشم	توس	شادی	نفرت	غم	تعجب
۲	خوش باشن من که یادم رفته دیگه :)	۲	۰	۱	۱	۱	۱
۳	باشد که #کرونا یک حقیقت را در گوش سنگین و مغز زنگ زده برخی فرو کند: هیچ کشوری جزیره‌ای جدا از دیگران نیست؛ اقتصاد، سلامت، امنیت، سیاست و محیط زیست همه به هم وابسته است.	۰	۰	۰	۲	۱	۰
۴	کرونا رو شکست میدیم؟ مرحله بعد چه کاری می‌خواهی بکنی؟!	۴	۳	۱	۳	۳	۴
۵	اگر در چند ماه اخیر تصمیم داشته اید وارد بورس شوید اما به هر دلیلی به تاخیر افتاده است، این چند روز را هم صبر کنید تا تکلیف مشخص شود. شاید لازم باشد ورود به بازار سرمایه را کلاً فراموش کنید. اگر برگشتی در کار بود کمی بالاتر ورود کنید. شاید آقای مددی همه واقعیت را نگوید.	۲	۴	۱	۲	۴	۲
۶	احساس می‌کنم اگر همین الان نخوابم، تمام مواد شیرینی که دارم توی خونه رو می‌خورم	۱	۱	۰	۰	۳	۲
۷	اینقدر گرفتار مسایل میشی که تخصص از دستت درد میره و به این فکر میکنی چطوری زنده بمونم	۲	۱	۰	۱	۰	۴
۸	ببینید این جلد های رنگارنگتان را که روی هیولاهای درونتان ساخته اید"	۱	۱	۰	۱	۱	۱

۴-۲ پیش پردازش

۱. باتوجه به اینکه تویتهای ما حاوی جملات زبان فارسی است و هدف ما در انتها تشخیص کلمات

تأثیر گذارتر در هر حس است، بنابراین بسیاری از کلمات در کلیه رشته‌ها تکراری است و فاقد ارزش است؛

بنابراین کلیه کلماتی که در سه فایل از پیش آماده شده را برای حذف از تویتهای مدنظر داریم.



جدول ۴-۲ فایل های پیش پردازش

۱	کوتاه	فعلی	غیر فعل
۲	و	کن	دیگران
۳	در	کرد	همچنان
۴	به	کردن	مدت
۵	از	باش	چیز
۶	که	بود	سایر
۷	می	بودن	جا
۸	این	شو	طی
۹	است	شد	کل
۱۰	را	شدن	کنونی
۱۱	با	دار	بیرون
۱۲	های	داشت	مثلا
۱۳	برای	داشتن	کامل
۱۴	آن	خواه	کاملا
۱۵	یک	خواست	آنکه
۱۶	شود	خواستن	موارد
۱۷	شده	گوی	واقعی
۱۸	خود	گفت	امور
۱۹	ها	گفتن	امورات

۲. حذف کردن کلیه کلمات غیر از الفبای فارسی

۳. حذف ستون نام احساسات که در اولین سطر داده های اصلی موجود است

۴. تعیین رفتار با مقدار وارد نشده مانند رکورد های خالی ، حذف رکورد هایی اعضای آن برابر با ۷ نباشد

۵. تبدیل دیتاها به لیست در فرمت جدید

۶. چک مجدد اینکه آیا پس از اعمال پیش پردازش ها متن رکورد خالی شده باشد

۷. بازگرداندن دیتاها به برنامه

۴-۳ رویکرد ما در عملیات بایس کردن استاپ وردها

سطر مربوط به هر داده ها را با معیار جداسازی (" ") از هم جدا می کنیم و به یک لیست جدید انتقال می دهیم و روی لیست جدید که حاوی کلمات است پیمایش انجام می دهیم و هر کجا به کلمه ای برخورد کردیم که در هر یک از سه فایل پیش پردازش موجود باشد، آن را حذف خواهیم کرد.

۴-۴ محاسبه TF-IDF

به منظور محاسبه TF-IDF کلمات نیاز به یک لیست یونیک از کلمات داخل دیتاست داریم تا از محاسبه TF-IDF تکراری برای کلمات جلوگیری شود.

