

افزایش داده و انتخاب مؤثر ویژگی در شبکه‌های مولد متخاصمی جهت تشخیص احساس از گفتار

آرش شیلاندری^{۱*}، حسین مروی^۲ و حسین خسروی^۳

اطلاعات مقاله	چکیده
<p>نوع مقاله: پژوهشی دریافت مقاله: ۱۴۰۰/۰۸/۰۳ بازنگری مقاله: ۱۴۰۰/۱۰/۲۸ پذیرش مقاله: ۱۴۰۰/۱۱/۰۴</p>	<p>تاکنون، یقینی مبتنی بر موفقیت و یا عدم موفقیت به‌کارگیری روش‌های کاهش ویژگی جهت افزایش کارایی سیستم‌های تشخیص احساس از گفتار حاصل نشده است. این مقاله با هدف افزایش داده‌ها در یک سیستم تشخیص احساس از گفتار، انتخاب ویژگی را مورد بحث و بررسی قرار می‌دهد. آزمایش‌ها بر روی چهار پایگاه داده متداول EMO-DB، eINTERFACE05، SAVEE و IEMOCAP در نرم‌افزار پایتون انجام گردیده و علاوه بر این، تجزیه و تحلیل داده‌ها بر روی هر چهار پایگاه داده برای چهار احساس غم، عصبانیت، خوشحالی و خنثی ارائه خواهد شد. یک شبکه افزایش داده متخاصمی جهت افزایش نمونه‌ها و دو شبکه انتخاب ویژگی ترکیبی معیار فیشر و الگوریتم جداساز خطی طی دو مرحله و با فیدبکی که از شبکه طبقه‌بند گرفته می‌شود، سیستم تشخیص احساس از گفتار را به نقطه بهینه‌ای از تعداد و ابعاد داده‌ها رسانیده و نشان می‌دهد آنالیز مؤلفه‌های اصلی روی داده‌های همبسته، مؤثرتر و الگوریتم جداساز خطی روی داده‌های با بعد کم بهتر عمل می‌کنند. همچنین که روش فیشر در کاهش سایز بهتر از آنالیز مؤلفه‌های اصلی عمل می‌کند. همچنین ماشین بردار پشتیبان جهت طبقه‌بندی احساسات مورد استفاده قرار گرفته است. نتایج نشان می‌دهد که استفاده از هر دو روش جداساز خطی و معیار فیشر به طور هم‌زمان در سیستم افزایش داده متخاصمی می‌تواند ویژگی‌ها را در ابعاد کمتر فیلتر نموده درحالی‌که اطلاعات احساسی را جهت طبقه‌بندی حفظ نماید. نتایج به‌دست آمده با تحقیقات اخیر مقایسه گردیده است که حاکی از دستیابی روش پیشنهادی به صحت ۸۶٫۳۲٪ در پایگاه داده برلین می‌باشد.</p>
<p>واژگان کلیدی: پردازش گفتار، انتخاب ویژگی، افزایش داده، تشخیص احساس از گفتار، شبکه‌های مولد متخاصمی.</p>	

۱- مقدمه

تشخیص احساس از گفتار تاکنون به دلیل نبود پایگاه داده با داده‌های فراوان و در دسترس با مسائل فراوانی روبه‌رو بوده است. عدم توانایی در انتخاب ویژگی‌های مهم و تأثیرگذار در شناسایی احساس، امکان استفاده از این سیستم‌ها را در کاربردهای برخط محدود می‌کند. همچنین وابسته بودن این سیستم به زبان، لهجه، سن، جنسیت و نوع حالت گوینده از عمده مسائل این روش بوده است.

روش‌های انتخاب ویژگی نقش مهمی در عملکرد بهینه مدل‌های یادگیری ماشین دارند. اگر ویژگی‌ها به‌درستی جهت آموزش شبکه عصبی انتخاب شوند این شبکه‌ها می‌توانند در زمانی کوتاه آموزش دیده و به‌درستی پاسخ دهند. تشخیص احساسات از گفتار شامل استخراج ویژگی و تجزیه و تحلیل بردارهای ویژگی استخراج شده از سیگنال گفتار می‌شود. سیستم تشخیص احساس از سیگنال گفتار از دیدگاه تشخیص الگو شامل پنج بخش می‌باشد:

* پست الکترونیک نویسنده مسئول: Shilandari@shahroodut.ac.ir

۱. دانشجوی دکتری تخصصی، دانشکده مهندسی برق، دانشگاه شاهرود

۲. دانشیار، دانشکده مهندسی برق، دانشگاه شاهرود

۳. دانشیار، دانشکده مهندسی برق، دانشگاه شاهرود

انتخاب ویژگی به‌منظور حذف ویژگی‌های اضافی و غیرمرتبط (نویزی) امری ضروری می‌باشد. ویژگی‌های اضافی ویژگی‌هایی هستند که اضافه‌شدن آن‌ها به بردارهای ویژگی پایگاه‌داده، اطلاعات جدیدی به مسئله اضافه نکند و بردارهای ویژگی غیرمرتبط نمونه‌هایی هستند که در طبقه‌بندی اهمیتی نداشته و حاوی اطلاعات مفیدی از داده‌های آن طبقه نباشند [۲].

به منظور انتخاب یک بردار ویژگی بهینه، باید همه زیرمجموعه‌های ممکن را بررسی نمود. اما از آن جا که این کار امری بسیار زمان‌بر و غیرعملی می‌باشد، ناگزیر به‌منظور انتخاب یک بردار ویژگی مناسب از الگوریتم‌های جستجو استفاده می‌شود. همچنین، استفاده از الگوریتم‌های کاهش ویژگی^۴ که با استفاده از یک تبدیل مناسب ویژگی را از فضای اصلی به فضایی با ابعاد پایین‌تر منتقل می‌کنند نیز متداول است [۱].

نکته جالب توجه دیگر، پدیدهٔ اوج^۵ می‌باشد که تعداد بهینه ویژگی‌ها را به‌صورت تابعی از تعداد نمونه‌ها و همبستگی بین ویژگی‌ها بیان می‌کند [۳]. با افزایش تعداد ویژگی‌های مورد استفاده، ابتدا نرخ تشخیص افزایش پیدا می‌کند اما پس از رسیدن به نقطه اوج، افزایش بیشتر تعداد ویژگی‌ها موجب کاهش نرخ تشخیص می‌شود.

هدف اصلی انتخاب ویژگی، تعیین خصوصیتی است که بهترین طبقه‌بندی را منجر می‌شوند. به نظر می‌رسد با انتخاب ویژگی‌های بهینه، سایز مجموعه ویژگی‌ها کاهش می‌یابد و عملکرد و دقت طبقه‌بندی افزایش می‌یابد. اگرچه تاکنون تعداد زیادی از روش‌های انتخاب ویژگی برای مطالعات تشخیص احساس از گفتار^۶ استفاده شده است ولی برخی از این روش‌ها منجر به کاهش دقت در تشخیص احساس می‌شوند چرا که این روش‌ها کلی بوده و مخصوص استفاده در سیگنال‌های صوت نمی‌باشند.

برای انتخاب بهترین ویژگی‌ها برای یک مدل یادگیری نظارت شده، روش‌های انتخاب ویژگی نظارت شده^۷ ارائه شده‌اند. هدف این دسته از الگوریتم‌ها، انتخاب بهترین زیرمجموعه از ویژگی‌ها برای تضمین عملکرد بهینه یک مدل نظارت شده است. این الگوریتم‌ها برای انتخاب بهترین ویژگی‌ها، از داده‌های برچسب زده شده^۸ استفاده می‌کنند.

۱- پیش‌پردازش ۲- استخراج ویژگی ۳- افزایش داده ۴- انتخاب ویژگی ۵- طبقه‌بندی. تعداد ویژگی‌های استخراج شده بسته به تعیین پارامترهای صوتی در نظر گرفته شده داشته و تغییرات آماری این پارامترها منجر به استخراج بردارهای ویژگی با ابعاد بالا می‌شوند. از این رو، از روش‌های انتخاب ویژگی جهت کاهش ابعاد بردارهای ویژگی استفاده می‌شود. همچنین ویژگی‌های نامرتبط یا تا حدودی مرتبط می‌توانند تأثیر منفی بر عملکرد سیستم داشته باشند. هرچند که معمولاً تعداد نمونه‌ها در پایگاه‌های داده احساسی اندک است ولی ابعاد بردارهای ویژگی استخراج شده در مقایسه با تعداد آنها بالاست. همه این ویژگی‌ها جهت شناخت احساسات مؤثر نیستند و از طرفی احساسات مختلف می‌توانند ویژگی‌های مختلف گفتار را تحت تأثیر قرار دهند [۱].

پیاده‌سازی روش‌های افزایش داده و انتخاب ویژگی و پاک‌سازی داده‌ها، اولین و مهم‌ترین گام در طراحی مدل‌های هوشمند یادگیری قلمداد می‌شوند. همچنین، هنگامی که ابعاد فضای ویژگی داده‌ها در مقایسه با تعداد نمونه‌ها بسیار زیاد است، با معضل نفرین ابعاد بالا^۲ مواجه خواهیم شد، استفاده از مجموعه ویژگی‌های مناسب، هزینه‌های محاسباتی (نظیر زمان لازم برای آموزش سیستم و همچنین استفاده از آنها در کاربردهای برخط) لازم برای آموزش بهینه سیستم را به‌شدت کاهش می‌دهد. محاسبه درجه اهمیت ویژگی‌ها و استفاده از آن‌ها در مرحله انتخاب ویژگی، گام مهمی در جهت تفسیرپذیری^۳ مدل‌های یادگیری ماشین بخصوص پس از افزایش داده خواهد بود [۱].

استفاده از تعداد ویژگی‌های زیاد به‌منظور طراحی سیستم تشخیص الگو باعث افت عملکرد الگوریتم طبقه‌بندی به دلیل تعدد پارامترهای تنظیم، پیچیدگی محاسباتی و آموزش ناکارآمد به دلیل ابعاد بالای مسئله می‌شود. تعداد بهینه ویژگی‌ها را می‌توان به‌عنوان تابعی از تعداد نمونه‌های پایگاه‌داده و همبستگی بین ویژگی‌ها بیان نمود. اگر تعداد ویژگی‌ها بیش از تعداد بهینه انتخاب شود، با پدیده نفرین ابعاد داده و یا بیش‌برازش مواجه شده و نرخ تشخیص کاهش پیدا می‌کند؛ بنابراین استفاده از الگوریتم‌های

⁵ Peaking Phenomenon

⁶ Speech Emotion Recognition (SER)

⁷ Supervised Feature Selection

⁸ Labelled Data

¹ Data Cleaning

² Curse of Dimensionality

³ Interpretability

⁴ Feature Reduction

روش‌ها، هر یک از ویژگی‌ها که در فضای جدید بیان می‌شوند شاید از نظر فیزیکی مفهوم خاصی نداشته باشند اما حاوی چکیده اطلاعات ویژگی‌های اصلی می‌باشند. بدین منظور برای حل یک مسئله با C ویژگی، ابعاد آن به $C-1$ ویژگی (فضای $C-1$ بعدی) کاهش می‌یابد.

