

اعتقادیابی متون فارسی بر اساس یادگیری عمیق با تفکیک احساس-کلمه

حسین علی کرمی*، امیرمسعود بیدگلی(نویسنده مسئول)**، حمید حاج سیدجوادی***
* گروه مهندسی کامپیوتر، واحد تهران شمال، دانشگاه آزاد اسلامی، تهران، ایران.
** گروه مهندسی کامپیوتر، واحد تهران شمال، دانشگاه آزاد اسلامی، تهران، ایران.
*** گروه ریاضی و علوم کامپیوتر، دانشگاه شاهد، تهران، ایران.

چکیده

اعتقادکاوی یا طبقه بندی متون بر اساس احساس و عقیده کاربران در وبسایت ها و رسانه های اجتماعی به مردم، شرکت ها و سازمان ها کمک می کند تا بتوانند تصمیم گیری های مهم را انجام دهند. اعتقادکاوی شامل یک سیستم برای تحلیل عقاید و احساسات مردم درباره یک موجودیت مانند محصولات، افراد، سازمان ها با توجه به نظرات، پیام ها و توییت های کاربران در رسانه های اجتماعی می باشد.

در این مقاله اعتقادکاوی متون فارسی بر اساس پیام ها، نظرات و توییت های کاربران در رسانه اجتماعی با استفاده از دو روش یادگیری عمیق **CNN** و **LSTM** با در نظر گرفتن شدت احساس کلمات و دیدگاه کاربران، در دو قطب مثبت و منفی با بازه ۲- و ۲+ طبقه بندی شده اند. در روش پیشنهادی ابتدا فرآیند پیش پردازش داده ها بر اساس تبدیل کاراکتر به عدد، حذف لیست واژه های اضافی و تحلیل چند واژه ای انجام می شود، در مرحله دوم جهت اعتقادکاوی و طبقه بندی متون فارسی از الگوریتم های یادگیری عمیق **CNN** و **LSTM** با تفکیک احساس بر روی کلمات (**WSD**) استفاده می شود تا شدت احساسات را با توجه به کلمات تشخیص دهد، در مرحله سوم دیدگاه کاربران با طبقه بندی متون در چهار دسته سیاسی، اجتماعی، اقتصادی و فرهنگی تشخیص داده می شود. ما مدل پیشنهادی را **CNN_WSD** و **LSTM_WSD** می نامیم. در این مقاله از مجموعه داده فارسی توییت برای ارزیابی و مقایسه روش پیشنهادی با سایر روش های یادگیری ماشین و یادگیری عمیق از جمله **DNN**, **CNN**, **LSTM** استفاده شده، پیاده سازی این روش با نرم افزار پایتون انجام شده است. میزان دقت روش پیشنهادی برای **CNN-WSD** و **LSTM-WSD** به ترتیب **95.8** و **94.3** درصد است.

واژگان کلیدی: اعتقادکاوی، پردازش زبان طبیعی (**NLP**)، یادگیری عمیق، متن کاوی

1. مقدمه

تحلیل احساسات که همچنین نظرکاوی نامیده می شود، بخشی از مطالعات است که به تحلیل عقاید، احساسات و نگرش های مردم

درباره موجودیت هایی همچون محصولات، سرویسها، سازمانها، افراد، رخدادها و موضوعات خاص میپردازد (Liu, 2012)

نویسنده مسئول: امیرمسعود بیدگلی (DrBidgoli@gmail.com),
am_bidgoli@iau-tnb.ac.ir

در فعالیت هایی که تاکنون بر روی متون فارسی انجام شده اغلب به طبقه بندی متون با توجه به زمینه و حوزه مورد مطالعه پرداخته شده است، در زمینه اعتقادکاوی و تحلیل احساسات بر روی نظرات و متون فارسی پیشرفت چشمگیری انجام نشده، در فعالیت هایی که تاکنون انجام شده اغلب بدون در نظر گرفتن معنا، شدت حساسیت کلمات و دیدگاه کاربران و کلمات اضافی فقط نظرات را در دو گروه مثبت و منفی طبقه بندی می کنند. وقتی در تحلیل احساسات کاربران نظرات در گروه های مثبت و منفی طبقه بندی می شود، اهمیت و ضرورت این موضوع احساس می شود که شدت احساسات نیز مورد ارزیابی قرار گیرد، در این مقاله با استفاده از روش پیشنهادی اعتقادکاوی با یادگیری عمیق، ابتدا نظرات کاربران در دو قطب مثبت و منفی دسته بندی شده و برای هر قطب شدت احساسات بین بازه -2 و $+2$ نسبت داده شده است، سپس دیدگاه کاربران از نظر سیاسی، اجتماعی، اقتصادی و فرهنگی تعیین می شود، اعتقادکاوی با تعیین شدت احساسات بر اساس دیدگاه کاربران باعث می شود که سازمان ها، شرکت ها و محققین به راحتی نظرات کاربران را در زمینه های مختلف تحلیل کنند و تصمیم گیری مناسب و تعیین اهداف را انجام دهند.

بیشتر پژوهش های انجام شده در حوزه پردازش اطلاعات متنی و تحلیل احساسات، بر روی داده کاوی و طبقه بندی اطلاعات، مانند تحلیل نظرات در دو قطب مثبت یا منفی تمرکز دارند. این در حالی است که در فرایند تصمیم گیری توسط مدیران سازمان ها و شرکت ها نیاز به اطلاع دقیق از دیدگاه و احساس کاربران می باشد، به عنوان مثال اگر درصد نظرات کاربران در دو قطب مثبت و منفی برابر باشد نمی توان تصمیم گیری انجام داد، اما در اعتقاد کاوی تعیین می کند که نیمی از کاربران با یک دیدگاه تعیین شده نظر مثبت و نیمی از آنها با دیدگاه دیگر نظر منفی دارند و یا حتی کاربرانی که از یک دیدگاه مشترک نظر داده اند چند درصد از آنها دیدگاه مثبت یا منفی دارند.

در اعتقادکاوی با تحلیل احساسات کاربران بر اساس شدت احساسات کلمات، زمینه ها و دیدگاه کاربران در دسته های سیاسی، اجتماعی، اقتصادی و فرهنگی مشخص می شود و نظرات کاربران را بر اساس دیدگاه آنها در دسته های مثبت و منفی بین بازه $+2$ و -2 تحلیل می کند.

در اعتقادکاوی و تحلیل احساسات بر اساس نظرات کاربران با یادگیری عمیق چالش ها و مشکلاتی مانند هزینه بالای آموزش بر اساس زمان یا حافظه استفاده شده، عدم وجود واژگان غنی و کامل، ابعاد بالای فضای ویژگی و ابهام در تشخیص مثبت یا منفی برخی از جملات وجود دارد. از سوی دیگر، علم یادگیری عمیق توانسته با پیشرفت خود به بسیاری از مسائل حوزه پردازش زبان طبیعی

پاسخ دهد و جایگزینی مناسب برای روش های سنتی باشد، یادگیری عمیق تاکنون از خود عملکرد بسیاری خوبی در بسیاری از شاخه های پردازش زبان طبیعی، خصوصاً تحلیل احساسات نشان داده است. مهمترین مزیت این روش، بی نیازی از استخراج دستی ویژگی ها است که به جای تخصص در حوزه زبان شناسی بر دسترسی به حجم بالای داده ها تکیه دارد.

برای مقابله با این مشکلات، در این مقاله از روش جدید یادگیری عمیق مبتنی بر شدت احساس کلمات بر اساس دیدگاه کاربران، به تحلیل احساسات و اعتقادکاوی نظرات پرداخته شده است. در این پژوهش روش پیشنهادی برای اعتقادکاوی نظرات فارسی کاربران در شبکه اجتماعی توییتر اهداف زیر را دنبال می کند:

- در پیش پردازش کلمات اضافی با توجه به لیست حذف می شود و داده ها به عدد و بصورت برداری تبدیل می شوند.

- در روش پیشنهادی بدلیل استفاده از روش یادگیری عمیق (CNN, LSTM) با تفکیک شدت احساس بر روی کلمه بر اساس دیدگاه کاربران در دقت رده بندی و تحلیل احساسات پیشرفت خوبی بدست آمده است.

- در روش پیشنهادی یک روش جدید با یادگیری عمیق با تفکیک و تحلیل احساسات در متون فارسی ارائه شده که متون در قسمت رمزنگاری به عدد تبدیل می شود تا متون به بردار تبدیل شود، این روش باعث افزایش دقت در اعتقادکاوی شود.

- ابتدا کلمات برگرفته از بدنه اصلی اسناد آموزشی با روش اسکپ گرام¹ مشخص می شود و با مقایسه کلمات با واژگان، شدت احساسات آن مشخص می شود و در دسته های مثبت و منفی در بازه $+2$ و -2 مشخص می شود.

- تحلیل احساسات کاربران بر اساس شدت احساس کلمات می تواند دید دقیق و صحیح تری نسبت به نظرات کاربران در زمینه های مختلف را برای تصمیم گیری به ما بدهد.

در سالهای اخیر، محققین مطالعات متعددی در زمینه تحلیل احساسات با یادگیری عمیق برای زبان انگلیسی ارائه داده اند اما در زمینه اعتقادکاوی و تحلیل دیدگاه های کاربران تحقیقات زیادی انجام نشده است. همچنین تحلیل احساسات در متون فارسی نیز پیشرفت چشمگیری نداشته است. به همین دلیل در این کار به اعتقادکاوی با تحلیل شدت احساس کلمات بر اساس دیدگاه کاربران پرداخته ایم. در این کار جدیدترین مدل های یادگیری عمیق، مانند شبکه های عصبی عمیق (DNN)، شبکه های عصبی مکرر (RNN) و شبکه های عصبی کانولوشن (CNN) و مدل های یادگیری ماشین و روش پیشنهادی جدید CNN_WSD, LSTM_WSD بررسی و ارزیابی شده است،

¹ Skip-gram

یادگیری عمیق شامل شبکه عصبی عمیق (DNN)، شبکه عصبی حلقوی (CNN)، شبکه عصبی مکرر (RNN) و LSTM است (کش دانگ و همکاران، ۲۰۲۰).

شبکه عصبی عمیق: یک شبکه عصبی با بیش از دو لایه است که برخی از آن‌ها لایه‌های پنهان هستند. شبکه‌های عصبی عمیق از مدل سازی ریاضی پیچیده برای پردازش داده‌ها به روش‌های مختلف استفاده می‌کنند.

شبکه عصبی حلقوی: یک نوع خاص از شبکه عصبی رو به جلو است که در اصل در مناطقی مانند بینایی ماشین، سیستم‌های توصیه کننده و پردازش زبان طبیعی به کار می‌رود. این یک معماری شبکه عصبی عمیق است، که به طور معمول از لایه‌های محکم و جمع با مخزن تشکیل شده است تا ورودی‌های یک لایه طبقه بندی کاملاً متصل را فراهم کند.

حافظه کوتاه مدت (LSTM): یک نوع خاص از شبکه عصبی مکرر (RNN) است که قادر به استفاده از حافظه طولانی به عنوان ورودی توابع فعال سازی در لایه پنهان است. سایر روش‌های یادگیری عمیق در (علی کرمی و همکاران، ۲۰۲۳) بیان شده است.

