

تشخیص لهجه‌های گفتار زبان فارسی از روی سیگنال گفتار با استفاده از روش‌های استخراج ویژگی کارآمد و ترکیب طبقه‌بندها

مجتبی شریف نوqابی^{*}، حسین مروی و دانیال دارابیان

گروه الکترونیک، دانشکده مهندسی برق و رباتیک، دانشگاه صنعتی شاهرود، شاهرود، ایران

چکیده

تشخیص لهجه از روی شکل موج گفتار یکی از شاخه‌های تاحدودی جدید در علم پردازش گفتار است. تشخیص لهجه می‌تواند تا حد زیادی باعث بهبود سامانه‌های بازشناسی گفتار شود. همانند هر سامانه بازشناسی، فرآیند تشخیص لهجه نیز شامل سه مرحله پیش‌پردازش، استخراج ویژگی و طبقه‌بندی است. در این مقاله سه روش کارآمد استخراج ویژگی شامل دامنه مرکزی طیفی (SCM)، مشتق نخست آن (ΔSCM) و تبدیل Zatk روش سیگنال گفتار اعمال شده‌اند و کارآبی این روش‌ها با روش‌های متداول استخراج ویژگی مانند ضرایب مل-کپستروم و مشتقات آن مقایسه شده است. علاوه بر این یک ویژگی جدید که با ایجاد تغییراتی در الگوریتم محاسبه ضرایب مل-کپستروم بدست می‌آید برای تشخیص لهجه‌ها در محیط‌های نویمای معروفی شده است. برای مرحله طبقه‌بندی از پنج طبقه‌بند مختلف، شامل MLP، KNN، SVM و RBF، PNN و ترکیب این طبقه‌بندها با یکدیگر استفاده شده است. نتایج آزمایش‌ها بیان گر بهبود نرخ بازشناسی لهجه‌ها با روش‌های پیشنهادی هستند.

وازگان کلیدی: دامنه مرکزی طیفی، ترکیب طبقه‌بندها، لهجه فارسی، ماشین بردار پشتیبان، ضرایب مل-کپستروم بهبود یافته.

۱- مقدمه

لهجه‌ای غیر از لهجه‌ای که با آن آموخته شده بخورد کند، شاهد کاهش چشم‌گیر بازدهی آن سامانه خواهیم بود. بنابراین اهمیت این مسئله ما را به پیاده‌سازی سامانه‌های برای تشخیص لهجه‌ها از روی شکل موج گفتار ترغیب می‌کند. به طور معمول برای تشخیص جنسیت، از مدل‌های مربوط به خود آن استفاده می‌شود. برای تشخیص لهجه نیز مدل‌های مختلفی ارائه شده است که روند کلی این مدل‌ها شامل سه مرحله استخراج ویژگی، طبقه‌بندی و تشخیص می‌باشد.

اگر نگاهی گذرا به تاریخچه و سابقه تشخیص لهجه‌ها از طریق گفتار داشته باشیم، با دسته‌بندی مشاهدات در دو حوزه لهجه‌های خارجی و فارسی به این نتیجه می‌رسیم که در حوزه لهجه‌های خارجی پژوهش‌های بیشتری انجام شده است که در این میان بر روی تشخیص لهجه‌های انگلیسی

در زبان‌شناسی به نحوه تلفظ کلمات در یک گروه خاص از افراد یک مکان یا ملت لهجه می‌گویند. مواردی که در لهجه یک گوینده مؤثر است برگرفته از سابقه اجتماعی- فرهنگی یکسان مردمی است که در یک موقعیت جغرافی قرار دارند و ممکن است الگوهای گوناگونی را به صورت مشترک در گفتارشان داشته باشند.

فرآیند بازشناسی خودکار گفتار یکی از زمینه‌های دانش هوش مصنوعی است. پژوهش‌ها نشان می‌دهند که تغییرات گفتار که مربوط به گوینده است، شامل تغییر در جنسیت، سن، لهجه، سلامتی و یا عدم سلامتی گوینده از دلایل تنزل کارآبی سامانه‌های بازشناسی گفتار هستند. که از میان این ویژگی‌ها دو ویژگی جنسیت و لهجه، بیشترین تأثیر را در این کاهش بازدهی دارند. هنگامی که یک سامانه تشخیص گفتار با یک لهجه خاص آموخت بییند، آن گاه با

۲- سامانه تشخیص لهجه

همان طور که در قبل بیان شد هر سامانه بازشناسی لهجه همانند هر سامانه خودکار بازشناسی گفتار شامل سه مرحله پیش پردازش، استخراج ویژگی و طبقه بندی است. در مرحله پیش پردازش، استخراج ویژگی و طبقه بندی این گفتار شامل حذف قسمت سکوت از گفتار انجام می شود. این کار با روش های مانند نرخ عبور از صفر و تعیین آستانه انرژی انجام می شود. درنهایت با اعمال یک فیلتر پیش تأکید این مرحله به پایان می رسد.

۳- روش های متداول استخراج ویژگی

مهم ترین بخش یک سامانه تشخیص لهجه استخراج ویژگی است. در مرحله استخراج ویژگی باید به دو نکته توجه بیشتری داشت. نخست این که باید قسمت هایی از سیگنال گفتار انتخاب شود که شامل بیشترین اطلاعات درباره لهجه آن باشد و نکته دوم این است که ویژگی استخراج شده باید قادر باشد اطلاعات مناسبی در مورد لهجه گفتار به ما ارائه کند و ما با استفاده از این ویژگی در مرحله طبقه بندی شاهد جداسازی بین لهجه های مختلف به بهترین نحو باشیم.

۳-۱- ضرایب مل-کپستروم (MFCC)

یکی از ویژگی های متداول که از سیگنال گفتار استخراج و در اغلب سیستم های تشخیص گفتار و گوینده از آن استفاده می شود، ضرایب مل-کپستروم ('MFCC') است. مراحل محاسبه این ضرایب در شکل (۱) نشان داده شده است. پس از مرحله پیش پردازش به دلیل غیر ایستان بودن سیگنال گفتار، آن را قاب بندی می کنیم تا ایستان فرض شود. طول قابها در این مرحله به طور تقریبی ۲۵ میلی ثانیه و بدون هم پوشانی در نظر گرفته شده است. دلیل عدم قراردادن هم پوشانی بین قابها، نخست این که نتایج آزمایش هاست که بیان گر نتیجه بهتر در حالت بدون هم پوشانی است و دوم این که میانگین گیری است که از هر چند فریم متوالی به عمل می آید. با توجه به این که ممکن است هر قاب به تنهایی شامل یک آواز کامل نباشد تا در تشخیص لهجه مفید باشد از بردارهای ویژگی مربوط به هر هفت قاب متوالی میانگین گرفته و آن را به عنوان یک بردار ویژگی شامل ضرایب مل-کپستروم و سایر ویژگی ها در نظر