تا کنون روش‌های مؤثر و متعددی ارائه شده است که با معرفی ویژگی‌های جدید و کارا دقت سیستم‌های بازشناسی احساس را افزایش می‌دهند. در این مقاله، نرم‌افزار `openSMILE [5]` جهت استخراج ویژگی‌ها استفاده گردیده و با روش‌های افزایش داده متخاصمی سعی در تولید بردارهای ویژگی جدید جهت افزایش نمونه‌های آموزش و متعادل نمودن پایگاه‌های داده شده است. همچنین، با استفاده از روش‌های انتخاب ویژگی، ویژگی‌های کمتر اثرگذار بر روی کارایی سیستم تشخیص احساس از گفتار حذف و ویژگی‌های مؤثر معرفی گردیدند. کمبود داده‌ها می‌تواند باعث شود که یک مدل یادگیری ماشین قادر به یادگیری توزیع واقعی داده‌ها نباشد و در نهایت منجر به مشکل بیش‌برازش می‌شود. به‌عنوان مثال، در هنگام آموزش یک شبکه عصبی عمیق با نمونه‌های آموزشی کم در هر طبقه، بیش‌برازش رخ می‌دهد، در صورتی که هر نمونه هزاران ویژگی دارد. برای حل این مشکل، می‌توان از قاعده‌مندسازی^{۱۰} برای ایجاد محدودیت در مدل استفاده کرد [6]. راه‌حل کلی دیگر، کاهش ابعاد با محدودیت پراکندگی است [7]. این رویکرد زمانی مؤثر خواهد بود که ویژگی‌های اضافی وجود داشته باشد. در غیر این صورت، اطلاعات مفید را از بین می‌برد و منجر به تخریب عملکرد می‌شود. همچنین، وجود بردارهای ویژگی کم‌اهمیت در آموزش شبکه عصبی، محاسبات را در این شبکه پیچیده و تأثیرگذاری نمونه‌ها را در آموزش بی‌اثر خواهد نمود؛ لذا انتخاب و حذف بردارهای ویژگی کم‌اهمیت موضوعی است که کمتر به آن در روش‌های افزایش داده به طور هم‌زمان پرداخته شده است. یکی از روش‌های مؤثر برای تقویت و افزایش داده‌ها، استفاده از شبکه‌های مولد متخاصمی است که توسط گودفلو^{۱۱} و همکاران در سال ۲۰۱۴ معرفی شده است [8]. در سال‌های اخیر،

بالین‌حال، در شرایطی که داده‌های برجسب زده در دسترس نیستند (یادگیری نظارت نشده)، روش‌هایی به نام روش‌های انتخاب ویژگی نظارت نشده^۱ پیاده‌سازی شده‌اند که ویژگی‌ها را بر اساس معیارهای مختلفی نظیر واریانس، آنتروپی، قابلیت ویژگی‌ها در حفظ اطلاعات مرتبط با مشابهت‌های محلی^۲ و سایر موارد امتیازبندی می‌کنند.

ویژگی‌های مرتبطی که از طریق فرایندهای مکاشفه‌ای نظارت نشده^۳ شناسایی شده‌اند، می‌توانند در مدل‌های یادگیری نظارت شده نیز مورداستفاده قرار بگیرند. چنین کاربردهایی از ویژگی‌های شناسایی شده، به سیستم یادگیری نظارت شده اجازه می‌دهد تا علاوه بر شناسایی میزان همبستگی ویژگی‌ها با برجسب طبقه داده‌ها، الگوهای دیگری نیز در داده‌های یادگیری شناسایی کنند. از دیدگاه طبقه‌بندی، روش‌های انتخاب ویژگی را می‌توان در چهار دسته طبقه‌بندی کرد: ۱- روش‌های فیلتر^۴ - روش‌های بسته‌بندی^۵ - روش‌های تعبیه شده^۶ - روش‌های ترکیبی^۷. گزینه دیگر برای انتخاب بهترین ویژگی‌ها، ترکیب روش‌های فیلتر و بسته‌بندی است. در این مقاله از یک فرایند دومرحله‌ای برای ترکیب دو روش فیلتر و بسته‌بندی استفاده شده است. در مرحله اول ویژگی‌ها بر اساس مشخصه‌های آماری فیلتر می‌شوند. در مرحله بعد، با استفاده از یک روش انتخاب ویژگی بسته‌بندی، بهترین ویژگی‌ها برای آموزش یک مدل یادگیری انتخاب می‌شوند و به طور مشخص از ترکیب معیار فیشر به‌عنوان روش انتخاب ویژگی فیلتر و از الگوریتم جداساز خطی به‌عنوان مکمل جهت انتخاب ویژگی‌ها استفاده شده است.

در الگوریتم‌های انتخاب ویژگی بر پایه فیلتر، بسته‌بندی و تعبیه شده، با حذف ویژگی‌های اضافی و غیرمرتبط یا انتخاب ویژگی‌های مؤثر، یک زیرمجموعه بهینه از ویژگی‌ها به دست می‌آید بطوریکه ویژگی‌های انتخابی بدون هیچ‌گونه تغییری از بین ویژگی‌های اصلی انتخاب می‌شوند. در مقابل این روش‌ها، الگوریتم‌های کاهش ویژگی نظیر الگوریتم جداسازی خطی^۸ و آنالیز مؤلفه‌های اصلی^۹ نیز وجود دارند که با اعمال یک تبدیل بر روی همه ویژگی‌ها آن‌ها را به یک فضای جدید با ابعاد کمتر انتقال می‌دهند [4]. در این

⁷ Hybrid

⁸ Linear Discriminant Analysis

⁹ Principal Component Analysis

¹⁰ Regularization

¹¹ Ian Goodfellow

¹ Unsupervised Feature Selection

² Local Similarity

³ Unsupervised Heuristics

⁴ Filter

⁵ Wrapper

⁶ Embedded

از گفتار را بهبود بخشند. به بیان دیگر، فرایند انتخاب ویژگی در داده‌های تولید شده سبب می‌شود که آنها هدفمند به داده‌های اصلی اضافه گردند و سبب آموزش هر چه بهتر شبکه طبقه‌بند احساس شوند و حجم محاسبات را کاهش و کارایی این شبکه را بیفزایند.

یک مشکل اساسی در آموزش شبکه‌های مولد متخاصمی، انتخاب ویژگی‌های مؤثر در سیستم تشخیص احساس از گفتار است. اگر ویژگی‌ها به صورت هوشمندانه افزوده نشوند، حجم محاسبات بالا رفته و سیستم تشخیص احساس کند و در استفاده‌های برخط غیر قابل استفاده خواهد شد. همچنین ممکن است بردارهای ویژگی‌های افزوده شده، مرز بین دو طبقه را از بین برده و یا داده‌های مرزی تولید شوند که کار آموزش شبکه طبقه‌بند را مشکل و یا راندمان طبقه‌بندی را پایین آورند. از این رو توجه به انتخاب ویژگی همیشه لازم خواهد بود.

بخش ۲ مروری بر راه‌حل‌های رایج برای مشکل کمبود داده‌ها و روش‌های استفاده شده جهت انتخاب ویژگی را ارائه می‌دهد. بخش ۳ طراحی شبکه پیشنهادی را توصیف می‌کند و تجزیه و تحلیل نظری را ارائه می‌دهد. بخش ۴ جزئیات آزمایش‌ها، از جمله توصیف داده‌ها، ویژگی‌ها، تنظیمات آزمایشی و پروتکل‌های ارزیابی را معرفی می‌کند. بخش ۵ نتایج تجربی را ارائه و تحلیل می‌کند. سرانجام، بخش ۶ نتیجه‌گیری و کارهای آینده را ارائه می‌دهد.

۲- مطالعات مرتبط

در سال ۲۰۱۱، بوژارت و همکارانش با وزن‌دهی به ضرایب مل کپستروپل باتوجه‌به موقعیت فرمت‌ها، ویژگی‌های جدیدی به نام WMFCC^۴ را ابداع نموده و با استفاده از این ویژگی‌ها و یک طبقه‌بند مبتنی بر مدل مخفی مارکوف به طبقه‌بندی چهار احساس خشم، آرامش، تأکید و عادی بر روی پایگاه داده AIBO پرداختند. آن‌ها به منظور ارزیابی روش خود از نرخ تشخیص بدون وزن استفاده نموده و به عدد ۷۰٪ دست یافتند [۹]. در یکی از تحقیقات شاخص صورت‌گرفته در این سال، هو و همکارانش ویژگی‌های طیفی جدیدی را معرفی نموده و از این ویژگی‌ها به منظور طبقه‌بندی هفت احساس موجود در پایگاه داده برلین^۵ و

شبکه‌های مولد متخاصمی به‌عنوان یکی از موفق‌ترین رویکردها برای تولید نمونه شناخته شده‌اند [۳۵]. با استفاده از یک بازی مخالف بین یک شبکه تشخیص‌دهنده و یک شبکه مولد، شبکه‌های مولد متخاصمی آموزش می‌بینند تا نمونه‌هایی تولید کنند که از داده‌های واقعی قابل تشخیص نیستند. علاوه بر این، آنها دارای سه خصوصیت عمده هستند [۸]: ۱- شبکه‌های مولد متخاصمی می‌توانند توزیع احتمال در مسائل پیچیده دنیای واقعی را یاد بگیرند.

۲- شبکه‌های مولد متخاصمی می‌توانند با داده‌های از دست‌رفته (داده‌های نویزی) آموزش داده شوند، حتی زمانی که برچسب‌های بسیاری از نمونه‌ها وجود ندارد و ۳- شبکه‌های مولد متخاصمی دارای خروجی‌های چند مدول هستند به این معنی که آنها می‌توانند چندین جواب صحیح متفاوت تولید کنند و تنوع نمونه‌های تولید شده را افزایش دهند ولی همچنان توجه به انتخاب ویژگی‌های ساخته شده و فیلتر نمودن نمونه‌ها امری ضروری است که در این مقاله به آن پرداخته شده است.

برخی مقالات از روش‌های GMM و HMM استفاده می‌کنند تا توزیع ویژگی‌های آکوستیکی را یاد بگیرند و سپس از طبقه‌بندی بیزین یا اصل ماکزیمم بخت^۱ برای تشخیص احساس استفاده می‌کنند. برخی از مدل‌های پس‌زمینه جامع^۲ استفاده می‌کنند تا از ویژگی‌های آکوستیکی، بردار ویژگی برای آموزش شبکه ماشین بردار پشتیبان بسازند که این روش بیشتر برای تشخیص گوینده استفاده می‌شود. برخی، روش‌های آماری را روی ویژگی‌های سطح پایین استفاده می‌کنند تا ویژگی‌های آماری سطح بالا را به دست آورند و سپس از ماشین بردار پشتیبان برای طبقه‌بندی ویژگی‌های کلی استفاده می‌شود. برخی از روش‌های KNN و یا درخت تصمیم‌گیری^۳ استفاده می‌کنند که نیاز به تعداد ویژگی‌های زیاد و ویژگی‌های ساخته شده به صورت دستی دارد.

این مقاله به کاربرد شبکه‌های مولد متخاصمی برای تولید داده‌های مصنوعی و انتخاب بهترین آنها پرداخته است که هدف از آن تولید داده‌هایی است که توزیع داده‌های واقعی را گسترش می‌دهند تا عملکرد سیستم تشخیص احساس

⁵ A database of German emotional speech

¹ Maximum Likelihood

² Universal Background Models (UBMs)

³ Decision Tree

⁴ Weighted MFCC

ANOVA بهره بردند. آنها نشان دادند که حتی با حذف ۲۲٪ از ویژگی‌ها، به طور متوسط دقت تشخیص احساسات را می‌توان تا ۲٫۲٪ بهبود داد. همچنین، شبکه ماشین بردار پشتیبان پیشنهادی دقت تشخیص را حداقل به میزان ۸٪ در مقایسه با طبقه‌بندی کننده یکپارچه شبیه‌سازی شده بهبود می‌دهد.

