

بازشناسی کد شناسائی شخصی و تصدیق هویت گوینده به منظور کنترل دسترسی از راه دور توسط تلفن

امیر نجاری
کارشناسی ارشد

محمد مهدی همایونپور
استادیار

دانشکده مهندسی کامپیوتر، دانشگاه صنعتی امیر کبیر

چکیده

این مقاله به بیان تحقیقات صورت گرفته در راستای طراحی و پیاده سازی یک سیستم تصدیق هویت گوینده از راه دور و از طریق خطوط تلفنی می پردازد. در این سیستم، گوینده هویت خود را توسط یک کد شناسائی ۷ رقمی اعلام می نماید. یک الگوریتم بازشناسی ارقام مبتنی بر شبکه عصبی پیشگو که تلفیقی از شبکه عصبی چند لایه پرسپترون و الگوریتم برنامه ریزی پویا می باشد اقدام به بازشناسی ارقام کد شناسائی کاربر می نماید. با اعلام هویت گوینده، مدل مرجع آن گوینده که با استفاده از الگوریتم تلفیقی درخت برآمدگی و الگوریتم ژنتیکی آموزش دیده استخراج و با مقایسه همان گفتار گوینده در بیان کد شناسائی وی با این مدل مرجع، عمل تصدیق هویت صورت می گیرد. در آزمایشات صورت گرفته برای بازشناسی ارقام تأثیر تعداد دفعات ارائه داده های آموزشی به شبکه عصبی، تأثیر برش یا عدم برش خروجی گره ها برای جلوگیری از حالت اشباع، تأثیر تعداد ویژگیها و نوع ویژگی از نقطه نظر ایستا و گذرا بودن و نیز تأثیر ایجاد تغییرات در پارامترهای یادگیری شبکه مورد بررسی قرار گرفت. همچنین طی آزمایشهایی کارائی سیستم تصدیق هویت پیشنهاد شده ارزیابی و روشهای مختلفی برای تعیین سطح آستانه تصمیم گیری مورد بررسی و آزمایش قرار گرفت و نتایج حاصل از بازشناسی ارقام و تصدیق هویت با استفاده از روشهای گفته شده با نتایج حاصل از روش کلاسیک چندی سازی برداری مقایسه گردید.

کلمات کلیدی

تصدیق هویت گوینده، بازشناسی گفتار، شبکه عصبی پیشگو، برنامه ریزی پویا، شبکه عصبی درخت برآمدگی، چندی سازی برداری، الگوریتم های ژنتیکی، نرمال سازی گروهی

Personal Identification Number (PIN) Recognition and Speaker Verification for Access Control Over PSTN

M. Mehdi Homayounpour
Assistant Professor

Amir Najari
M.Sc. Student

Computer Engineering Department
Amirkabir University of Technology

Abstract

An implementation of a speaker verification system is presented in which the speaker claims his/her identity by a seven digits PIN over telephone. PIN digits are recognized by a hybrid predictive neural network system using multi-layer perceptron and dynamic programming algorithm. The speaker verification was done by comparing a speaker's model to his/her speech when uttering the seven digits PIN. The speaker's reference model was trained using Bumb tree and Genetic algorithm. The total number of epochs used to train the neural network, feature types and the dimension of feature vectors, and the effect of different neural net parameters were studied on Farsi digit recognition performance. Some other experiments were conducted to assess the performance of the proposed speaker verification method and to evaluate the effect of different types of decision threshold calculation on the speaker verification error rate.

Keywords

Speech Recognition, Speaker Recognition, Speaker Verification, Predictive Neural Network, Bumb Tree, Genetic Algorithm, Dynamic Programming, Cohort Normalization

