

مدیریت وجوه گم‌شده مبتنی بر یادگیری تقابلی در ادغام داده‌های چندوجهی برای تشخیص افسردگی در شبکه‌های اجتماعی

حامد مروی^۱، ابوالفضل نادی^{۲*}، محمد مهدی کیخا^۳

^۱ دانشجوی کارشناسی ارشد علوم کامپیوتر، دانشکده ریاضی، آمار و علوم کامپیوتر، دانشگاه تهران، تهران، ایران

^۲ استادیار گروه علوم کامپیوتر، دانشکده ریاضی، آمار و علوم کامپیوتر، دانشگاه تهران، تهران، ایران

^۳ استادیار گروه علوم کامپیوتر، دانشگاه سیستان و بلوچستان، زاهدان، ایران

چکیده

مقاله پژوهشی

تاریخ دریافت:

۱۴۰۴/۶/۳

تاریخ پذیرش:

۱۴۰۴/۸/۷

کلیدواژه‌ها:

مدل‌های چندوجهی، وجوه گم‌شده، ادغام داده‌های چندوجهی، یادگیری تقابلی، شبکه‌های اجتماعی

نویسنده مسئول:

a.nadi@ut.ac.ir

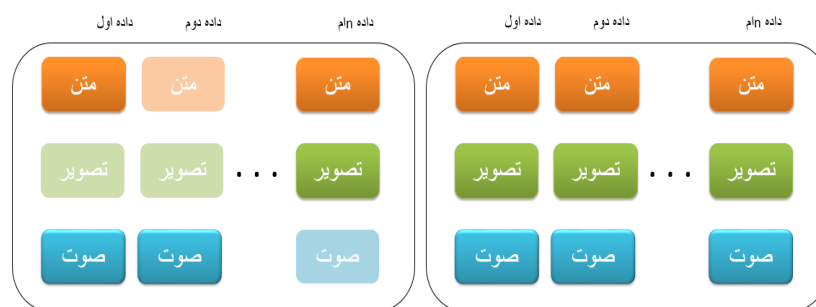
تحلیل داده‌های شبکه‌های اجتماعی اهمیت بنیادینی در استخراج الگوهای رفتاری کاربران دارد. مدل‌های چندوجهی که ترکیبی از اطلاعات متنی، تصویری و سایر منابع را به کار می‌گیرند، ابزارهای مناسبی برای این دست از تحلیل‌ها هستند. با این حال، چالشی اساسی در این مدل‌ها، فقدان برخی وجوه در بخشی از نمونه‌های داده است؛ برای نمونه، کاربری ممکن است تنها متن منتشر کند و هیچ تصویری به اشتراک نگذارد. این مسئله سبب می‌شود مدل‌های چندوجهی نتوانند از تمامی اطلاعات موجود به‌طور کامل بهره‌برداری کنند. در این مقاله، روشی برای بهره‌گیری از داده‌های ناقص در مدل‌های چندوجهی ارائه شده است. ابتدا، مدل‌های تک‌وجهی برای پردازش هر وجه به‌صورت مستقل آموزش داده شدند. سپس، یک رمزگذار مبتنی بر یادگیری تقابلی طراحی و آموزش داده شد که هدف آن، برآورد بردار ویژگی وجوه گم‌شده با اتکا به ویژگی‌های وجوه موجود است. در نهایت، داده‌های متنی و تصویری (واقعی یا بازسازی شده) در یک مدل چندوجهی ادغام شده و برای تحلیل رفتار کاربران مورد استفاده قرار گرفت. نتایج آزمایش‌ها بر روی دادگان MDDL و بر پایه‌ی معیارهای دقت و F1 نشان می‌دهند که مدل چندوجهی پیشنهادی، با دقت ۹۰/۱۷ درصد و امتیاز F1 برابر با ۹۰/۶۴ درصد، عملکرد بهتری نسبت به مدل‌های تک‌وجهی مبتنی بر متن (۸۷/۸۷ درصد دقت) و تصویر (۷۳/۳۷ درصد دقت) دارد. این نتایج مؤید آن است که مدل پیشنهادی، بدون نیاز به حذف داده‌های ناقص، از اطلاعات موجود به‌طور مؤثر بهره‌برداری می‌کند.

doi : 10.22034/ABMIR.2025.23578.1160

۱- مقدمه

بردارها با هم ترکیب می‌شوند [۶][۵][۴]. با این حال، یکی از چالش‌های اصلی در این مدل‌ها که استفاده از برخی از وجه‌ها را دچار مشکل می‌کند این است که در بعضی از مواقع یک یا چند وجه در دسترس نبوده و یا از دست رفته است؛ به این معنا که ممکن است یک یا چند وجه در برخی از نمونه‌های داده در دسترس نباشد. برای مثال، در سیستم‌های تشخیص احساسات یا تحلیل رفتار کاربران در شبکه‌های اجتماعی، ترکیب اطلاعات متنی، تصویری و صوتی می‌تواند به تشخیص دقیق‌تر و کامل‌تر کمک کند، اما لزوماً پست‌های کاربران همه‌ی وجه‌ها را ندارند و ممکن است کاربر فقط متنی بنویسد و تصویر یا صوتی در پستش ارسال نکند. در دسترس نبودن هر یک از این وجوه یعنی بخشی از اطلاعاتی که می‌توانست به تصمیم‌گیری بهتر کمک کند، از دست می‌رود. به‌عنوان نمونه، در دسترس نبودن تصویر ممکن است منجر به از دست رفتن ویژگی‌های بصری مهمی شود که به‌طور غیرمستقیم بر دقت تشخیص یا پیش‌بینی مدل تأثیرگذار باشد. در روش‌های گذشته، برای این‌که در زمان اجرا بتوان از نمونه داده‌هایی که یک یا چند وجه آن‌ها گم شده است استفاده کرد، آن وجه‌ها در مدل نادیده گرفته می‌شد و در نتیجه اطلاعات زیادی از دست می‌رفت. در شکل ۱ می‌بینیم که در شرایطی با سه وجه، نمونه‌های داده می‌توانند به دو دسته تقسیم شوند: داده‌های کامل (که اطلاعات را از هر سه وجه دارند) و داده‌های ناقص (که اطلاعات از یک یا چند وجه در دسترس نیستند)

افسردگی یکی از شایع‌ترین اختلالات روانی در جهان است که حدود ۵/۷ درصد از بزرگسالان را درگیر می‌کند و در مواردی می‌تواند منجر به خودکشی شود [۱]. با گسترش شبکه‌های اجتماعی، پژوهش‌های متعددی برای شناسایی نشانه‌های افسردگی از رفتار و محتوای کاربران انجام شده است. در این میان، مدل‌های چندوجهی^۱ به دلیل توانایی ترکیب داده‌های متنی و تصویری، نقش مهمی در بهبود دقت و جامعیت تحلیل احساسات و پیش‌بینی وضعیت روانی کاربران دارند [۳][۲]. استفاده از داده‌های چندوجهی می‌تواند دقت پیش‌بینی و درک جامع از اطلاعات را بهبود بخشد، زیرا ترکیب وجه‌های مختلف به نمایش بهتری از یک پدیده منجر می‌شود. در ادغام داده‌های^۲ چندوجهی، سه روش اصلی برای ترکیب داده‌ها وجود دارد: ادغام اولیه^۳، ادغام نهایی^۴ و ادغام ترکیبی^۵. در ادغام اولیه، داده‌های خام پیش از هرگونه استخراج ویژگی با هم ترکیب می‌شوند. به‌عنوان مثال، اگر وجه‌های متن و تصویر داشته باشیم، جریان‌های ورودی متن و تصویر پس از هم‌ترازی و پیش‌پردازش اولیه به یک معماری مشترک داده می‌شوند تا مدل ویژگی‌ها را به‌صورت مشترک بیاموزد. در روش ادغام نهایی، بردارهای تعبیه‌شده‌ی هر وجه به‌دست می‌آیند و روی هر کدام مدل مستقلی اجرا می‌شود؛ سپس در انتها نتایج این مدل‌ها با هم ترکیب می‌شود تا به پاسخ نهایی برسیم. در روش ادغام ترکیبی، ابتدا بردار تعبیه‌شده‌ی هر وجه استخراج می‌شود. سپس روی آن‌ها پردازش‌های مستقلی انجام می‌گیرد و درنهایت این