در سال ۱۳۹۶، حریمی و همکارش تشخیص احساسات از گفتار را با استفاده از مدل و ویژگی‌های دینامیکی غیرخطی بررسی نموده و احساس‌های با سطح برانگیختگی یکسان را با استفاده از ویژگی‌های دینامیکی غیرخطی از یکدیگر جدا نمودند. در روش پیشنهادی که بر روی پایگاه‌داده برلین با استفاده از تکنیک ۱۰ تکه برابر ارزیابی شد، نرخ بازشناسی ۹۶٫۳۵٪ و ۸۷٫۱۸٪ برای زنان و مردان به دست آمد [۱۳].

در مرجع [۱۴] ابتدا برای هر سگمنت یعنی مجموعه‌ای از فریم‌های پشت‌سرهم، با استفاده از شبکه عصبی عمیق یک توزیع احتمالی حالت احساسی مبتنی بر سگمنت تولید شده و سپس از روی این توزیع‌های احتمالی، ویژگی‌های در سطح کلی به دست آمده است. این ویژگی‌های کلی به یک ELM داده شده (یک شبکه عصبی خاص تک‌لایه مخفی ساده و مؤثر) تا احساسات کلی به دست آید. علت استفاده از ELM در گام آخر این است که ELM ساده‌تر از DNN است و داده‌های کمتری برای آموزش نیاز دارد و در اینجا عملکرد خیلی بهتری از ماشین بردار پشتیبان دارد. نتایج عملی نشان می‌دهند که رویکرد پیشنهاد شده اطلاعات احساسی را به طرز مؤثری از ویژگی‌های سطح پایین به دست می‌آورد و منجر به افزایش ۲۰٪ دقت در مقایسه با رویکردهای جدید شده است.

در مرجع [۱۵] از روش اصلاح شده VQ برای کاهش ویژگی‌های آکوستیکی استفاده شده است. به این صورت که روش VQ و همچنین روش VQ تفاضلی روی ویژگی‌های آکوستیکی در سطح فریم استفاده شده است. استفاده از VQ تفاضلی پایداری ویژگی‌ها را با اضافه کردن دینامیک زمانی که حین آنالیز آماری از بین می‌رود تقویت می‌کند. مکانیزم استفاده شده در اظهارات مختلف ارزیابی شده است به منظور رسیدن به این هدف یک پایگاه‌داده محلی هم‌زمان

تخمین احساس در فضای پیوسته بر روی پایگاه‌داده VAM استفاده نمودند. آن‌ها با استفاده از الگوریتم انتخاب ویژگی دومرحله‌ای فیلتر - بسته‌بند مبتنی بر معیار فیشر، الگوریتم جداساز خطی و جستجوی روبه‌جلو و طبقه‌بند مبتنی بر ماشین بردار پشتیبان به بهترین نرخ تشخیص ۸۵٫۶٪ برای پایگاه‌داده برلین دست پیدا کردند. همچنین با استفاده از رگرسیون مبتنی بر ماشین‌های بردار پشتیبان متوسط ضریب همبستگی ۷۳٪ را برای پارامترهای تخمین زده شده در فضای پیوسته احساس به دست آوردند [۱۰]. در این سال لوکا و همکارانش سیستم خود را که به منظور بازشناسی احساس بر روی یک پایگاه‌داده طبیعی که با استفاده از مکالمات تلفنی ضبط شده طراحی نموده بودند آزمایش کرده و نتایج را با استفاده از ماتریس‌های تداخل گزارش نمودند [۱۱]. در سال ۲۰۱۲، طیف وسیعی از ویژگی‌های عروزی و طیفی به منظور جداسازی چهار احساس خشم، خوشحالی، ناراحتی و عادی بر روی پایگاه‌داده IEMOCAP^۱ مورد آزمایش قرار گرفتند. این تحقیقات با استفاده از الگوریتم‌های انتخاب ویژگی جستجوی روبه‌جلو و طبقه‌بندی مبتنی بر ماشین بردار پشتیبان، کارایی ویژگی‌های مختلف را در طبقه‌بندی هر احساس به خوبی با یکدیگر مقایسه نموده‌اند [۱۲].

در سال ۲۰۱۳، مروی و همکارش از ترکیب ویژگی‌های طیفی (ویژگی‌هایی که از طیف سیگنال دست می‌آیند مثل فرمنت‌ها، MFCC و PLP) و ویژگی‌های عروزی (ویژگی‌هایی که از آنالیز سیگنال در حوزه زمان دست می‌آیند و اغلب از منحنی فرکانس گام و انرژی سیگنال استخراج می‌شوند) استفاده کردند. آنها از الگوریتم دومرحله‌ای شامل معیار فیشر و الگوریتم جداساز خطی به منظور کاهش ویژگی‌ها و همچنین به منظور تشخیص احساس از گفتار استفاده نمودند و نشان دادند که ترکیب ویژگی‌های عروزی و طیفی باعث افزایش متوسط نرخ تشخیص می‌شود و نرخ تشخیص در پایگاه‌داده درام را برای گویندگان مرد به ۴۷٫۲۸٪ و نرخ تشخیص در گویندگان زن را به ۵۵٫۷۴٪ افزایش داده است. در همین سال، منصور شیخان و همکارانش از شبکه ماشین بردار پشتیبان برای تشخیص احساسات گفتار با استفاده از روش انتخاب ویژگی

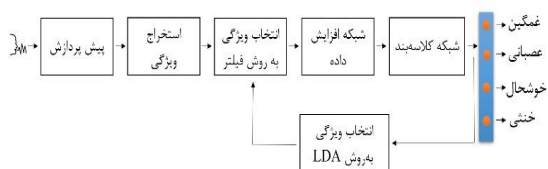
^۱ Interactive emotional dyadic motion capture database

عصبانیت، آرامش، ترس و استرس [۲۴]. در سال‌های اخیر، از شبکه‌های مولد متخاصمی برای تشخیص احساس از گفتار استفاده شده است. به‌عنوان مثال، چانگ و شفر از یک شبکه مولد متخاصمی عمیق^۱ برای یادگیری نمایش متمایز از گفتار احساسی به روشی نیمه نظارت استفاده کردند [۲۵].

۳- روش پیشنهادی

۳-۱- ساختار کلی پیشنهاد شده

شکل (۱)، ساختار کلی شبکه طبقه‌بندی کننده احساس پیشنهاد شده را برای طبقه‌بندی چهار احساس نمایش می‌دهد. از روش ترکیبی، جهت انتخاب ویژگی و از روش افزایش داده متخاصمی جهت افزایش داده‌ها و از ماشین بردار پشتیبان جهت طبقه‌بندی استفاده گردیده است.



شکل ۱- ساختار کلی شبکه تشخیص احساس پیشنهادی

تاکنون تحقیقات بسیاری در زمینه افزایش داده‌ها و یا انتخاب ویژگی‌های بهینه در سیستم‌های تشخیص احساس از گفتار انجام شده است ولی به‌ندرت به مسئله انتخاب ویژگی در نمونه‌های مصنوعی تولید شده و یا موضوع افزایش داده و انتخاب ویژگی به‌صورت هم‌زمان پرداخته شده است. همان‌طور که در شکل (۱) و ساختار کلی بررسی شده در این مقاله نمایش داده شده است طی دو مرحله ویژگی‌های مؤثر انتخاب و با فیدبکی که از شبکه طبقه‌بند گرفته می‌شود تکنیک مذکور بهینه خواهد شد.

۳-۲- پیش‌پردازش و استخراج ویژگی

در این مقاله، پس از انجام پیش‌پردازش و حذف نویز از سیگنال گفتار، بردارهای ویژگی توسط نرم‌افزار openSMILE استخراج گردید [۵]. از ویژگی‌های سطح پایین مثل انرژی فریم میانگین مربع^۲، ضرایب کپسترال مل فرکانسی^۳، نرخ عبور از صفر^۴، احتمال صدا^۵ و فرکانس‌های اساسی^۶ استفاده شد.

با پایگاه داده EMO-DB الگوریتم پیشنهاد شده استفاده شده است. همچنین پارامترهای طراحی طبقه‌بندی کننده و ویژگی‌های کاهش داده شده برای عملکرد بهینه فراهم شده است و در نهایت ویژگی‌های دیگری در سطح فریم و در سطح کل سخن اضافه شده تا بسته پیشنهادی تکمیل گردد. در مرجع [۱۶] ابعاد ویژگی را از ۲۷۶ به ۷۵ توسط روش SFFS کاهش دادند و تشخیص احساس به میزان ۲٫۷٪ بهتر شد. در یک مطالعه از روش‌های مختلفی برای کاهش ویژگی از ۵۸ به ۱۸، ۲۸، ۳۱ و ۳۳ استفاده شد و دقت تشخیص احساس به ترتیب به میزان ۴٫۵٪ و ۳٪ و ۲٫۲٪ و ۰٫۷٪ افزایش پیدا کرد.

در مرجع [۱۷] عملکرد تشخیص هنگامی که ابعاد ویژگی توسط روش mRMR از ۳۸۰ به ۱۲۱ کاهش پیدا کرد، به میزان ۱٫۵٪ کم شد. در مرجع [۱۸] هنگامی که ابعاد ویژگی از ۵۵ به ۴۹ و ۴۵ و ۲۴ و ۸ توسط روش FCBF کاهش پیدا کرد تغییر در میزان تشخیص به ترتیب به میزان ۰٫۹٪ و ۱٫۱٪- و ۲٫۳٪- و ۳٫۴٪- به دست آمد. در مرجع [۱۹] نرخ تشخیص با کاهش ویژگی از ۲۰۴ به ۸۷ توسط روش FCBF و طبقه‌بند ماشین بردار پشتیبان به میزان ۱٫۵٪ افزایش پیدا کرد. هو و همکاران [۲۰] از شبکه عصبی کانولوشن بسیار عمیق برای تولید ویژگی‌های اضافی برای آموزش مدل‌های صوتی استفاده کردند و دریافتند که افزایش داده به ساخت سیستم‌های تشخیص گفتار کمک بسیاری می‌کند. همچنین مقالاتی در خصوص بهبود نرخ تشخیص احساس با استفاده از تفکیک جنسیتی در سال ۲۰۱۷ [۲۱]، استخراج ویژگی‌های مقاوم گفتار در سال ۲۰۲۰ [۲۲] و تولید سیگنال‌های مصنوعی به کمک شبکه‌های عصبی مصنوعی در سال ۱۳۹۵ [۲۳] نیز منتشر گردیده‌اند. در سال ۲۰۲۱، چوراسیا و همکارانش یک سیستم تشخیص احساس از گفتار را بر اساس ویژگی‌های استخراج شده و به‌دست آمده از طیف‌نگارهای فرکانس مل (MFCC) پیشنهاد کردند. آن‌ها از شبکه‌های عصبی عمیق یک‌بعدی و پایگاه داده RAVDESS و از شش احساس بر اساس جنسیتشان (مرد و زن) استفاده نموده و توانستند به دقت ۸۲٫۳ درصد در طبقه‌بندی برسند. احساسات طبقه‌بندی شده توسط آن‌ها عبارت‌اند از: شادی، غم،

⁴ Zero-Crossing Rates

⁵ Voice Probabilities

⁶ Fundamental Frequencies

¹ Deep Convolutional GAN (DCGAN)

² Root-Mean-Square (RMS) Frame Energies

³ Mel-Frequency Cepstrum Coefficients (MFCCs)

(C=۴) استفاده گردیده است. در این مرحله، واریانس بین کلاسی و درون کلاسی مطابق رابطه فیشر محاسبه و مقدار به دست آمده میزان اهمیت ویژگی را مشخص می کند.

$$FDR(u) = \frac{2}{c(c-1)} \sum_{c1} \sum_{c2} \frac{(\mu_{c1,u} - \mu_{c2,u})^2}{\sigma_{c1,u}^2 + \sigma_{c2,u}^2} \quad (1)$$

$$1 \leq c_1 < c_2 \leq c$$

در این رابطه $\sigma_{ci,u}^2$ و $\mu_{ci,u}$ مقدار میانگین و واریانس طبقه i ام $(i=1,2,3,4)$ را برای ویژگی u نشان می دهند. امتیاز ویژگی i ام بوسیله رابطه ۱ محاسبه گردیده و تمامی ویژگی ها بر اساس امتیازشان از بیشتر به کمتر مرتب می شوند.