رویکرد ما در تشخیص لیست کلمات یونیک، تبدیل دیتاها به فرمی است که در ادامه قابل استفاده باشد. با حلقه روی هر سطر از دیتاها پیش پردازش شده تشکیل یک لیست جدید و جاگذاری متن تویت در اولین عنصر آن، و تشکیل یک لیست موقت برای جای گذاری هر یک از شش احساس دیگر و در انتها قراردادن لیست موقت در دومین عنصر لیست اصلی و تکرار این عملیات که در انتها منجر به دستیابی به لیستی است که حاوی دیتاها با فرمت موردنظر ما است.

در انتها این مرحله در صورت موفقیت پیام موفق بودن عملیات چاپ می شود.

۴-۵ اعلان متغیر

باتوجه به این موضوع که در صورتی که TF-IDF از مقداری پایین تر شود و ورودی متنی بلند باشد بر روی عملکرد برنامه، تأثیر منفی می گذارد پس نیاز است یک معیاری بگذاریم و در صوت مواجه شدن TF-IDF های کمتر از آن آنها را در نظر نگیریم و از آن گذر کنیم.

باتوجه به ماهیت دیتاست و وجود امتیاز برای هر یک از احساسات شش گانه در صورتی که با متنی مواجه شدیم که امتیاز یک احساس آن زیر یک عدد خاص باشد در محاسبه TF-IDF آن کلمه در جمله ای که احساس مربوطه آن زیر یک مقدار خاص باشد، آن حس و TF-IDF اش را نادیده می گیریم. در جهت بهبود کارایی این عمل و کنترل راحت تر آن دو متغیر تعریف می کنیم و آن را در متن اصلی برنامه مقداردهی می کنیم.

۴-۶ مدل ذخیره سازی TF-IDF

روند کلی به این صورت است که یک پیمایش روی لیست کلمات یونیک می زنیم و TF-IDF آن کلمه در لیست کلیه توییت ها و با توجه امتیاز آن توییت که وزن آن TF-IDF محسوب می شود، حساب می کنیم.

به طور دقیق تر تابع محاسبه کننده TF-IDF سه مقدار کلمه، مجموعه توییت و امتیاز احساسات هر جمله را دریافت می کند و شش متغیر به ازای هر یک از احساسات با مقدار اولیه ۰ اعلان می کند و پس از آن اقدام به محاسبه IDF آن کلمه در مجموع توییت ها می کند پس از آن با اعمال یک پیمایش روی کلیه توییت ها و محاسبه tf کلمه ورودی در توییت پیمایش شده چک می کنیم که مقدار TF آن صفر نباشد تا از عملیات اضافه جلوگیری کنیم زیرا در صورتی که هر یک از متغیرهای TF-IDF صفر باشد مقدار TF-IDF صفر خواهد بود. در ادامه حلقه ای که توییت ها را پیمایش می کند و مقدار TF را محاسبه می کنیم، مقدار بدست آمده برای TF-IDF را که در وزن آن



توییت ضرب شده است را به مقدار مربوط به آن احساس اضافه می کنیم و در انتها یک لیستی حاوی شش متغیر مقدار TF-IDF برای هریک از شش احساس است باز می گردانیم.

مقادیر باز گردانده شده نمایش دهنده امتیازهای مختلف آن کلمه در شش احساس هستند که این مقادیر را در یک دیکشنری با ساختمان داده ای به این شکل که کلید هر مقدار، کلمه خواهد بود و مقدار آن برابر با همان لیستی است که از تابع محاسبه کننده TF-IDF دریافت کرده ایم

. به علت حجم بالای پردازش اطلاعات دیکشنری نهایی را ذخیره می کنیم تا در ادامه مسیر و انجام تست روی داده از همین مدل تولید شده استفاده کنیم و مجبور نباشیم هر بار مقدار TF-IDF کلمات را مجدداً محاسبه کنیم. تا اینجای کار ما موفق به تولید مدل شده ایم در ادامه می توانیم جزئیات مدل را مشاهده کنیم.