تمرکز بیشتری وجود داشته است. به عنوان نمونه در سال ۱۹۹۶/رسلان و هانسن با ارائه یک سامانه طبقه بندی لهجه برای زبان انگلیسی یکی از برجسته ترین و نخستین پژوهش های انجام شده در این زمینه را انجام دادند. (رسلان، ۱۹۹۶). فاریا در سال ۲۰۰۵ تلاش دیگری برای طبقه بندی لهجه ها با استفاده از روش های SVM و GMM انجام داد. در این پژوهش علاوه بر ویژگی های صوتی گفتار از ویژگی های کلامی آن نیز استفاده شده است. (فاریا، ۲۰۰۵). پدرسن و دیدریچ در سال ۲۰۰۷ طبقه بندی لهجه را با کمک SVM و یادگیری درخت تصمیم برای دو لهجه عربی و هندی از زبان انگلیسی انجام داد. (پدرسن، ۲۰۰۷)، علاوه بر این پژوهش های دیگری نیز در حوزه تشخیص لهجه صورت گرفته است. (چن، ۲۰۰۱)، (کامف، ۱۹۹۶) و (زنگ، ۲۰۰۵). در حوزه لهجه های فارسی، نتایج مشاهدات ما نشان می دهد که تاکنون تعداد سه پژوهش در این زمینه انجام شده است، در نخستین پژوهش پنج لهجه اصفهانی، تهرانی، آذری، کردی و مازندرانی انتخاب شده و با استخراج ویژگی های ضرایب مل-کپستروم و مشتقات اول و دوم آن و همچنین انرژی هر قاب در مرحله طبقه بندی از شبکه عصبی احتمالاتی و ماشین بردار پشتیان استفاده شده است. (قلی پور، ۱۳۹۱) در دومین پژوهش صورت گرفته سه لهجه تهرانی، اصفهانی و کرمانشاهی برای انجام آزمایش های مختلف انتخاب شده اند. همانند بیشتر مقالات، این مقاله نیز در مرحله استخراج ویژگی از ویژگی های متداولی همچون ضرایب مل-کپستروم، فرکانس های فرمنت و انرژی استفاده کرده است. در مرحله طبقه بندی سه طبقه بند مختلف شامل KNN، SVM و MLP به کار گرفته شده است (ربیعی، ۲۰۱۰). در آخرین پژوهش بررسی شده پنج لهجه اصفهانی، تهرانی، آذری، شمالی و جنوبی برای انجام آزمایش ها برگزیده شده است. در این مقاله ویژگی های ضرایب مل-کپستروم و مشتق نخست و دوم آن، انرژی هر قاب، سه فرکانس فرمنت نخست و SDC از سیگنال گفتار دارای لهجه استخراج شده است. هدف اصلی این مقاله ترکیب دو طبقه بند GMM و PRLM است که با انجام این کار نرخ تشخیص لهجه های مختلف نسبت به حالتی که این طبقه بند ها به طور تکی استفاده می شوند، به اندازه مطلوبی افزایش می یابد. (جلالوند، ۲۰۱۲)

فصل نهم



۳-۳- فرکانس‌های فرمنت

یکی دیگر از ویژگی‌های حوزه فرکانس که از سیگنال گفتار استخراج می‌شود، فرکانس‌های فرمنت هست. فرمنت‌ها پارامترهای مربوط به لوله صوتی هستند که برای هر آوازی مقدار فرکانس فرمنت متفاوتی در مقایسه با آواهای دیگر وجود دارد که مشخصه همان آوا به شمار می‌رود. به طور معمول از سه یا دو فرکانس فرمنت نخست (2F یا 3F) به عنوان ویژگی در سامانه‌های تشخیص گفتار و لهجه استفاده می‌شوند. برای محاسبه این ویژگی روش‌های مختلفی وجود دارد که در اینجا با استفاده از رابطه بازگشتی ضرایب پیش‌بینی خطی (LPC^۱) فرکانس‌های فرمنت محاسبه شده‌اند.

۳-۴- پیش‌گویی خطی مبتنی بر درک انسان (PLP^۲)

یکی از معایب آنالیز LPC آن است که فرآیند و نحوه عملکرد دستگاه شنوازی انسان را در محاسبه ویژگی‌ها منظور نمی‌کند. به عبارت دیگر LPC در تمامی فرکانس‌ها سیگنال گفتار را به یک صورت تخمین می‌زند که این مطابق با دستگاه شنوازی انسان نیست. برای حل این مشکل و به منظور هماهنگ کردن روش پیش‌گویی خطی با دستگاه شنوازی انسان روش پیش‌گویی خطی مبتنی بر درک انسان پیشنهاد شد. (گریگر ، ۲۰۱۱) برای محاسبه این ویژگی به جای مقیاس مل از یک مقیاس جدید به نام بارک استفاده شده است. برخی از آزمایش‌های انجام‌شده بیان گر این مسئله هست که روش MFCC نسبت به PLP دارای نتایج بهتری در تشخیص گفتار است؛ اما روش PLP نسبت به تغییرات تعداد ضرایب و تعداد فیلترهای مورد استفاده در محاسبه ضرایب دارای نتایج پایدارتری نسبت به MFCC است. علاوه‌بر این ویژگی PLP در برابر نویه مقاومت بیشتری از خود نشان می‌دهد و مصون‌تر است. (فان، ۲۰۰۰) به طور کلی این آنالیز به دلیل کم‌بودن پیچیدگی محاسباتی آن یک ویژگی کاراست و سیگنال گفتار را با بعدی کمتر از خود سیگنال نمایش می‌دهد. و این ویژگی‌ها باعث شده است که از این آنالیز در تشخیص گفتار مستقل از گوینده بیشتر استفاده شود.



(شکل-۱) : مراحل محاسبه MFCC

۲-۳- مشتق نخست و دوم ضرایب مل-کپستروم

با اعمال سایر مراحل شکل (۱) ضرایب مل-کپستروم حاصل می‌شود. علاوه بر ضرایب مل-کپستروم از مشتق نخست ($\Delta MFCC$) و مشتق دوم آن ($\Delta\Delta MFCC$) نیز به عنوان ویژگی استفاده می‌شود که از روابط (۱) و (۲) به دست می‌آیند. (بهروزا، ۲۰۱۲).

$$\Delta MFCC = D[n] = C[n+m] - C[n-m] \quad (1)$$

$$\Delta\Delta MFCC = DD[n] = D[n+m] - D[n-m] \quad (2)$$

که در این روابط C بردار ضرایب مل-کپستروم و n شماره قاب است. m را در عمل به طور تقریبی ۲ یا ۳ در نظر می‌گیرند.

^۱ Linear Prediction Cepstral

^۲ Perceptual Linear Prediction

۲-۴- اعمال تبدیل Zak به سیگنال گفتار

ویژگی دیگری که در این مقاله از آن به عنوان یک ویژگی جدید استفاده شده است، تبدیل Zak است. این تبدیل قادر است ضرایبی از سیگنال گفتار، هم از حوزه زمان و هم از حوزه فرکانس استخراج کند. رابطه^(۵) نحوه محاسبه این تبدیل را نشان می‌دهد (اوسلندر، ۱۹۹۱).

$$Z_{sig}(t, v) = \sum_{n=-\infty}^{n=+\infty} S(t+n)e^{-j2\pi nv} \quad (5)$$

در این رابطه منظور از S همان سیگنال گفتار است، v دوره تناوب را نشان می‌دهد و t نشان دهنده قراردادشتن سیگنال در حوزه زمان است. با اعمال این تبدیل به هر قاب سیگنال گفتار یک ماتریس در خروجی خواهیم داشت که شامل N سطر و M ستون خواهد بود. N تعداد ضرایب حاصل از تبدیل در حوزه زمان و M تعداد ضرایب حاصل در حوزه فرکانس است.

با پایان یافتن مرحله استخراج ویژگی و تشکیل بردار ویژگی‌های مختلف که از ترکیب ویژگی‌های استخراج شده به دست می‌آیند، وارد مرحله طبقه‌بندی می‌شویم.

۵- طبقه‌بندی

پس از تشکیل بردار ویژگی‌ها به منظور دسته‌بندی و تشخیص آن‌ها باید از یک طبقه‌بند مناسب استفاده کرد. در این مقاله ابتدا از پنج طبقه‌بند مختلف شامل شبکه عصبی پرسپترون چندلایه (MLP)، شبکه عصبی احتمالاتی (PNN)، طبقه‌بند k نزدیک‌ترین همسایه (KNN)، ماشین بردار پشتیبان (SVM)، و شبکه توابع بنیادی شعاعی (RBF) استفاده شده است. در مرحله بعد بعضی از طبقه‌بندهای به کار گرفته شده به دو روش مختلف و به صورت دوتایی و سه‌تایی با یکدیگر ترکیب شده‌اند. هدف این است که با مقایسه نتایج مربوط به هر کدام از این طبقه‌بندها و ترکیب آن‌ها، یک حالت را که دارای عملکرد بهتر و پیوینه‌تری است، انتخاب کنیم. ابتدا توضیح مختصری درباره هر کدام از طبقه‌بندها ارائه می‌شود و سپس به بررسی نتایج خواهیم پرداخت.

۵-۱- شبکه MLP^۲

این شبکه شامل سه لایه به نام‌های ورودی، مخفی و خروجی است که تعداد سلول‌های هر لایه به روش سعی و

² Multi-Layer Perceptron

ویژگی‌هایی که بیان شد متدائل‌ترین ویژگی‌های هستند که در یک سامانه تشخیص گفتار و لهجه به کار می‌روند.

۴- روش‌های پیشنهادی استخراج ویژگی

در این قسمت از مقاله به بیان توضیحات و جزئیات بیشتری از روش‌های پیشنهادی برای مرحله استخراج ویژگی می‌پردازیم.