همان طور که می دانیم، ویژگی هایی برای طبقه بند احساسی مناسب هستند که واریانس درون کلاسی حداقل و واریانس بین کلاسی حداکثری داشته باشد [۲۵]. این مرحله از انتخاب ویژگی، بسیار سریع خواهد بود و در مدت زمان کوتاهی به ویژگی ها امتیاز داده می شود و امتیازهای پایین حذف می گردند؛ این آستانه از قبل مشخص شده است. در این مقاله این آستانه مقادیر ۷۰٪ کمتر از میانگین بالاترین ده امتیاز فیشر در نظر گرفته شد؛ لذا از ۲۱۸۵ ویژگی به ۱۵۸۲ ویژگی برای هر نمونه خواهیم رسید. البته می توان به جای انتخاب مقدار آستانه، تعداد ویژگی های مورد نظر را به عنوان مقدار آستانه تعیین نمود (مثلاً ۱۵۰۰ ویژگی با بالاترین امتیاز فیشر). اشکال اساسی این مرحله این است که هدف، انتخاب ویژگی هایی است که بتوانند عملکرد شبکه طبقه بند را بالاتر از قبل برده و راندمان آن را افزایش دهند ولی این خطر همیشه وجود دارد که ویژگی های انتخاب شده ترکیب خوبی را تشکیل ندهند. به بیان دیگر ممکن است ویژگی هایی انتخاب شوند که به تنهایی خوب عمل کنند ولی در کنار هم عالی عمل نکنند.

۳-۴- افزایش داده

شکل (۲)، شبکه مولد متخاصمی و اهداف شبکه های مولد و تشخیص دهنده را نشان می دهد. هدف از شبکه مولد، فریب شبکه تشخیص دهنده است؛ بنابراین شبکه عصبی مولد آموزش داده شده است تا خطای طبقه بندی نهایی (بین داده های واقعی و تولید شده) را به حداکثر برساند. در مقابل هدف شبکه تشخیص دهنده این است که داده های تولید شده جعلی را کشف کنند؛ بنابراین شبکه عصبی

برای استخراج این ویژگی ها، فریم های ۲۵ میلی ثانیه با هم پوشانی ۱۰ میلی ثانیه از سیگنال گفتار در نظر گرفته شدند. سپس ویژگی های سطح پایین توسط عملگرهای آماری مانند حداکثر، حداقل، دامنه، انحراف معیار و ... برای استخراج ویژگی های سطح بالا پردازش شدند. با استفاده از این عملگرهای آماری به ویژگی های سطح پایین، ویژگی های فریم، فریم شده به ویژگی های بیانی^۱ تعمیم داده شد و برای هر چهار مجموعه پایگاه داده، ویژگی های با واریانس صفر حذف شدند و بقیه ویژگی ها به طور مستقل توسط z-norm نرمال گردید. تعداد و نوع این ویژگی ها برای هر نمونه، ۲۱۸۵ ویژگی مطابق جدول ۱ می باشد.

جدول ۱- تعداد و نوع ویژگی های استخراج شده

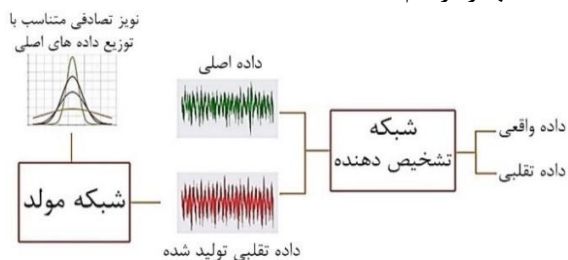
مشخصات ویژگی	نوع ویژگی
میانگین تعداد الگوها در هر قاب برای هر باند فرکانسی تعداد نسبی هر الگو در باندهای فرکانسی و تعداد نسبی هر الگو در باندهای فرکانسی	الگوهای طیفی (۲۰۴ ویژگی):
اعمال ۲۰ تابع آماری به E1 و E2 و مشتقات اول و دوم آن ها	انرژی هارمونیک ها (۷۸۰ ویژگی):
اعمال ۲۰ تابع آماری به منحنی فرکانس گام، منحنی انرژی، منحنی نرخ عبور از صفر، منحنی اپراتور انرژی تیگر و همچنین نسبت طول زمانی مصوت ها به صامت ها و کلیه مشتقات اول و دوم آن ها	ویژگی های عروزی (۲۴۱ ویژگی):
اعمال ۲۰ تابع آماری به ضریب اول MFCC، ۴ فرمنت اول و مشتقات اول و دوم آن ها	ویژگی های طیفی (۹۶۰ ویژگی):
۲۰ تابع آماری شامل: مقدار کمینه، بیشینه، برد، میانگین، میانه، ۱۰٪ و ۲۵٪، صدک های اول، پنجم، دهم، بیست و پنجم، هفتاد و پنجم، نودم، نود و پنجم و نود و نهم، برد میان چارکی، واریانس، انحراف معیار، چولگی و کشیدگی.	

۳-۳- انتخاب ویژگی به روش فیلتر

پس از استخراج بردارهای ویژگی احساسی و به منظور فیلتر نمودن آن ها از معیار فیشر جهت انتخاب ویژگی مقدماتی به روش فیلتر مطابق با رابطه ۱ و با چهار طبقه احساسی

¹ Utterance-Level Features

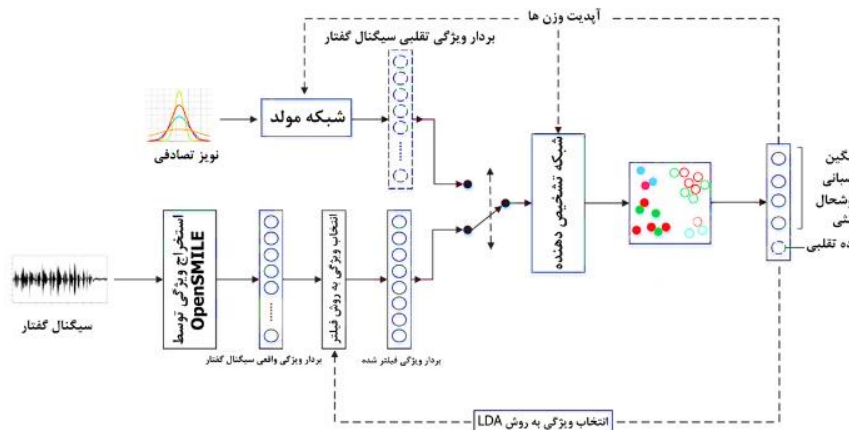
می‌شود که شبکه تشخیص‌دهنده قادر به تشخیص نمونه‌های تولید شده از نمونه‌های اصلی نبوده و با احتمال ۵۰٪ آنها را از هم جدا نماید.



شکل ۲- نمایش شبکه مولد متخاصمی

از آنجاکه مدل پیشنهادی در شکل (۳)، بردار ویژگی داده‌ها را به صورت یک‌به‌یک بین یک منبع نویز و یک نمونه هدف می‌آموزد، برای یک مجموعه داده گفتار دارای N برچسب طبقه‌های احساسی، نیاز به ایجاد انتقال برداری بین هر جفت از آنها است، یعنی $N(N-1)/2$ انتقال، که بسیار زیاد است.

تشخیص‌دهنده آموزش‌دیده تا خطای طبقه‌بندی نهایی را به حداقل برساند. مطابق شکل (۲)، در هر تکرار از روند آموزش، وزن‌های شبکه مولد به منظور افزایش خطای طبقه‌بندی (گرادیان صعودی خطا بر روی پارامترهای مولد) به‌روز می‌شود درحالی‌که وزن‌های شبکه تشخیص‌دهنده برای کاهش این خطا (گرادیان نزولی خطا بر پارامترهای متمایزکننده) به‌روز می‌شود. هر دو شبکه سعی می‌کنند شوند. رقابت بین آنها باعث می‌شود این دو شبکه باتوجه‌به اهداف مربوطه خود پیشرفت کنند. از نقطه نظر تئوری بازی، می‌توان این روش را به‌عنوان یک بازی دونفره در نظر گرفت که وضعیت تعادل در آن وضعیتی است که شبکه مولد، داده‌ها را با توزیع هدفمند دقیق تولید کند درحالی‌که شبکه تشخیص‌دهنده، پیش‌بینی "واقعی" یا "تولید شده" بودن آنها را با احتمال $1/2$ انجام دهد. در واقع بازی زمانی تمام



شکل ۳- ساختار شبکه افزایش داده متخاصمی و انتخاب ویژگی پیشنهادی

بر این، روش پیشنهاد شده در این مقاله، نمونه‌های مصنوعی را در فضای ویژگی ایجاد می‌کند. به این معنی که ورودی شبکه افزایش داده متخاصمی، یک سیگنال گفتار خام نیست بلکه بردارهای ویژگی مورد استفاده برای شبکه طبقه‌بند است که توسط نرم افزار OpenSMILE که یک نرم افزار منبع باز برای استخراج خودکار ویژگی‌ها از سیگنال‌های صوتی و طبقه‌بندی سیگنال‌های گفتاری و موسیقی است استخراج می‌شوند [۲۶]. دلیل این امر این است که هدف اصلی این مقاله انتخاب ویژگی‌های مهم به جای ترکیب گفتار است. باین‌حال، تولید نمونه در فضای ویژگی‌ها هم مزایا و هم مضراتی دارد. مزیت آن این است که می‌توان بدون در نظر گرفتن سنتز گفتار، بر تحقیقات شبیه سازی توزیع داده‌ها از طریق شبکه های مولد

در روش پیشنهاد شده در این مقاله، از داده‌های برچسب زده شده از هر نوع احساسات به‌عنوان دامنه هدف استفاده شده است درحالی‌که داده‌ها در دامنه منبع، بدون برچسب هستند. بنابراین، یک مجموعه داده هدف مصنوعی تولید می‌شود که به اندازه داده‌های منبع واقعی باشند و دارای همان احساسات باشند. به‌عنوان مجموعه داده هدف واقعی، داده‌های مصنوعی نیز پس از جداسازی توسط روش انتخاب ویژگی، برای آموزش شبکه طبقه‌بند مورد استفاده قرار می‌گیرند. درحالی‌که منبع مصنوعی در دامنه مجموعه داده بدون برچسب قرار دارد. به جای آموزش N شبکه مولد متخاصمی به طور جداگانه، شبکه های افزایش داده متخاصمی در یک چارچوب کامل قرار گرفتند تا نمونه‌های تولید شده هر احساس هدف با یکدیگر مرتبط شوند. علاوه