جدول ۴-۳ کلمه بر حسب TF-IDF برای هر شش احساس

۱	کلید	خشم	ترس	شادی	نفرت	غم	شگفتی
۲	مرحله	۱۳.۱۰	۶.۰۷	۱۱.۸۹	۹.۰۱	۸.۵۶	۷.۳۵
۳	میخواهی	۱۵.۰۶	۵.۱۲	۲.۱۲	۱۲.۶۴	۷.۴۱	۱۸.۲۹
۴	بخوری	۱۳.۱۱	۵.۶۲	۰	۷.۴۲	۸.۷۶	۹.۹۶
۵	ماه	۱۳.۳۹	۱.۸۱	۷.۸۶	۱۳.۲۵	۲۷.۴۸	۱۰.۱۴
۶	تصمیم	۲.۸۶	۱.۹۲	۳.۵۰	۸.۸۳	۴.۰۸	۴.۳۷
۷	اید	۲.۸۸	۰.۸۲	۰	۳.۲۹	۲.۹۵	۱.۸۵
۸	وارد	۱۱.۰۹	۵.۲۹	۱۰.۵۹	۸.۱۲	۶.۴۰	۱۳.۰۱
۹	بورس	۶۶.۶۴	۴۶.۰۰	۶.۳۲	۳۹.۴۷	۸۶.۱۶	۵۶.۸۸
۱۰	شوید	۲.۴۵	۰.۸۹	۰	۲.۲۳	۱۳.۸۹	۰
۱۱	دلیلی	۹.۰۹	۶.۱۳	۰.۶۳	۰	۱.۴۸	۰.۶۳
۱۲	تاخیر	۱.۹۱	۰.۹۶	۰	۰	۰.۹۶	۰
۱۳	افتاده	۰.۸۹	۱.۲۵	۳.۵۶	۱.۹۶	۷.۰۰	۱.۰۷



۴-۷ تست داده های جدید

گفته شد محاسبه جدول TF-IDF در لحظه و بصورت بلادرنگ ممکن نیست و بار محاسباتی دارد بدین منظور دیتاها را ذخیره کردیم و جهت استفاده مجدد مدل تولید شده را فقط به برنامه اضافه می کنیم و دوباره آن را تولید نمی کنیم. به جهت سهولت در کدنویسی و مشخص بودن دیتاها، شش لیست خالی تعریف می کنیم و در این لیست ها داده هایمان را به شکل یک لیست دوتایی شامل کلمه و مقدار TF-IDF آن کلمه در یک حس خواهد بود. در انتها شش لیست خواهیم داشت که در هر لیست کلمه کلمات یک حس همراه با ارزش آن کلمه در همان حس موجود می باشد

به جهت پیش بینی از دیکشنری قبلی استفاده می کنیم بدین صورت که یک حلقه روی تمام کلیدها و مقدار کلیدها می زنیم و پیمایش انجام می دهیم.

۴-۸ روش های تست ورودی

این قسمت از برنامه به نحوی نوشته شده است که در دو حالت تست داده های دیگر و اجرای به تنهایی می تواند استفاده شود در صورتی که این برنامه منحصر اجرا شود اقدام به دریافت ورودی از کاربر می کند و عملیات تشخیص حس را انجام می دهد، حالت دیگر به این صورت است که در یک برنامه دیگر که دیتا های تست را آماده کرده ایم این تابع را صدا می زنیم و خروجی ما به ازای یک رشته تنها حس تشخیص داده شده خواهد بود در صورتی که در حالت اجرای انحصاری کد بدون فراخوانی در کد دیگر جزئیات بیشتری نمایش داده می شود.

۹-۴ نحوه تشخیص احساس از TF-IDF

رویکرد انتخابی ما در تشخیص و انتخاب حس به عنوان احساس غالب متن بدین شکل خواهد بود که با دریافت رشته ورودی شش متغیر برای نگهداری مقدار امتیاز متن به ازای هر شش حس به صورت جدا از هم تعریف می کنیم.

رشته ورودی را با معیار قراردادن " " از هم جدا می کنیم و در یک لیست قرار می دهیم و پیمایش را روی آن انجام می دهیم. وضعیت هر کلمه در هر شش حس را حساب می کنیم. بدین صورت که اگر کلمه پیمایش شده در هر لیست تولید شده هر حس وجود داشته باشد، مقدار TF-IDF آن بازگردانده می شود و در غیر این صورت مقدار منطقی نادرست^۱ بازگردانده می شود سپس با بررسی وضعیت کلمه در هر حس در صورتی که مقدار مخالف نادرست باشد TF-IDF بازگردانده شده را به متغیرهای اصلی که نگهدارنده مقدار کل TF-IDF های هر حس هستند اضافه می کنیم. در انتها از میان شش متغیر حسی که بیشترین مقدار عددی را دارا باشد را به عنوان احساس قالب متن معرفی می کنیم.