شکل (۱): نمونه‌های با وجه‌های کامل (راست) و نمونه‌های با وجه‌های گم‌شده (چپ)



۲- پیشینه پژوهش

پژوهش‌های پیشین در زمینه تشخیص افسردگی از داده‌های شبکه‌های اجتماعی را می‌توان به سه دسته اصلی تقسیم کرد: (۱) روش‌های تک‌وجهی که معمولاً بر تحلیل داده‌های متنی متمرکز هستند، (۲) روش‌های چندوجهی که از ترکیب داده‌های متنی، تصویری و صوتی برای بهبود عملکرد استفاده می‌کنند، و (۳) روش‌های مدیریت داده‌های دارای وجه گم‌شده که اگرچه الزاماً در حوزه تشخیص افسردگی نیستند، اما از نظر فنی در زمینه یادگیری چندوجهی و مقابله با نبود داده در یک یا چند وجه مرتبط هستند. در ادامه، هر یک از این دسته‌ها به صورت جداگانه بررسی می‌شوند تا روند تکامل پژوهش‌ها و شکاف‌های موجود در این حوزه مشخص شود.

۱-۲ روش‌های تک‌وجهی

امروزه کاربران احساسات و فعالیت‌های خود را از طریق پست‌هایشان در شبکه‌های اجتماعی به اشتراک می‌گذارند. این امر فرصتی را برای ارزیابی وضعیت سلامت روان، از جمله افسردگی، فراهم کرده. پس می‌توانیم با استفاده از پست‌های کاربران، وضعیت روانی آن‌ها را ارزیابی کنیم و نشانه‌های مرتبط با سلامت روان آن‌ها را شناسایی کنیم. توئیتر به‌عنوان منبعی ارزشمند برای درک افکار و احساسات افراد مطرح شده است، زیرا کاربران معمولاً تجربیات شخصی خود را به صورت عمومی به اشتراک می‌گذارند. این ویژگی، امکان تحلیل متن پیام‌های عمومی آن‌ها را برای شناسایی الگوهای زبانی مرتبط با افسردگی فراهم می‌کند [۸] [۷].

روش‌های اولیه تحلیل احساسات تک‌وجهی بودند و عمدتاً بر روی متن تمرکز داشتند. با این حال، تکیه صرف بر داده‌های متنی برای استخراج احساسات بیان‌شده توسط انسان‌ها کافی نیست، زیرا معنای آنچه یک کاربر می‌گوید اغلب بر اساس رفتارهای غیرکلامی تغییر می‌کند [۹].

در مطالعه‌ای که توسط کوپراسمیت و همکاران [۱۰] انجام شد، داده‌های شبکه‌های اجتماعی تحلیل شد و به همبستگی معناداری بین زبان استفاده‌شده در توئیترها و شیوع علائم افسردگی پی برده شد. آن‌ها دریافتند که افرادی با سطح بالاتر افسردگی بیشتر تمایل دارند احساسات منفی را در توئیترهای خود بیان کنند.

دشپاند و راتو [۱۱] روشی مبتنی بر پردازش زبان طبیعی پیشنهاد کردند تا اطلاعات متنی موجود در توئیتر را مدل‌سازی کنند. برخلاف روش‌های پیشین، آن‌ها از کیسه واژگان^۱ برای نمایش متن توئیتر به‌عنوان یک بردار تنک استفاده کردند، که این روش به دسته‌بندی امکان می‌دهد ویژگی‌های نهان را به صورت خودکار یاد بگیرد. طبقه‌بندی بیز ساده^۲ آموزش‌دیده آن‌ها برای معیار F1 به امتیاز ۰/۸۳ دست یافت، درحالی که ماشین بردار پشتیبان^۳ به امتیاز ۰/۷۹ رسید. در مطالعه‌ای توسط دی‌چودری و همکاران [۱۲]، ارتباط بین شاخص‌های زبانی و افسردگی در پلتفرم‌های شبکه‌های اجتماعی بررسی شد. آن‌ها دریافتند که افراد دارای علائم افسردگی بیشتر از زبان منفی استفاده می‌کنند که نشان‌دهنده همبستگی بین بیان منفی و افسردگی است.

مورارکا و همکاران یک مدل مبتنی بر روبرتا^۴ برای دسته‌بندی بیماری‌های روانی در شبکه‌های اجتماعی توسعه دادند. دادگان آن‌ها شامل ۱۵۹/۱۷ پست از سایت ردیت^۵ است. نتایج پژوهش آن‌ها نشان داد که مدل‌های مبتنی بر روبرتا نسبت به روش‌های سنتی پردازش زبان طبیعی مانند حافظه طولانی مدت کوتاه‌مدت^۶ یا ماشین بردار پشتیبان دقت بالاتری در تشخیص بیماری‌های روانی از محتوای متنی شبکه‌های اجتماعی دارند [۱۳].

در مجموع، روش‌های تک‌وجهی اگرچه در استخراج نشانه‌های زبانی مؤثرند، اما قادر به بهره‌گیری از اطلاعات غیرکلامی نیستند و در نتیجه درک جامع از وضعیت روانی کاربر ارائه نمی‌دهند.

۲-۲ روش‌های چندوجهی

با وجود پیشرفت‌هایی که در تشخیص افسردگی با استفاده از داده‌های تک‌وجهی حاصل شده، داده‌های چندوجهی دارای مزایای بیشتری هستند. پژوهش در زمینه تشخیص افسردگی به‌کمک

⁴ RoBERTa

⁵ Reddit

⁶ Long Short-Term Memory

¹ Bag of Words

² Naive Bayes

³ Support Vector Machine



ویژگی‌های تصویری نیز بدون بهره‌گیری از شبکه‌های عصبی عمیق و صرفاً براساس توزیع رنگ‌های سرد و گرم انجام گرفت. علاوه بر این، ویژگی‌های رفتاری کاربران از جمله تعداد پست‌ها، فعالیت شبانه و میزان تعاملات اجتماعی مورد بررسی قرار گرفتند. مدل پیشنهادی آن‌ها، از Bi-GRU همراه با مکانیزم توجه برای پردازش متن و پرسپترون چندلایه^۷ برای ادغام داده‌های چندوجهی استفاده کرد. نتایج این مدل نشان داد که ترکیب داده‌های متنی، تصویری و رفتاری، نسبت به روش‌های تک‌وجهی، دقت بیشتری در شناسایی افسردگی در شبکه‌های اجتماعی دارد.

با وجود بهبود عملکرد مدل‌های چندوجهی، این مطالعات فرض کرده‌اند که تمام وجوه داده در دسترس هستند و به مسئله داده‌های ناقص توجه نکرده‌اند.