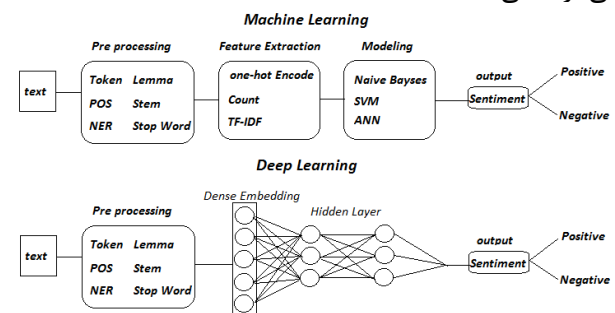
تا مشکلات مختلف در رابطه با تجزیه و تحلیل احساسات را حل کند. در روش پیشنهادی برای اعتقادکاوی و رده‌بندی نظرات از یادگیری عمیق با تفکیک شدت احساسات کلمات بر روی مجموعه داده‌های فارسی استفاده شده، این روش سبب استخراج ویژگی‌های مناسب و افزایش دقت و تشخیص شدت احساسات در بازه +۲ و -۲ می‌شود، بنابراین روش پیشنهادی برای بهبود تصمیم‌گیری و نظرکاوی در متون فارسی اهمیت بالایی دارد.

در این مقاله در بخش اول مقدمه ای در زمینه اعتقادکاوی و رده بندی متون فارسی، در بخش دوم ادبیات تحقیق و پیش زمینه ای در مورد روش‌های یادگیری عمیق، در بخش سوم به کارهای گذشته در زمینه اعتقادکاوی و تحلیل احساسات پرداخته شده، در بخش چهارم پیش پردازش متون فارسی و اعتقادکاوی نظرات با روش پیشنهادی CNN_WSD و LSTM_WSD ارائه شده، در بخش پنجم پردازش داده‌ها، در بخش ششم نتایج و آزمایشات روش پیشنهادی و مقایسه با سایر روش‌ها و سپس به نتیجه گیری پرداخته شده است.

2. ادبیات تحقیق

در روش پیشنهادی اعتقادکاوی مبتنی بر الگوریتم یادگیری عمیق می‌باشد، قبل از بررسی و تحلیل مدل بهینه شبکه عصبی عمیق در روش پیشنهادی مقدماتی از یادگیری عمیق بیان می‌شود.

یادگیری عمیق یک رویکرد چند لایه را با لایه‌های پنهان شبکه عصبی تطبیق می‌دهد. در روش‌های یادگیری سنتی ماشین، ویژگی‌ها به صورت دستی یا با استفاده از روش‌های انتخاب ویژگی مشخص و استخراج می‌شوند. اما در مدل‌های یادگیری عمیق، ویژگی‌ها به طور خودکار استخراج می‌شوند و به دقت و عملکرد بهتری دست می‌یابند. شکل ۱ اختلاف‌های طبقه بندی و تحلیل احساسات بین دو رویکرد یادگیری ماشین و یادگیری عمیق را نشان می‌دهد. شبکه‌های عصبی مصنوعی و یادگیری عمیق در حال حاضر بهترین راه حل‌ها را برای بسیاری از مشکلات در زمینه های اعتقادکاوی، تحلیل احساسات، تشخیص گفتار و پردازش زبان طبیعی ارائه می‌دهند.



شکل ۱: تفاوت بین دو رویکرد تحلیل احساسات با روش‌های یادگیری ماشین (بالا) و یادگیری عمیق (پایین).

3. کارهای گذشته

در اوایل دهه ۱۹۶۰ مردم شروع به مطالعه طبقه بندی متن کردند. در آن زمان قوانین طبقه بندی با توجه به پدیده‌ها و قوانین زبان نوشته می‌شد. تا دهه ۱۹۹۰، مردم شروع به مطالعه تکنولوژی طبقه بندی خودکار مبتنی بر کامپیوتر کردند. در این روشها ابتدا به پیش برچسب گذاری داده‌ها، یادگیری قوانین و سپس به طبقه بندی و آموزش نمونه‌های جدید از دسته‌های ناشناخته به طور خودکار پرداخته شد. نتایج نشان می‌دهد که در زمینه حجم داده‌های بزرگ، دقت طبقه بندی آن بسیار بهتر از تعریف قوانین است. بنابراین تحقیقات کنونی بر طبقه بندی خودکار متون با الگوریتم‌های هوش مصنوعی متمرکز است (یاو و همکاران، ۲۰۱۱).

طبقه بندی متون یکی از وظایف اصلی یادگیری ماشین است. هدف آن طراحی الگوریتم‌های مناسب برای استخراج ویژگی‌ها و طبقه بندی متون به صورت خودکار است. در گذشته، اساساً از طبقه بندی کلمات کلیدی و طبقه بندی سنتز معنایی با شبکه عصبی استفاده می‌شد.

محققان در طول سال‌های ۲۰۰۰-۲۰۱۵ به ارزیابی و مقایسه روش‌های تجزیه و تحلیل احساسات و متن کاوی پرداخته اند (پیراری و همکاران، ۲۰۱۷). در زمینه تحلیل احساسات مانند استخراج احساسات، نظر کاوی، استخراج نظرات و اعتقادکاوی،

روش‌هایی برای تجزیه و تحلیل احساسات و اندیشه کاوی در کلمه، جمله و سطح سند، برای نظرات مصرف کنندگان در عبارات بیان شده است. با توجه به فازی بودن شخصیت، تکنیک های یادگیری ماشین سنتی نمی‌توانند نظرات را به خوبی نشان دهند. برای رفع این مشکل روش تحلیل احساسات با فاز معنایی برای حل مسئله پیشنهاد شده است (فانگ و همکاران، ۲۰۱۸).

مطالعه افکار عمومی می‌تواند اطلاعات ارزشمندی را در اختیار ما قرار دهد. تجزیه و تحلیل احساسات در شبکه‌های اجتماعی، مانند توییتر یا فیس بوک، به ابزاری قدرتمند برای تحلیل نظرات کاربران تبدیل شده و کاربردهای گسترده ای دارد. با این حال، چالش‌هایی در صحت تحلیل احساسات در پردازش زبان طبیعی (NLP) پیش آمده. در سال‌های اخیر نشان داده شده است که مدل‌های یادگیری عمیق یک راه حل امیدوار کننده برای چالش‌های تحلیل احساسات و اعتقادکاوی است. در مطالعاتی مانند قطبیت احساسات که از یادگیری عمیق برای حل مشکلات تحلیل احساسات استفاده شده است، مدل‌هایی با استفاده از فرکانس سند، فرکانس معکوس (TF-IDF) و تعبیه کلمه بر روی یک سری مجموعه داده‌ها اعمال شده اند (کچ دانگ و همکاران، ۲۰۲۰).

با توجه به ویژگی‌های مورد استفاده برای بهره برداری از نظرات مصرف کننده، مدل کیف واژه‌ها بیانگر سند سنتی است که در آن فرکانس‌های کلمه برای هر کلمه (عبارت) در واژگان محاسبه می‌شود (جانسون و همکاران، ۲۰۱۵). با این حال، این رویکرد منجر به بازنمایی اسناد پراکنده در ابعاد بالا می‌شود، علاوه بر این، این روش معنی کلمه را نادیده می‌گیرد. برای غلبه بر این مشکلات، روش تعبیه کلمات به جای تک کلمات معرفی شده اند، که یک زمینه کوتاه در نظر گرفته شده است (تانگ و همکاران، ۲۰۱۵). در مقایسه روش تعبیه کلمات با مدل کیف واژه^۲، روش تعبیه کلمات در مدل سازی متن کلمه و معنی کلمه نیز مؤثر است. بعد از تولید اسناد مناسب، می‌توان از مدل‌های مختلف شبکه عصبی و سایر روشهای یادگیری ماشین، مانند ماشین‌های بردار پشتیبان، برای تولید کلمات تعبیه شده و طبقه بندی عقاید استفاده کرد (دو و همکاران، ۲۰۱۹).

یکی از اولین پژوهش‌های صورت گرفته در زمینه نظرکاوی برای زبان فارسی مربوط به گردآوری مجموعه داده‌ای با نام PersianClues است. این پژوهش با استفاده از یک روش ابتکاری بدون ناظر به تحلیل احساسات می‌پردازد. در واقع تغییری که در این روش صورت گرفته اضافه کردن مجموعه کلمات حاوی بار معنایی به عنوان بردار ویژگی‌ها در مرحله یادگیری است

(شمس و همکاران، ۲۰۱۲)، پژوهش دیگری نیز تحت عنوان ایجاد یک سیستم نظرکاوی با استفاده از الگوریتم‌های ناظر انجام گرفته است. در گام نخست آن، یک لغتنامه احساس برای زبان فارسی به کمک شبکه واژگانی فارسی موجود، فارسی نت، گسترش داده شده است. این پژوهش با استفاده از سه الگوریتم یادگیری ماشین، شامل ماشین بردار پشتیبانی، بیز ساده و رگرسیون منطقی به ارزیابی روش پیشنهادی خود پرداخته است (بصیری و همکاران، ۲۰۱۴).

یکی از روش‌های اعتقادکاوی و تحلیل احساسات بر روی متون فارسی، تحلیل احساسات کاربران بر اساس دیدگاه آنها است که با روش‌های یادگیری عمیق CNN، DNN و LSTM و تکنیک تعبیه کلمات^۳ و TF-IDF به آن پرداخته شده بر اساس نتایج این مقاله روش تعبیه کلمات^۴ با LSTM برای تحلیل احساسات مبتنی بر دیدگاه کاربران نتایج بهتری را داشته است (علی کرمی و همکاران، ۲۰۲۳).

دو و هانگ^۵ مقاله ای تحت عنوان تحقیقات طبقه بندی متن با شبکه‌های عصبی مکرر مبتنی بر توجه را ارائه دادند که روش پیشنهاد شده در این مقاله مزایای هر دو روش را در نظر می‌گیرد. با استفاده از یک مکانیزم توجه بر روی یادگیری وزن برای هر کلمه استفاده می‌شود. در این روش، کلمات کلیدی وزن بیشتری خواهند داشت و کلمات رایج وزن کمتری خواهند داشت. بنابراین، در نظرکاوی متون و اعتقادیابی نه تنها همه کلمات را در نظر می‌گیرد بلکه توجه بیشتری به کلمات کلیدی نیز می‌کند (دو و هانگ، ۲۰۱۸).

یکی دیگر از پژوهش‌های صورت گرفته در زبان فارسی تحت عنوان بهره‌برداری از یادگیری عمیق در تحلیل احساسات است. در این پژوهش از مدل یادگیری عمیق شامل شبکه عصبی کانولوشن استفاده شده است. و در نهایت مدل یادگیری عمیق معرفی شده خود را با روش‌های کم عمق یادگیری ماشین همچون پرسپترون چندلایه مقایسه نموده اند (دشتی پور و همکاران، ۲۰۱۸).

بصیری و کبیری مقاله‌ای تحت عنوان بهبود تجزیه و تحلیل احساسات و اعتقادیابی در زبان فارسی با استفاده از اصلاح واژگان را ارائه دادند، که هدف از تجزیه و تحلیل احساسات^۶ مبتنی بر اصلاح واژگان این است که مشکل استخراج افکار مردم از نظرات آنها در وب را با استفاده از واژگان کلمات از پیش تعریف شده رفع کنند. با این حال، اعتقادکاوی و نظرکاوی برای زبان فارسی در مقایسه با

³ Word Embedding

⁴ Word Embedding

⁵C. Du, L. Huang

⁶ Sentiment Analysis

¹ N. cach dang et al

² bag-of-words

انگلیسی متفاوت است. استفاده از روش مبتنی بر واژگان در فارسی، یک رشته جدید است. منابع محدودی برای تحلیل احساسات و نظرکاوی در زبان فارسی وجود دارد که دقت روش‌های موجود مبتنی بر واژگان پایین‌تر از زبان‌های دیگر است (بصیری و کبیری، ۲۰۱۸).