۴-۱- دامنه مرکزی طیفی و مشتق نخست آن (ΔSCM و SCM)

ویژگی دامنه مرکزی طیفی نخستین بار در سال ۲۰۱۰ به عنوان یک ویژگی در تشخیص گوینده استفاده شده است. این ویژگی و مشتق نخست آن با اعمال تغییراتی در مراحل محاسبه ضرایب مل-کپستروم به دست می‌آید. برای محاسبه^(۱) SCM با توجه به شکل (۱) به جای مرحله محاسبه انرژی از رابطه (۳) دامنه‌های مرکزی طیفی را محاسبه کرده و سپس لگاریتم و تبدیل کسینوسی ضرایب حاصل را به دست می‌آوریم (کارن، ۲۰۱۰).

$$M_k = \frac{\sum_{f=l_k}^{u_k} f |S[f]| \omega_k[f] |}{\sum_{f=l_k}^{u_k} f} \quad (3)$$

در این رابطه u_k فرکانس لبه پایین و l_k فرکانس لبه بالای هر فیلتربانک، مثلثی شکل است. $|S[f]|$ تبدیل یافته هر قاب از سیگنال گفتار از حوزه زمان به حوزه فرکانس و $\omega_k[f]$ نشان‌دهنده هر کدام از فیلتربانک‌های مثلثی شکل هستند. همانند ضرایب مل-کپستروم برای ویژگی اخیر نیز مشتقاتی وجود دارد که می‌توان آن‌ها به عنوان یک ویژگی در کنار سایر ویژگی‌ها استفاده کرد. در این مقاله از مشتق نخست ویژگی SCM استفاده شده است که با توجه به رابطه (۴) به دست می‌آید.

$$\Delta SCM = D[n] = S[n+m] - S[n-m] \quad (4)$$

که در این رابطه S بردار ضرایب دامنه مرکزی طیفی و n شماره قاب است. همانند رابطه محاسبه $\Delta MFCC$ ، m را در عمل به طور تقریبی ۲ یا ۳ در نظر می‌گیرند.

فصل نهم

¹ Spectral Centroid Magnitude



این نمونه‌ها k_n نزدیک‌ترین همسایه‌های x هستند. در حالت کلی k را به صورت k_n در نظر می‌گیریم که k_n تابعی تعریف شده از n است. اگر چگالی نقاط آموزش اطراف x زیاد باشد، سلول کوچک می‌شود و بنابراین نتیجه به دست آمده نتیجه بهتری است و در صورتی که چگالی نقاط آموزش اطراف x کم باشد، سلول بزرگ می‌شود.

۴-۴- طبقه‌بند^۳ SVM

SVM درواقع یک طبقه‌بندی کننده دودویی است که دو طبقه را با استفاده از یک مرز خطی از هم جدا می‌کند. در این روش با استفاده از تمامی باندها و یک الگوریتم بهینه‌سازی، نمونه‌هایی که مرز طبقه‌ها را تشکیل می‌دهند، به دست می‌آورند. این نمونه‌ها را بردارهای پشتیبان گویند. تعدادی از نقاط آموزشی را که کمترین فاصله تا مرز تصمیم‌گیری دارند، به عنوان زیرمجموعه‌ای برای تعریف مرزهای تصمیم‌گیری و به عنوان بردار پشتیبان می‌توان در نظر گرفت (واپنایک، ۱۹۹۱). البته قبل از تقسیم خطی باید دقت کرد که برای این که ماشین بتواند داده‌های با پیچیدگی بالا را دسته‌بندی کند، داده‌ها را با تابعی مشخص به فضایی با ابعاد بالاتر منتقل می‌کنیم. برای محاسبه مرز تصمیم‌گیری دو طبقه به طور کامل جدا از هم از روش حاسیه بهینه استفاده می‌شود.

۴-۵- شبکه^۴ RBF

RBF روشی برای تقریب توابع است و یادگیری با آن ارتباط نزدیکی با شبکه‌های عصبی مصنوعی دارد. در این روش فرضیه یادگرفته شده به صورت رابطه (۶) است.

$$\hat{f}(x) = w_0 + \sum_{u=1}^k w_u K_u(d(x_u, x)) \quad (6)$$

w_0 نشان‌دهنده بردار وزن نخستین و w_u بردار وزن‌های بروزشده است. x خروجی مطلوب، x_u ورودی تابع و d فاصله بین ورودی و خروجی مطلوب است. و درنهایت $(x) \hat{f}$ خروجی شبکه RBF را به‌ازای هر ورودی نشان می‌دهد. در این روش از تعداد k تابع کرنل برای تقریب تابع استفاده می‌شود. تابع کرنل به‌طور معمول به صورت تابع گاوسی رابطه (۷) با واریانس σ_u^2 انتخاب می‌شود.

³ Support Vector Machine

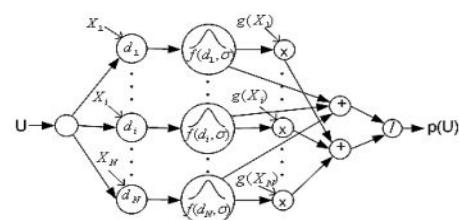
⁴ Radial Basis Function

خطا مشخص می‌شود. این شبکه بر مبنای الگوریتم پس انتشار خطای آموزش می‌بیند. بدین ترتیب که خروجی‌های واقعی با خروجی‌های دلخواه مقایسه و وزن‌ها با الگوریتم پس انتشار به صورت تحت نظارت تنظیم می‌شوند تا الگوی مناسب به وجود آید. روند یادگیری هنگامی متوقف می‌شود که مجموع کل خطای از مقدار آستانه تعیین شده کمتر شود و یا تعداد کل دوره آموزش به پایان برسد (هیکین، ۲۰۰).

۲-۵- شبکه^۱ PNN

این طبقه‌بند از نوع شبکه‌های عصبی RBF است و برپایه تابع چگالی احتمال نمایی و قانون تصمیم‌گیری Bayes عمل می‌کند. PNN یک شبکه پیش‌رونده سه‌لایه بوده و از یادگیری بانظارت استفاده می‌کند. لایه نخست ورودی‌ها را دریافت، لایه میانی بردار احتمال را بر پایه توزیع هر طبقه تعیین و لایه آخر مقدار بیشینه احتمال را انتخاب کرده و طبقه مربوطه را تعیین می‌کند. این شبکه به صورت موازی، محاسبات را انجام می‌دهد، آموزش ساده‌ای دارد و به‌دلیل سرعت بالا در کاربردهای بی‌درنگ مورد استفاده قرار می‌گیرد (قليپور على، ۱۳۹۱).

شکل (۲) تصویری از این شبکه را نشان می‌دهد.



شکل -۲: نحوه عملکرد شبکه PNN (قليپور على، ۱۳۹۱)

۳-۵- شبکه^۲ KNN

KNN یک الگوریتم آموزش با سربرستی است. در حالت کلی از این الگوریتم به دو منظور تخمین تابع چگالی توزیع داده‌های آموزش و طبقه‌بندی داده‌های آزمون بر اساس الگوهای آموزش استفاده می‌شود. برای تخمین $(x) p$ از روی n نمونه آموزش توسط الگوریتم KNN می‌توانیم یک سلول به مرکزیت x ایجاد کرده و اجازه دهیم شاع این سلول تا حدی گسترش بیندا کند که k_n نمونه آموزش را در برگیرد.

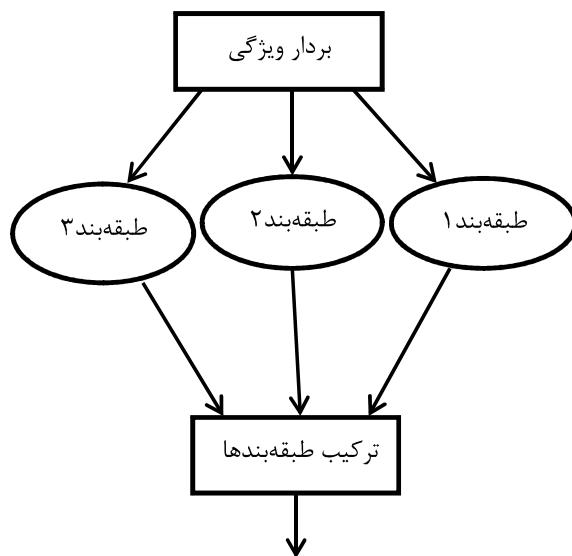
¹ Probabilistic Neural Network

² K-Nearest Neighbor

نتیجه حاصل را با هم جمع می‌کنیم و سپس نرخ بازشناسی را با ترکیب جدید به دست می‌آوریم.