۲-۳ روش‌های مدیریت وجه گم‌شده

بسیاری از دادگان چندوجهی شامل نمونه‌هایی هستند که همه وجوه را در اختیار ندارند. این مشکل در داده‌های جمع‌آوری شده از شبکه‌های اجتماعی بسیار اتفاق می‌افتد زیرا ممکن است کاربر، متنی را پست کند اما تصویری را به آن پست الحاق نکند و این نمونه داده دارای وجه تصویر نیست. روش‌هایی که در ادامه بررسی می‌کنیم مستقیماً در حوزه تشخیص افسردگی متمرکز نیستند، اما از نظر فنی در زمینه‌ی یادگیری چندوجهی و مدیریت داده‌های ناقص کاربرد دارند. در پژوهش‌های قبلی، چند روش برای مدیریت داده‌های ناقص ارائه شده است. یکی از روش‌ها، حذف داده‌های ناقص است. در این روش، نمونه‌هایی که یک یا چند وجه را ندارند کنار گذاشته می‌شوند. این روش اگرچه ساده‌ترین رویکرد است، اما باعث از دست رفتن اطلاعات مفید و کاهش اندازه مجموعه داده‌ها می‌شود [۱۹]. روش دیگر، نادیده گرفتن وجوه گم‌شده است. در این روش تمام داده‌های موجود در هر نمونه استفاده می‌شود، بدون آنکه برای وجوه‌های غایب جایگزینی ایجاد کند. به عنوان مثال، اگر یک نمونه داده دارای دو وجه از سه وجه باشد، مدل تنها از همان دو وجه استفاده خواهد کرد. این رویکرد در

داده‌های چندوجهی بسیار مورد توجه قرار گرفته شده و تعداد قابل توجهی از پژوهشگران را به خود جذب کرده است [۱۴]. مدل‌های چندوجهی با چالش‌هایی روبه‌رو هستند که در سه دسته‌ی اصلی جای می‌گیرند: «ترجمه و هم‌ترازی داده‌ها»، «پیچیدگی محاسباتی»، و «ادغام مؤثر وجوه مختلف». هم‌ترازی داده‌های چندوجهی به معنی پیدا کردن ارتباط‌های معنادار بین انواع داده‌های مختلف مانند تصویر و متن است که به دلیل تفاوت‌های ساختاری، یکی از چالش‌های بزرگ محسوب می‌شود. ترجمه چندوجهی نیز به تبدیل داده‌ها از یک وجه به وجه دیگر اشاره دارد. همچنین، پردازش داده‌های متنوع به‌طور هم‌زمان نیازمند منابع محاسباتی زیاد است، به‌خصوص در مدل‌های عمیق و مبدل‌ها^۸ که با داده‌های بزرگ سروکار دارند. در نهایت، چگونگی ادغام مؤثر داده‌های چندوجهی به‌گونه‌ای که اطلاعات کلیدی حفظ شوند، یکی از مسائل مهمی است که می‌تواند بر دقت نهایی مدل تأثیرگذار باشد [۱۶] [۱۵].

فانگ و همکاران [۱۷] مدلی چندوجهی برای تشخیص افسردگی پیشنهاد دادند که آر مکانیزم توجه چندسطحی^۳ استفاده کرده که از سه وجه داده صوتی، بصری و متنی استفاده می‌کند. این مدل برای استخراج ویژگی‌های هر وجه، از حافظه طولانی‌مدت کوتاه‌مدت و نوع دو طرفه آن^۴ همراه با مکانیزم توجه استفاده می‌کند. در مرحله بعدی، خروجی‌ها را از طریق شبکه ادغام توجه^۵ ترکیب می‌کند تا اطلاعات معنادارتر استخراج شود. نتایج آزمایش‌ها روی مجموعه داده DAIC-WOZ نشان داد که این مدل نسبت به روش‌های قبلی عملکرد بهتری دارد و توانسته با کمک مکانیزم توجه چندسطحی، ویژگی‌های کلیدی افسردگی را با دقت بیشتری شناسایی کند.

وانگ و همکاران [۱۸] روشی چندوجهی برای شناسایی افسردگی در شبکه‌های اجتماعی پیشنهاد دادند که از متن، تصاویر و ویژگی‌های رفتاری کاربران در پلتفرم ویبو^۶ استفاده می‌کند. در این پژوهش، برای پردازش متن از مدل XLNet جهت استخراج بردارهای معنایی پست‌ها و پروفایل کاربران استفاده شد. تحلیل

⁵ Attention Fusion Network

⁶ Weibo

⁷ Multi-Layer Perceptron

⁸ Missing Modality

¹ Translation and Alignment

² Transformer

³ Multi-Level Attention Mechanism

⁴ Bidirectional Long Short-Term Memory



خود را حفظ کند. نتایج این پژوهش نشان می‌دهد که این روش علاوه بر کاهش پیچیدگی مدل، امکان استفاده مؤثر از داده‌های ناقص را بدون نیاز به بازسازی مستقیم فراهم می‌کند [۲۲].

رضا و همکاران برای مدیریت وجه‌های گم‌شده در مدل‌های چندوجهی، روشی براساس حذف تصادفی وجه‌ها در مرحله آموزش ارائه دادند. در این روش، در هر مرحله آموزش، به صورت تصادفی یک یا چند وجه از داده‌ها حذف می‌شود. این کار مدل را مجبور می‌کند تا بازنمایی‌های مقاومی یاد بگیرد که بتوانند در زمان تست، حتی در صورت ناقص بودن داده‌های ورودی، عملکرد قابل قبولی داشته باشند. برخلاف روش‌هایی که برای بازسازی وجه گم‌شده از شبکه‌های مولد تخصصی یا خودرمزگذارها استفاده می‌کنند، این روش تلاشی برای بازسازی ندارد. بلکه به مدل آموزش می‌دهد تا صرفاً با استفاده از داده‌های در دسترس، پیش‌بینی دقیقی انجام دهد. این رویکرد نه تنها هزینه محاسباتی را کاهش می‌دهد، بلکه در شرایط واقعی که برخی از وجوه داده ممکن است همواره در دسترس نباشند عملکرد مطلوبی از خود نشان می‌دهد [۲۳].

وانگ و همکاران روشی برای یادگیری چندوجهی در شرایطی که برخی وجوه در داده‌های ورودی ناقص هستند، ارائه داده‌اند. ابتدا برای هر وجه، مدل‌های معلم^۵ تک‌وجهی آموزش داده می‌شوند که ویژگی‌های مربوط به آن وجه را استخراج می‌کنند. سپس از طریق تقطیر دانش^۶ برای انتقال اطلاعات از این مدل‌های معلم به یک مدل دانش‌آموز^۷ چندوجهی استفاده می‌شود. مدل دانش‌آموز علاوه بر ادغام داده‌های چندوجهی، یاد می‌گیرد که در صورت نبود برخی وجوه، ویژگی‌های وجوه گم‌شده را تخمین بزند. این یادگیری از طریق حذف مصنوعی وجود در مرحله آموزش انجام می‌شود. در نتیجه مدل دانش‌آموز قادر است با استفاده از وجوه موجود، بردار ویژگی وجوه گم‌شده را بازسازی کند و از آن‌ها هنگام ادغام داده استفاده کند [۲۴].

مقایسه با حذف کامل داده‌ها، اطلاعات بیشتری را حفظ می‌کند. با این حال، مشکل اصلی یعنی بهره‌نبردن از وجه‌های گم‌شده همچنان پابرجاست [۱۹].

کای و همکاران [۲۰] در پژوهش خود از مدل‌های مولد^۱ مانند خودرمزگذارها^۲ و شبکه‌های مولد تخصصی^۳ برای بازسازی وجه گم‌شده استفاده کرده‌اند. نتایج آن‌ها نشان می‌دهد که این روش می‌تواند امکان استفاده از تمام داده‌ها را فراهم کند، اما در عین حال، پیچیدگی محاسباتی بالایی دارد و کیفیت داده‌های بازسازی شده همواره ثابت نیست [۲۰].

در پژوهشی دیگر روشی برای بازسازی بردار ویژگی^۴ وجه گم‌شده ارائه شد که در آن، به جای بازسازی کامل وجه، تنها بردار ویژگی آن تخمین زده می‌شود. در این روش شبکه‌ای طراحی شد که با یادگیری ویژگی‌های مشترک بین وجوه مختلف، قادر است ویژگی‌های گم‌شده را از طریق اطلاعات موجود بازسازی کند. مدل ابتدا با داده‌هایی که تمام وجوه را دارند آموزش داده می‌شود و سپس، در هر مرحله یک یا چند وجه به صورت تصادفی حذف می‌شوند تا مدل یاد بگیرد ویژگی‌های از دست‌رفته را بازسازی کند. نتایج این پژوهش نشان می‌دهد که این روش در مقایسه با بازسازی مستقیم داده‌ها، بهینه‌تر بوده و امکان تخمین ویژگی‌های از دست‌رفته را با دقت بیشتری فراهم می‌کند [۲۱].