برای هر شبکه مولد با چگالی احتمال $p(g)$ ، بهترین شبکه تشخیص دهنده آن است که بتواند این تابع را مینیمم کند.

$$\mathbb{E}_{x \sim p_t}[1 - D(x)] + \mathbb{E}_{x \sim p_g}[D(x)] = \int_{\mathcal{R}} (1 - D(x))p_t(x) + D(x)p_g(x)dx \quad (4)$$

اگر به منظور مینیمم کردن نسبت به D در این انتگرال، تابع داخل انتگرال به ازای هر مقدار x مینیمم شود به این ترتیب بهترین شبکه تشخیص دهنده ممکن برای شبکه مولد مشخص خواهد شد.

$$(p_t(x) \geq p_g(x)) \quad (5)$$

در حقیقت، یکی از بهترین انتخاب‌ها این است که $p_t(x) = p_g(x)$ شود. سپس G مورد نیاز است که این تابع را ماکزیمم کند.

$$\int_{\mathcal{R}} (1 - D_G^*(x))p_t(x) + D_G^*(x)p_g(x)dx = \int_{\mathcal{R}} \min(p_t(x), p_g(x))dx \quad (6)$$

به طور مشابه، برای به حداکثر رساندن این انتگرال نسبت به g ، این تابع باید برای هر مقدار x ماکزیمم شود. از آنجاکه چگالی $p(t)$ مستقل از مولد G است، G نمی‌تواند به گونه‌ای بهتر از حالت زیر تنظیم شود:

$$p_g(x) \geq p_t(x) \quad (7)$$

البته، از آنجاکه $p(g)$ یک چگالی احتمال است که باید به ۱ برسد، لزوماً برای بهترین G داریم:

$$p_g(x) = p_t(x) \quad (8)$$

بنابراین، نشان داده شده است که در یک مورد ایدئال با ظرفیت‌های نامحدود شبکه‌های مولد و تشخیص دهنده، نقطه بهینه تنظیم شبکه تخصی به گونه‌ای است که شبکه مولد چگالی‌ای مانند همان چگالی واقعی را تولید کند و شبکه تشخیص دهنده مجبور است در یکی از دو حالت، صحیح بودن را انتخاب کند. در آخر نیز G این معادله را حداکثر می‌کند:

$$\frac{1}{2} \int_{\mathcal{R}} \min(p_t(x), p_g(x))dx = \int_{\mathcal{R}} \frac{\min(p_t(x), p_g(x)) p_t(x) + p_g(x)}{p_t(x) + p_g(x)} dx \quad (9)$$

به این شکل، انتظار می‌رود که G می‌خواهد احتمال شبکه تشخیص دهنده را در تشخیص اشتباه به حداکثر برساند. رایانه‌ها اساساً می‌توانند متغیرهای شبه تصادفی ساده تولید

متخاصمی متمرکز شد و ایراد آن در این است که احساسات نمونه‌های تولید شده فاقد یک شکل واضح برای ارزیابی ادراکی انسان است اما هنوز هم می‌توان ویژگی‌های عاطفی آنها را با مقایسه شباهت آنها با نمونه‌های داده واقعی آزمایش نمود.

الگوسازی شبکه‌های مولد متخاصمی در اصل نیاز به تعریف دو چیز دارد: ۱- معماری ۲- تابع خطا. شبکه‌های مولد متخاصمی همان‌طور که گفته شد از دو شبکه اصلی تشکیل شده‌اند:

یک شبکه مولد که از ورودی تصادفی z با چگالی $p(z)$ استفاده می‌کند و خروجی $G(z)$ که باید (پس از آموزش) دنبال کند، توزیع احتمال هدفمند را تولید می‌کند.

یک شبکه تشخیص دهنده که ورودی x را دریافت می‌کند و می‌تواند یک داده واقعی $X(t)$ باشد که چگالی آن با $p(x)$ نشان داده شده است و احتمال $D(x)$ را به عنوان واقعی بودن داده x گزارش می‌کند.

اگر داده‌های "واقعی" و "تولید شده" به یک نسبت به شبکه تشخیص دهنده ارسال شوند، خطای مطلق مورد انتظار از شبکه تشخیص دهنده به شرح زیر خواهد بود:

$$E(G, D) = \frac{1}{2} \mathbb{E}_{x \sim p_t}[1 - D(x)] + \frac{1}{2} \mathbb{E}_{z \sim p_z}[D(G(z))] = \frac{1}{2} (\mathbb{E}_{x \sim p_t}[1 - D(x)] + \mathbb{E}_{x \sim p_g}[D(x)]) \quad (2)$$

در این رابطه، $D(x)$ احتمال واقعی بودن نمونه x توسط شبکه تشخیص دهنده، $D(G(z))$ احتمال واقعی بودن نمونه x تولید شده توسط شبکه مولد و p_z توزیع نویز ورودی است. هدف شبکه مولد این است که شبکه تشخیص دهنده را که هدف آن، توانایی تفکیک داده‌های واقعی و تولید شده است، فریب دهد؛ بنابراین، هنگام آموزش شبکه مولد، هدف این است که این خطا حداکثر شود، در حالی که سعی می‌شود این خطا برای شبکه تشخیص دهنده به حداقل برسد و نمونه‌های تقلبی به خوبی از نمونه‌های اصلی تشخیص داده شوند تا شبکه مولد بتواند نمونه‌های بهتری تولید کند که توسط شبکه تشخیص دهنده قابل شناسایی نباشد و داریم:

$$\max_G (\min_D E(G, D)) = \mathbb{E}_{x \sim p_{data}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))] \quad (3)$$

در شبکه‌های مولد متخاصمی فرض بر آن است که یک مجموعه داده ورودی $(s_i, y_i)_{i=1}^N$ که نیمی از داده‌ها از نمونه‌های واقعی X و نیمی از نمونه‌های مصنوعی $G(Z)$ می‌باشند موجود است. هر نمونه آموزش s_i با یک برچسب y_i مطابقت دارد. تمام نمونه‌های واقعی به‌عنوان یک و همه نمونه‌های مصنوعی به‌عنوان صفر برچسب گذاری شده‌اند. از آنجاکه شبکه تشخیص دهنده را می‌توان یک طبقه‌بند دودویی در نظر گرفت، عملکرد این خطا را می‌توان به‌عنوان یک خطای متقاطع باینری تعریف کرد [۲۷].

$$J^{(D)}(D, G) = H((s_i, y_i)_{i=1}^N, D) = \quad (10)$$

$$-\frac{1}{N} \sum_{i=1}^N [y_i \log D(s_i) + (1 - y_i) \log(1 - D(s_i))]$$

اگر y_i با عدد یک برای $s_i = x$ و صفر برای $s_i = G(z)$ جایگزین شود و همچنین میانگین‌ها با انتظارات جایگزین گردند، تابع خطای شبکه تشخیص دهنده به شرح زیر خواهد بود [۲۷]:

$$J^{(D)}(D, G) = \quad (11)$$

$$-\mathbb{E}_{x \sim p_{data}(x)} [\log D(x)] - \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))]$$

که در آن " p " توزیع داده از نمونه واقعی x است. در یک بازی مینیماکس (همچنین به آن جمع صفر نیز گفته می‌شود)، کل خطای بازیکنان همیشه صفر است. به این معنی که خطای شبکه مولد برعکس $J(D)$ است. با این حال، هنگامی که گرادیان نزولی با توجه به G محاسبه می‌شود، تنها قسمت دوم در رابطه ۱۲ مهم است. بنابراین، تابع خطای شبکه مولد در یک بازی مینیماکس می‌تواند به شرح زیر است [۲۷]:

$$J^{(G)}(G) = \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))] \quad (12)$$

کل بازی را می‌توان با یک عبارت اعتبارسنجی خلاصه کرد. در واقع رابطه ۱۴، هدف اصلی در شبکه‌های مولد متخاصمی است و داریم:

$$\min_G \max_D V(D, G) = \quad (13)$$

$$\mathbb{E}_{x \sim p_{data}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))]$$

کنند. به‌عنوان مثال می‌توانند متغیرهایی تولید کنند که بسیار نزدیک به توزیع یکنواخت عمل می‌کنند. روش‌های مختلفی برای تولید متغیرهای تصادفی پیچیده‌تر وجود دارد از جمله مفهوم "روش تبدیل" که شامل بیان یک متغیر تصادفی به‌عنوان تابعی از متغیرهای تصادفی ساده‌تر است. در یادگیری ماشین، مدل‌های مولد سعی می‌کنند داده‌هایی از توزیع احتمالی خاص (پیچیده) تولید کنند. مدل‌های مولد شبکه‌های یادگیری عمیق به‌عنوان شبکه‌های عصبی (توابع بسیار پیچیده) مدل شده‌اند که به‌عنوان ورودی یک متغیر تصادفی ساده را گرفته و یک متغیر تصادفی را تولید می‌کنند که از توزیع احتمال هدفمند پیروی می‌کنند. این شبکه‌های مولد را می‌توان "مستقیم" آموزش داد (با مقایسه توزیع داده‌های تولید شده با توزیع داده‌های واقعی). این ایده شبکه‌های مولد تطبیق دهنده است. این شبکه‌های مولد همچنین می‌توانند به‌صورت "غیرمستقیم" با تلاش برای فریب شبکه دیگری که در همان زمان آموزش داده‌شده است آموزش داده شوند تا داده‌های تولید شده را از داده‌های واقعی متمایز کند. این ایده شبکه‌های مولد تخصصی است. این روش برای تغییر دادن تابع خطا از مقایسه مستقیم به غیرمستقیم واقعاً کاری است که می‌تواند برای دستاوردهای بعدی در حوزه یادگیری عمیق بسیار الهام‌بخش باشد.

الگوریتم ۱- نحوه آموزش شبکه مولد متخاصمی با روش نزول گرادیان. تعداد گام‌ها برای آموزش شبکه تشخیص دهنده برابر k و به‌عنوان پارامتر اولیه در نظر گرفته شده است.

برای تعداد تکرارهای آموزش تکرار شود:
 برای تعداد k تکرار شود:
 به تعداد m از فضای نویز اولیه $\mathbf{p}_g(\mathbf{z})$ نمونه‌برداری می‌شود.
 $\mathbf{z} = \{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$
 به تعداد m از توزیع اولیه داده‌ها \mathbf{p}_{data} نمونه‌برداری می‌شود.
 $\mathbf{x} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$
 با محاسبه گرادیان زیر خطای شبکه متمایزکننده محاسبه شود:

$$\nabla_{\theta_d} \frac{1}{m} \sum_i [\log D(\mathbf{x}^{(i)}) + \log(1 - D(G(\mathbf{z}^{(i)})))]$$

 پایان حلقه دوم
 به تعداد m از فضای نویز اولیه $\mathbf{p}_g(\mathbf{z})$ نمونه‌برداری می‌شود.
 $\mathbf{z} = \{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$
 وزن‌های شبکه مولد با روش کاهش گرادیان به‌صورت زیر به‌روزرسانی می‌شود:

$$\nabla_{\theta_g} \frac{1}{m} \sum_i [\log(1 - D(G(\mathbf{z}^{(i)})))]$$

 اتمام حلقه اول

می‌دهد بستگی دارد. هر چه احتمال واقعی بودن بالاتر باشد، خطای کمتری نیز در شبکه مولد دریافت می‌کند. مشاهده می‌شود که منحنی مینماکس در انتهای سمت چپ بسیار صاف است، به این معنی که شبکه مولد دارای شیب بسیار کمی است. با استفاده از گرادیان نزولی شیب شبکه مولد قبلاً متوقف شده است تا در مرحله اولیه عملکرد خود را بهبود بخشد و وزن‌های آن آپدیت گردند. با استفاده از محدودیت‌های تابع خطای در نظر گرفته شده، منحنی بازی اکتشافی غیر اشباع، شیب خود را در انتهای کار از دست می‌دهد. در این زمان، فرایند آموزش به حد مطلوب رسیده و نمونه‌های تولید شده قادر به فریب دادن شبکه تشخیص دهنده هستند. بنابراین، تابع خطای شبکه مولد در بازی غیر اشباع اغلب در عمل مورد استفاده قرار می‌گیرد درحالی‌که رابطه ۱۴ برای تحلیل نظری استفاده می‌گردد.