در این پژوهش از هر دو روش مذکور استفاده شده است. روش میانگین‌گیری برای زمانی که سه طبقه‌بند را با هم ترکیب کردایم و روش تعیین وزن هنگامی که دو طبقه‌بند با یکدیگر ترکیب شده است، استفاده می‌شود.

نکته دیگری که در مورد ترکیب طبقه‌بندها باید در نظر گرفت، این است که بردار ویژگی که هر کدام از طبقه‌بندها با آن آموزش دیده‌اند هم می‌تواند متفاوت باشد و هم به‌طور کامل مشابه هم باشند که در اینجا ویژگی استخراج شده برای آموزش هر کدام از طبقه‌بندها با هم یکسان است. شکل (۳) روش استفاده شده را نشان می‌دهد.



(شکل-۳) : روش استفاده شده برای ترکیب طبقه‌بندها

روش ترکیب طبقه‌بندها ، طبق مشاهدات ما ، تاکنون یکبار در مورد تشخیص لهجه‌های فارسی به کار گرفته شده است که در آن دو طبقه‌بند GMM و PRLM با یکدیگر به روش موازی ترکیب و باعث بهبود نتایج تا حد مطلوبی (حدود ۲۰ درصد) شده‌اند. (جلال‌وند، ۱۴۰۲)

۶- تشخیص لهجه‌ها در محیط نوفه‌ای

یکی از عواملی که ممکن است در هر سامانه تشخیص گفتاری باعث کاهش کارآیی و همچنین نرخ بازشناسی آن سامانه شود، وجود نوфе در محیط است. حال اگر دو عامل

$$K_u(d(x_u, x)) = e^{\frac{1}{2\sigma_u^2}d(x_u, x)^2} \quad (7)$$

نشان داده شده است در صورتی که تعداد کافی تابع کرنل گاؤسی انتخاب شود، با استفاده از این شبکه می‌توان هر تابعی را با خطای تاحدودی کم تقریب زد. از مزایای این شبکه می‌توان به آموزش آسان‌تر آن نسبت به شبکه‌های عصبی معمولی، که از روش پسانشار استفاده می‌کنند، نام برد.

۶-۵- ترکیب طبقه‌بندها^۱

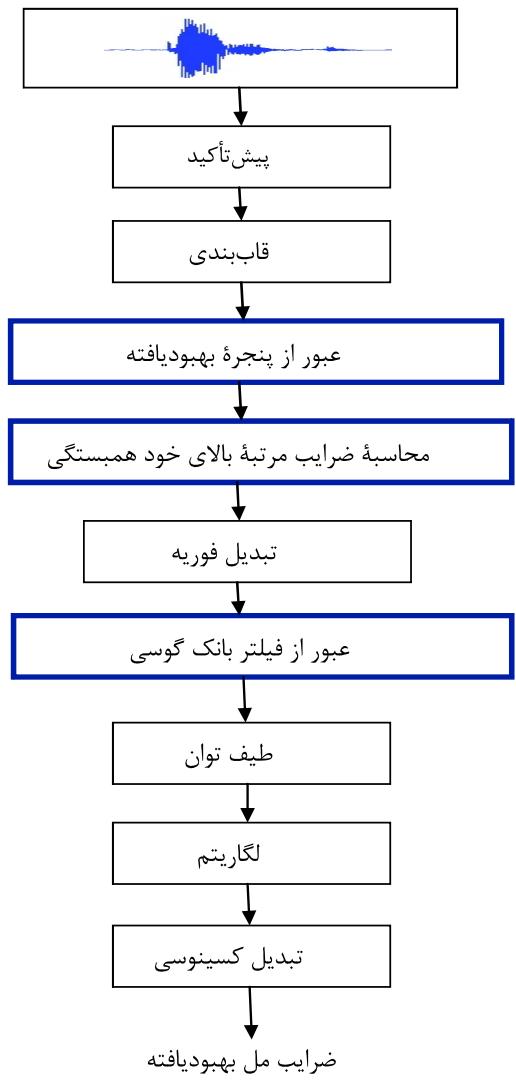
هدف نهایی طراحی یک سامانه شناسایی الگو رسیدن به بهترین عملکرد طبقه‌بندی برای مساله موجود است. (هیکین، ۱۹۹۹) از آنجایی که هیچ طبقه‌بندی به‌طور کامل قادر به حل تمام مسائل نیست، ترکیب طبقه‌بندها به عنوان روشی برای کاهش خطای طبقه‌بندی و افزایش کارآیی سامانه‌ها پیشنهاد شده است.

ایده اصلی این طرح، ایجاد طبقه‌بندهای گوناگون با ناحیه‌های خطای متفاوت است تا سامانه ترکیبی با بهره‌گیری از نقاط قوت تک‌تک طبقه‌بندها و رفع خطای ایجاد شده یک طبقه‌بند، با اطلاعات به دست آمده از طریق طبقه‌بندهای دیگر، خطای کل سامانه را کاهش دهد (دایتنیج و فروند، ۱۹۹۷ و ۱۹۹۵).

برای ترکیب طبقه‌بندها، روش‌ها و طرح‌های مختلفی پیشنهاد شده است که از آن جمله می‌توان روش‌های سری، موازی و ترکیبی را نام برد. در این مقاله از روش موازی استفاده شده است. در این روش همه طبقه‌بندها به‌طور همزمان با هم ترکیب می‌شوند. یکی از ساده‌ترین فرآیندهایی که در مرحله ترکیب استفاده می‌شود، روش میانگین‌گیری است؛ بدین صورت که نتیجه حاصل از تمام طبقه‌بندهای موردنظر برای ترکیب کردن با هم جمع و سپس بر تعداد طبقه‌بندهای ترکیب شده تقسیم می‌شوند.

روش دیگر این است که برای هر کدام از طبقه‌بندها یک وزن خاص تعیین و سپس این وزن را در نتیجه نهایی هر کدام از طبقه‌بندها ضرب کنیم. وزن‌ها باید به گونه‌ای باشند که مجموع وزن‌های به کار رفته برابر یک شود. بهینه‌ترین وزن ممکن برای هر طبقه‌بند را از روش سعی و خطابه دست می‌آوریم. پس از ضرب کردن وزن در نظر گرفته شده برای هر طبقه‌بند در طبقه‌بند موردنظر،

^۱ Combination classifier



(شکل-۴): روش پیشنهادی برای محیط نوشهای (ضرایب مل-کپستروم بهبودیافته)

ویرگی این نوشه با فرض ناهمبسته‌بودن نسبت به سیگنال اصلی این است که تابع خودهمبستگی مربوط به آن تا حدود زیادی نسبت به زمان بدون تغییر و نزدیک به صفر است. بنابراین تابع خودهمبستگی به صورت رابطه (۱۰) بیان خواهد شد.

$$R_{XX}(m, k) = R_{SS}(m, k) + R_{dd}(m) \quad (10)$$

در این رابطه k شماره هر قاب است. رابطه (۱۰) نشان می‌دهد که تابع خودهمبستگی نوشه ناهمبسته به سیگنال، مستقل از فریم است و با توجه به ناچیز بودن آن

لهجه و نوشه در یک گفتار کنار هم قرار گیرند، شرایط دشوارتری برای تشخیص آن ایجاد خواهد شد. بنابراین نیاز است تا روشی برای حذف نوشه از یک سیگنال گفتار لهجه‌دار قبل از تشخیص نوع لهجه آن پیشنهاد شود. برای حذف نوشه از یک سیگنال گفتار، تاکنون روش‌های متعددی پیشنهاد شده است که اغلب آن‌ها با ایجاد تغییراتی در مراحل محاسبه ضرایب مل-کپستروم به دست آمدان. این تغییرات را به سه دسته مختلف تقسیم‌بندی می‌توان کرد:

۱- مدل‌های بهبودیافته شامل تغییر در بلوك‌های پایه این الگوریتم

۲- مدل‌های بهبودیافته شامل یک بلوك تکمیل‌کننده که به الگوریتم پایه اضافه شده است.