وانگ و همکاران روشی برای مدل‌سازی هم‌زمان ویژگی‌های مشترک و اختصاصی در داده‌های چندوجهی پیشنهاد داده‌اند. در این رویکرد، ویژگی‌های هر وجه به دو دسته ویژگی‌های مشترک (که بین تمام وجوه مشابه هستند) و ویژگی‌های اختصاصی (که مختص همان وجه‌اند) تقسیم می‌شوند. در صورت گم‌شدن یک وجه، این مدل به جای بازسازی مستقیم بردار ویژگی، از ویژگی‌های مشترک موجود در سایر وجوه برای استخراج اطلاعات استفاده می‌کند. به عبارت دیگر، مدل فرض می‌کند ویژگی‌های مشترک بین وجوه مختلف حاوی اطلاعات کافی برای پردازش هستند و در نتیجه، می‌تواند بدون نیاز به بازسازی مستقیم وجه گم‌شده، عملکرد

⁵ Teacher Model

⁶ Knowledge Distillation

⁷ Student Model

¹ Generative Model

² Auto Encoder

³ Generative Adversarial Network

⁴ Feature Vector



آموزش داده می‌شوند تا بتوانند به‌طور مستقل ویژگی‌های مربوط به هر وجه را استخراج کرده و پیش‌بینی انجام دهند. در مرحله بعد، یک رمزگذار مبتنی بر یادگیری تقابلی به کار گرفته می‌شود تا بردار ویژگی مربوط به وجه گم‌شده را برآورد کند. در نهایت، مدل چندوجهی با ترکیب ویژگی‌های متنی و تصویری-چه به صورت واقعی و چه تخمین‌زده-فرآیند دسته‌بندی را انجام می‌دهد. در ادامه، هر یک از این مراحل به تفصیل توضیح داده خواهند شد.

ابتدا مدل‌های تک‌وجهی برای پردازش متن و تصویر آموزش داده شده‌اند. این مدل‌ها به‌طور مستقل ویژگی‌های مربوط به هر وجه را استخراج کرده و پیش‌بینی می‌کنند که آیا یک توییت مربوط به یک کاربر افسرده است یا خیر. ابتدا، مدل‌های تک‌وجهی آموزش داده شدند. برای پردازش متن، از مدل XLM-RoBERTa استفاده شده است. این مدل به دلیل آموزش بر روی داده‌های چندزبانه، قادر به پردازش توییت‌هایی به زبان‌های مختلف است. متن خام مستقیماً به مدل داده می‌شود. خروجی توکن [CLS] برای نمایش ویژگی‌های کلی متن استخراج شده و از آن برای تشخیص افسرده یا غیرافسرده بودن توییت استفاده می‌شود.

برای پردازش تصویر، از مدل ViT^۳ استفاده شده است. ورودی این مدل باید تصاویری با اندازه ۲۲۴×۲۲۴ باشد؛ به همین منظور، ابتدا تصویر مقیاس‌بندی می‌شود تا یکی از ابعاد (طول یا عرض) به ۲۲۴ برسد، سپس-در صورتی که بعد دیگر کوچک‌تر از ۲۲۴ باشد-نواحی خالی با میانگین رنگ پیکسل‌های تصویر پر می‌شود.^۴ تصویر ورودی در مدل ViT به ۱۹۶ بخش ۱۶×۱۶ پیکسلی تقسیم می‌شود و از هر بخش، یک بردار ویژگی استخراج می‌گردد. این بردارها وارد مبدل می‌شوند و مدل بر اساس آن‌ها، نتیجه‌گیری نهایی را انجام می‌دهد.

در تصویر ۲ معماری کلی مدل‌های متنی و تصویری نشان داده شده است: بخش بالا مربوط به پردازش متن توسط مدل XLM-RoBERTa است-ابتدا، ورودی به واحدهای زبانی شکسته می‌شود^۵ و سپس پردازش شده و ویژگی‌های آن استخراج می‌گردد؛ بخش پایین نیز به پردازش تصویر اختصاص دارد که با استفاده از

در مجموع، پژوهش‌های پیشین در زمینه‌ی مدیریت وجه گم‌شده به دو رویکرد اصلی تقسیم می‌شوند. گروه اول، روش‌های بازسازی محور هستند که با استفاده از شبکه‌های مولد مانند خودرمزگذارها یا شبکه‌های مولد تخصصی تلاش می‌کنند داده ازدست‌رفته را بازسازی کنند. این روش‌ها هرچند دقیق‌اند، اما از نظر محاسباتی سنگین بوده و نیاز به داده آموزشی زیاد و تنظیم‌گرهای متعدد دارند و در داده‌های واقعی پایداری محدودی دارند. گروه دوم شامل روش‌هایی است که سعی می‌کنند مدل را نسبت به فقدان وجه مقاوم کنند، بدون آن‌که وجه گم‌شده را بازسازی کنند؛ مانند حذف تصادفی وجه‌ها در آموزش یا مدل‌سازی ویژگی‌های مشترک میان وجوه. این روش‌ها سبک‌تر هستند، اما معمولاً فقط از وجوه موجود استفاده می‌کنند و بازنمایی دقیقی از وجه غایب به دست نمی‌دهند. بنابراین هنوز نیاز به روشی وجود دارد که در عین سادگی و کارایی، بتواند در داده‌های شبکه‌های اجتماعی، ویژگی‌های وجه گم‌شده را با استفاده از هم‌ترازی تقابلی بین وجوه مختلف تخمین بزند.

در بخش بعدی، روش پیشنهادی برای مدیریت وجه گم‌شده توضیح داده خواهد شد. ابتدا مدل‌های تک‌وجهی معرفی می‌شوند و سپس نحوه استفاده از یادگیری تقابلی برای تخمین بردار ویژگی وجه گم‌شده بررسی خواهد شد.

۳- روش پیشنهادی

در سال‌های اخیر، برخی پژوهش‌ها به جنبه‌هایی از داده‌های ناقص در مدل‌های چندوجهی اشاره کرده‌اند، اما عموماً این مسئله به صورت محدود یا در قالب حذف داده‌های ناقص بررسی شده است. پژوهش حاضر با تمرکز ویژه بر شرایطی که یکی از وجوه (به‌ویژه تصویر) در دسترس نیست، تلاش می‌کند با بهره‌گیری از یادگیری تقابلی میان وجوه متنی و تصویری، بازنمایی مؤثری از وجه گم‌شده به دست آورد و این چالش را حل کند.