اورولاگین^۱ مقاله‌ای تحت عنوان یک رویکرد جدید برای تجزیه و تحلیل احساسات و تجسم فکری و طبقه‌بندی خلاصه اخبار ارائه داد که در این کار روش‌های موثر برای استخراج داده‌ها را پیشنهاد می‌کند. در این روش بصورت یک مرور کلی و خلاصه‌ای از متن و تحلیل احساسات می‌تواند احساسات بیان شده در متن را به صورت محاسبات به دست آورد. خلاصه‌سازی متن و تحلیل احساسات بر روی اخبار بی‌بی‌سی انجام شده است. روش خلاصه‌سازی با متن جایگزینی برای تجزیه و تحلیل احساسات مورد استفاده قرار می‌گیرد و طرح‌های تجسم سه بعدی برای نشان دادن اطلاعات احساسات ارائه شده است (اورولاگین، ۲۰۱۸).

کیویو و لی^۲ مقاله‌ای تحت عنوان تجزیه و تحلیل احساسات متن کوتاه در میکرو بلاگ براساس تجزیه وابستگی را ارائه دادند، این روش برای حل مشکلات ارتباط بین کلمات عاطفی و اصلاح کننده و احساسات متن کوتاه از طریق ساختار احساسات و قوانین محاسبه احساسات پیشنهاد شده، که به تحلیل احساسات متن کوتاه کمک می‌کند. احساسات متن کوتاه با توجه به تأثیرات مختلف روابط میان جملات و سهم هر جمله به محاسبه اندیشه متن کوتاه پرداخته می‌شود (کیویو و لی، ۲۰۱۸). همچنین با در نظر گرفتن معنای کلمات در نظر کاوی، با توجه به اینکه برخی از ویژگی‌های انتخابی مناسب نیستند و منجر به افزایش خطاها در طبقه‌بندی می‌شوند، ویژگی‌های بهینه انتخاب می‌شوند و سپس این ویژگی‌ها به ماشین یادگیری داده می‌شود (علی کرمی و همکاران، ۲۰۱۹).

در زمینه تحلیل احساسات، عنان روشی را برای اعتقاد یابی و تفکر کاوی در سطح کلمه، جمله و سند برای نظرات مصرف کنندگان به زبان چینی ارائه کرد. (عنان^۳، ۲۰۲۰)، او یک رویکرد مجموعه‌ای برای انتخاب ویژگی ارائه کرد، که چندین لیست ویژگی‌های فردی را که با روش‌های مختلف انتخاب ویژگی به دست آمده اند جمع می‌کند تا زیر مجموعه ویژگی‌های قوی‌تر و کارآمدتری به دست آید. تا تجزیه و تحلیل احساسات در مورد بررسی محصول بر اساس جاسازی کلمات وزن دار و شبکه‌های عصبی عمیق انجام دهد (عنان، ۲۰۲۲).

یک معماری شبکه عصبی کانولوشن دو طرفه توسط وانگ پیشنهاد شد که از دو لایه دو طرفه LSTM و GRU استفاده می‌کند تا با اتصال دو لایه پنهان از جهت‌های مخالف به یک زمینه، هم زمینه‌های گذشته و هم آینده را استخراج کند (وانگ و همکاران، ۲۰۱۹).

وانگ و همکاران یک روش افزایش بازنمایی خلاف واقع (CRA) را برای آموزش و طبقه‌بندی احساسات در اعتقاد یابی نظرات کاربران ارائه کردند تا نتیجه عملکرد تعمیم دامنه هدف را بهبود بخشند (وانگ و وان^۴، ۲۰۲۲).

سیجی مای و همکاران یک چارچوب جدید با نام HyCon برای یادگیری متضاد ترکیبی نمایش سه وجهی برای احساسات چندوجهی ارائه کردند (مای و همکاران^۵، ۲۰۲۲).

ترکی و روی یک روش تشخیص سخنان منفی را با استفاده از تجسم ابری کلمه و یادگیری گروهی با بردار شمارنده ارائه کردند، آنها از یک چارچوب محاسباتی استفاده کردند که تکنیک‌های تقویت داده‌ها را برای بازنمایی و عملکرد بهتر استفاده کردند (ترکی و روی^۶، ۲۰۲۲).

4. روش پیشنهادی یادگیری عمیق با تفکیک

شدت احساس کلمات

در روش پیشنهادی مجموعه داده توییت‌های فارسی برای تحلیل احساسات و اعتقاد کاوی نظرات کاربران استفاده شده است. تحلیل احساسات براساس قطبیت انجام شده است، قطبیت جملات در این پیکره به صورت عددی بین +۲ و -۲ نمایش داده شده اند که عدد کوچکتر نشانگر قطبیت کمتر (بار منفی بیشتر) است.

هدف اصلی در متن کاوی، تحلیل احساسات کاربران بر اساس قطبیت، تشخیص شدت احساسات با تفکیک احساسات-کلمه است. در این مقاله تحلیل احساسات با روش پیشنهادی یادگیری عمیق مبتنی بر احساسات-کلمه انجام می‌شود، که برای آنالیز دقیق و صحیح تر نظرات در دو مرحله زیر انجام می‌شود:

در مرحله اول احساسات کاربران در دسته‌های مثبت و منفی در باز +۲ و -۲ با در نظر گرفتن کلمه-احساس مشخص می‌شود. به صورتی که کلمات با استفاده از مدل اسکپ گرام آموزش داده می‌شود.

در مرحله دوم واژگان بدست آمده از بدنه اصلی اسناد با چند واژگان مقایسه می‌شود تا قطبیت و شدت احساسات مبتنی بر واژگان را اضافه کند.

⁴ K. Wang and X. Wan

⁵ S. Mai, Y. Zeng, S. Zheng and H. Hu

⁶ T. Turki and S.S. Roy

¹ Siddhaling Urologin

² Lirong Qiu, Jie Li

³ Onan

برای به دست آوردن بازنمایی سند برای لایه بعدی در معماری CNN_WSD , $LSTM_WSD$ ، میانگین مقادیر بردارها از ماتریس وزن تعبیه شده محاسبه شد.

برای تکمیل بازنمایی احساس-کلمه با قطبیت و شدت احساسات، از چندین واژگان احساساتی از پیش تعریف شده استفاده کردیم. برای به دست آوردن یک ارزیابی احساساتی معتبر، پیشنهاد می شود به یک واژگان واحد اعتماد نکنید. علاوه بر این، ترکیبی از شاخص های احساسات مبتنی بر واژگان، مشکل حساسیت به عقاید غیرمستقیم را که معمولاً در مدلهای یادگیری ماشین وجود دارد، غلبه می کند. برای محاسبه قطبیت احساسات، از دو واژگان دست ساز کلمات مثبت و منفی استفاده کردیم (بصیری و همکاران، ۲۰۱۸).

یکی از کمبودهای این واژگان این است که وزن مساوی بدون توجه به شدت احساس آنها به همه کلمات اختصاص می یابد. برای پرداختن به این مسئله، ما شاخص های شدت احساسات به دست آمده از واژگان زیر را با نقاط قوت احساساتی از پیش آموزش دیده درج کردیم: بنابراین، نمرات کلی مثبت و منفی را می توان برای هر واژگان محاسبه کرد. علاوه بر این، ترکیب چندین واژگان، پوشش واژگان بالاتری را تضمین می کند (هاجک و همکاران، ۲۰۲۰).

برای به دست آوردن نمایندگی n -gram، وزن هر n -gram به شرح زیر محاسبه می شود:

$$\omega_{ij} = (1 + \log(tf_{ij})) \times \log(N/df_i) \quad (3)$$

جایی که ω_{ij} وزن نام n -gram را در سند j -ام (بررسی) بیان می کند. df_i و tf_{ij} ؛ $j = 1, 2, \dots, N$ ؛ به ترتیب بیانگر فرکانس اصطلاحات و اسناد^۴ هستند. بنابراین، طول بررسی در نظر گرفته می شود، و وزن نسبتاً بالاتری به n -gram نادر اختصاص می یابد. برای پردازش بیشتر، n -gram ها با توجه به وزن آنها رتبه بندی می شوند، و n -gram های برتر برای ورود به لایه بازنویسی اسناد در معماری CNN_WSD , $LSTM_WSD$ انتخاب می شوند.

برای آنالیز دادهها پیش پردازش بر روی اسناد آموزشی انجام شده و داده های گسیخته^۱ حذف شده، و بردارهای مبتنی بر $tf-idf^2$ از متن ورودی به دست آمده و بردارهای حاصله، نرمال سازی می شوند. در مرحله پیش پردازش کاهش ابعاد بر روی متون انجام می گردد (هاجک و همکاران^۳، ۲۰۲۰). در نهایت بعد از پیش پردازش که در زیر جزئیات آن شرح داده شده، داده های پیش پردازش شده به عنوان داده های آموزشی برای اعتقاد کاوی و رده بندی نظرات روش پیشنهادی اعمال می گردد.

در تحلیل احساسات و دسته بندی خودکار نظرات از یادگیری عمیق مبتنی بر احساس-کلمه CNN_WSD , $LSTM_WSD$ بر روی اسناد انجام می شود. CNN (Convolutional Neural Networks), $LSTM$ (Long-Short-Term Memory) با دو لایه پنهان مترام برای پردازش در ویژگی های ورودی استفاده می شود. این دو نمایندگی شامل احساس-کلمه و n -gram است.

بازنمایی کلمه- احساس در دو مرحله تولید می شود: در مرحله اول، تعبیه کلمات با استفاده از مدل اسکپ گرام آموزش داده می شود. مرحله دوم، واژگان به دست آمده از بدنه اصلی اسناد ورودی با چند واژگان مقایسه شده است تا قطبیت و شدت احساسات مبتنی بر واژگان را اضافه کند (علی کرمی و همکاران، ۲۰۲۳).

برای محاسبه ماتریس وزن تعبیه شده، عملکرد تعبیه شده برای هر کلمه w_t در واژگان اعمال می شود. تابع تعبیه برای دنباله $W = \{w_1, w_2, \dots, w_t, \dots, w_T\}$ تا عملکرد تابع هدف زیر به حداکثر برسد:

$$E = \frac{1}{T} \sum_{t=1}^T \sum_{-c \leq j \leq c} \text{Logp}(W_t + j | W_t) \quad (1)$$

c نمایانگر شعاع پنجره متن است (چند کلمه اطراف را در نظر می گیریم) و $p(W_t + 1 | W_t)$ احتمال کلمه خروجی با توجه به کلمات ورودی محاسبه شده با استفاده از الگوریتم سلسله مراتبی (softmax) است.

$$\rho(\omega | \omega I) = \prod_{j=1}^{L(\omega)-1} \sigma([n(\omega, j+1) = ch(n(\omega, j))] v_{n(\omega, j)}^T v_{\omega I}) \quad (2)$$

جایی که ωI و ωO به ترتیب کلمات ورودی و خروجی هستند. v^w و v^w به ترتیب نشانگرهای بردار کلمات ورودی و خروجی را نشان می دهند. $n(w, j)$ گره j ام در درخت باینری است. $L(w)$ طول مسیر در درخت است. $ch(n)$ یک گره فرزند را نشان می دهد. $\sigma(x)$ یک تابع عملکرد سیگموئیدی را مشخص می کند، در صورتی که اگر x صحیح باشد، $[[x]] = 1$ ؛ در غیر این صورت $[[x]] = -1$.

¹Outlier

²term frequency-inverse document frequency

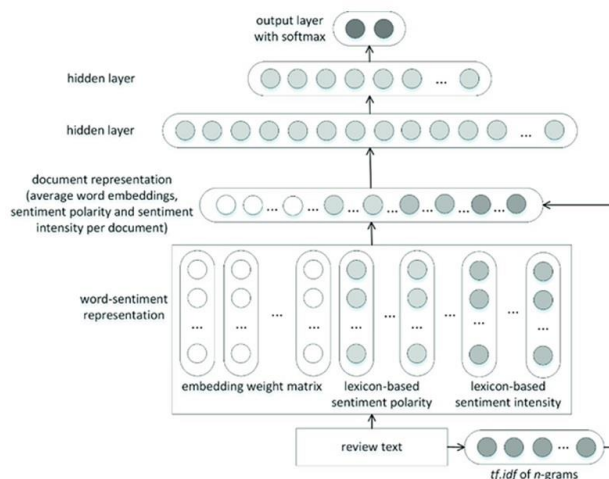
³Hajek P., Barushka A., Munk M.