۳- بهبود در پیاده‌سازی سخت‌افزاری این الگوریتم توسط کاستن محاسبات ضرب و جمع در الگوریتم پایه

در این مقاله با ترکیب چند روش از روش‌های استفاده شده به نتیجه مطلوبی در کاهش نوشه از یک سیگنال گفتار لهجه‌دار رسیده‌ایم. شکل (۴) نحوه به دست آمدن ویرگی جدید را که شامل سه تغییر نسبت به الگوریتم نشان داده شده در شکل است، نشان می‌دهد.

تغییرات ایجاد شده به شرح زیر است.

الف) در بیشتر الگوریتم‌های بررسی شده برای مرحله عبور از پنجره از یک پنجره همینگ ساده استفاده شده است؛ اما پنجره‌ای که ما در اینجا استفاده می‌کنیم، با رابطه (۸) نشان داده می‌شود. (ساهیدولا، ۲۰۱۲)

W(n) یک پنجره همینگ ساده با طول مناسب است. در پنجره پیشنهادشده دو پارامتر پراکندگی طیفی و عرض بخش اصلی پنجره افزایش می‌یابد که یک تغییر مطلوب است؛ اما هم‌گرایی بخش‌های جانبی کاهش پیدا می‌کند که این یک تغییر نامطلوب است؛ اما در مجموع تغییرات ایجاد شده حاکی از بهبود نتایج است و بتاراین این تغییر را اعمال می‌کنیم. (ساهیدولا، ۲۰۱۲)

ب) نوشه اضافه شده به سیگنال را می‌توان با رابطه (۹) بیان کرد.

$$w_{new}(n) = nw(n) \quad (8)$$

$$X(n) = S(n) + d(n) \quad (9)$$

که در آن $s(n)$ سیگنال ورودی و $d(n)$ نوشه اضافی شونده به سیگنال است.

جملات لهجه‌دار مربوط به این مقاله از پایگاه داده فارس‌دات پژوهشگاه توسعه فناوری‌های پیشرفته خواجه نصیرالدین طوسی که معتبرترین پایگاه دادگان گفتار زبان فارسی است، استخراج شده‌اند. این پایگاه داده شامل ۶۰۸۰ جمله است که توسط ۳۰۴ نفر گوینده ایرانی با ده لهجه مختلف بیان شده است. از هر گوینده ۲۰ جمله ضبط گردیده است. در این پایگاه داده گویندگان از لحاظ سن، جنس، سطح تحصیلات و لهجه‌هایشان متمایز شده‌اند که شرایط خوبی را برای انجام پژوهش‌های مختلف فراهم کرده است. در این پژوهش برای شناسایی لهجه‌ها پنج لهجه مختلف انتخاب شده است. در انتخاب نوع لهجه‌ها سعی شده است، لهجه بیشتر مردم مناطق ایران را لحاظ گرفتار شامل شود. این لهجه‌ها شامل تهرانی، ترکی، اصفهانی، شمالی و جنوبی هستند.

۲-۷ آزمایش نخست: مقایسه لهجه‌ها به صورت دوتایی

در آزمایش نخست با انتخاب ویژگی MFCC از میان ویژگی‌های متداول و ویژگی SCM از بین ویژگی‌های پیشنهادی و با بهکاربردن تمامی طبقه‌بندی‌های معروفی شده در این مقاله، نرخ بازنگاری لهجه‌ها به صورت دوتایی به دست آمده‌اند. هدف از این آزمایش مقایسه عملکرد ویژگی پیشنهادی با ویژگی‌های متداول در حالی که دو طبقه داریم، است.

جدول (۱) نتایج را برای ویژگی MFCC نشان می‌دهد.

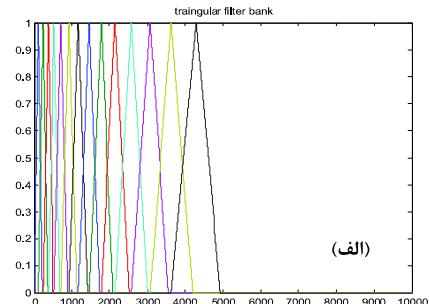
(جدول -۱): نرخ تشخیص لهجه‌ها به صورت دوتایی با

طبقه‌بندی‌های مختلف و ویژگی متداول MFCC

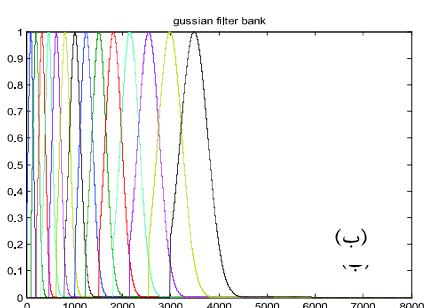
میانگین	RBF	SVM	KNN	PNN	MLP	طبقه‌بند
69.23	72.21	82.62	30.46	80.4	80.05	اصفهانی و تهرانی
76.95	74.95	83.03	62.1	83.3	81.4	تهرانی و ترکی
69.13	70.59	80.55	32.04	81.65	80.86	ترکی و اصفهانی
74.07	76.57	81.65	46.26	82.34	83.55	شمالی و جنوبی
70.59	71.24	82.75	35.91	81.79	81.26	تهرانی و جنوبی
70.15	69.79	79.72	40.68	80.55	80.05	تهرانی و شمالی
70.28	74.31	80.13	35.03	80.96	80.99	ترکی و جنوبی
69.81	72.21	83.58	34.9	78.62	79.78	ترکی و شمالی
75.15	75.12	82.89	55.27	81.37	81.13	اصفهانی و شمالی
74.18	71.08	80.82	55.03	82.34	81.67	اصفهانی و جنوبی

می‌توان با این روش اثر نویه را کاهش داد. تخریب سیگنال توسط نویه در ضرایب مرتبه پایین تابع خودهمبستگی، بیشتر از ضرایب مرتبه بالا است، بنابراین با حذف آن‌ها، از اثر نویه کاسته خواهد شد.

ج) آخرین تغییر ایجادشده در محاسبه ضرایب مل-کپستروم، تغییر نوع فیلتربانک از مثلثی به گوسی است. در فیلتربانک مثلثی اطلاعات بخش‌هایی از قاب که در نقاط ابتدایی و انتهایی و خارج از زیربخش‌ها، قرار می‌گیرند از دست می‌روند؛ زیرا مثلثها در خارج از زیرباندها وزنی ندارند؛ اما اگر به جای این فیلتربانک از یک فیلتربانک گوسی که متقارن نیز می‌باشد، استفاده کنیم، به دلیل وجود وزن در خارج از زیرباندهای آن، مانع از دست‌رفتن اطلاعات در این بخش‌ها می‌شود. مزیت دیگر آن نسبت به فیلتربانک مثلثی، شروع و اتمام آن با شبکه کمتر و ملایم‌تر است. این فیلتربانک با ایجاد همبستگی بیشتر بین زیرباندها اطلاعات از دست‌رفته در مزهای را کاهش داده و در بهبود الگوریتم و به دنبال آن بالابردن نرخ بازنگاری سامانه خودکار پردازش گفتار مؤثر خواهد بود. شکل (۵) این دو نوع فیلتربانک را نشان می‌دهد.



(الف)



(ب)

(شکل -۶): (الف): فیلتربانک مثلثی، (ب): فیلتربانک گوسی

فصل نهم



۷-۳- آزمایش دوم: مقایسه همه لهجه‌ها در حضور ویژگی‌های مختلف

در آزمایش دوم، عملکرد بردار ویژگی‌های مختلف متداول و پیشنهادی در حضور همه لهجه‌های منتخب این مقاله سنجیده شده است.

ابتدا با تشکیل بردار ویژگی‌هایی از MFCC و مشتقات آن و همچنین فرکانس‌های فرمت، نرخ بازناسی را در حضور طبقه‌بندهای مختلف به دست می‌آوریم. جدول (۳) نتایج حاصل از این آزمایش را نشان می‌دهد.