۳-۵- انتخاب ویژگی به روش ترکیبی

برای مواجهه با مشکل ابعاد بیش از حد داده‌ها، کاهش ابعاد با ترکیب ویژگی‌ها انجام شد. ترکیبات خطی به جهت محاسباتی و تحلیل ساده‌تر می‌باشند. در واقع روش خطی، داده‌های با ابعاد زیاد را بر روی فضای با ابعاد کمتر منعکس می‌کند. این روش بر اساس آنالیز تفکیک‌پذیری خطی می‌باشد و ماتریس انتقال بهینه را به‌عنوان ماتریسی که کوواریانس بین طبقات را حداکثر و کوواریانس درون طبقات را حداقل می‌کند، تعریف می‌نماید. در این روش تعداد پایه‌های خطی به‌وسیله تعداد طبقه‌ها محدود می‌شود (۴ طبقه) و پایه‌ها همیشه متعامد نیستند؛ لذا از چهار پایه خطی جهت پیاده‌سازی استفاده گردیده است. همچنین، ماتریس‌های کوواریانس درون کلاسی و برون کلاسی به فرم زیر محاسبه می‌شوند [۲۹].

$$S_w = \frac{1}{N} \sum_{i=1}^I \sum_{n_i=1}^{N_i} (X_{n_i} - \mu_i)(X_{n_i} - \mu_i)^T \quad (15)$$

$$S_b = \frac{1}{N} \sum_{l=1}^I N_l (\mu_l - \mu)(\mu_l - \mu)^T \quad (16)$$

که در آن N تعداد کل نمونه‌ها، N_i تعداد نمونه‌های طبقه i -ام، μ_i میانگین طبقه i -ام، I تعداد طبقه‌ها و X_{n_i} نمونه n_i -ام می‌باشند. S_w همبستگی درون کلاسی و S_b تمایز بین کلاسی را اندازه‌گیری می‌کنند. هدف به دست آوردن ماتریس تبدیل W است به نحوی که رابطه فیشر را ماکزیمم کند.

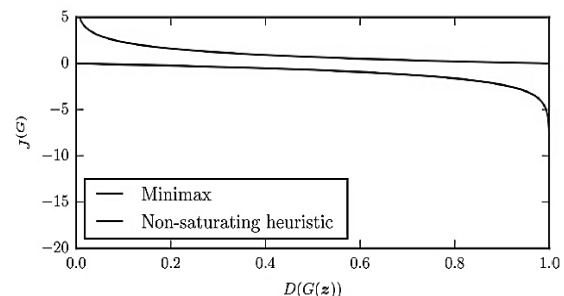
$$J(W) = \frac{W^T \text{trace}(S_b) W}{W^T \text{trace}(S_w) W} \quad (17)$$

شبکه مولد سعی می‌کند آن را به حداکثر برساند درحالی‌که شبکه تشخیص‌دهنده سعی در به حداقل رساندن آن دارد. این بدان معناست که به دنبال کمترین G هستیم که بیشترین D را نتیجه دهد.

اما بروز رسانی خطا در شبکه مولد در یک بازی مینماکس در عمل به‌خوبی انجام نمی‌شود و هنگامی که وزن‌های شبکه مولد هنوز به‌صورت مناسب انتخاب نشده‌اند، شبکه تشخیص‌دهنده به‌راحتی می‌تواند تشخیص دهد که نمونه‌های تولید شده جعلی هستند و با اطمینان بالا آنها را رد می‌کند [۲۷]. در نتیجه، شیب شبکه مولد جهت اصلاح وزن از بین می‌رود و گرادیان آن کوچک می‌شود. تابع خطا در شبکه مولد در یک بازی اکتشافی غیر اشباع به شرح زیر است [۲۷]:

$$J^{(G)}(G) = - \mathbb{E}_{z \sim p_z(z)} [\log(D(G(z)))] \quad (14)$$

شکل (۴)، تفاوت عملکرد خطای شبکه مولد در بازی اکتشافی مینماکس و غیر اشباع را نشان می‌دهد. محور افقی احتمال توصیف یک نمونه مصنوعی از نظر شبکه تشخیص‌دهنده را نشان می‌دهد. هرچه این مقدار بالاتر باشد، خطا در شبکه مولد کمتر است. به‌عبارت‌دیگر، هرچه شبکه تشخیص‌دهنده به عدد یک نزدیک می‌شود، نمونه تولید شده توسط شبکه مولد واقعی‌تر است. در قسمت سمت چپ نمودار، $D(G(z))$ نزدیک به صفر است و اغلب در شروع یک فرایند آموزش ظاهر می‌شود. در این زمان، شبکه تشخیص‌دهنده می‌تواند به‌راحتی تشخیص دهد که یک نمونه واقعی یا ساختگی است، زیرا شبکه مولد شروع به نمونه‌گیری از توزیع نویز تصادفی $p_z(z)$ با پارامترهای تصادفی می‌کند.



شکل ۴- مقایسه خطای شبکه مولد در یک بازی مینماکس و یک بازی اکتشافی غیر اشباع [۲۸]

شبکه مولد $J(G)$ برای تولید نمونه $G(z)$ به احتمال $D(G(z))$ که شبکه تشخیص‌دهنده به نمونه اختصاص

که علاوه بر صحت طبقه‌بندی برای مسائل طبقه‌بندی یک مدل شبکه عصبی مورد استفاده قرار می‌گیرند. سه سنجه اضافه دیگری که کمتر متداول اما محبوب هستند عبارت‌اند از: ضریب کاپای کوهن^۷، MCC^۸ و ماتریس درهم ریختگی^۹. از معیار دقت و صحت جهت معیار ارزیابی آزمایش‌ها استفاده گردیده است. شبیه‌سازی‌ها با استفاده از Python و Keras و Tensorflow2.2 در نرم‌افزار پایتون انجام گردیدند. (3.8 (64-bit)

۵- نتایج

به‌منظور مقایسه توزیع‌های پیچیده در فضای با ابعاد بالا، ابتدا اندازه‌گیری هم‌پوشانی مقادیر ویژگی‌های منحصر به فرد را معرفی کرده و داریم:

$$f = \frac{(\mu_1 - \mu_2)^2}{\sigma_1^2 + \sigma_2^2} \quad (19)$$

که در آن $\mu_1, \mu_2, \sigma_1, \sigma_2$ میانگین و انحراف معیار از دو مجموعه داده برای یک بعد از ویژگی خاص هستند. هرچه این مقدار بزرگتر باشد، منطقه همپوشانی کمتری بین دو مجموعه داده وجود دارد.

در این مقاله، تلاش شده است ویژگی‌هایی که بیشترین ارتباط با احساسات دارند انتخاب و سپس بر اساس آن ویژگی‌های انتخاب شده، شبکه طبقه‌بند جهت جدا نمودن چهار احساس آموزش داده شود. اگر اطلاعات اضافی که با عواطف مرتبط نیست فیلتر شوند، می‌توان عملکرد شبکه طبقه‌بند را برای تشخیص احساس از سیگنال گفتار بهبود بخشید و یا دست‌کم سرعت آن را بالا برد. داده‌های مصنوعی تولید شده توسط روش پیشنهادی احتمالاً در مناطقی رخ می‌دهند که مطمئناً می‌توانند به احساساتشان اختصاص یابند؛ بنابراین، تفاوت بین داده‌های مصنوعی و داده‌های هدف نشانگر تغییرات مورد نیاز برای تقویت بخشیدن عواطف در داده‌ها است.

بردار تفاوت بین مجموعه داده مصنوعی تولید شده و داده‌های موجود در پایگاه داده برلین محاسبه گردیدند، سپس بزرگ‌ترین ابعاد ویژگی n بردار تفاوت استخراج شد. این اختلاف به‌عنوان مقدار همپوشانی تعریف شده در رابطه ۱۹ محاسبه گردید. علاوه بر این، به ترتیب دو شبکه ماشین بردار پشتیبان بر اساس ویژگی‌های انتخاب شده ۱- فقط

شبکه طبقه‌بند از نوع ماشین بردار پشتیبان و از نوع چهار طبقه می‌باشد. همچنین از کرنل تابع پایه شعاعی لاپلاس^۱ استفاده گردیده است.

$$k(x, y) = \exp\left(-\frac{\|x-y\|}{\sigma}\right) \quad (18)$$

در این مرحله از انتخاب ویژگی، شبکه طبقه‌بند (ماشین بردار پشتیبان) در جریان انتخاب ویژگی است و عملکرد آن را به‌صورت مستقیم کنترل می‌نماید. ترکیب این دو مرحله انتخاب ویژگی و همچنین افزایش هوشمندانه داده‌ها به روش افزایش داده متخاصمی، سیستم تشخیص احساس از گفتار را به نقطه بهینه‌ای از تعداد داده و ابعاد داده‌ها رسانید.

استفاده از روش‌های انتخاب ویژگی همیشه موفقیت طبقه‌بندی را افزایش نمی‌دهد. نتایج در ادامه نشان می‌دهد که استفاده از هر دو روش جداساز خطی و آنالیز مؤلفه‌های اصلی به طور هم‌زمان، منجر به نتیجه بهتری نسبت به هر کدام از آنها می‌شود. زیرا آنالیز مؤلفه‌های اصلی روی داده‌های همبسته مؤثرتر عمل می‌کند و الگوریتم جداساز خطی روی داده‌های با بعد کم بهتر عمل می‌کند و روش فیشر در کاهش ساینز بهتر از آنالیز مؤلفه‌های اصلی عمل می‌کند.

۴- آزمایش‌ها

۴-۱- پایگاه داده^۲

آزمایش‌ها این مقاله بر روی چهار پایگاه داده متداول EMO-DB، eINTERFACE05، SAVEE و IEMOCAP انجام گردیده و علاوه بر این، تجزیه و تحلیل داده‌ها بر روی هر چهار پایگاه داده برای چهار احساس غمگین، عصبانی، خوشحال و خنثی ارائه شده است.

معروف‌ترین این پایگاه داده‌ها، پایگاه داده گفتار احساسی برلین^۳ [۳۰] است. پایگاه داده گفتار احساسی برلین، یک مجموعه داده کوچک است که شامل ۸۰۰ جمله است که به هفت طبقه احساس تقسیم شده‌اند. تمام گفته‌ها توسط ده بازیگر حرفه‌ای ضبط شده است.