۱۰-۴ دقت مدل

برای محاسبه دقت مدل از یک دیتاستی استفاده می کنیم که حاوی متن و نوع احساس آن می باشد [4].

جدول ۴-۴ نمونه از داده های تست

مقدار	کلید	1
غمگین	خیلی کوچک هستن و سائزشون بدرد نمیخوره میخوام پس بدم	۲
نفرت	از صدای پرنده دم دمای صبح متنفرم متنفرم متنفرم	۳
غمگین	این وضع به طرز خنده داری گریه داره...	۴
نفرت	خب من رسماً از یک نفر متنفرم، چون از گربه بدش میاد از صبح شروع کرده رو مخ من راه رفتن؛ شپش میگری، کز از میگیری، کک داره، درد داره، مرض داره	۵

¹ False



دقت مدل تولید شده با توجه به مطالب گفته شده روی دیتاهای تست به شکل زیر خواهد بود

جدول ۴-۵ کد نمونه‌ها

نمونه‌ها	پیش بینی های درست	دقت	ترس	نفرت	غمگین	خشمگین	خوشحال	تعجب
۱۲۴۰	۱۰۹۰	۸۷.۹۰۳۲۲۵	۵۵۱	۵۱	۵۹	۴۱۲	۶۸	۹۹



فصل ۵

نتیجه گیری و پیشنهادات

۵-۱ نتیجه گیری

همانطور که گفته شد هدف از این پروژه ارائه ی ساز و کاری نوین در تشخیص احساسات با استفاده از روش های بازیابی اطلاعات بوده است. آنچه اهمیت این پروژه را دوچندان کرده بود عدم وجود اقدامات مشابه در زبان فارسی بود که با وجود منابع بسیار اندک توانستیم به مقدار زیادی این پروژه را در حالت بهینه ارائه کنیم. در انتهای این راه ما می توانیم با دریافت یک متن از ورودی احساس آن را تشخیص دهیم، این امر می تواند در زمینه های بازاریابی و فروش بسیار کاربردی باشد و می توانیم آن را مسائل مختلفی استفاده کنیم.



در ادامه بیشتر به اهمیت متن کاوی در متون ادبیات فارسی اشاره می کنیم.

متن کاوی در متون ادبیات فارسی اهمیت زیادی دارد و می تواند به درک عمیق تر از ساختارها، مضامین، و احساسات

مختلف متون ادبی کمک کند. در زیر توضیحاتی در مورد اهمیت این موضوع ارائه شده است:

تفکر علمی و فرهنگی : متون ادبیات فارسی غنی از تفکر علمی و فرهنگی هستند. متن کاوی در این متون به ما امکان

می دهد تا درک عمیقی از این تفکرات را پیدا کنیم و برای شناخت بهتر فرهنگ و ادبیات ایرانیان اقدام نماییم.

شناخت احساسات و عواطف : ادبیات فارسی غنی از عواطف و احساسات انسانی است. متن کاوی در این حوزه به ما

این امکان را می دهد تا عمیق ترین احساسات نویسندگان را درک کرده و با زبان ادبی انتقال دهیم.

شناسایی نقش و تأثیر نویسنده : متون ادبیات فارسی در بسیاری از موارد تحت تأثیر تجربیات و دیدگاه های نویسنده

قرار دارند. با استفاده از متن کاوی، ما می توانیم نقش و تأثیر نویسنده در آثارش را بهترین طور ممکن شناسایی کنیم.

تحلیل ساختار و اندیشه : ادبیات فارسی اغلب دارای ساختارها و اندیشه های پیچیده ای هستند. متن کاوی به ما این

امکان را می دهد تا اجزا و عناصر مختلف یک متن را تجزیه و تحلیل کنیم و به فهم بهتری از ساختار ادبی و اندیشه

نویسنده برسیم.

تفسیر ساختارهای زبانی : زبان ادبی فارسی دارای ساختارهای زبانی خاصی است. متن کاوی در این حوزه به ما این

امکان را می دهد که ساختارهای زبانی پیچیده را تجزیه و تحلیل کرده و به درک عمیق تری از زبان ادبی فارسی

برسیم.