⁴ term and document frequency

```

while i < n do
if P(Wo|Wi<
vw is the vector marker of the input words and v'w
represent
the vector marker of the output words.
x= Softmax(Wt + j) // Equation (2)
If x is integer, [[x]] = 1; Otherwise [[x]] = -1.
If x is correct, [[x]] = 1; Otherwise [[x]] = -1.
end if where (Wij weight of each n-gram in the document)
// Equation(3)
.Repeating number or using the frequency product of each
word
tf_idf // Equation (4)
i= i + 1
Obtain well-trained sentiment- word embeddings(WSD):
using
synonym Word and the words of the context around it.
end while
return.

```



شکل ۲- معماری پیشنهادی برای استخراج نظرات و اعتقادکاوی

دو لایه پنهان بعدی برای پردازش رابطه پیچیده بین نمایش اسناد و کلاسهای مثبت / منفی از خروجی استفاده می شود. برای جلوگیری از اضافه کردن اتصالات و اثربخشی آموزش، به ترتیب از یکپارچه سازی ترک خوردگی 0.2 و 0.5 برای ورودی و دو لایه پنهان استفاده کردیم. الگوریتم نزول شیب مینی دسته ای با $b = 100$ مینی دسته^۱، میزان یادگیری 0.1 و 1000 تکرار، رفتار همگرایی خوب و پایدار را برای ما فراهم کرده است. اعداد مختلف nh_2 و nh_1 در دو لایه پنهان $\{27, 26, 25, 24\}$ برای به دست آوردن معماری مطلوب مورد آزمایش قرار گرفت. همانطور که در زیر ارائه می شود، بهترین نتایج برای $nh_1 = 25$ و $nh_2 = 24$ نوروں بدست آمد. توجه داشته باشید که ما همچنین با یک لایه پنهان آزمایش کرده ایم اما بدون پیشرفت. عملکرد هدف با از دست دادن آنتروپی متقاطع نشان داده شد. پیچیدگی کلی مدل ارائه شده می تواند به صورت:

$O(b \times I \times (m \times nh_1 + nh_1 \times nh_2 + nh_2 \times nO))$ بیان شود، که I تعداد تکرارها را نشان می دهد. m تعداد ویژگی های موجود در لایه ارائه اسناد را نشان می دهد. nh_1, nh_2 و nO به ترتیب تعداد نوروں ها در لایه های پنهان اول و دوم و لایه خروجی را نشان می دهند.

الگوریتم ۱. الگوریتم روش پیشنهادی تحلیل احساسات

Parameters: $n(w, j)$ is the j th node in the binary tree. $L(w)$ is the path length in the tree. $ch(n)$ is a child node. $\sigma(x)$ specifies a sigmoid function.

Input: Set of words in a sentiment lexicon $w = \{w_1, w_2, \dots, w_n\}$

Output: Well-trained sentiment- word embeddings Et Initialization

$w \text{ embed}(w) // \text{Equation (1)}$

$i = 0, n = \text{total number of documents}$

1 mini-batches

مراحل طراحی روش پیشنهادی در زیر بیان شده:

1.4 پیش پردازش و خلاصه سازی اولیه

در حقیقت، پیش پردازش وظیفه نگاشت متن داده شده به یک نمای منطقی را بر عهده دارد. به عبارت دیگر استخراج ویژگی و وزندهی و کاهش ابعاد در این قسمت انجام می گیرد. بسته به کاربرد استخراج ویژگی می تواند بسیار ساده و یا بسیار مفصل باشد. تحلیل واژگانی شامل عملیات مربوط به یکسان سازی متن، قواعد مربوط به نشانه گذاریها و مرزبندی بین کلمات می باشد. بعد از این مرحله عموماً دسته ای از کلمات بی ارزش که متناوباً تکرار می شوند و بار معنایی خاصی ندارند: مانند حرف ربط ("و"، "که"، "تا"، "وقتی که"، "اگر"، "اما"، "اینکه")، حرف اضافه ("به"، "با"، "از"، "در")، فعل ربطی ("است"، "بود"، "شد") و حرف تعریف ("یک" در "یک دانشجوی نمونه کسی است که ...") از متن داده شده حذف می شوند. سپس با استفاده از الگوریتم های ریشه یابی، به منظور بهینه سازی ویژگی های استخراج شده، کلمات ریشه یابی می شوند. در نهایت با استفاده از گروه های اسمی کلمات دسته بندی می گردند.

الگوریتم ۲. نحوه تخصیص شدت احساس کلمات

- 1: A set of words in Vocabulary = $(W_1, W_2, W_3, \dots, W_V)$
- 2: A dictionary of input words and their corresponding sentiment polarity
- 3: // Get sentiment polarity of words based on SentiWordNet and Sentipers lexicons
- 4: $\text{Sentiment_class} = \{ \} // \text{Dictionary of words and their corresponding polarity class}$
- 5: for each W in Vocabulary:

کلاس‌های کمتر و تعداد رخداد بیشتر آن در یک کلاس، کلماتی را به عنوان ویژگی نهایی انتخاب کنند که هر چه بیشتر نماینده یک کلاس خاص باشند.

2.1.4 بازنمایی متون با روش تعبیه واژه های عصبی

برای وزن‌دهی به ویژگی‌ها می‌توان از رویکردهای متفاوتی بهره برد. در ساده‌ترین حالت این وزن‌دهی می‌تواند به صورت باینری انجام شود. انتخاب دیگر وزن‌دهی به هر کلمه با توجه به تعداد تکرار هر کلمه و استفاده از حاصلضرب فرکانس هر کلمه در معکوس فرکانس سند که معمولاً به صورت زیر تعریف می‌شود:

$$TF - IDF(ti, dj) = tf(ti, dj) \times \log(N/(N(ti))) \quad (5)$$

ti تعداد اسنادی از مجموعه $N(ti)$ و N نماینده تعداد کل اسناد می باشد.

در نظر کاوی داده ها با یادگیری عمیق ابتدا داده های متنی باید با روش های رمزنگاری به عدد تبدیل شود، زیرا هر مدل پایه ریاضی دارد. با وجود روش های برداری فوق که اغلب در الگوریتم های یادگیری ماشین سنتی استفاده می شود، ما از تعبیه واژه های عصبی در مدل های یادگیری عمیق بهره مند می شویم که هر کلمه به یک بردار کم ابعاد به نام ویژگی کلمه تبدیل می شود (تورین و همکاران ، 2010). تعریف کلمات در این نمایش باعث می شود که کلمات مشابه مستقیماً پیدا شوند. بنابراین می توانیم از این مزیت برای کاوش نظر نیز استفاده کنیم. به طور خلاصه ، هر جمله رمزگذاری می شود، کلمات بعدی برداری می شوند. در نتیجه ، هرچه بردارهای مشابه بیشتری در میان جملات قرار بگیرند ، شباهت بیشتری خواهند داشت. شایان ذکر است که تعبیه کلمه عصبی نه تنها یک ردیاب مترادف است بلکه یک روش برای یافتن کلمات از یک خانواده (به عنوان مثال گربه ، سگ) است. این بردارها به دو صورت لایه جاسازی آنلاین و جاسازی کلمات از قبل آموزش داده شده انجام می شود، روش لایه جاسازی آنلاین به مجموعه داده موجود متکی است و در فرایند یادگیری عصبی عملی خواهد شد. در حقیقت ، بردارهای خروجی از ورودی با استفاده از هیچگونه عمل ریاضی محاسبه نمی شوند. بنابراین ، هر کلمه در جملات همانطور که ظاهر می شود با یک عدد صحیح رمزگذاری می شود. در این حالت ، Vs تعداد کلمات مجموعه واژگان را نشان دهد و Ev بعد تعبیه بردارها را نشان می دهد. سپس ، پس از آموزش شبکه عصبی ، انتظار داریم یک بردار تعبیه شده در اندازه به شرح زیر باشد. در این میان ، شکل ۳ نحوه آنلاین بودن را نشان می دهد.

```

6: Score = 0
7: Word = Lemmatization(W) // Lemmatize the word
8: Synsets = GetSynsets_SentiWordNet(Word) // Obtain the
  synsets of a word from SentiWordNet
9: if length (Synsets) > 0 do // That means the word exists in
  SentiWordNet lexicon
10: Score = Average( Synsets.positive_scores) - Average(
  Synsets.negative_scores)
11: else
12: Score = getPolarity_SenticNet(Word)
13: if Not Score do // That means the word does not exist in the
  SenticNet lexicon
14: Score = getPolarity_VADER(Word)
15: end if
16: end if
17: // Assign sentiment class based on the score obtained before
18: append (W: "Strong Negative") to Sentiment_class ifScore <= -1
19: append (W: "Negative") to Sentiment_class if-1 < Score < 0
20: append (W: "Neutral") to Sentiment_class ifScore == 0
21: append (W: "Positive") to Sentiment_class if0 < Score < 1
22: append (W: "Strong Positive") to Sentiment_class ifScore >= 1
23: end for return Sentiment_class

```

همانطور که در الگوریتم 2 خلاصه شده است، کلمه مورد نظر با ارزش امتیاز احساسات آن از مجموعه داده جایگزین شده است. از طرف دیگر، از آنجایی که کلمه ممکن است در قطبیت های مختلف در SentiWordNet و Sentpres ظاهر شود، از تفاوت میانگین امتیازات مثبت و منفی استفاده کردیم. این با فرمول زیر نشان داده می شود:

$$Score = \sum_{i=1}^k Synsetp(i)k - \sum_{i=1}^k Synsetn(i)k \quad (4)$$

در فرمول 4 مقدار k به تعداد ظاهر کلمه اشاره دارد، $Synsetp$ نشان دهنده امتیاز مثبت و $Synsetn$ نشان دهنده امتیاز منفی است.

1.1.4 انتخاب ویژگی اولیه

انتخاب ویژگی، تکنیکی است که برای مواجهه با داده‌های با ابعاد بالا استفاده می‌شود. مسئله اصلی این است که داده‌هایی که دارای ابعاد بالا هستند زمان بیشتری را برای پردازش صرف می‌کنند. یکی از راه‌های کم کردن زمان محاسبات، انتخاب ویژگی‌هایی از فضای مسئله است که در تعیین جوابها موثر هستند و از باقی ویژگی‌ها صرف نظر می‌شود. بدین صورت داده‌هایی با ابعاد کمتر به وجود می‌آیند که بعد از انجام اعمالی نظیر رده‌بندی، جواب‌هایی مشابه داده‌های اولیه تولید می‌کنند، در روش پیشنهاد شده برای حل این مسئله از پارامترهای معکوس فرکانس مطابقت (ICF) استفاده می‌کند و مربوط به هر کلمه uni و یکنواختی (ICF) سپس کلمه‌ای را انتخاب می‌کند که حد آستانه کوچکتر از ICF مشخص و بزرگتر از حد آستانه تعیین شده داشته باشد. در uni حقیقت این پارامترها سعی می‌کنند با امتیاز دادن به تکرار کلمه در

در این مقاله از مجموعه داده توئیتر فارسی استفاده شده، نظرات فارسی در مجموعه داده شبکه اجتماعی توئیتر شامل ۱۰۰۰۰۰۰ سند است، برای هر دسته نظرکاوی با مثبت و منفی بودن نظر کاربران نمایش داده می‌شود. مجموعه داده توئیتر شامل فیلدهای زیر می باشد:



شکل ۳: لایه جاسازی آنلاین

2.4 اعتقادکاوی

در اعتقادکاوی، هم بازخوردهای مثبت و هم منفی مجموعه داده استفاده شده که نظرات فارسی کاربران در توئیتر می باشد را با استفاده از روش پیشنهادی CNN_WSD, LSTM_WSD رده بندی و درصد بازخورد مثبت و منفی نظرات کاربران در مورد هر یک از موضوعات استخراج می‌شود تا مشخص شود که آیا از نظر کاربران درصد بازخورد مثبت هریک از موضوعات بیشتر است یا درصد بازخورد منفی، درصد مثبت و منفی هر یک از موضوعات با توجه به مجموع فراوانی نسبی تکرار هر یک از لغات مثبت(زیبا، عالی، دوست داشتنی و غیره) و منفی(بد، زشت، نامناسب، افتضاح و غیره) مربوط به هر موضوع محاسبه می‌شود. در مجموعه داده آموزشی هر جمله با ویژگی‌هایی مانند قطب، کلمات کلیدی و اهداف شرح داده می‌شود.