مدل پیشنهادی بر پایه ساختاری چندمرحله‌ای طراحی شده است که هدف آن تحلیل داده‌های چندوجهی و مدیریت وجوه گم‌شده است. در گام اول، دو مدل تک‌وجهی^۱ برای پردازش متن و تصویر

⁴ Padding

⁵ Tokenize

¹ Single Modal Model

² Encoder

³ Vision Transformer

نمونه برآورد کند [۲۹]. برای آموزش مدل T2V از یادگیری تقابلی استفاده می‌شود تا تصویر مجازی (تخمین زده) تا حد امکان به بردار ویژگی تصویر واقعی نزدیک باشد. معیار نزدیکی، تشابه کسینوسی^۵ بین دو بردار است. در یادگیری تقابلی، هدف این است که بردار ویژگی مجازی هر نمونه به تصویر خودش نزدیک باشد، و با تصاویر سایر نمونه‌ها فاصله داشته باشد. ابتدا با رمزگذارهای متن و تصویر که به صورت تک‌وجهی آموزش دیدند، بردار ویژگی‌های نمونه داده‌هایی که در آن‌ها هر دو وجه موجودند محاسبه می‌شوند. فرض می‌کنیم که در هر گروه^۶، مدل n نمونه به‌عنوان ورودی دریافت می‌کند. در این حالت، n بردار ویژگی متن با نمادهای T_1, \dots, T_n و n بردار ویژگی تصویر با نمادهای I_1, \dots, I_n در اختیار خواهیم داشت. تابع هزینه^۷ به صورت زیر است:

$$L = \frac{1}{n} \sum_{i=1}^n (1 - \cos(V_i, I_i) + \max(0, x_i - \text{margin})) \quad (1)$$

$$x_i = \frac{1}{(n-1)} \sum_{j \neq i} \cos(V_i, I_j) \quad (2)$$

که در آن:

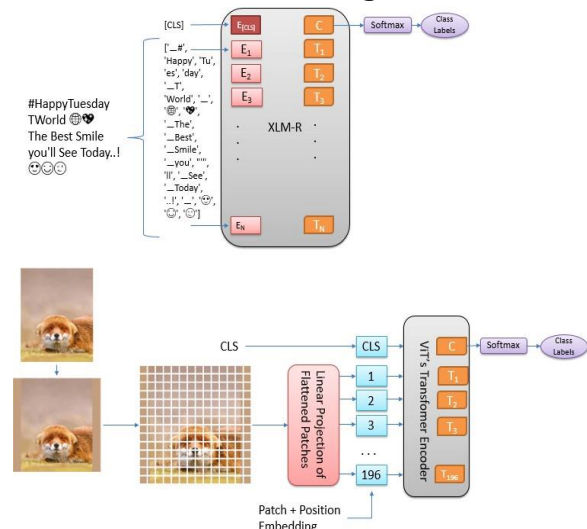
- V_i بردار ویژگی مجازی (تخمین زده) برای نمونه i است.
- I_i بردار ویژگی تصویر برای نمونه i است.
- $\cos(V_i, I_i)$ میزان تشابه کسینوسی بین بردار ویژگی‌های مجازی و واقعی نمونه i است. پس برای نمونه i، تابع هزینه از جمع دو بخش تشکیل شده است:

$$\max(0, x_i - \text{margin})$$

$$1 - \cos(V_i, I_i)$$

بخش $1 - \cos(V_i, I_i)$ موجب می‌شود که بردار ویژگی مجازی i تا حد امکان شبیه به بردار ویژگی تصویر i شود؛ بخش دوم نیز باعث می‌شود بردار ویژگی مجازی i از بردار ویژگی تصویر نمونه‌های دیگر دور شود. پارامتر margin تضمین می‌کند اگر میزان تشابه کسینوسی بین بردار ویژگی مجازی i با بردار ویژگی تصویر i بیشتر از مقدار margin از هم دور باشند، دیگر نیازی به «دورتر کردن» این دو بردار نباشد و قسمت دوم برابر با صفر شود. مدل T2V

مدل ViT انجام می‌شود؛ این مدل ابتدا تصویر را به قطعات کوچک‌تر تقسیم می‌کند و سپس با بهره‌گیری از مکانیسم توجه^۱، ویژگی‌های آن را استخراج می‌کند.



شکل (۲): معماری مدل‌های تک‌وجهی: تصویر بالا مدل متنی

ViT و XLM-RoBERTa و تصویر پایین مدل تصویری

در این مرحله، مدل‌های تک‌وجهی برای متن و تصویر توانسته‌اند به صورت مستقل ویژگی‌های خود را استخراج کرده و بر اساس آن تصمیم‌گیری کنند. این مدل‌ها در مرحله بعدی به‌عنوان بخش‌های پایه برای مدل چندوجهی مورد استفاده قرار خواهند گرفت. راه‌حل ما برای مسئله وجه از دست‌رفته، تخمین بردار ویژگی آن وجه است. برای این کار، از یادگیری تقابلی^۲ استفاده می‌کنیم؛ روشی که در مدل‌های شناخته‌شده‌ای مانند SimCLR [۲۵] [۲۵] و CLIP [۲۶] [۲۶] به‌کار رفته است. در یادگیری تقابلی، تلاش می‌شود نمونه‌هایی که به هم شباهت دارند، بردارهای ویژگی نزدیکی داشته باشند و در مقابل، نمونه‌های غیرمشابه دارای بردارهای دور از هم باشند [۲۷]. در مواردی که وجه تصویری در دسترس نیست، می‌کوشیم با استفاده از بردار ویژگی متن، بردار متناظر تصویر را تخمین بزنیم. برای این کار، یک رمزگذار^۳ به نام T2V^۴ طراحی و آموزش داده می‌شود؛ این مدل، بردار ویژگی متن را به‌عنوان ورودی دریافت کرده و می‌کوشد بردار ویژگی متناظر تصویر را برای همان

⁵ Cosine Similarity

⁶ Batch

⁷ Loss Function

¹ Attention Mechanism

² Contrastive Learning

³ Encoder

⁴ Text-toVisual Encoder

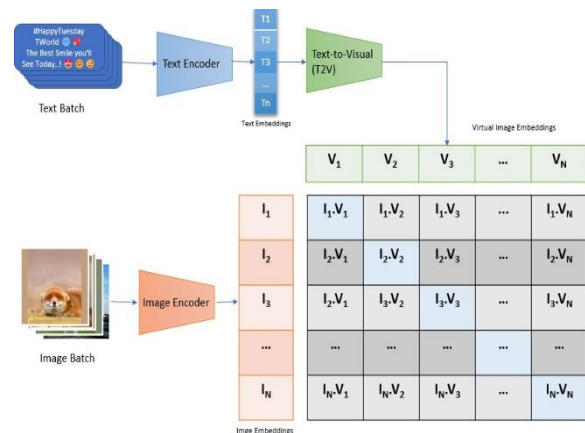
ساختار در شکل ۳ دیده می‌شود. در ادامه، بردار ویژگی متن و تصویر (واقعی یا تخمینی) با هم ترکیب شده^۳ شده و به یک شبکه پرسپترون چندلایه داده می‌شود تا عمل دسته‌بندی انجام گیرد. ورودی این لایه، بردارهای ادغام‌شده‌ی متنی و تصویری است که دارای اندازه‌ی $768+768$ برابر ۱۵۳۶ بعد است. در این معماری، لایه مخفی اول بعد را به ۵۱۲ کاهش می‌دهد و در لایه بعدی به ۲۵۶؛ سپس لایه خروجی -۲ بعدی است که به کمک کلاس افسرده/غیرافسرده معلوم می‌شود.

جدول (۱): پارامترهای مدل‌های مورد استفاده

توضیحات	پارامتر	مدل
Xlm-roberta-base 128 768	مدل متن حداکثر طول دنباله ابعاد خروجی	XLN RoBERTa
google/vit-base-patch16-224 $3 \times 224 \times 224$ Normalize(mean=[0.5, 0.5, 0.5],std=[0.5, 0.5, 0.5]) 768	مدل تصویر شکل ورودی نرمال‌سازی ابعاد خروجی	ViT
T2V Encoder 768 768 Contrastive Learning 0.004 5 256	مدل متن به تصویر ورودی خروجی تابع هزینه نرخ یادگیری تعداد تکرار اندازه دسته	T2V
$768 + 768 = 1536$	ابعاد ادغام	Fusion
Dense $1536 \rightarrow 512 \rightarrow 256$ $\rightarrow 2$	معماری شبکه	Multimodal Classification

در جدول ۱ توضیحاتی درباره مدل‌های تک‌وجهی متن و تصویر، مدل T2V و مدل چندوجهی آمده است. همچنین نرخ یادگیری^۴، تعداد تکرار^۵ و اندازه دسته^۶ مدل T2V گزارش شده. پس از معرفی روش پیشنهادی، لازم است عملکرد آن را در شرایط مختلف ارزیابی کنیم. در بخش بعدی، نتایج آزمایش‌ها ارائه شده و نتایج بررسی می‌شوند.