(جدول -۳): نرخ بازناسی لهجه‌ها با ویژگی‌های متداول

ویژگی	طبقه‌بند	MLP	PNN	KNN	SVM	RBF
MFCC		55.63	72.33	22.07	75.3	25.73
MFCC + ΔMFCC		65.95	71.3	21.11	73.73	36.04
MFCC + ΔMFCC + ΔΔMFCC		68.45	70.67	19.53	74.63	68.70
MFCC + ΔMFCC + ΔΔMFCC + 2F		65.86	72.37	21.22	74.09	63.46

در مرحله بعد بردارهایی از ویژگی‌های مختلف پیشنهادی را تشکیل می‌دهیم و آن‌ها را به طبقه‌بندهای مختلف اعمال می‌کنیم تا نرخ بازناسی به ازای پنج لهجه مختلف به دست آید.

جدول (۴) نتایج حاصل را نشان می‌دهد.

(جدول -۴): نرخ بازناسی لهجه‌ها با ویژگی‌های پیشنهادی

ویژگی	طبقه‌بند	MLP	PNN	KNN	SVM	RBF
SCM		53	71.83	26.93	75.16	36.75
SCM + ΔSCM		58.13	71.70	26.33	74.54	33.31
Zak Transform		48.4	71.65	24.56	74.11	52.61

در جدول (۵) میانگین نرخ بازناسی به ازای هر کدام از ویژگی‌ها با یکدیگر مقایسه شده‌اند.

(جدول -۵): مقایسه میانگین نرخ بازناسی لهجه‌ها با ویژگی‌های مختلف

ویژگی	میانگین نرخ تشخیص
MFCC	50.21
MFCC + ΔMFCC	53.62
MFCC + ΔMFCC + ΔΔMFCC	60.39
MFCC + ΔMFCC + ΔΔMFCC + 2F	59.4
SCM	52.73
SCM + ΔSCM	52.8
Zak Transform	54.26

همان‌طور که نتایج جدول (۱) نشان می‌دهد، سه طبقه‌بند PNN و SVM بـه طور میانگین دارای عملکرد بهتری نسبت به سایر طبقه‌بندها هستند.

در جدول (۲) نتایج به ازای ویژگی پیشنهادی نشان داده شده است.

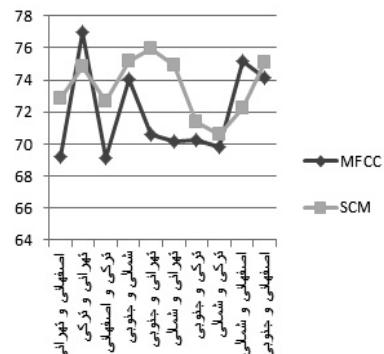
(جدول -۲): نرخ تشخیص لهجه‌ها به صورت دوتایی با

طبقه‌بندهای مختلف و ویژگی پیشنهادی

میانگین	RBF	SVM	KNN	PNN	MLP	طبقه‌بند
72.86	75.44	83.40	44.16	83.11	78.22	اصفهانی و تهرانی
74.81	69.79	81.77	62.38	80.88	79.25	تهرانی و ترکی
72.71	74.31	79.11	45.17	78.96	76	ترکی و اصفهانی
75.18	70.27	82.96	58.98	82.96	80.74	شمالي و جنوبي
75.97	76.73	84.14	57.69	84	77.33	تهرانی و جنوبي
74.89	74.63	82.51	57.17	82.07	78.07	تهرانی و شمالي
71.35	66.55	82.51	46.26	81.92	79.55	ترکی و جنوبي
70.57	73.18	80.44	41.29	80	78.96	ترکی و شمالي
72.28	64.62	80.44	61.69	80.59	74.07	اصفهانی و شمالي
75.09	71.08	81.62	63.67	81.33	77.77	اصفهانی و جنوبي

همان‌طور که نتایج دو جدول اخیر نشان می‌دهد، ویژگی پیشنهادی SCM دارای میانگین عملکرد بهتری نسبت به ویژگی MFCC در تشخیص لهجه‌ها به صورت دوتایی است. همچنین طبقه‌بند KNN نسبت به سایر طبقه‌بندها بازدهی کمتر دارد. شکل (۶) مقایسه میانگین عملکرد این دو ویژگی را نشان می‌دهد.

لازم به ذکر است که به دلیل افزایش حجم مقاله از به دست آوردن نرخ بازناسی لهجه‌ها به صورت دوتایی و در حضور دو ویژگی پیشنهادی دیگر، یعنی ΔSCM و تبدیل Zak صرف نظر می‌کنیم.



(شکل -۷): مقایسه میانگین عملکرد دو ویژگی و SCM در تشخیص لهجه‌ها به صورت دوتایی

طبقه‌بندها استفاده می‌کنیم. به طور معمول وزنی که به طبقه‌بندی با نرخ بازشناسی کمتر تعلق می‌گیرد از وزنی که به طبقه‌بندی با نرخ بازشناسی بیشتر تعلق می‌گیرد، بیشتر است. به عنوان مثال اگر شبکه MLP دارای نرخ بازشناسی ۴۲/۷ درصد و شبکه SVM دارای نرخ بازشناسی ۶۳/۲۴ درصد باشد، وزنی که به MLP تعلق می‌گیرد، برابر ۰/۶ و وزنی که به PNN تعلق می‌گیرد، برابر ۰/۴ است. جدول (۶) نتایج حاصل را نشان می‌دهد.

(جدول - ۶): نرخ بازشناسی لهجه‌ها با ترکیب طبقه‌بندها

به صورت دوتایی

ویژگی	طبقه‌بند	MLP	PNN	SVM	MLP + PNN	PNN + SVM	SVM + MLP
MFCC + ΔMFCC + ΔΔMFCC	64.62	71.52	71.23	84.6	87.71	86.59	
Zak Transform	42.7	70.03	63.24	76.57	85.58	83.93	

همان‌طور که نتایج جدول نشان می‌دهد، با انجام عمل ترکیب طبقه‌بندها شاهد بهبود نرخ بازشناسی تا حد مطلوبی خواهیم بود. لازم به ذکر است که دلیل تفاوت نرخ بازشناسی به‌ازای هر کدام از طبقه‌بندها در جدول اخیر با مقادیر جدول (۳) و جدول (۴) این است که با هر بار راه اندازی طبقه‌بندها، وزن‌های نخستین به صورت تصادفی تعیین می‌شود و این باعث تغییر نرخ بازشناسی می‌شود. لازم به ذکر است که بردار وزن نخستین برای هر بار آموزش شبکه، توسط نرم افزار به صورت خود کار تعیین که باعث پایابی نتایج می‌شود. در جدول (۷) میانگین عملکرد روش ترکیب به صورت دوتایی و میزان بهبود نرخ بازشناسی نشان داده شده است.

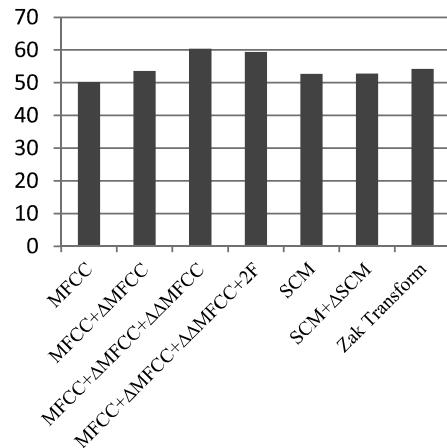
(جدول - ۷): مقایسه میانگین عملکرد طبقه‌بندها و ترکیب آنها

ویژگی	طبقه‌بند	میانگین بدون ترکیب	میانگین ترکیب به صورت دوتایی	میزان بهبود
MFCC + ΔMFCC + ΔΔMFCC	69.12	86.3	17.18%	
Zak Transform	58.65	82.02	23.37	

در مرحله بعد به‌ازای همان دو ویژگی انتخابی، هر سه طبقه‌بند را با یکدیگر ترکیب می‌کنیم. در این مرحله ترکیب کردن را با روش میانگین‌گیری انجام می‌دهیم، به این

همان‌طور که مقادیر این جدول نشان می‌دهد، در میان ویژگی‌های متداول استفاده شده در این مقاله، بردار ویژگی $MFCC + \Delta MFCC + \Delta \Delta MFCC$ بالترين بازدهی است که دليل اين امر می‌تواند وجود مشتقات ضرایب مل-کپستروم باشد. در میان ویژگی‌های پیشنهادشده تبدیل Zak بیشترین بازدهی را دارد که این امر نیز می‌تواند ناشی از حضور همزمان ویژگی‌های حوزه زمان و فرکانسی در این تبدیل باشد. نمودار میله‌ای شکل (۷) مقایسه این بازدهی را نشان می‌دهد.