جدول ۱: مجموعه ای از احساسات رتبه بندی شده بر اساس قطبیت

Ranked	Emotion Class
-2	Furious
-1	Angry
0	Neutral
+1	Happy
+2	Delighted

قطبیت هر جمله از یک مجموعه E بین بازه +2 و -2 انتخاب شده است. $E = \{-2, -1, 0, +1, +2\}$ که در جدول ۱ نمایش داده شده.

3.4 پردازش داده

عملکرد یادگیری ماشین و یادگیری عمیق اغلب به اندازه و کیفیت داده های آموزشی بستگی دارد که جمع آوری آنها اغلب خسته کننده است (وی و زو ، 2019). بنابراین، از چندین مجموعه داده متعادل استفاده شده است.

5. نتایج و آزمایشات

در اعتقادکاوی هر کلمه به عنوان یک ویژگی در داده‌های آموزشی برای رده‌بندی در نظر گرفته می‌شود و کلمات در متون فارسی بر اساس نرخ تاثیر در نظرکاوی در رده‌بندی اندازه‌گیری می‌شود و بر اساس میزان دقت Acc و زمان رده‌بندی اندازه‌گیری می‌شود.

این آزمایشات بر روی مجموعه داده توئیتر با ۱۰۰۰۰۰۰ سند انجام شده است. ما از مدل‌های یادگیری عمیق مبتنی بر شدت احساسات کلمات با TF-IDF و تعبیه کلمات بر روی مجموعه داده های توئیتر استفاده کرده ایم. برای یادگیری تعبیه کلمات (word embeddings)، از مدل Skip-Gram برای آموزش داده ها استفاده شده است.

1.5 مجموعه داده توئیتر

نظرات فارسی در مجموعه داده شبکه اجتماعی توئیتر شامل ۱۰۰۰۰۰۰ سند است، برای هر دسته نظرکاوی با مثبت و منفی بودن نظر کاربران نمایش داده می‌شود. مجموعه داده توئیتر شامل فیلدهای زیر می باشد:

- “target” is the polarity of the tweet;
- “id” is the unique ID of each tweet;
- “date” is the date of the tweet;
- “query_string” indicates whether the tweet has been collected with any particular query keyword (for this column, 100% of the entries labeled are with the value “NO_QUERY”);
- “user” is the Twitter handle name of the user who tweeted;
- “text” is the verbatim text of the tweet.

شکل ۴: نمونه مجموعه داده توئیتر

ما از فیلدهای target, text برای انجام آزمایشات استفاده می‌کنیم.

جدول ۲. رتبه بندی جملات مجموعه داده توئیتر

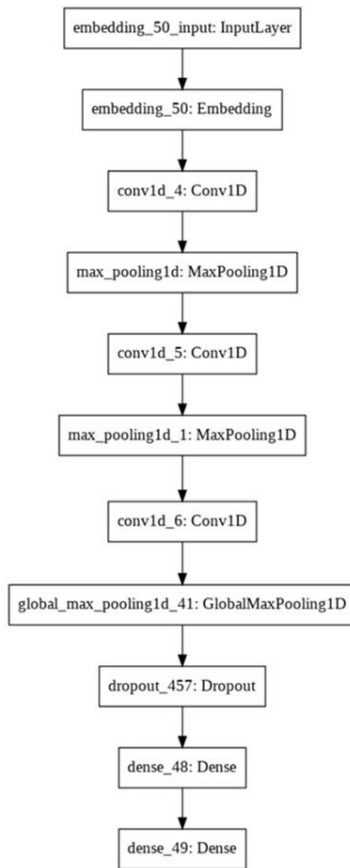
Ranked	-2	-1	0	+1	+2
Sentences	5394	93998	425084	294538	180986

تعداد جمله های تمام دسته ها ۱۰۰۰۰۰۰ می باشد که تعداد داده های هر دسته نیز در جدول ۲ مشخص شده است. (علی کرمی و همکاران، ۲۰۲۳)

1.5 پارامترها

از آنجا که هر متن به صورت مجموعه ای از کلمات تعبیه شده در آمده، در لایه اول مدل های شبکه عصبی خود، تعداد نرون ها به اندازه بیشترین طول جملات بر حسب کلمه است. در این مجموعه داده، طولانی ترین متن موجود شامل ۲۵۷ کلمه بوده و بنابراین در لایه اول ۲۵۷ نرون خواهیم داشت.

همانطور که در شکل ۵ نشان داده شده است، ما با تنظیمات مختلف مدل را آزمایش کرده ایم. در تعبیه کلمات به کمک Keras، در مدل های از قبل تعیین شده بهترین عملکرد با تعبیه کلمه $embed_size = 300$ بدست آمد. از آنجا که یادگیری تعبیه در خود شبکه عصبی صورت میگیرد این ابعاد $num_words = 2000$ در نظر گرفته شده است. مدل اسکپ گرام در محیط پایتون پایتون با استفاده از سیستم یادگیری عمیق کلب گوگل^۱ آموزش داده شد.

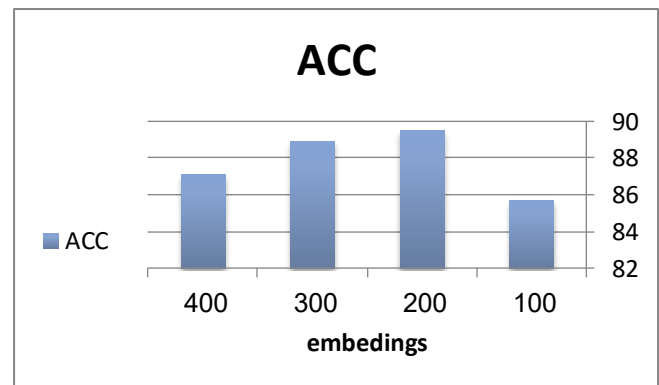


شکل ۶ ساختار لایه های در مدل های پیشنهادی LSTM-WSD

در این قسمت ساختار CNN مبتنی بر تعبیه کلمات (WSD) برای لایه های شبکه عصبی استفاده شده است، در روش پیشنهادی شبکه عصبی کانولوشن (CNN) است که یکی از موفق ترین ساختارهای شبکه های عصبی است این مدل برای داده های متنی خصوصاً در مسائل طبقه بندی متن نیز به خوبی عمل میکند تنظیمات CNN در روش استفاده شده در شکل ۷، قابل مشاهده است.

Layer (type)	Output Shape	Param #
embedding_1 (Embedding)	(None, 40, 300)	4500300
conv1d_1 (Conv1D)	(None, 40, 64)	57664
conv1d_2 (Conv1D)	(None, 40, 32)	6176
max_pooling1d_1 (MaxPooling1D)	(None, 13, 32)	0
conv1d_3 (Conv1D)	(None, 13, 16)	1552
conv1d_4 (Conv1D)	(None, 13, 8)	264
global_average_pooling1d_1 (GlobalAveragePooling1D)	(None, 8)	0
dense_1 (Dense)	(None, 1)	9
Total params: 4,565,965		
Trainable params: 65,665		
Non-trainable params: 4,500,300		

شکل ۷. ساختار لایه های در مدل های پیشنهادی CNN



شکل ۵. تأثیر تعداد تعبیه کلمات بر عملکرد مدل CNN_WSD، LSTM_WSD با دو لایه پنهان از $nh1 = 25$ و $nh2 = 24$ نرون.

2.5 مدل ها

در این قسمت از ساختار LSTM مبتنی بر تعبیه کلمات برای لایه های شبکه عصبی استفاده شده است. نخستین ساختار حافظه طولانی کوتاه مدت دوطرفه (LSTM) است که بر پایه شبکه های عصبی بازگشتی طراحی شده و دوطرفه بودن آن امکان دریافت اطلاعات توسط گذشته و آینده را به لایه خروجی آن اضافه میکند. در شکل ۶ لایه های تعبیه شده برای این ساختار قابل مشاهده است.

3.5 نتایج

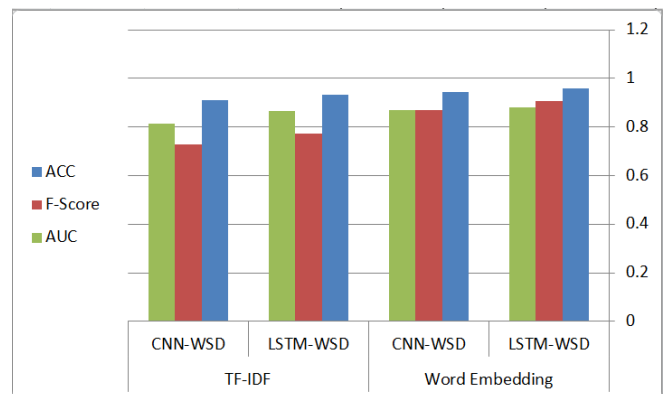
در آزمایشات ما، سه روش ارزیابی بر روی مجموعه داده توپیتر در نظر گرفته می شود: دقت (Acc)، مساحت زیر منحنی (AUC) و F-Score. برای ارزیابی عملکرد مدل پیشنهادی، از روش k-fold cross validation استفاده شده که k برابر با ۱۰ می باشد.

جدول ۳. نتایج آزمایش های ACC، F-Score و AUC برای توپیتر مجموعه داده ها.

Experiments	TF-IDF		Word Embedding	
	CNN-WSD	LSTM-WSD	CNN-WSD	LSTM-WSD
ACC	0.910036	0.934178	0.943248	0.958863
F-Score	0.727000	0.772000	0.871000	0.907000
AUC	0.814408	0.864623	0.868279	0.881397

نتایج تعبیه کلمات برای روش LSTM_WSD : ACC برابر با 95.8 درصد و F-Score برابر با 0.907 و AUC برابر با 0.864 است و روش CNN_WSD با ACC برابر با 94.3 درصد و F-Score برابر با 0.871 و AUC برابر با 0.889 است.

نتایج جدول ۳ نشان میدهد که استفاده از تعبیه کلمات برای تحلیل احساسات با یادگیری عمیق نتایج بهتری نسبت به TF-IDF برای ارزیابی دقت (Acc)، مساحت زیر منحنی (AUC) و F-Score دارد.



شکل ۸- نمودار ارزیابی روش های یادگیری عمیق در مجموعه داده توپیتر

نتایج شکل ۸ نشان میدهد، بهترین رفتار با ترکیب LSTM و تعبیه کلمات (Word embedding) نشان داده می شود و بالاترین مقادیر ACC، F_Score و AUC با تعبیه کلمات LSTM + داده شد. ما می توانیم تأیید کنیم که تعبیه کلمات یک تکنیک مناسب تر از TF-IDF برای انجام تجزیه و تحلیل احساسات است، پس از تجزیه و تحلیل نتایج مربوط به کیفیت پیش بینی ها، لازم است اطلاعاتی در مورد هزینه محاسباتی مرتبط با القاء مدل ها به دست آوریم.