یک پرسپترون چندلایه ساده بدون لایه پنهان^۱ است که اندازه لایه‌ی ورودی ۷۶۸ و اندازه لایه‌ی خروجی نیز ۷۶۸ است. در شکل ۳ نیز نشان داده شده است: اگر ماتریسی بسازیم که در سطرها آن بردار ویژگی مجازی تصویر را قرار دهیم و در ستون‌ها بردار ویژگی تصویر واقعی را، آن‌گاه درایه‌های ماتریس میزان تشابه کسینوسی بین بردارهای مجازی و واقعی تصویر خواهند بود؛ می‌کشیم قطر اصلی ماتریس (که با رنگ آبی مشخص شده) به ۱ نزدیک شود و سایر درایه‌ها نزدیک به صفر شوند. معیار ارزیابی T2V را میانگین تشابه کسینوسی بین جفت بردارهای تصویر مجازی با تصویر واقعی‌اش قرار داده‌ایم که در بخش نتایج بررسی می‌شود.



شکل (۳): آموزش مدل T2V

پس از آموزش مدل T2V، گام بعدی ادغام^۲ ویژگی‌های متنی و تصویری برای انجام دسته‌بندی چندوجهی است. در این ساختار، اگر تصویر موجود باشد، بردار ویژگی واقعی آن به کار می‌رود؛ در غیر این صورت، بردار ویژگی تصویر با کمک مدل T2V و براساس بردار متنی تخمین زده شده و جایگزین می‌شود. در این مرحله، ابتدا بردار ویژگی‌های متنی برای تمام داده‌ها استخراج می‌شود؛ سپس برای نمونه‌هایی که تصویر در دسترس دارند، بردار ویژگی تصویری واقعی نیز به دست می‌آید؛ اگر در داده‌ای تصویرش موجود نبود، بردار ویژگی‌اش تخمین زده می‌شود-این

⁴ Learning Rate

⁵ Epoch

⁶ Batch Size

¹ Hidden Layer

² Fusion

³ Concatenate



۴- آزمایش و نتایج

در این پژوهش، از دادگاه MDDL که در مقاله [۲۸] معرفی شده است، استفاده شده است. این دادگان شامل مجموعه‌ای از توییت‌های کاربران شبکه اجتماعی توییتر است که برای تحلیل و تشخیص افسردگی طراحی شده است. داده‌های این مجموعه شامل اطلاعات متنی و تصویری است که از پست‌های کاربران استخراج شده‌اند. این مجموعه داده شامل دو گروه «کاربر افسرده» و «کاربر غیر افسرده» را در بر می‌گیرد. کاربران افسرده کسانی هستند که در محتوای منتشرشده خود به افسردگی اشاره کرده‌اند؛ مثلاً عباراتی مانند «من دچار افسردگی شده‌ام» یا «پزشک به من گفته که افسرده‌ام» را در توییت‌هایشان نوشته‌اند. در مقابل، کاربران غیر افسرده کسانی هستند که در محتوای متنی خود اشاره‌ای به افسردگی نداشته‌اند و از کلمه «افسرده» استفاده نکرده‌اند.

از هر کاربر، مجموعه‌ای از ویژگی‌ها در سه حوزه متن، تصویر و تعاملات اجتماعی استخراج شده است. این ویژگی‌ها شامل اطلاعاتی درباره ساختار شبکه اجتماعی کاربر (مانند تعداد دنبال‌کنندگان و تعاملات)، محتوای متنی پست‌ها، تصویر به اشتراک گذاشته شده و همچنین برچسب افسرده بودن یا نبودن کاربر است. در مجموع، ۳۵۷/۲۱۵ نمونه داده در اختیار داریم که از این تعداد، ۳۴/۳۹۳ نمونه دارای برچسب مثبت (افسرده) هستند. این آمار بیانگر نامتوازن بودن مجموعه داده است، به طوری که تنها حدود ۶/۹ درصد از نمونه‌ها مربوط به کاربران افسرده بود و باقی موارد به کاربران غیرافسرده تعلق دارد. برای مدیریت این نامتوانی، در مرحله آموزش مدل‌های تک‌وجهی و همچنین هنگام آموزش مدل چندوجهی، داده‌ها را متوازن کرده‌ایم (۵۰ درصد افسرده، ۵۰ درصد غیرافسرده) تا مدل دچار سوگیری^۱ نسبت به کلاس اکثریت نشود. داده‌ها به سه بخش آموزش، اعتبارسنجی و آزمون با نسبت ۷۰، ۱۵ و ۱۵ درصد تقسیم شدند تا از تعمیم‌پذیری مناسب مدل اطمینان حاصل شود.

در مراحل ابتدایی، تنها از داده‌هایی استفاده شده که تمامی وجوه اطلاعاتی-از جمله تصویر- را دارا بودند. برای شبیه‌سازی شرایط نبود وجه تصویری، ۳۰ درصد از نمونه‌ها به صورت تصادفی از

میان داده‌های کامل انتخاب و به عنوان نمونه‌های فاقد تصویر علامت‌گذاری شدند. به دلیل تصادفی بودن انتخاب، نسبت نمونه‌ها در هر کلاس تقریباً ثابت باقی ماند. این نمونه‌ها به گونه‌ای تنظیم شدند که مدل در حین آموزش و در فرایند ادغام، بتواند شرایط فقدان وجه تصویری را نیز تجربه کرده و آن را در یادگیری خود لحاظ کند. این مجموعه داده برای آموزش یک مدل چندوجهی باهدف تشخیص افسردگی در شبکه‌های اجتماعی به کار گرفته شده است. افزون بر آن، مدل برای سنجش توانایی خود در شرایطی که یکی از وجوه (تصویر) در دسترس نیست، نیز ارزیابی شد تا مشخص شود آیا می‌تواند در چنین موقعیت‌هایی نیز عملکرد مؤثری ارائه دهد یا نه.

آموزش مدل با متوازن‌سازی داده‌ها آغاز می‌شود تا تعداد نمونه‌های افسرده و غیرافسرده برابر شود. بررسی‌ها نشان داد که انجام پیش‌پردازش^۲ بر روی توییت‌ها اثر معناداری بر عملکرد مدل ندارد بنابراین، متن خام توییت‌ها بدون پیش‌پردازش وارد مدل XLM-RoBERTa می‌شود تا مدل متنی آموزش ببیند. در بخش تصویری، چند مدل مختلف از جمله ViT, ResNet18 و CLIP مورد بررسی قرار گرفتند و در نهایت ViT به دلیل دقت بهتر انتخاب شد. پس از آموزش مدل‌های تک‌وجهی، بردارهای ویژگی متن و تصویر استخراج می‌شوند و با استفاده از آن‌ها، رمزگذار T2V آموزش داده می‌شود. پارامتر margin در مرحله آموزش به صورت تجربی تعیین شد؛ چند مقدار مختلف در بازه ۰/۱ و ۰/۵ بررسی شد و مقدار ۰/۳ بر اساس عملکرد روی داده اعتبارسنجی انتخاب شد. میانگین تشابه کسینوسی در این مرحله ۶۱ درصد به دست آمد. در ادامه، ادغام داده‌ها انجام می‌شود و نتایج آن در ادامه ارائه خواهد شد.