کمک به پژوهش های بیشتر : نتایج متن کاوی می توانند به پژوهش های بیشتر در زمینه ادبیات فارسی و متون مشابه

کمک کنند و محققان را به سمت پرسش ها و چالش های جدید هدایت کنند.

این نکات تنها یک آغاز برای شروع بحث در مورد اهمیت متن کاوی در متون ادبیات فارسی هستند. حال که جزئیات

بیشتری از اهمیت این موضوع می دانیم می توانیم بهتر در مورد بررسی تحقیقات انجام شده اظهار نظر نمود.



در تحقیقات این پروژه ما ابتدا سعی کردیم دیتاست را بر اساس اصول اولیه نرمال کردن آنرا به یک حداقل برای اجرا کردن روش های امتیاز دهی تبدیل کنیم. در این مسیر ما می توانیم از گفایل های آماده که برای مثال کلمات توقف را گسترش داده اند کمک بگیریم اما با توجه به ماهیت فایل های ما در این پروژه و وجود ادبیات گفتاری و نوشتاری خاص که این روزه در شبکه های اجتماعی رواج دارند امکان استفاده مستقیم از این فایل ها وجود ندارد و نیاز است برای استفاده در پروژه شخصی سازی شود، در این مسیر ابتدا با بررسی تک تک توییت ها که در این پروژه آنها را یک سند در نظر می گیریم به نحوه نوشتار آن دقت می کنیم و فایل های پیش پردازش را بر طبق آن تغییر می دهیم یا از ابتدا آن ها را می نویسیم.

در قدم بعدی با توجه به اینکه پروژه عمیقاً در زبان فارسی پژوهش انجام می دهد اکیدا نیاز است دیتاست را از هرگونه الفبای غیر فارسی پاکسازی کنیم که در این موضوع نیز با توجه به دیتاست نیاز است بررسی کنیم در توییت هایی که الفبای غیر فارسی در آن است با حذف تمامی کاراکترهای غیر فارسی مفهوم جمله چطور تغییر می کند؟ پس در این موضوع نیز نیاز است با بررسی دقیق تر اقدام به حذف الفبای غیر فارسی نمود.

در دیتاست اصلی با در سطر اول نام ویژگی های هر سطر آورده شده است، نیاز است در اولین اجرا و اولین اقدامی که صورت می گیرد آنها را از دیتاست خارج کنیم، راهکارهای دیگری که وجود دارد شامل نادیده گرفتن سطر اول است، یا برنامه به نوعی بهینه شود که در صورت مواجه شدن با اطلاعاتی که در قالب خواسته شده نباشد آنرا تشخیص دهد و آنرا در محاسبات لحاظ نکند. قطعاً برنامه ای که بتواند در همه شرایط بدون خطا اجرا شود و توانایی برخورد با انواع اطلاعات وارد نشده را داشته باشد از نظر عملکرد و پایداری در سطح بالاتری قرار می گیرد و در محیط های عملیاتی بار هزینه مالی و زمانی کمتری را بر دوش متقاضی برنامه می گذارد و همچنین برای برنامه نویس نیز این مسئله حائز اهمیت است که تا بالاترین حد امکان برنامه بهینه اجرا شود.

مورد دیگری که در این پروژه وجود دارد برخورد با اطلاعاتی هست که پس از اعمال فیلترها دیگر قابل استفاده نیست، تصور کنید با اعمال کلمات توقف و حذف حروف غیر فارسی به جمله ای برسیم که دارای تنها یک یا دو



کلمه باشد، مسلماً اطلاعاتی که می توانیم از این متن دریافت کنیم کم و ناچیز خواهد بود و در ابعاد بزرگتر می تواند کارایی مدل ما را پایین بیاورد.

بدین جهت در متن برنامه تنها توییت هایی را مورد بررسی قرار می دهیم که دارای کلمه بیشتر از سه باشند و در صورت مواجهه با توییت هایی که که کمتر از مقدار گفته شده باشد آنرا نادیده خواهیم گرفت.