در این مرحله از آزمایش‌ها به این نکته دست یافتیم که در حضور پنج لهجه، برخی از ویژگی‌های متداول استفاده شده کارآیی بهتری نسبت به ویژگی‌های پیشنهادی دارند و در میان طبقه‌بندها، ماشین بردار پشتیبان، کارآیی بالاتری دارد.



(شکل - ۸): نمودار میله‌ای مقایسه بازدهی ویژگی‌های مختلف در تشخیص لهجه‌ها

۴-۷- آزمایش سوم: ترکیب طبقه‌بندها

در این بخش از مقاله قصد داریم با استفاده از روش ترکیب طبقه‌بندها بازشناسی لهجه‌ها را انجام دهیم. از میان ویژگی‌های متداول و پیشنهادی، ویژگی را انتخاب می‌کنیم که دارای بالاترین بازدهی باشد. این ویژگی‌ها شامل بردار ویژگی $MFCC + \Delta MFCC + \Delta \Delta MFCC$ و ضرایب تبدیل Zak است.

به منظور انجام عمل ترکیب کردن، طبقه‌بندهای SVM، PNN و MLP که بالاترین بازدهی را داشته‌اند، انتخاب شده‌اند.

ابتدا طبقه‌بندها را به صورت دوتایی با یکدیگر ترکیب

می‌کنیم. برای انجام این کار از روش وزن دهنی به خروجی

سال ۱۳۹۵ شماره ۲ پیاپی ۲۸

۷-۵- آزمایش چهارم : تشخیص لجههای در محیط نویه‌ای

ویژگی پیشنهادی در این مقاله را که با ایجاد سه تغییر در الگوریتم پایه محاسبه ضرایب مل-کپستروم به دستمی آید با نام AGMFCC نام‌گذاری می‌کنیم. برای انجام آزمایش از نسبت سیگنال به نویه‌های مختلف استفاده شده است. در مرحله نخست چهار ویژگی مختلف در کنار طبقه‌بند PNN به کار گرفته شده است. جدول (۹) نرخ بازناسی لجههای در محیط نویه‌ای با طبقه‌بند PNN و ویژگی‌های مختلف نتایج حاصل را نشان می‌دهد.

(جدول -۹): نرخ بازناسی لجههای در محیط نویه‌ای با طبقه‌بند PNN و ویژگی‌های مختلف

SNR	ویژگی	MFCC	SCM	Zak T	AGMFCC
Without noise		71.41	69.13	71.54	72.14
25dB		71.05	71.38	70.69	71.65
10dB		70.6	72.08	70.72	70.33
5dB		70.76	71.16	69.12	72.30
0dB		72.23	69.06	67.83	70.49
-5dB		71.25	63.17	67.05	72.41
میانگین		71.21	69.17	69.49	71.55

برای این‌که نشان داده شود ویژگی AGMFCC مستقل از طبقه‌بند در محیط نویه‌ای بهترین عملکرد را دارد، همان آزمایش بالا با طبقه‌بند SVM انجام شده است. (جدول (۱۰) نرخ بازناسی لجههای در محیط نویه‌ای با طبقه‌بند SVM و ویژگی‌های مختلف نتایج را نشان می‌دهد.

(جدول -۱۰): نرخ بازناسی لجههای در محیط نویه‌ای با طبقه‌بند SVM و ویژگی‌های مختلف

SNR	ویژگی	MFCC	SCM	Zak T	AGMFCC
Without noise		74.58	73.95	74.29	74.63
25dB		74.13	75.75	74.84	75.21
10dB		73.73	74.18	72.14	73.69
5dB		71.54	75.16	74.23	73.95
0dB		66.97	71.94	71.10	75.92
-5dB		58.52	69.97	68.58	72.12
میانگین		69.95	73.49	72.53	74.25

همان‌طورکه نتایج جدول (۹) و جدول (۱۰) نشان می‌دهد، ویژگی پیشنهادی AGMFCC که برای محیط‌های نویه‌ای استفاده می‌شود، بدون توجه به نوع طبقه‌بند بهترین

صورت که نتیجه تمام طبقه‌بندهای موردنظر را با هم جمع کرده و بر تعداد آن‌ها تقسیم می‌کنیم و سپس با تعیین یک آستانه خروجی نهایی این ترکیب را به دست می‌آوریم. جدول (۸) نتایج حاصل را نشان می‌دهد.

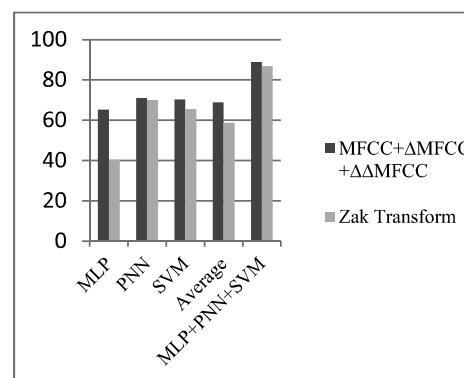
(جدول -۸): نرخ بازناسی لجههای با ترکیب طبقه‌بندها به صورت سه‌تایی

ویژگی	طبقه‌بند	MLP	PNN	SVM	MLP+SVM + PNN
MFCC + ΔMFCC + ΔΔMFCC		65.22	71.05	70.33	88.85
Zak Transform		40.27	70.03	62.52	86.93

همان‌طورکه جدول بالا نشان می‌دهد با انجام این روش ترکیبی نرخ بازناسی افزایش پیدا می‌کند. برای ویژگی متداول استفاده شده نسبت به میانگین خروجی طبقه‌بندها ۱۹/۹۸درصد بهبود و برای ویژگی پیشنهادی ۲۸/۲درصد بهبود را شاهد هستیم.

در این جدول نیز بنا به دلیلی که در قبل بیان شد، نرخ بازناسی حاصل از هر کدام از طبقه‌بندها با مقادیر جدول‌های قبلی مقداری تفاوت دارد.

نمودار شکل (۸) مقایسه نتایج را نشان می‌دهد.



(شکل -۹): نمودار مقایسه عملکرد طبقه‌بندها و ترکیب آنها به صورت سه‌تایی

نتایج بخش ۴-۷ بیان گر بهبود نتایج با بهکاربردن روش ترکیب طبقه‌بندهای است و این روش را می‌توان به عنوان شیوه‌ای کارآمد در تشخیص گفتار و لجه استفاده کرد.

۹- مراجع

قلی پور علی، صداقی محمدحسین و شمسی موسی، "طبقه‌بندی برخی از لهجه‌های زبان فارسی با استفاده از شبکه عصبی احتمالاتی" بیستمین کنفرانس مهندسی برق ایران، دانشگاه تهران، ۱۳۹۱.

A. Rabiee and S. Setayeshi, "Persian Accents Identification Using an Adaptive Neural Network", 2th .Int. Conf. on Education Tech-nology and Computer Science, pp. 7-11, 2010.

A.Faria "Accent Classification for Speech Recognition", Machine Learning for Multimodal Interaction, pp.285-293, Springer Berlin Hei-delberg, 2005.

C.Pedersen and J.Diederich , "Accent Class-ification Using Support Vector Machines", 6th ICIS ,2007.

F.Phan, M.T.Evangelia.sideman "Speaker Ident-ification using Nerula Network and Wavelets" IEEE Engineering in medicin and Biology Magazine, vol.191, pp.92-101,2000.

H.Behrava, Master's thesis, "Dialect and Accent Recognition". School of Computing, Eastern Finland university, 2012.

J.M.Karen Kua, T.Thiruvaran, M.Nosratighods, E.Ambikairajah and J.Epps. "Investigation of Spectral Centroid Magnitude and Frequency for Speaker Recognition", In Odyssey-2010, pp 007.

K.Kumpf and R.W.King "Automatic accent classif-ication of foreign accented australian English speech" Fourth international conference onspoken language, pp 1740-1743 vol.3 , Philadelphia,1996.