جدول ۴. زمان آزمایشات مدل های مختلف روش پیشنهادی با پردازنده

گرافیکی GPU

Experiments	TF-IDF		Word Embedding	
	CNN-WSD	LSTM-WSD	CNN-WSD	LSTM-WSD
Tweeter Dataset	2min 48s	23min 52s	8min 01s	10min 39s

جدول ۴ زمان پردازش مورد نیاز برای القاء مدل ها از مجموعه داده توپیتر را نشان می دهد و شامل زمان CPU مورد نیاز در آزمایش ها است.

جدول ۴ نشان می دهد که استفاده از TF-IDF، به زمان محاسبه طولانی تری نسبت به استفاده از تعبیه کلمات (Word Embedding) نیاز دارد. این یکی دیگر از دلایلی است که تکنیک تعبیه کلمات را توصیه می کنیم. با این حال، LSTM با TF-IDF و هم با تعبیه کلمات وقت گیرترین الگوریتم است. در نهایت، خلاصه کلی از نتایج بدست آمده در آزمایشات را در زیر توضیح می دهیم:

دو مدل یادگیری عمیق (CNN و LSTM) با تفکیک شدت شدت احساسات کلمات برای انجام آزمایش های تجزیه و تحلیل احساسات استفاده شد. مشخص شد که مدل LSTM بهترین نتیجه بین زمان پردازش و دقت نتایج را ارائه می دهد. اگرچه مدل LSTM در هنگام استفاده از تعبیه کلمات از بالاترین درجه دقت برخوردار بود، اما زمان پردازش آن بیشتر از مدل CNN بود. مدل LSTM هنگام استفاده از تکنیک TF-IDF موثر نیست و زمان پردازش بسیار بیشتر آن منجر به نتایجی می شود که بهتر نیستند.

جدول ۵. مقایسه و ارزیابی نتایج روش پیشنهادی با سایر مدل های

یادگیری عمیق

Model	AUC	Acc [%]	F-score	Testing time [s]
LSTM_WSD (this study)	88.13 ± 0.91	0.958 ± 0.011	0.907 ± 0.009	10.398 ± 0.226
CNN_WSD, (this study)	86.82 ± 0.91	0.943 ± 0.005	0.871 ± 0.006	8.012 ± 0.131
LSTM Alikarami et al.	84.05 ± 0.28	0.917 ± 0.003	0.841 ± 0.002	2.042 ± 0.128
CNN Alikarami et al.	84.29 ± 0.17	0.921 ± 0.001	0.844 ± 0.003	8.139 ± 0.286
DNN Alharbi et al.	84.14 ± 0.71	0.842 ± 0.002	0.828 ± 0.005	8.14 ± 0. 31
Tree-LSTM	86.06 ±	0.867 ±	0.851 ±	4.15 ± 0.

برای مدل های تعبیه کلمات و TF-IDF را با روش DNN-WSD (حساس به کلمه) مقایسه می کنیم. نتایج جدول ۶ نشان میدهد در رفتار هر دو مدل تعبیه کلمات و TF-IDF مدل مورد ارزیابی DNN-WSD برای تحلیل احساسات بر روی مجموعه داده توپیتر عملکرد

جدول ۶. نتایج آزمایش های ACC، F-Score و AUC برای توپیتر مجموعه داده ها.

Experiments	TF-IDF		Word Embedding	
	DNN-WSD	LSTM-WSD	DNN-WSD	LSTM-WSD
ACC	0.757757	0.934178	0.790962	0.958863
F-Score	0.763832	0.772000	0.788766	0.907000
AUC	0.764996	0.864623	0.788166	0.881397

مناسبی ندارد و مدل های LSTM و CNN عملکرد بهتری نسبت به DNN عملکرد بهتری دارند.

جدول ۴ و ۶ نشان می دهد که استفاده از TF-IDF، به زمان محاسبه طولانی تری نسبت به استفاده از تعبیه کلمات نیاز دارد. این یکی دیگر از دلایلی است که تکنیک تعبیه کلمات را توصیه می کنیم. با این حال، LSTM وقت گیرترین الگوریتم است، هم با TF-IDF و هم با Word Embedding. با توجه به اینکه بهبود LSTM نسبت به DNN و CNN در مورد اخیر چندان قابل توجه نیست، استفاده از این دو روش را می توان از نظر کاهش زمان و هزینه محاسباتی مناسب تر دانست. ما میدانم که انواع مختلف مجموعه داده ها بر نتایج تجزیه و تحلیل احساسات متفاوت تأثیر می گذارد.

4.5 نتایج

در آزمایشات ما، سه روش ارزیابی بر روی مجموعه داده توپیتر در نظر گرفته می شود: دقت (Acc)، مساحت زیر منحنی (AUC) و F-Score. برای ارزیابی عملکرد مدل پیشنهادی، از روش k-fold cross validation استفاده شده که k برابر با ۱۰ می باشد. LSTM، CNN و DNN برای به دست آوردن بازنمایی سطح معنایی جمله استفاده شدند. در این روش، ابعاد حالت های پنهان / سلول روی ۳۰۰ تنظیم شده است که مربوط به تعداد کلمات تعبیه شده است. معماری CNN از لایه حلقوی (convolutional layer) با پنج فیلتر و حداکثر لایه مخلوط max pooling layer تشکیل شده است. نمایش اسناد برای هر دو مدل به عنوان ترکیب بیان جملات با استفاده از GRU ساخته شد (تانگ و همکاران، ۲۰۱۵).

Model	AUC	Acc [%]	F-score	Testing time [s]
and Discourse-LSTM Kraus et al.	0.28	0.006	0.003	11
GRU, and hybrid approaches Do et al.	85.29 ± 0.16	0.851 ± 0.002	0.854 ± 0.006	11.2 ± 0.31
RNN Abid et al.	84.74 ± 0.91	0.850 ± 0.003	0.838 ± 0.008	9.16 ± 0.21
Coattention-LSTM + Location Yang et al.	87.52 ± 0.19	0.891 ± 0.002	0.871 ± 0.003	5.032 ± 0.18
Bi-LSTM Wu et al.	84.85 ± 0.06	0.871 ± 0.002	0.841 ± 0.003	3.042 ± 0.45

برای ارزیابی جامع اثربخشی مدل پیشنهادی، عملکرد آن را در مقابل مدل های موجود زیر مقایسه کردیم:

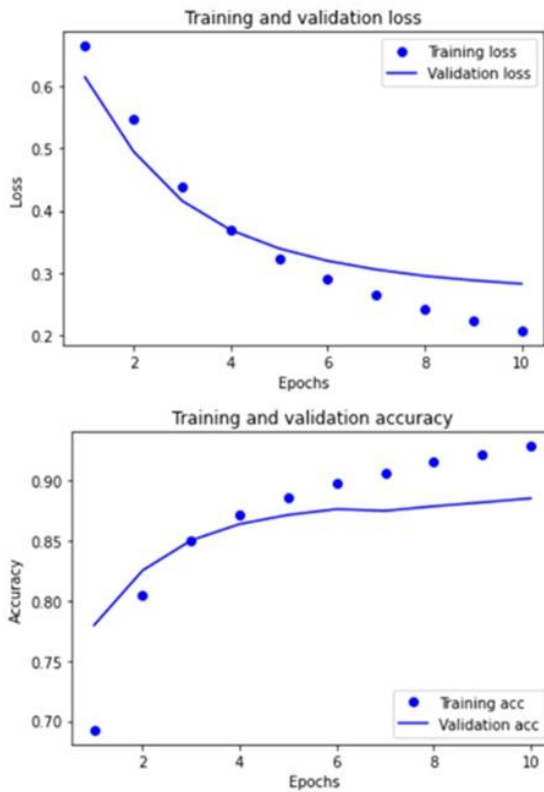
در جدول ۵ روش های مختلف یادگیری عمیق از جمله LSTM، CNN، DNN، Tree-LSTM، GRU، Bi-LSTM و Coattention-LSTM که از سال ۲۰۱۸ تا ۲۰۲۳ توسط محققین برای تحلیل احساسات بر روی مجموعه داده توپیتر ارائه شده است با استفاده از نرم افزار پایتون برای مجموعه داده فارسی توپیتر پیاده سازی و شبیه سازی شده اند و در حالت برابر ارزیابی شده است، با توجه به اینکه روش های یادگیری عمیق عملکرد بهتری نسبت به روش های یادگیری ماشین دارد ما فقط روش های پر کاربرد یادگیری عمیق را مورد ارزیابی و مقایسه قرار داده ایم.

نتایج جدول ۵ نشان میدهد، بهترین رفتار با ترکیب LSTM-WSD و تعبیه کلمات نشان داده می شود و بالاترین مقادیر ACC، F-Score و AUC با تعبیه کلمات LSTM + داده شد. همچنین روش CNN-WSD با تعبیه کلمات نیز نسبت به اکثر روش های مورد مقایسه نتایج بهتری داشته و زمان اجرای مدل CNN نسبت به LSTM بهتر می باشد، نتایج بدست آمده نشان می دهد روش های LSTM، Tree LSTM و Bi LSTM زمان اجرای پایین تری نسبت به مدل های تعبیه کلمات و TF-IDF برای انجام تجزیه و تحلیل احساسات است، اما میزان دقت طبقه بندی این سه مدل به ترتیب ۹۱٪، ۸۶٪ و ۸۴٪ درصد است.

پس از تجزیه و تحلیل نتایج مربوط به کیفیت پیش بینی روش های مورد ارزیابی ما در جدول ۶ نتایج روش LSTM-WSD را

¹ Tang, D., Qin, B., Liu, T.

آزمون‌های معنی‌داری آماری که ما انجام دادیم راهی برای ارزیابی نتایج آزمون‌ها در اختیار ما قرار داد. سطح قابل اعتمادی که معمولاً برای نتایج استفاده می‌شود 95٪ است، همچنین به صورت $p = 0.05$ نوشته می‌شود و به عنوان سطح p شناخته می‌شود. همانطور که در شکل 13 نشان داده شده است، آزمون‌های معناداری آماری را برای آموزش و اعتبار سنجی اضافه کردیم.



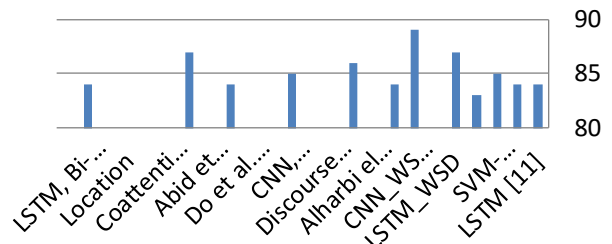
شکل ۱۲. نمودار دقت بر اساس دوره‌ها.

شکل ۱۲ عملکرد مدل پیشنهادی LSTM+WSD را بر اساس ۱۰ دوره نشان می‌دهد. حداکثر میزان دقت برای ۱۰ تکرار 93٪ و میزان خطا 0.189 برای LSTM آموزش مدل تعبیه کلمه است، اما در Validation میزان دقت 88٪ و میزان خطا برای ۱۰ تکرار 0.281 است.