از چالش هایی که در این مسیر می توان به آن اشاره کرد بار محاسباتی حساب کردن مقدار tf-idf است، از جمله اقداماتی که در جهت کنترل کردن این مورد در برنامه لحاظ شده است ایجاد کردن لیست کلمات است، بدین صورت که از اولین خط خواندن دیتاست تمامی کلماتی که به آن برخورد می کنیم را به یک لیست اضافه می کنیم، این نکته قابل توجه است که تنها زمانی کلمه را به لیست اضافه می کنیم که آن کلمه در لیست ما وجود نداشته باشد و بدین ترتیب به لیستی دسترسی خواهیم داشت که حاوی تمام کلمات داخل دیتاست بدون در نظر گرفتن تکرار است، این امر باعث می شود که ما در طول برنامه و محاسبه مقادیر برای کل دیتاست، هیچ گاه اقدام به بررسی مقداری تکراری نکنیم، با توجه به بار محاسباتی ای که گفته شد همین اقدام می تواند پیچیدگی زمانی پروژه را تا حد بسیار بالایی بکاهد و سرعت عمل کد ما را بالا ببرد.

با این حال پیچیدگی زمانی بالا کماکان در برنامه وجود خواهد داشت. اقدامی که در این پروژه انجام شد بررسی انواع روش های ممکن به جهت تسریع این فرایند ها بود که بهترین اقدام ممکن ذخیره مدل تولید شده در فایل می باشد.

بدین صورت که با اجرای برنامه و پس از انجام پیش پردازش ها برنامه دیکشنری مربوط به مقدار tf-idf کلمات به تفکیک شش احساس، آن را ذخیره می کنیم.

با توجه به اینکه در این مرحله از پروژه با یک دیکشنری مواجه هستیم باید روش صحیحی به جهت ذخیره آن انتخاب کنیم.

این روش انتخابی باید ویژگی های حداقلی را دارا باشد که می توان به موارد ذیل اشاره نمود.



نخست اینکه دلیلی که ما به ذخیره دیکشنری روی آوردیم کم کردن پیچیدگی زمانی این فرایند است، بنابراین با توجه به اینکه حجم داده های ما بسیار بالا می باشد باید روشی اتخاذ شود که پیچیدگی زمانی حداقلی را نیز دارا باشد. از طرف دیگر برنامه به خودی خود حجم بالایی از رم را اشغال می کند بنابراین نباید روش ذخیره سازی ما با فدا کردن پیچیدگی مکانی تنها پیچیدگی زمانی را بهبود ببخشد. روشی که ما در این پروژه از بهره جویی کردیم استفاده از توابع داخلی پایتون به جهت نوشتن و خواندن ساختمان داده های خود پایتون است که این کار را به نحو احسن انجام می دهند.

زمانی که ما در برنامه از توابع آماده خود پایتون استفاده می کنیم به این معنا است که این برنامه در بهترین حالت ممکن اجرا می شود و این موضوع می تواند اطمینان حاصل کند که بهینه تری روش ممکن استفاده شده است.

۵-۲-۲ پیشنهادات

اقداماتی که در جهت بهتر شدن این مدل که در این پروژه استفاده شده است می توان انجام داد به بهینه کردن الگوریتم محاسبه TF-IDF می توان اشاره کرد. این معیار امتیاز دهی اصلی ترین معیاری است که ما در این پروژه استفاده کرده ایم. مشخص است که هرچه این الگوریتم بهینه تر عمل کند جوابی که در انتها خواهیم داشت نیز بهینه خواهد بود، برای مثال می توانیم مقداری وزن جملات را نرمال کنیم که یکباره یک توییت با وزن بالا تاثیر کمتری در کلمات غیر مرتبط بگذارد.

با توجه به مورد گفته شده در بخش قبلی می توان در راستای بهبود عملکرد همین مدل انتخاب شده پنی استفاده از TF-IDF وزن دار اقدام به بالابردن کارایی برنامه کنیم، مواردی که می توان بطور مستقیم به آن اشاره کرد کلماتی است که به عنوان کلمه توقف از متن حذف می شوند.

با بهینه کردن هرچه بیشتر این کلمات می توانیم بازدهی بیشتری در برنامه داشته باشیم.



باید توجه داشت که مدل استفاده شده در این پروژه تا حد معقولی نتیجه صحیح را باز می گرداند و با تغییر جزئیات در پیش پردازش و نوع الگوریتم بهبود زیادی مشاهده نخواهد شد و مورد انتظار است زیر ده درصد بهبود حاصل شود.