L. Auslander, I. Gertner, R. Tolimieri, "The Discrete Zak Transform Application to Time-Frequency Analysis and Synthesis of Nonstationary Signals" IEEE Trans. onSignal Proc., Vol. 39, No. 4, pp. 825-835, April 1991.

M. Sahidullah and G. Saha , "A novel wind-owing technique for efficient computation of mfcc for spea-ker recognition"Arxiv, pp 1206-2437 v1.,2012.

M.L.Arslan and H.L.Hansen , "Language Accent Classification in American English", Speech Communication 18(4): 353-367,1996.

O.Grigeor, C.Grigeor and V.Velican," Impaired Spe-ech Evaluation using Mel-Cepstrum Ana-lysis" International Journal Of Circuit, Systems And Signal Processing, 2014.

S.Haykin, "Neural Networks", Macmillan Coll-ege Publishing Company, 1999.

S.Jalalvand, A.Akbari, and B.Nasersharif, "A classifier combination approach for Farsi accents recog-nition" 20th Iranian Conference on Electric Engineering, pp.716-720, Tehran, Iran,2012.

عملکرد را دارد. همان‌طور که نتایج دو جدول اخیر نشان می‌دهد، در شرایط نوشهای ویژگی AGMFCC نسبت به سایر ویژگی‌ها میانگین عملکرد بهتری دارد.

۸- نتیجه‌گیری کلی

در این مقاله هدف ما تشخیص لهجه‌های مختلف زبان فارسی از روی شکل موج گفتار بود. برای این منظور سه مرحله پیش‌پردازش، استخراج ویژگی و طبقه‌بندی را اجرا کردیم. در مرحله استخراج ویژگی از تعدادی ویژگی متداول Zak و تبدیل ΔSCM ، SCM و RBF استفاده شد. درنهایت در مرحله طبقه‌بندی، پنج طبقه‌بند مختلف شامل MLP ، RBF و KNN ، SVM ، PNN و ترکیبی از آن‌هایی که دارای بالاترین بازدهی بودند، بهصورت دوتایی و سه‌تایی به کار گرفته شد. البته لازم به ذکر است که علاوه‌بر طبقه‌بندی‌های مورد استفاده در این مقاله، طبقه‌بندی‌های دیگری نیز وجود دارند که از آن‌ها می‌توان برای تشخیص گفتار، گوینده و لهجه استفاده کرد که از میان آن‌ها می‌توان به طبقه‌بند $I\text{-vector}$ / $PLDA$ اشاره کرد. $I\text{-vector}$ روش تحلیل تفکیک خطی احتمالاتی ($PLDA$) و اغلب برای تشخیص گوینده مستقل از متن استفاده می‌شود و با ایجاد تغییراتی می‌توان آن را برای تشخیص گوینده وابسته به متن نیز استفاده کرد. (استافلاکیس)

نتایج آزمایش‌هایی که در سه مرحله مختلف انجام شد نشان داد که در بین ویژگی‌های متداول، بردار ویژگی شامل ضرایب مل-کپستروم و مشتق نخست و دوم آن و در بین ویژگی‌های پیشنهادی ویژگی تبدیل Zak دارای بیشترین بازدهی هستند. در بین طبقه‌بندی‌های مختلف نیز طبقه‌بند SVM دارای میانگین عملکرد بهتری نسبت به سایر طبقه‌بندی‌ها است. علاوه‌بر این انجام عمل ترکیب طبقه‌بندی نیز باعث شد که نرخ بازنگاری پنج لهجه مختلف، بین ۱۷ تا ۲۸ درصد بهبود پیدا کند.

علاوه‌بر این یک ویژگی جدید که با ایجاد سه تغییر در الگوریتم اصلی محاسبه ضرایب مل-کپستروم به دست می‌آید، و آن را AGMFCC نامیدیم، برای تشخیص لهجه‌ها در محیط نوشهای استفاده شد؛ که نتایج آزمایش‌ها نشان دهنده بهبود عملکرد سامانه تشخیص لهجه‌ها در محیط نوشهای با استفاده از این ویژگی است.

فصل نهم





حسین مروی مدرک کارشناسی خود را در رشته مهندسی الکترونیک از دانشگاه فردوسی اخذ کرد. ایشان مقطع کارشناسی ارشد را در دانشگاه شیراز و در رشته مهندسی برق گرایش مخابرات به پایان رسانیدند. مدرک دکترای تخصصی خود را در کشور انگلستان، دانشگاه Surrey، مرکز CVSSP و در زمینه پردازش گفتار اخذ کردند. نامبرده هم‌اکنون به عنوان هیئت‌علمی با رتبه دانشیاری در دانشگاه صنعتی شهرورد، دانشکده مهندسی برق و راتیک مشغول تدریس است. زمینه‌های موردن علاقه ایشان، پردازش سیگنال‌ها، پردازش صوت، تشخیص گفتار، تشخیص گوینده و بهسازی گفتار است.

نشانی رایانمۀ ایشان عبارت است از:

h.marvi@shahroodut.ac.ir



دانیال دارابیان در سال ۱۳۸۶ در رشته دبیر فنی الکترونیک در دانشگاه حکیم سبزواری پذیرفته و با گذراندن واحدهای درسی لازم در سال ۱۳۹۰ در مقطع کارشناسی فارغ‌التحصیل شد. در مهر همان سال در رشته مهندسی الکترونیک گرایش سیستم در مقطع کارشناسی ارشد دانشگاه صنعتی شهرورد پذیرفته شد؛ و سپس در شهریور ۱۳۹۲ با دفاع از پایان‌نامه خود در زمینه بازنگاری مقاوم گفتار با استفاده از ضرایب مل-کیستروم بهبودیافته، تحصیلات در مقطع کارشناسی ارشد را به پایان برد. هم‌اکنون به عنوان هنرآموز در آموزش و پرورش و دانشگاه‌های غیرانتفاعی مشغول به تدریس است.

نشانی رایانمۀ ایشان عبارت است از:

Danial.darabian1@gmail.com

Simon Haykin, Neural Networks, Macmillan College Publishing Company, 1999.

T. G. Dietterich and G. Bakiri(1995), "Solving multiclass learning problems via error correcting output codes", J. of Artificial Intelligence Research 2, pp263-286.

T.Chen, C.Huang, E.Chang and J.Wang, "Automatic Accent Identification Using Gaussian Mixture Models" Conference on Automatic Speech Recognition and Understanding, pp.343-346, Italy, 2001.

T.Stafylakis, P. Kenny, P. Ouellet, J. Perez, M. Kockmann, and P. Dumouchel. "I-Vector-/PLDA Variants for Text-Dependent Speaker Recognition." Preprint submitted to Computer, Speech and Language, 2013.

V. Vapnik and A. Chervonenkis, "The necessary and sufficient conditions for consistency in the empirical risk minimization method," Pattern Recognition and Image Analysis, vol. 1, no. 3, pp. 283-305, 1991.

Y. Freund and R. Schapire.(1997) "A decision-theoretic generalization of on-line learning and an application to boosting". Journal of Computer and System Sciences, 55, 1, pp. 119-139.

Y.Zheng, R.Sproat, L.Gu, I.Shafran, H.Zohu, Y.Su, D.Jurafsky, R.Star, S-Y.Yoon, "Accent detection and speech recognition for Shanghai-Accented Mandarin", Interspeech, pp. 217-220, 2005.



مجتبی شریف نوقابی در سال ۱۳۸۶ در رشته دبیر فنی الکترونیک در دانشگاه حکیم سبزواری پذیرفته شد و با گذراندن واحدهای درسی لازم در سال ۱۳۹۰ با رتبه چهارم در بین

دانشجویان ورودی خود در مقطع کارشناسی فارغ‌التحصیل شد. در مهر همان سال در رشته مهندسی الکترونیک گرایش سیستم در مقطع کارشناسی ارشد دانشگاه صنعتی شهرورد پذیرفته شد؛ و سپس در بهمن ۱۳۹۲ با دفاع از پایان‌نامه خود در زمینه کسب رتبه دوم در بین دانشجویان هنرآموزی خود، تحصیلات در مقطع کارشناسی ارشد را به پایان برد. هم‌اکنون به عنوان هنرآموز در آموزش و پرورش و دانشکده‌های فنی مشغول به تدریس است.

نشانی رایانمۀ ایشان عبارت است از:

mojtabasharif@chmail.ir