5.5 بحث

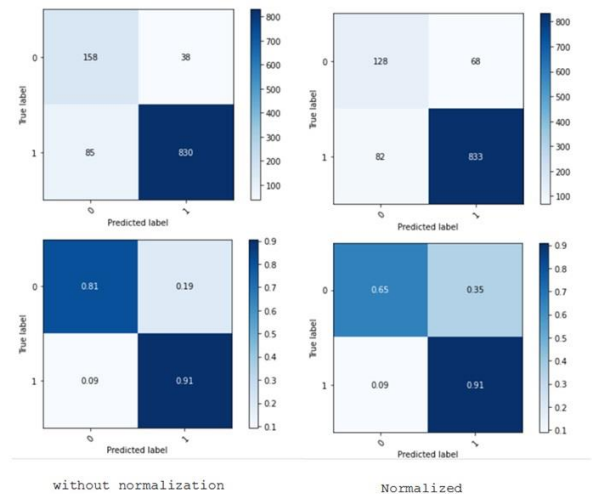
با توجه به جداول 3 تا 6 می‌توان نتیجه گرفت که Recall، دقت، AUC و F-score روش پیشنهادی LSTM و تعبیه کلمه بهتر از سایر روش‌های ارزیابی شده است. این نتایج برای ترکیب LSTM و TF-IDF به ترتیب 0.776، 0.774، 0.641 و 0.569 است. نتایج به دست آمده در مقایسه با تعبیه کلمات و TF-IDF برای مجموعه داده توییتر و سایر مجموعه‌های داده مبتنی بر روش LSTM نشان می‌دهد که نتایج تعبیه کلمات در LSTM بهتر از TF-IDF است.

AUC



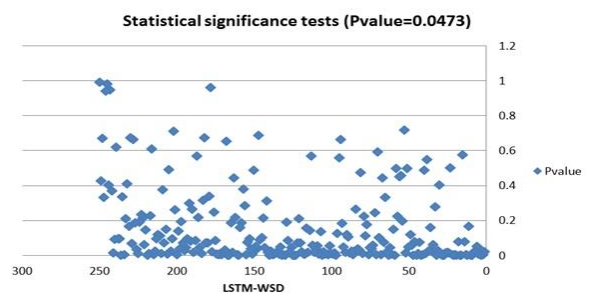
شکل ۹. مقایسه میزان AUC روش پیشنهادی نسبت به سایر روش‌ها

نتایج شکل ۹ نشان می‌دهد که میزان مساحت زیر منحنی AUC روش LSTM_WSD برابر با 86.49، CNN_WSD برابر با 88.92 است.



شکل ۱۰. نمودار سردرگمی برای LSTM+WSD

نمودار سردرگمی (Confusion) بر اساس سطح دقت مدل پیشنهادی LSTM و جاسازی کلمه برای مجموعه داده SentiPers طراحی شده است. این نمودار به دو صورت نرمال شده و بدون نرمال سازی طراحی شده است. به منظور ارزیابی عملکرد مدل LSTM و جاسازی کلمه، به صورت گرافیکی بر اساس 10 دوره از مجموعه داده SentiPers نمایش داده می‌شود.



شکل ۱۱. نمودار Pvalue برای آزمون‌های معناداری آماری.

همچنین برای ترکیب CNN + جاسازی کلمه و DNN + جاسازی کلمه، 0.791 و 0.802 Recall، 0.791 و 0.802 Dقت، 0.791 FScore، 0.801 و 0.788 است.

نتایج به دست آمده نشان می‌دهد که روش تعبیه کلمه و LSTM با دقت 93 درصد و درصد خطای 0.189 نسبت به سایر روش‌های ارزیابی شده عملکرد بهتری دارند. اما زمان اجرای LSTM و تعبیه کلمه بیشتر از CNN و DNN است.

در این مقاله محدودیت‌های پژوهشی وجود دارد. ابتدا، ما کلیدواژه‌های اصلی با فرکانس بالا را برای باور کاوی انتخاب کردیم و ممکن است برخی از کلیدواژه‌های مهم با فرکانس پایین نادیده گرفته شده باشند.

دوم: نتایج نشان می‌دهد که باور کاوی شامل رشته‌های بسیاری مانند علوم کامپیوتر، زبان‌شناسی و مهندسی برق است که نشان‌دهنده روند تحقیقات بین‌رشته‌ای است.

بنابراین، کارهای آینده باید از اسناد متنوع برای کشف ماهیت بین رشته ای استفاده کنند.

6.5 پیشنهادات و کاربردها

با توجه به اینکه حجم داده در شبکه‌های اجتماعی و وبسایت‌ها بالا است و بر اساس قطبیت کلمات تحلیل احساسات و اعتقادکاوی بر روی جملات انجام می‌شود پیشنهاد می‌شود از روش‌های یادگیری عمیق برای افزایش دقت طبقه بندی در تحلیل احساسات و اعتقادکاوی استفاده شود.

نتایج بدست آمده از این پژوهش نشان می‌دهد که عملکرد مدل یادگیری LSTM حساس به کلمه (WSD) نسبت به سایر مدل‌های یادگیری عمیق DNN، RNN و CNN بهتر است. همچنین مدل تعبیه کلمات عملکرد بهتری نسبت به TF-IDF برای تجزیه و تحلیل نظرات دارد.

به محققان پیشنهاد می‌شود بر روی مدل‌های یادگیری عمیق LSTM و کانولوشن + تعبیه کلمات و مدل‌های ترکیبی بر اساس حساسیت به کلمات و دیدگاه کاربران فعالیت کنند.

تحلیل احساسات و اعتقادکاوی داده‌های شبکه‌های اجتماعی باعث درک نظرات کاربران شده و کاربردهای گسترده‌ای در تصمیم‌گیری و سیاست‌گذاری‌ها دارد و در زمینه‌های مختلف سیاسی، اجتماعی، اقتصادی، فرهنگی و ورزشی کاربرد دارد.

به عنوان مثال در زمینه اقتصادی تحلیل نظرات کاربران بر روی محصولات می‌تواند به فروشگاه‌ها و تولیدکنندگان کمک کند تا نظرات مشتریان را تحلیل کند تا سیاست‌گذاری و تصمیم‌های مناسب را جهت افزایش جذب مشتریان انجام دهند.

در زمینه سیاسی به دولت‌ها کمک می‌کند در زمینه‌های مختلف نظرات کاربران را در شبکه‌های اجتماعی تحلیل کنند و سیاست‌گذاری‌های خود را بر اساس دیدگاه و نظرات مردم انجام دهد.

در تحقیقات آینده، با بررسی انجمن‌های احساسات-کلمه (word-sentiment) و همچنین دیدگاه‌های کاربران در حوزه‌های مختلف می‌توان تجزیه و تحلیل کامل‌تری انجام داد. یکی از محدودیت‌های مدل پیشنهادی این است که در آن فقط ویژگی‌های محلی ضبط شده‌اند. استخراج ویژگی n-gram استفاده شده در این تحقیق، شباهت معنایی یا توانایی تمایز کلمات را در نظر نمی‌گیرد. بنابراین، نمایش‌های n-gram افزایش یافته (چن و همکاران، ۲۰۱۹) برای کاهش ابعاد و کمی بودن داده‌ها توصیه می‌شود. استفاده از یک روش انتخاب ویژگی مؤثر نیز ممکن است منجر به کاهش پیچیدگی محاسباتی و بهبود راندمان زمانی شود (باروشکا و هاجک، ۲۰۲۰). از طرح‌های جایگزین مبتنی بر تعبیه نیز می‌توان استفاده کرد (عنان، ۲۰۱۹).

6. نتیجه‌گیری

در این مطالعه، ما مدل پیشنهادی LSTM، CNN، کارآمد با تفکیک شدت احساسات بر روی کلمه WSD برای استخراج نظر و اعتقادکاوی ارائه داده‌ایم. ما با انجام آزمایشات گسترده در مجموعه داده فارسی توییت، عملکرد مدل را در مقایسه با نمایش‌های اولیه کلمه ثابت کردیم. ما مدل پیشنهادی را با روش‌های موجود یادگیری عمیق و سایر روش‌های یادگیری ماشین مقایسه کردیم. از این رو، اثربخشی مدل ارائه شده نشان داده شد.

نتایج آزمایشات نشان می‌دهد که تفکیک شدت احساسات کلمات فقط بر اساس تعبیه کلمات مؤثر هستند. ادغام WSD با N-gram، پیشرفت بیشتری را ارائه می‌دهد.

در این پژوهش، علاوه بر تحلیل احساسات بر روی مجموعه داده، احساسات کاربران در مقوله‌های مختلف سیاسی، اقتصادی، فرهنگی و علمی مورد تجزیه و تحلیل قرار گرفت تا تحلیل‌گران بتوانند جهت‌گیری فکری کاربران در آن زمینه و دیدگاه‌های مثبت و منفی را نشان دهند. این مطالعه از WSD کارآمد با تمایز شدت احساسات کلمات برای عقیده کاوی و تحلیل باور استفاده کرد. مدل پیشنهادی با روش‌های یادگیری عمیق موجود و سایر روش‌های یادگیری ماشین برای نشان دادن اثربخشی مقایسه شد.

- ACM Transactions on Asian and Low-Resource Language Information Processing (TALLIP), Vol 17(4), pp. 1145-1154.
14. Cach Dang, N. , Moreno-García, M.N. and De la Prieta, F., (2020), *Sentiment Analysis Based on Deep Learning: A Comparative Study*, *Electronics* 2020, 9, 483; doi:10.3390/electronics9030483.
 15. Catal, C., Nangir, M.: *A sentiment classification model based on multiple classifiers*. *Appl. Soft Comput.* 50, 135–141 (2017)
 16. Chen, X., Xue, Y., Zhao, H., Lu, X., Hu, X., Ma, Z.: *A novel feature extraction methodology for sentiment analysis of product reviews*. *Neural Comput. Appl.* 31(10), 6625–6642 (2019)
 17. Chen, Z.; Liu, B. *Lifelong machine learning*. *Synth. Lect. Artif. Intell. Mach. Learn.* 2018, 12, 1–207. [CrossRef]
 18. Dashtipour, K. et al., (2018). *Exploiting Deep Learning for Persian Sentiment Analysis*. s.l., s.n.
 19. Dastgheib, M.B. and Koleini, S., (2019), *Persian Text Classification Enhancement by Latent Semantic Space*, *International Journal of Information Science and Management*, Vol 17(1), pp. 33-46.
 20. Do, H.H., Prasad, P.W.C., Maag, A., Alsadoon, A.: *Deep learning for aspect-based sentiment analysis: a comparative review*. *Expert Syst. Appl.* 118, 272–299 (2019)
 21. Do, H.H.; Prasad, P.; Maag, A.; Alsadoon, A.J. *Deep Learning for Aspect-Based Sentiment Analysis: A Comparative Review*. *Expert Syst. Appl.* 2019, 118, 272–299. [CrossRef]
 22. Du, C. and Huang, L., (2018), *Text Classification Research with Attention-based Recurrent Neural Networks*, *International Journal of Computers Communications & Control*, ISSN 1841-9836, 13(1),pp. 50-61.
 23. Fang, Y., Tan, H. and Zhang, J., (2018), *Multi-Strategy Sentiment Analysis of Consumer Reviews Based on Semantic Fuzziness*, *IEEE. Transactions and content mining are permitted for academic research only*, Vol 6, pp.20625-20631.
 24. Ferrara, E., Varol, O., Davis, C., Menczer, F., and Flammini, A., (2016), “The rise of social bots,” *Commun. ACM*, vol. 59, no. 7, pp. 96–104.
 25. H. Alikarami, A.M. Bidgoli and M. Sadeghzadeh, *Text mining of Persian texts based on Cellular Learning Automata and optimizing parameters of SVM*, 4th international congress on engineering,
 1. A. Onan, *Sentiment analysis on product reviews based on weighted word embeddings and deep neural networks*. *Concurrency and Computation: Practice and Experience*, 33(23), e5909, 2020.
 2. A. Onan, *Bidirectional convolutional recurrent neural network architecture with group-wise enhancement mechanism for text sentiment classification*, *Journal of King Saud University - Computer and Information Sciences*, Volume 34, Issue 5, Pages 2098-2117, 2022.
 3. Abid, F.; Alam, M.; Yasir, M.; Li, C.J. *Sentiment analysis through recurrent variants latterly on convolutional neural network of Twitter*. *Future Gener. Comput. Syst.* 2019, 95, 292–308.
 4. Alharbi, A.S.M.; de Doncker, E. *Twitter sentiment analysis with a deep neural network: An enhanced approach using user behavioral information*. *Cogn. Syst. Res.* 2019, 54, 50–61.
 5. Alikarami, H. and Khadem, F., (2016), *Data Mining Using Genetic Algorithms and Cellular Learning Automata Based on Factor Analysis and Cluster Analysis*, *1st International Conference on New Research Achievements in Electrical and Computer Engineering*, Tehran, Iran.
 6. Available online: <http://alt.qcri.org/semeval2017/> (accessed on 12 March 2020).
 7. Available online: <http://help.sentiment140.com/site-functionality> (accessed on 12 March 2020).
 8. Available online: <http://www.cs.cornell.edu/people/pabo/movie-review-data/> (accessed on 12 March 2020).
 9. Available online: <https://www.kaggle.com/c/word2vec-nlp-tutorial/data> (accessed on 12 March 2020).
 10. Available online: <https://www.kaggle.com/crowdflower/twitter-airline-sentiment> (accessed on 12 March 2020).
 11. Barushka, A., Hajek, P.: *Spam detection on social networks using cost-sensitive feature selection and ensemble-based regularized deep neural networks*. *Neural Comput. Appl.* 1–19 (2020)
 12. Basiri, M. E., Nilchi, A. R. N. & Ghassem-aghvae, N., (2014). *A Framework for Sentiment Analysis in Persian*.
 13. Basiri, M.E. and kabiri, A., (2018), *Words Are Important: Improving Sentiment Analysis in the Persian Language by Lexicon Refining*,