به جهت رسیدن به کارایی بهتر نیاز است دیتاست اولیه ای که روی آن عملیات محاسبه انجام می دهیم به درستی وزن دهی شده باشد و در حالت پیش فرض و فرم اولیه اش نیز نیاز به پیش پردازش خاصی نداشته باشد، هرکجا که پیش پردازش انجام دهیم ممکن است مقداری از کارایی را از دست بدهیم.

استفاده از روش های ماشین لرنینگ برای تشخیص احساسات در متون

در عصر اطلاعات و فناوری، تحلیل احساسات در متون به ویژه در زمینه ادبیات فارسی، به یکی از چالش های اصلی پژوهشی تبدیل شده است. استفاده از روش های ماشین لرنینگ، به خصوص شبکه های عصبی و ماشین های بردار پشتیبان، در این حوزه امکانات گسترده تری برای تشخیص احساسات و اندازه گیری میزان هر احساس ارائه می دهد.

استفاده از شبکه های عصبی

شبکه های عصبی به دلیل توانایی بالا در یادگیری الگوهای پیچیده، برای تحلیل احساسات در متون ادبیات فارسی بسیار مناسب هستند. این شبکه ها می توانند اطلاعات احساساتی متن ها را به صورت خودکار و بدون نیاز به الگوهای دقیق، استخراج کنند. با آموزش این شبکه ها بر روی دیتاستی که شامل متون با برچسب های احساسات مختلف است، می توان به تشخیص و اندازه گیری احساسات در متون پرداخت.



استفاده از ماشین‌های بردار پشتیبان

ماشین‌های بردار پشتیبان نیز از قدرتمندترین روش‌های ماشین‌لرنینگ برای تحلیل احساسات متنی هستند. این روش با توجه به توانایی خود در دسته‌بندی داده‌ها، می‌تواند برای تفکیک احساسات مثبت، منفی، یا حتی احساسات میانه در متون ادبیات فارسی به کار گرفته شود. با تعیین یک مدل بهینه و استفاده از ویژگی‌های مناسب، می‌توان به تشخیص احساسات با دقت بالا دست یافت.

مقایسه با روش‌های بازیابی اطلاعات و TF-IDF

استفاده از شبکه‌های عصبی و ماشین‌های بردار پشتیبان برای تشخیص احساسات در مقایسه با روش‌های بازیابی اطلاعات و TF-IDF، اغلب نتایج بهتری ارائه می‌دهد. زیرا این روش‌ها قابلیت یادگیری و تفکیک بهتری را دارند و می‌توانند الگوهای پیچیده‌تری را در متون ادبیات فارسی تشخیص دهند. در حالی که روش‌های سنتی معمولاً بر اساس ویژگی‌های ساختاری یا تکراری متن عمل می‌کنند، ماشین‌لرنینگ با توانایی یادگیری و استفاده از ویژگی‌های پیچیده‌تر به دقت بیشتری دست می‌یابد.

استفاده از ماشین‌لرنینگ در تشخیص احساسات متن‌ها، نه تنها به افزایش دقت در تحلیل احساسات ادبیات فارسی منجر می‌شود، بلکه امکانات جدیدی برای تحقیقات آینده در زمینه متن‌کاوی و ادبیات فارسی ارائه می‌دهد.

مراجع

- [1] N. Sabri, R. Akhavan, and B. Bahrak, “Emopars: A collection of 30k emotion-annotated persian social media texts,” in *Proceedings of the Student Research Workshop Associated with RANLP 2021*, 2021, pp. 167–173.
- [2] S. S. Sadeghi, H. Khotanlou, and M. Rasekh Mahand, “Automatic Persian text emotion detection using cognitive linguistic and deep learning,” *J. AI Data Min.*, vol. 9, no. 2, pp. 169–179, 2021.
- [3] C. D. Manning, P. Raghavan, and H. Schütze, *Introduction to Information Retrieval*. USA: Cambridge University Press, 2008.
- [4] H. Mirzaee, J. Peymanfard, H. H. Moshtaghin, and H. Zeinali, “ArmanEmo: A Persian Dataset for Text-based Emotion Detection,” *arXiv Prepr. arXiv2207.11808*, 2022.