با توجه به نتایج به دست آمده، مدل چندوجهی عملکرد بهتری نسبت به مدل‌های تک‌وجهی متن و تصویر ارائه داده است. این مدل توانسته دقت و امتیاز FI بالاتری نسبت به مدل‌های تک‌وجهی به دست آورد. همچنین، در میان مدل‌های تک‌وجهی، مدل متنی عملکرد بهتری نسبت به مدل تصویری داشته است. دقت و فراخوانی بالاتر این مدل نشان می‌دهد که روش‌های پردازش زبان

² Preprocessing

¹ Bias



رمزگذار T2V امکان جایگزینی وجه گم‌شده را فراهم آورد و سبب شد که حتی نمونه‌های ناقص نیز در فرآیند آموزش نقش داشته باشند. این موضوع نه تنها مانع کاهش حجم داده شد، بلکه باعث افزایش پایداری مدل در شرایط واقعی گردید. به طور کلی، یافته‌های این پژوهش نشان می‌دهد که ادغام چندوجهی همراه با مدیریت صحیح وجوه گم‌شده می‌تواند موجب بهبود معنادار عملکرد مدل‌های تشخیص افسردگی شود و زمینه را برای تحلیل دقیق‌تر وضعیت روان‌شناختی کاربران در شبکه‌های اجتماعی فراهم کند. در ادامه، چند مسیر پیشنهادی برای بهبود مدل مطرح شده. همچنین استفاده از مکانیزم توجه متقابل^۱ می‌تواند به یادگیری مؤثرتر ارتباط بین ویژگی‌های متنی و تصویری کمک کند. وجه‌های دیگری-شامل زمان ارسال توییت، اطلاعات کاربری (تعداد دنبال‌کنندگان و دنبال‌شوندگان) و میزان تعامل کاربران (تعداد لایک‌ها، بازشرها و پاسخ‌ها)- نیز در دسترس هستند. این اطلاعات می‌تواند بازتابی از میزان اثرگذاری پست‌ها و نوع تعاملات کاربر باشد و در تحلیل بهتر وضعیت احساسی او نقش مهمی ایفا کند. مدل کنونی هر توییت را به طور مستقل به عنوان افسرده یا غیرافسرده دسته‌بندی می‌کند؛ اما هدف اصلی، تشخیص وضعیت روانی کاربر است. در ادامه مسیر، مدل به گونه‌ای توسعه خواهد یافت که به جای تحلیل تک توییت‌ها، بتواند وضعیت روانی کاربران را بر اساس مجموعه‌ای از توییت‌های آن‌ها ارزیابی کند. این بهبودها می‌تواند عملکرد مدل را در تشخیص افسردگی افزایش داده و تحلیلی جامع‌تر از وضعیت روان‌شناختی کاربران در شبکه‌های اجتماعی ارائه دهند.

اقرار

در اینجا لازم می‌دانیم تا سپاسگزاری خود را از حمایت‌های بی‌دریغ اپراتور هوش مصنوعی (هورا) شرکت ایران جی‌پی‌یو (در راستای تأمین سخت‌افزار محاسباتی گرافیکی مورد نیاز این تحقیق ذیل تفاهم‌نامه به شماره ۰۳۰۹۱۳/۱۰۰ در جهت تعالی حوزه هوش مصنوعی کشور اعلام کنیم.

طبیعی در تحلیل محتوای متنی و شناسایی افسردگی مؤثرند. این نتایج-که در جدول ۲ قابل مشاهده است-حاکی از آن است که مدل‌های چندوجهی با ترکیب متن و تصویر می‌توانند ویژگی‌های مکمل را از هر دو منبع استخراج کنند. در نتیجه، استفاده از چنین رویکردی می‌تواند به بهبود عملکرد مدل و دقت بیشتر در تحلیل احساسات و تشخیص افسردگی در شبکه‌های اجتماعی منجر شود.

جدول (۲): نتایج مدل‌های تک‌وجهی و مدل چندوجهی در

تشخیص افسردگی در سطح توییت

مدل	دقت روی داده آموزش	دقت روی داده تست	معیار F1 روی داده تست
متن	۹۶/۷۹	۸۷/۸۷	۸۷/۹۴
تصویر	۸۱/۱۵	۷۳/۳۷	۷۴/۳۲
چند وجهی	۹۵/۳۲	۹۰/۱۷	۹۰/۶۴

۵- نتیجه‌گیری

در این پژوهش، از روشی مشابه مقاله [۲۹] به منظور مدیریت وجه گم‌شده در مدل‌های چندوجهی جهت تشخیص افسردگی در شبکه‌های اجتماعی استفاده شد. این روش به جای حذف نمونه‌های ناقص، می‌کوشد ویژگی‌های وجه گم‌شده را از اطلاعات موجود پیش‌بینی برآورد کند و از این طریق، از داده‌های ناقص نیز در فرآیند آموزش بهره‌بردارد. نتایج نشان داد که مدل متنی نسبت به مدل تصویری عملکرد بهتری دارد. دلیل این امر آن است که در داده‌های شبکه‌های اجتماعی، محتوای متنی مستقیماً احساسات و افکار کاربران را منعکس می‌کند و بنابراین سرنخ‌های قوی‌تری برای شناسایی افسردگی فراهم می‌آورد، در حالی که تصاویر اغلب متنوع، غیرساختاریافته یا کم‌ارتباط با وضعیت روانی کاربر هستند. این یافته با مطالعات پیشین [۱۰، ۱۲، ۱۳] نیز هم‌راستا است که تأکید داشته‌اند زبان کاربران افسرده ویژگی‌های متمایزی دارد. همچنین، مدل چندوجهی عملکرد بالاتری نسبت به مدل‌های تک‌وجهی کسب کرد. این برتری ناشی از ترکیب اطلاعات متنی و تصویری و ایجاد هم‌افزایی میان آن‌ها است؛ به گونه‌ای که ضعف هر وجه توسط دیگری جبران می‌شود. علاوه بر این، بهره‌گیری از