۴-۲- معیار ارزیابی داده‌ها و محیط شبیه‌سازی

دقت^۴، صحت^۵ و امتیاز اف یک^۶ سه سنجه متداولی هستند

^۶ F1-Score

^۷ Cohen's Kappa

^۸ Matthews Correlation Coefficient

^۹ Confusion Matrix

^۱ Laplacian Radial Basis Function

^۲ Dataset

^۳ Berlin Database of Emotional Speech (EmoDB)

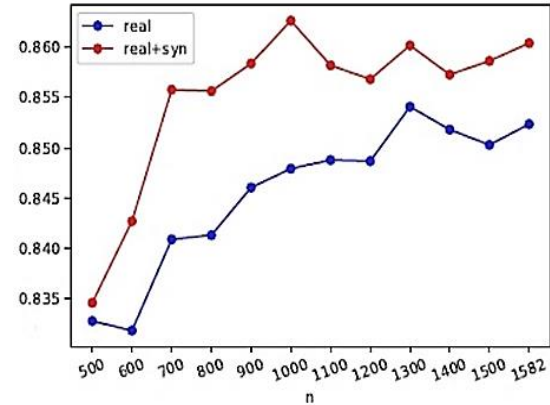
^۴ Precision

^۵ Recall

پایگاه داده و روش ارزیابی دارد) عمومیت مسئله ازدست رفته و الگوریتم به جای یادگیری قاعده طبقه بندی، به صورت خاص داده های آموزشی را فرامی گیرد که این مسئله کاهش دقت طبقه بندی را در مواجهه با داده های جدید (داده های تست) در پی دارد.

جدول ۲، نتایج طبقه بندی چهار احساس از سیگنال گفتار به وسیله ماشین بردار پشتیبان و برای چهار پایگاه داده نشان می دهد. هنگامی که فقط از انتخاب ویژگی به روش فیلتر (معیار فیشر) و فقط بر روی پایگاه داده اصلی کار کنیم، نتایج مطابق ستون سوم خواهد بود. تفاوت ستون سوم و چهارم افزودن داده توسط شبکه افزایش داده متخاصمی به پایگاه داده اصلی است. همان طور که مشاهده می شود در تمامی پایگاه داده های بکار گرفته شده معیار صحت افزایش یافته است چرا که داده های متنوع تری در اختیار طبقه بند جهت آموزش بوده است. ستون پنجم، شش و هفتم نتایج همین روند را به ترتیب به وسیله سه روش انتخاب ویژگی PCA، SFS [۳۱] و FCBF [۳۲] نشان می دهد. همان طور که مشاهده می شود در اکثر داده ها کاهش معیار صحت گزارش گردیده است. علت اصلی آن این است که این روش های انتخاب ویژگی نتوانسته اند تمامی ویژگی های اصلی را جهت آموزش ماشین بردار پشتیبان در خود جای دهند و با کاهش ویژگی ها سبب کاهش معیار صحت گردیده اند. در ستون آخر، نتایج روش پیشنهادی گزارش گردیده است. همان طور که ملاحظه می شود، در پایگاه داده برلین صحت ۸۶٫۳۲٪ گزارش گردیده است که در مقایسه با سایر روش های بکار گرفته شده در [۳۱] و [۳۲] بالاتر است. نمودار میله ای ۱ همچنین نمایشی از مقایسه هر شش روش و موفقیت هر کدام به صورت مقایسه با دیگری می باشد. بالاترین نرخ دقت مربوط به پایگاه داده برلین و روش پیشنهادی است. در سایر پایگاه داده های مورد آزمایش بیشترین نرخ دقت مربوط به پایگاه داده SAVEE و با استفاده از روش پیشنهادی است اگرچه که در دو پایگاه داده دیگر روش PCA به موفقیت بالاتری دست یافته است چرا که به نظر می رسد آنالیز مؤلفه های اصلی روی داده های همبسته مؤثرتر عمل می کند و الگوریتم جداساز خطی روی داده های با بعد کم بهتر عمل می کند. همچنان که روش فیشر در کاهش سایز بهتر از آنالیز مؤلفه های اصلی عمل می کند.

داده های واقعی و ۲- ترکیب داده های واقعی و مصنوعی آموزش داده شد. مقادیر مختلف n برای مشاهده رابطه بین عملکرد شبکه طبقه بند و تعداد ویژگی های انتخاب شده بررسی گردید.



شکل ۵- نتایج طبقه بندی بر اساس تعداد ویژگی ها

شکل (۵)، معیار صحت هر دو شبکه طبقه بند را نشان می دهد. همان طور که ملاحظه می شود، شبکه طبقه بند بر اساس ترکیبی از داده های واقعی و مصنوعی وقتی $n = 1000$ است به حداکثر خود یعنی ۸۶٫۳۲٪ می رسد و شبکه طبقه بند مبتنی بر داده های واقعی وقتی $n = 1300$ می باشد به حداکثر خود یعنی ۸۵٫۲۰٪ خواهد رسید. می توان استنباط کرد که آخرین ۵۸۲ ویژگی طبقه بند با داده های "واقعی + مصنوعی" و ۲۸۲ ویژگی آخر برای طبقه بندی کننده با داده های "واقعی" حاوی اطلاعات بسیار کمی در مورد احساسات هستند. بنابراین، داده های مصنوعی نه تنها می توانند جهت افزایش داده ها بلکه برای انتخاب ویژگیها به منظور بهبود عملکرد شبکه طبقه بند مورد استفاده قرار گیرند و دقت و سرعت این شبکه را در طبقه بندی افزایش دهد. همچنین، پدیده اوج که پیشتر توضیح داده شد به وضوح در شکل (۵) قابل ملاحظه است و همان طور که مشاهده می شود با افزایش تعداد ویژگی های مورد استفاده، ابتدا نرخ تشخیص افزایش پیدا می کند اما پس از رسیدن به نقطه اوج، افزایش بیشتر تعداد ویژگی ها موجب کاهش نرخ تشخیص می شود. علت اصلی این پدیده افت کارایی الگوریتم طبقه بند به دلیل پیچیدگی و ابعاد زیاد فضای ویژگی در مقابل تعداد داده های آموزشی می باشد. به عبارت دیگر، در چنین شرایطی به دلیل تعدد پارامترهای تنظیم طبقه بند (به دلیل ابعاد بالای فضا) و محدود بودن تعداد داده های آموزشی (که بستگی به

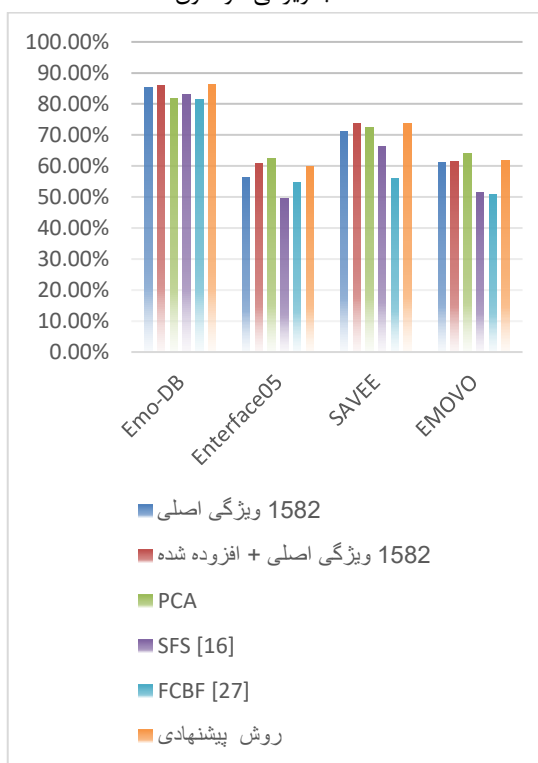
جدول ۲- مقایسه نرخ موفقیت طبقه‌بندی در پایگاه‌های داده

متفاوت و به‌وسیله روش‌های انتخاب ویژگی گوناگون

روش پیشنهادی	FCBF[۳۲]	SFS[۳۱]	PCA	۱۵۸۲ ویژگی اصلی + افزوده	۱۵۸۲ ویژگی اصلی	طبقه‌بند	پایگاه‌داده
۸۶,۳۲	۸۱,۳۲	۸۲,۹۳	۸۱,۷۱	۸۶,۰۰	۸۵,۲۰	SVM	EMO-DB
۵۹,۷۷	۵۴,۶۲	۴۹,۴۹	۶۲,۳۱	۶۰,۶۵	۵۶,۳۳	SVM	Enterface05
۷۳,۸۲	۵۵,۸۳	۶۶,۲۶	۷۲,۳۹	۷۳,۶۵	۷۰,۹۸	SVM	SAVEE
۶۱,۶۷	۵۰,۸۹	۵۱,۴۸	۶۳,۹۱	۶۱,۳۸	۶۱,۱۰	SVM	IEMOCAP

نمودار ۱- مقایسه نرخ موفقیت طبقه‌بندی به‌وسیله روش‌های

انتخاب ویژگی گوناگون



برای ارزیابی عملکرد روش پیشنهاد شده، نتایج به‌دست‌آمده با یک مطالعه اخیر [۳۳] که در آن نویسندگان

از الگوریتم‌های مختلف کاهش ابعاد غیرخطی برای استخراج نمایش ویژگی‌های سطح پایین برای تشخیص احساس از سیگنال گفتار استفاده کرده‌اند مقایسه شده است. سه الگوریتم کاهش ابعاد با عملکرد بالا در [۳۳] انتخاب گردیده که این روش‌ها شامل مقیاس‌بندی چندبعدی^۱، تجزیه و تحلیل مؤلفه‌های اصلی^۲ و رمزگذار خودکار است. نتایج در جدول ۳ ارائه شده است.

جدول ۳- مقایسه نتایج با استفاده از الگوریتم‌های مختلف

کاهش ابعاد ویژگی در پایگاه‌داده IEMOCAP

روش	تکنیک افزایش داده	UAR (%)
SMACOF MDS [33]	ندارد	۵۸,۸
PCA [33]	ندارد	۵۷,۷
AUTOENCODER [33]	ندارد	۵۷,۸
SMACOF MDS [33]	دارد	۵۸,۹
PCA [33]	دارد	۵۸,۳
AUTOENCODER [33]	دارد	۵۸,۵
روش پیشنهادی	دارد	۶۱,۶۷

در [۳۳]، نویسندگان از ماشین بردار پشتیبان برای طبقه‌بندی ویژگی‌هایی که توسط هر الگوریتم کاهش ابعاد آموخته شده است، استفاده کردند. از جدول ۳ می‌توان دریافت که مدل پیشنهادی عملکرد بهتری نسبت به سایر تکنیک‌های کاهش ابعاد غیرخطی دارد و نشان می‌دهد که مدل پیشنهادی به طور کارآمد ویژگی‌ها را در ابعاد کمتر فیلتر می‌کند درحالی‌که اطلاعات احساسی را جهت طبقه‌بندی حفظ می‌نماید.