36. Kim, Y., 2014. *Convolutional Neural Networks for Sentence Classification*. Doha, Qatar, s.n.
37. Kraus, M.; Feuerriegel, S. *Sentiment analysis based on rhetorical structure theory: Learning deep neural networks from discourse trees*. *Expert Syst. Appl.* 2019, 118, 65–79.
38. Kumar, S.; Gahalawat, M.; Roy, P.P.; Dogra, D.P.; Kim, B.-G.J.E. *Exploring Impact of Age and Gender on Sentiment Analysis Using Machine Learning*. *Electronics* 2020, 9, 374.
39. LeCun, Y., Bengio, Y. & Hinton, G., 2015. *Deep learning*. *Nature*, Volume 521, pp. 436–444.
40. Li, L.; Goh, T.-T.; Jin, D. *How textual quality of online reviews affect classification performance: A case of deep learning sentiment analysis*. *Neural Comput. Appl.* 2018, 1–29.
41. Liu, B., 2012. *Sentiment Analysis and Opinion Mining*. *Synthesis lectures on human language technologies*, pp. 1-167.
42. Maas, A.L.; Daly, R.E.; Pham, P.T.; Huang, D.; Ng, A.Y.; Potts, C. *Learning word vectors for sentiment analysis*. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1*, Portland, OR, USA, 19–24 June 2011; pp. 142–150.
43. Marzieh Fadaee, Arianna Bisazza, and Christof Monz. 2017. *Data Augmentation for Low-Resource Neural Machine Translation*. *arXiv e-prints*, page arXiv:1705.00440.
44. Mousavirad, S.J. and Ebrahimpour-Komleh, H., (2014), *Wrapper Feature Selection using Discrete Cuckoo Optimization Algorithm*, *Austrian E-Journals of Universal Scientific Organization*, Vol. 4(11), Apr, pp. 709-721.
45. Onan, A.: *Deep learning based sentiment analysis on product reviews on Twitter*. In: Younas, M., Awan, I., Benbernou, S. (eds.) *Innovate-Data 2019*. CCIS, vol. 1054, pp. 80–91. Springer, Cham (2019). https://doi.org/10.1007/978-3-030-27355-2_6
46. Piryani, R., Madhavi, D. and Singh, V.K., (2017), “*Analytical mapping of opinion mining and sentiment analysis research during 2000–2015*,” *Information Processing & Management*, vol. 53, no. 1, pp. 122–150.
47. Qiu, L. and Li, J., (2018), *Sentiment analysis of short texts in microblog based on dependency parsing*, *springer: Cluster Computing*, Volume 21, Issue 1, pp 985-995.
48. Roustaei, A. and Rastegari, H., (2018), *Persian question classification using headword and semantic features*, *IEEE, technology & applied science, New Zealand-Auckland*, 2019.
26. H. Alikarami, A. M. Bidgoli and H. H. S. Javadi, (2023), "Belief Mining in Persian Texts Based on Deep Learning and Users' Opinions (revised December 2022)," in *IEEE Transactions on Affective Computing*, doi: 10.1109/TAFFC.2023.3288407.
27. Hajek P., Barushka A., Munk M. (2020) *Opinion Mining of Consumer Reviews Using Deep Neural Networks with Word-Sentiment Associations*. In: Maglogiannis I., Iliadis L., Pimenidis E. (eds) *Artificial Intelligence Applications and Innovations*. AIAI 2020. IFIP Advances in Information and Communication Technology, vol 583. Springer, Cham. https://doi.org/10.1007/978-3-030-49161-1_35.
28. Hassan, A. and Mahmood, A., (2018), *Convolutional Recurrent Deep Learning Model for Sentence Classification*, *IEEE*, Vol 6, pp. 13949 – 13957.
29. Hosseini, P. et al., 2018. *SentiPers: A Sentiment Analysis Corpus for Persian*. *arXiv*.
30. Jason Wei and Kai Zou. 2019. *EDA: Easy Data Augmentation Techniques for Boosting Performance on Text Classification Tasks*. *arXiv e-prints*, page arXiv:1901.11196.
31. Jeong, B.; Yoon, J.; Lee, J.-M. *Social media mining for product planning: A product opportunity mining approach based on topic modeling and sentiment analysis*. *Int. J. Inf. Manag.* 2019, 48, 280–290. [CrossRef]
32. Johnson, R., Zhang, T.: *Effective use of word order for text categorization with convolutional neural networks*. In: *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 103–112 (2015)
33. Joseph Turian, Lev-Arie Ratinov, and Yoshua Bengio. 2010. *Word representations: A simple and general method for semi-supervised learning*. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, pages 384–394, Uppsala, Sweden. Association for Computational Linguistics.
34. K. Wang and X. Wan, *Counterfactual Representation Augmentation for Cross-Domain Sentiment Analysis*. *IEEE Transactions on Affective Computing*, 2022.
35. Kausar, S., Huahu, X., Shabir, M.Y., Ahmad, W.: *A sentiment polarity categorization technique for online product reviews*. *IEEE Access* 8, 3594–3605 (2019)

59. Urologin, S., (2018), *Sentiment Analysis Visualization and Classification of Summarized News Articles: A Novel Approach*, (IJACSA) *International Journal of Advanced Computer Science and Applications*, Vol. 9, No. 8, pp. 616-624.
60. Wang, Y.; Wang, M.; Xu, W. A sentiment-enhanced hybrid recommender system for movie recommendation A big data analytics framework. *Wirel. Commun. Mob. Comput.* 2018, 2018. [CrossRef]
61. Woolley, S.C., (2016), "Automating power: Social bot interference in global politics," *First Monday*, vol. 21, no. 4.
62. Wu, C.; Wu, F.; Wu, S.; Yuan, Z.; Liu, J.; Huang, Y. Semi-supervised dimensional sentiment analysis with variational autoencoder. *Knowl. Based Syst.* 2019, 165, 30–39.
63. Y. Wang, M. Wang, W. Xu, A sentiment-enhanced hybrid recommender system for movie recommendation A big data analytics framework. *Wirel. Commun. Mob. Comput.* 2018, 2018.
64. Yang, C.; Zhang, H.; Jiang, B.; Li, K.J. Aspect-based sentiment analysis with alternating coattention networks. *Inf. Process. Manag.* 2019, 56, 463–478. [CrossRef]
65. Yao, Q.Z., Song, Z.L. and Peng, C., (2011), *Research on text categorization based on LDA*, *Computer Engineering and Applications*, Vol 47(13), pp. 150–153. *Journal of Theoretical and Applied Information Technology*, Vol 96(21), pp. 7206-7214.
49. S. Mai, Y. Zeng, S. Zheng and H. Hu, (2022). Hybrid contrastive learning of tri-modal representation for multimodal sentiment analysis. *IEEE Transactions on Affective Computing*.
50. Schmitt, M.; Steinheber, S.; Schreiber, K.; Roth, B. Joint Aspect and Polarity Classification for Aspect-based Sentiment Analysis with End-to-End Neural Networks. *arXiv* 2018, arXiv:1808.09238.
51. Shams, M., Shakery, A. & Faili, H., (2012). A non-parametric LDA-based induction method for sentiment analysis. Shiraz, Iran, s.n.
52. Shayaa, S. and et al., (2018), *Sentiment Analysis of Big Data: Methods, Applications, and Open Challenges*, IEEE. *Translations and content mining are permitted for academic research only*, Vol 6, pp. 37807-37827.
53. Singh, V.K.; Mukherjee, M.; Mehta, G.K. Combining collaborative filtering and sentiment classification for improved movie recommendations. In *Proceedings of the International Workshop on Multi-disciplinary Trends in Artificial Intelligence*, Hyderabad, India, 7–9 December 2011; pp. 38–50.
54. Singhal, P.; Bhattacharyya, P. *Sentiment Analysis and Deep Learning: A Survey*; Center for Indian Language Technology, Indian Institute of Technology: Bombay, Indian, 2016.
55. Sohrabi, M.K. and Roshani, R., (2017), *Frequent itemset mining using cellular learning automata*, *Computers in Human Behavior*, Vol 68, pp. 244-253.
56. Stai, E.; Kafetzoglou, S.; Tsiropoulou, E.E.; Papavassiliou, S.J. A holistic approach for personalization, relevance feedback & recommendation in enriched multimedia content. *Multimed. Tools Appl.* 2018, 77, 283–326.
57. T. Turki and S.S. Roy, *Novel Hate Speech Detection Using Word Cloud Visualization and Ensemble Learning Coupled with Count Vectorizer*. *Applied Sciences*, 12(13), 6611, 2022.
58. Tang, D., Qin, B., Liu, T.: Document modelling with gated recurrent neural network for sentiment classification. In: *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pp. 1422–1432 (2015)

The belief of Persian text mining based on deep learning with emotion-word separation

Abstract:

Belief analysis or the classification of texts based on the feelings and opinions of users on websites and social media helps people, companies and organizations to make important decisions. Belief mining includes a system for analyzing people's opinions and feelings about an entity such as products, people, organizations, according to the opinions, messages and tweets of users in social media.

In this article, the belief analysis of Persian texts based on the messages, comments and tweets of users in social media and websites of 4 datasets using two deep learning methods, CNN, LSTM, taking into account the sense of the word, in two poles, positive and negative with intervals. 2- and 2+ are classified. In the proposed method, first the process of data pre-processing based on character to number conversion, removing the list of extra words and multi-word analysis is done, then for belief analysis and classification of Persian texts CNN, LSTM machine learning algorithm with word sense separation (WSD) is used to Recognize the intensity of emotions according to the words. We call the proposed model CNN_WSD and LSTM_WSD.

In the proposed method, the Persian Twitter dataset is used for evaluation and then it is compared with other machine learning and deep learning methods, DNN, CNN, LSTM, in the implementation of this method, python software is used. The accuracy rate of the proposed method for LSTM-WSD and CNN-WSD is 95.8 and 94.3%, respectively.

Keywords: Belief mining, natural language processing (NLP), deep learning, text mining.