¹ Cross-Attention



References

- [1] World Health Organization (WHO), "Depressive disorder (depression)," WHO Fact Sheets, Aug. 29, 2025. [Online]. Available: <https://www.who.int/news-room/fact-sheets/detail/depression>
- [2] W. Guo, J. Wang, and S. Wang, "Deep multimodal representation learning: a survey," *IEEE Access*, vol. 7, pp. 63373–63394, 2019, doi: 10.1109/ACCESS.2019.2916887.
- [3] P. Xu, X. Zhu, and D. A. Clifton, "Multimodal learning with transformers: a survey," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 45, no. 10, pp. 12113–12132, 2023, doi: 10.1109/TPAMI.2023.3275156.
- [4] D. Ailyn, "Multimodal data fusion techniques," 2024. [Online]. Available: https://www.researchgate.net/publication/383887675_Multimodal_Data_Fusion_Techniques
- [5] T. Meng, X. Jing, Z. Yan, and W. Pedrycz, "A survey on machine learning for data fusion," *Information Fusion*, vol. 57, pp. 115–129, 2020, doi: 10.1016/j.inffus.2019.12.001.
- [6] M. Pawłowski, A. Wróblewska, and S. Sysko-Romańczuk, "Effective techniques for multimodal data fusion: a comparative analysis," *Sensors*, vol. 23, p. 2381, 2023, doi: 10.3390/s23052381.
- [7] D. Lee, S. Park, J. Kang, D. Choi, and J. Han, "Cross-lingual suicidal-oriented word embedding toward suicide prevention," in *Findings of the Association for Computational Linguistics: EMNLP 2020*, R. Cotterell, S. Eger, and S. Wiseman, Eds., Online, Nov. 2020, pp. 2208–2217, doi: 10.18653/v1/2020.findings-emnlp.200.
- [8] B. G. Bokolo and Q. Liu, "Deep learning-based depression detection from social media: comparative evaluation of ML and transformer techniques," **Electronics**, vol. 12, no. 21, p. 4396, 2023, doi: 10.3390/electronics12214396.
- [9] L. Zhu, Z. Zhu, C. Zhang, Y. Xu, and X. Kong, "Multimodal sentiment analysis based on fusion methods: a survey," **Information Fusion**, vol. 95, pp. 306–325, 2023, doi: 10.1016/j.inffus.2023.02.028.
- [10] G. Coppersmith, M. Dredze, and C. Harman, "Quantifying mental health signals in Twitter," in *Proc. Workshop on Computational Linguistics and Clinical Psychology: From Linguistic Signal to Clinical Reality*, Baltimore, MD, USA, Jun. 27, 2014, pp. 51–60, doi: 10.3115/v1/W14-3207.
- [11] M. Deshpande and V. Rao, "Depression detection using emotion artificial intelligence," in **Proc. 19th IEEE Int. Conf. Intelligent Sustainable Systems**, Palladam, Tirupur, India, Dec. 2017, pp. 858–862, doi: 10.1109/ISS1.2017.8389299.
- [12] M. De Choudhury, M. Gamon, S. Counts, and E. Horvitz, "Predicting depression via social media," in *Proc. Int. AAAI Conf. Web and Social Media*, Limassol, Cyprus, Jun. 5–8, 2013, vol. 7, pp. 128–137, doi: 10.1609/icwsm.v7i1.14432.
- [13] A. Murarka, B. Radhakrishnan, and S. Ravichandran, "Classification of mental illnesses on social media using RoBERTa," in *Proc. 12th Int. Workshop on Health Text Mining and Information Analysis*, E. Holderness et al., Eds., Online, 2021, pp. 59–68. [Online]. Available: <https://aclanthology.org/2021.louhi-1.7/>
- [14] S. Yang, L. Cui, L. Wang, T. Wang, and J. You, "Enhancing multimodal depression diagnosis through representation learning and knowledge transfer," **Heliyon**, vol. 10, no. 4, p. e25959, Feb. 2024, doi: 10.1016/j.heliyon.2024.e25959.
- [15] T. Baltrušaitis, C. Ahuja, and L. Morency, "Multimodal machine learning: a survey and taxonomy," **IEEE Trans. Pattern Anal. Mach. Intell.**, vol. 41, no. 2, pp. 423–443, 2019, doi: 10.1109/TPAMI.2018.2798607.
- [16] S. Yang, L. Cui, L. Wang, T. Wang, and J. You, "Cross-modal contrastive learning for multimodal sentiment recognition," **Applied Intelligence**, vol. 54, pp. 4260–4276, 2024, doi: 10.1007/s10489-024-05355-8.
- [17] M. Fang, S. Peng, Y. Liang, C.-C. Hung, and S. Liu, "A multimodal fusion model with multi-level attention mechanism for depression detection," **Biomedical Signal Processing and Control**, vol. 82, p. 104561, Apr. 2023, doi: 10.1016/j.bspc.2022.104561.
- [18] Y. Wang, Z. Wang, C. Li, Y. Zhang, and H. Wang, "Online social network individual depression detection using a multitask heterogeneous modality fusion approach," **Information Sciences**, vol. 609, pp. 727–749, 2022, doi: 10.1016/j.ins.2022.07.109.
- [19] R. Wu, H. Wang, H.-T. Chen, and G. Carneiro, "Deep multimodal learning with missing modality: a survey," **arXiv preprint* arXiv:2409.07825*, 2024.
- [20] L. Cai, Z. Wang, H. Gao, D. Shen, and S. Ji, "Deep adversarial learning for multi-modality missing data completion," in **Proc. 24th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining**, pp. 1158–1166, 2018, doi: 10.1145/3219819.3219963.
- [21] M. Ma, J. Ren, L. Zhao, S. Tulyakov, C. Wu, and X. Peng, "SMIL: Multimodal learning with



- severely missing modality,” in *Proc. 35th AAAI Conf. Artificial Intelligence (AAAI '21)*, vol. 35, no. 3, pp. 2302–2310, 2021, doi: 10.1609/AAAI.V35I3.16330.
- [22] H. Wang, Y. Chen, C. Ma, J. Avery, L. Hull, and G. Carneiro, “Multi-modal learning with missing modality via shared-specific feature modelling,” in Proc. 2023 IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR), pp. 15878–15887, 2023, doi: 10.1109/CVPR56347.2023.
- [23] M. K. Reza, A. Prater-Bennette, and M. S. Asif, “Robust multimodal learning with missing modalities via parameter-efficient adaptation,” *arXiv preprint* arXiv:2310.03986, 2024.
- [24] Q. Wang, L. Zhan, P. Thompson, and J. Zhou, “Multimodal learning with incomplete modalities by knowledge distillation,” in *Proc. 26th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining*, 2020, pp. 1828–1838, doi: 10.1145/3394486.3403234.
- [25] X. Chen, S. Kornblith, M. Noroozi, and G. E. Hinton, “A simple framework for contrastive learning of visual representations,” in *Proc. 37th Int. Conf. Machine Learning (ICML)*, 2020, pp. 1597–1607, doi: 10.48550/arXiv.2002.05709.
- [26] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever, “Learning transferable visual models from natural language supervision,” in *Proc. Int. Conf. Mach. Learn. (ICML)*, vol. 139, pp. 8748–8763, 2021, doi: 10.48550/arXiv.2103.00020.
- [27] A. Jaiswal, A. R. Babu, M. Z. Zadeh, D. Banerjee, and F. Makedon, “A survey on contrastive self-supervised learning,” *arXiv preprint* arXiv:2010.05113, 2021.
- [28] G. Shen, J. Jia, L. Nie, F. Feng, C. Zhang, T. Hu, T.-S. Chua, and W. Zhu, “Depression detection via harvesting social media: a multimodal dictionary learning solution,” in *Proc. 26th Int. Joint Conf. Artif. Intell. (IJCAI)*, 2017, pp. 3838–3844, doi: 10.24963/ijcai.2017/535.
- [29] C. Zhang, F. Nian, and J. Lee, “Toward robust multimodal learning using multimodal foundational models,” *arXiv preprint* arXiv:2401.13697, 2024, doi: 10.48550/arXiv.2401.13697.

Contrastive-Learning-Based Handling of Missing Modalities in Multimodal Data Fusion for Depression Detection on Social Networks

Hamed Marvi¹, Abolfazl Nadi^{2*}, Mohammad Mehdi Keikha³

¹ MS.c Student, Department of Computer Science, School of Mathematics, Statistics and Computer Science, College of Science, University of Tehran, Tehran, Iran

² Assistant Professor, Department of Computer Science, School of Mathematics, Statistics and Computer Science, College of Science, University of Tehran, Tehran, Iran

³ Assistant Professor, Department of Computer Science, University of Sistan and Baluchestan, Zahedan, Iran

Article Information

Original Research Paper

Received:

2025 August 25

Accepted:

2025 October 29

Keywords:

Multimodal Models, Missing Modality, Multimodal Data Fusion, Contrastive Learning, Social Media

Corresponding Author*:

a.nadi@ut.ac.ir

Abstract

The analysis of social network data plays a fundamental role in uncovering users' behavioral patterns. Multimodal models that combine textual, visual, and other information sources are effective tools for such analyses. However, a major challenge in these models is the absence of certain modalities in parts of the dataset; for instance, a user may post only text without sharing any images. This issue prevents multimodal models from fully exploiting all available information. In this paper, a method is proposed for leveraging incomplete data within multimodal models. First, unimodal models are trained independently to process each modality. Then, a contrastive learning-based encoder is designed and trained to estimate the feature vector of the missing modality using the features of the available modalities. Finally, textual and visual data—either real or reconstructed—are fused in a multimodal model to analyze user behavior. Experimental results on the MDDL dataset, based on accuracy and F1-score metrics, demonstrate that the proposed multimodal model achieves superior performance, reaching an accuracy of 90.17% and an F1-score of 90.64%, outperforming unimodal text-based (87.87% accuracy) and image-based (73.37% accuracy) models. These findings confirm that the proposed model effectively utilizes available information without discarding incomplete data.

 : 10.22034/ABMIR.2025.23578.1160

E-ISSN: [2821-2037](https://doi.org/10.22034/ABMIR.2025.23578.1160) /The Author 2025. Published by Yazd University This is an open access article under the CC BY 4.0 License (<https://creativecommons.org/licenses/by/4.0/>).