برای بررسی بهتر مدل پیشنهادی در یک محیط متقاطع، آزمایش‌هایی با استفاده از داده‌های واقعی، مصنوعی، واقعی + مصنوعی، واقعی + مصنوعی (کاهش یافته) انجام شد و نتایج آنها با دو منبع [۳۴] و [۲۸] مقایسه گردید. همان‌طور که در جدول ۴ گزارش شده است، روش پیشنهادی به نتایج بهتری در مقایسه با دو منبع ذکر شده دست‌یافت. همچنین، انتخاب ویژگی به روش پیشنهادی توانست نتایج را اندکی بهبود دهد. این نشان می‌دهد که مدل پیشنهادی عملکرد سیستم تشخیص احساس از گفتار را در یک محیط متقابل با استفاده از ترکیب داده‌های واقعی و مصنوعی و همچنین هنگامی که داده‌های آموزشی با این نمونه‌های مصنوعی افزایش می‌یابد و سپس فیلتر می‌گردند،

^۲ PCA

^۱ SMACOF (MDS)

آموزش می‌بیند، معیار صحت وقتی $n=1000$ است به حداکثر خود یعنی ۸۶٫۳۲٪ رسیده و شبکه طبقه‌بند آموزش داده شده مبتنی بر داده‌های واقعی وقتی $n=1300$ می‌باشد به حداکثر خود یعنی ۸۵٫۲۰٪ خواهد رسید. می‌توان استنباط کرد که آخرین ۵۸۲ ویژگی طبقه‌بند با داده‌های "واقعی + مصنوعی" و ۲۸۲ ویژگی آخر برای طبقه‌بندی‌کننده با داده‌های "واقعی" حاوی اطلاعات بسیار کمی در مورد احساسات هستند و می‌توان آن‌ها را حذف نمود. از ترکیب دو روش معیار فیشر و الگوریتم جداساز خطی جهت کاهش این ویژگی‌ها استفاده و معیار صحت ۸۶٫۳۲٪ در پایگاه داده احساسی برلین گزارش شد. بنابراین، داده‌های مصنوعی نه تنها می‌توانند جهت افزایش داده‌ها بلکه برای انتخاب ویژگی‌ها به منظور بهبود عملکرد شبکه طبقه‌بند مورد استفاده قرار گیرند و دقت و سرعت این شبکه را در طبقه‌بندی افزایش دهد. با توجه به تحقیقات انجام شده می‌توان گفت که انتخاب ویژگی عملکرد تشخیص احساس را بهتر می‌کند و به طبقه‌بند، داده، نسبت کاهش سائز و حتی به نحوه جمع‌آوری داده‌ها بستگی دارد.

بهبود می‌بخشد و استفاده از روش انتخاب ویژگی آن را بهینه می‌نماید.

جدول ۴- نتایج ارزیابی بین کلاسی

روش	داده‌های واقعی	داده‌های مصنوعی	داده‌های واقعی + مصنوعی (کاهش یافته)	داده‌های واقعی + مصنوعی (کاهش یافته)
Sahu et al. [34]	۴۵٫۱۴	۳۳٫۹۶	۴۵٫۴۰	-
Bao et al. [28]	۴۵٫۵۸	۴۱٫۵۸	۴۶٫۵۲	-
روش پیشنهادی	۴۶٫۱	۴۲٫۲۱	۴۶٫۶۵	۴۶٫۹۵

۶- نتیجه‌گیری

انتخاب ویژگی‌های مؤثر و افزایش داده‌ها از طریق شبکه‌های مولد متخاصمی در پایگاه داده‌های کوچک یک موضوع مهم در تشخیص احساس از سیگنال گفتار است. مدل پیشنهادی در این مقاله از ترکیب دو شبکه انتخاب ویژگی، یک شبکه افزایش داده متخاصمی و یک ماشین بردار پشتیبان تشکیل شد. نشان داده شد هنگامی که شبکه طبقه‌بند بر اساس ترکیبی از داده‌های واقعی و مصنوعی

مراجع

- [1] Eva Lieskovaš, Maroš Jakubec, Roman Jarina and Michal Chmulík, "A Review on Speech Emotion Recognition Using Deep Learning and Attention Mechanism," in Electronics, May. 2021.
- [2] Akçay, M.B.; Oğuz, K. Speech emotion recognition: Emotional models, databases, features, preprocessing methods, supporting modalities, and classifiers. Speech Communication, 2020, pp.116, 56–76.
- [3] C. Sima and E. R. Dougherty, "The peaking phenomenon in the presence of feature-selection," Pattern Recognition Letter, vol. 29, no. 11, 2008, pp. 1667–1674.
- [4] J. Rong, G. Li, and Y.-P. P. Chen, "Acoustic feature selection for automatic emotion recognition from speech," Information Processing & Management, vol. 45, no. 3, 2009, pp. 315–328.
- [5] F. Eyben, F. Weninger, F. Gross, and B. Schuller, "Recent developments in openSMILE, the Munich open-source multimedia feature extractor," in Proc. 21st ACM International Conference on Multimedia, 2013, pp. 835–838.
- [6] S. N. Negahban, P. Ravikumar, M. J. Wainwright, and B. Yu, "A unified framework for high-dimensional analysis of M-estimators with decomposable regularizers," statistical science, vol. 27, no. 4, Nov 2012, pp. 538–557.
- [7] H. Zou, T. Hastie, and R. Tibshirani, "Sparse principal component analysis," Journal of Computational and Graphical Statistics, vol. 15, no. 2, 2006, pp. 265–286.
- [8] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, Y. Bengio. "Generative adversarial nets". In: Advances in neural information processing systems, 2014.
- [9] E. Bozkurt, E. Erzin, Ç. E. Erdem, and A. T. Erdem, "Formant position based weighted spectral features for emotion recognition," Speech Communication., vol. 53, no. 9, 2011, pp. 1186–1197.
- [10] S. Wu, T. H. Falk, and W.-Y. Chan, "Automatic speech emotion recognition using modulation spectral features," Speech Communication, vol. 53, no. 5, 2011, pp. 768–785.

- [11] P. Laukka, D. Neiberg, M. Forsell, I. Karlsson, and K. Elenius, "Expression of effect in spontaneous speech: Acoustic correlates and automatic detection of irritation and resignation," *Computer Speech & Language*, vol. 25, no. 1, 2011, pp. 84–104.
- [12] H. Pérez-Espinosa, C. A. Reyes-García, and L. Villaseñor-Pineda, "Acoustic feature selection and classification of emotions in speech using a 3D continuous emotion model," *Biomed. Signal Process. Control*, vol. 7, no. 1, 2012, pp. 79–87.
- [۱۳] علی حریمی، علیرضا احمدی فرد، علی شهزادی و خشایار یغمایی، "تشخیص احساس از روی گفتار با استفاده از طبقه‌بند مبتنی بر مدل و ویژگی‌های دینامیکی غیر خطی"، نشریه مهندسی برق و مهندسی کامپیوتر ایران، دوره ۱۵، شماره ۲، تابستان ۱۳۹۶، صفحه ۱۵۲-۱۴۵.
- [14] K. Han, D. Yu, and I. Tashev, *Speech Emotion Recognition Using Deep Neural Network and Extreme Learning Machine.*, vol. 3, 2014, pp. 232–243.
- [15] H. Palo and M. Mohanty, "Modified-VQ Features for Speech Emotion Recognition," *Journal of Mathematical Sciences*, vol. 16, Sep 2016, pp. 406–418.
- [16] B. Schuller, R. Müller, M. Lang, and G. Rigoll, *Speaker independent emotion recognition by early fusion of acoustic and linguistic features within ensembles.*, vol. 2, 2005, pp. 565–572.
- [17] I. Luengo, E. Navas, and I. Hernáez, "Feature Analysis and Evaluation for Automatic Emotion Identification in Speech," *Multimedia, IEEE transaction*, vol. 12, Nov 2010, pp. 490–501.
- [18] D. Gharavian, M. Sheikhan, and F. Ashohtedel, "Emotion recognition improvement using normalized formant supplementary features by a hybrid of DTW-MLP-GMM model," *Neural Computing & Applications*, vol. 22, no. 6, 2013, pp. 1181–1191.
- [19] X. Zhao, S. Zhang, and B. Lei, "Robust emotion recognition in a noisy speech via sparse representation," *Neural Computing & Applications*, vol. 24, Jun. 2013.
- [20] H. Hu, T. Tan, and Y. Qian, "Generative adversarial network-based data augmentation for noise-robust speech recognition," in *Proc. The international Conference on Acoustics, Speech, & Signal Processing (ICASSP)*, Apr 2018, pp. 5044–5048.
- [21] A. Harimi and Kh. Yaghmaie, "improving speech emotion recognition via gender classification," in *Journal of Modeling in Engineering.*, vol.48, 2017, pp. 184–200.
- [22] M. Sadeghi, H. Marvi and A. Ahmadifard "A New and Efficient Feature Extraction Method for Robust Speech Recognition Based on Fractional Fourier Transform and Differential Evolution Optimizer," in *Journal of Modeling in Engineering.*, vol.61, 2020, pp. 86–96.
- [۲۳] سیدعلی سلیمانی ابوری، محمدرضا فدوی امیری و حسین مروی، "تولید سیگنال مصنوعی زلزله به کمک مدلی جدید در فشرده‌سازی و آموزش شبکه‌های عصبی مصنوعی"، نشریه مدل‌سازی در مهندسی، دوره ۱۴، شماره ۴۶، پائیز ۱۳۹۵، صفحه ۷۵-۸۵.
- [24] M. Chourasia, S. Haral, S. Bhatkar, and S. Kulkarni, "Emotion Recognition from Speech Signal Using Deep Learning." In *Intelligent Data Communication Technologies and Internet of Things*, 2021, pp. 471–481.
- [25] J. Chang, S. Scherer. "Learning representations of emotional speech with deep convolutional generative adversarial networks". In: *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2017.
- [26] F. Eyben, F. Wenginger, F. Gross and B. Schuller. "Recent developments in openSMILE, the Munich open-source multimedia feature extractor". In: *Proc. 21st ACM international conference on Multimedia*. ACM, vol. 5, 2013, pp. 232–240.
- [27] I. Goodfellow. "NIPS 2016 tutorial: Generative adversarial networks". In: *arXiv preprint arXiv:1701.00160*, 2016.
- [28] F. Bao, M. Neumann, and N. T. Vu, "CycleGAN-based emotion style transfer as data augmentation for speech emotion recognition." In *Interspeech*, 2019, pp. 2828–2832.
- [29] W. Y. Zhao, "Discriminant component analysis for face recognition," in *Proceeding's 15th International Conference on Pattern Recognition. ICPR-2000*, vol. 2, 2000, pp. 818–821.
- [30] F. Burkhardt, A. Paeschke, M. Rolfes, W. F. Sendlmeier, and B. Weiss, "A database of German emotional speech," in *Proc. 9th European Conference on Speech Communication and Technology*, 2005, pp 1–4.

- [31] B. Yang and M. Lugger, "Emotion recognition from speech signals using new harmony features," *Signal Processing*, vol. 90, no. 5, 2010, pp. 1415–1423.
- [32] I. Luengo, E. Navas, and I. Hernandez, "Feature Analysis and Evaluation for Automatic Emotion Identification in Speech," *Multimedia, IEEE transaction*, vol. 12, Nov 2010, pp. 490–501.
- [33] G. Paraskevopoulos, E. Tzinis, N. Ellinas, T. Giannakopoulos, and A. Potamianos, "Unsupervised low-rank representations for speech emotion recognition," *Proc. Interspeech 2019*, 2019, pp. 939–943.
- [34] S. Sahu, R. Gupta, and C. Espy-Wilson, "On enhancing speech emotion recognition using generative adversarial networks," *Proc. Interspeech 2018*, 2018, pp. 3693–3697.
- [35] S. Latif , M. Asim, R. Rana, S. Khalifa, R. Jurdak and B.W. Schuller, "Augmenting Generative Adversarial Networks for Speech Emotion Recognition, " *Proc. Interspeech*, 521-525, doi: 10.21437/Interspeech.2020-3194, 2020.