امروزه شاهد فراگیر شدن فناوری پردازش خودکار گفتار در امور صنعتی، پزشکی، اداری و عمومی جامعه می‌باشیم. یکی از شاخه‌های مهم این فناوری تصدیق^۱ هویت گوینده می‌باشد. تصدیق هویت گوینده چنانچه از راه دور و از طریق سیستم‌های تلفنی، شبکه‌های مخابراتی، کامپیوتری و اینترنت صورت گیرد ضمن ایجاد تسهیلات فراوان، از جاذبه بالایی برخوردار می‌باشد. در یک مثال میتوان به موردی اشاره نمود که در آن فردی از محل کار خود یا از منزل با بانک محل حساب خود تماس گرفته و پس از آنکه هویت وی از طریق صدای وی تأیید شد، امکان انجام عملیات بانکی به وی داده شود. مثالی دیگر از اینگونه موارد، دسترسی به اطلاعات و بانکهای اطلاعاتی محرمانه از راه دور می‌باشد. تحقیقات بین‌المللی گسترده‌ای در حال حاضر در این زمینه در حال انجام میباشد که هدف از این تحقیقات عملیاتی و کاربردی نمودن این سیستم هاست. در مواردی نیز این سیستم‌ها به مرحله عملیاتی رسیده و در حال حاضر در مراکزی از قبیل بانکها و شرکتهای بیمه، مراکز کنترل پالایشگاه‌ها و نیروگاه‌ها، آزمایشگاه‌های تحقیقات استراتژیک و حتی بیمارستانها از سیستم‌های تصدیق هویت در امن نمودن دسترسی افراد مجاز به اطلاعات و اماکن خاص استفاده می‌شود. یقیناً امن نمودن دسترسی به اطلاعات و اماکن به گونه‌ای که نیاز به استفاده از کارت‌های مغناطیسی، کدهای محرمانه و استفاده از کلمات رمز که امکان گم شدن، جعل شدن و حتی فراموشی آنها وجود دارد را مرتفع نماید، در کشور ما نیز مد نظر بوده و خیلی زود میتواند همگانی شده و در فرهنگ مردم کشور ما وارد گردیده و موجبات راحتی و تسهیلات بیشتر در زندگی مردم را فراهم آورد.

تصدیق هویت گوینده شاخه‌ای از مبحث بازشناسی گوینده است که هدف از آن شناسایی افراد از طریق صدای آنها می‌باشد. سابقه تحقیقات در این زمینه به بیش از دو دهه قبل میرسد. بازشناسی گوینده میتواند با شنیدن، مشاهده اسپکتروگرام توسط افراد متخصص و نیز بصورت خودکار صورت گیرد. روش اول و دوم بیشتر جنبه‌های پلیسی و قانونی داشته و به منظور شناخت مجرمین و یا رفع سوء ظن از متهمین مورد استفاده قرار می‌گیرد. آنچه در اینجا مد نظر است استفاده از روشهای خودکار در

بازشناسی گوینده می‌باشد. در این روشها با دریافت صدای یک فرد از طریق میکروفون و رقمی نمودن آن اقدام به استخراج ویژگیهای^۲ از گفتار می‌نمائیم که می‌توانند در متمایز نمودن این فرد از سایر گویندگان مفید واقع شود. به کمک این ویژگیها، در مرحله آموزش مدلی از گوینده ساخته شده و در مراحل بازشناسی با مقایسه صدای گوینده با این مدل و تعیین میزان شباهت بین آندو، عمل بازشناسی صورت می‌گیرد. ویژگیهای متعددی چون چگونگی تغییرات فرکانس ارتعاش تارهای صوتی که اصطلاحاً گام^۳ نامیده میشود، فرکانسهای تشدید مجرای گفتار یا باصطلاح فرمانتها^۴، ضرائب طیف فوریه سیگنال گفتار، ضرائب پیشگویی خطی^۵ و مشتقات آن چون ضرائب انعکاسی^۶، ضرائب نسبت سطوح مقطع^۷ و ضرائب زوج خطوط طیفی^۸، ضرائب کپستروم حاصل از آنالیز فوریه در معیار مل^۹، ضرائب کپستروم حاصل از آنالیز پیشگویی خطی^{۱۰} و بسیاری ویژگیهای دیگر برای مشخص نمودن و بیان خصوصیات وابسته به گوینده موجود در گفتار مورد استفاده قرار گرفته و کارائی آنها ارزیابی شده است. این ویژگیها عموماً سعی در مدل نمودن خصوصیات مجرای گفتار و یا خصوصیات سیگنال تحریک ناشی از ارتعاش تارهای صوتی گوینده نموده و به تنهایی یا ترکیبی از آنها برای بازشناسی گوینده بکار رفته اند. کارائی این ویژگیها می‌تواند در محیط‌های نویزی و در مواردی که صدا از طریق خطوط تلفنی و مخابراتی منتقل میشود تحت الشعاع قرار گیرد [۱]. روشهایی چون لیفتر^{۱۱} نمودن میانگذر و نیز ویژگیهای دیگری چون RASTA و PLP پیشنهاد شده اند که تا اندازه‌ای با تاثیرات نویز و محدودیتهای حاصل از خطوط مخابراتی چون محدودیت پهنای باند، اضافه شدن نویز و مانند آن مقابله می‌نمایند. به منظور مدل نمودن گویندگان نیز روشهای متعددی پیشنهاد گردیده است که از آن جمله میتوان به روشهای در هم پیچیدن زمانی^{۱۲}، چندی سازی برداری^{۱۳}، مدل مخفی مارکوف^{۱۴}، روشهای آماری مرتبه د و^{۱۵} روشهای شبکه عصبی^{۱۶} [۲] و سیستم‌های فازی [۳] اشاره نمود.

در این مقاله سیستمی برای تصدیق هویت گویندگان از طریق خطوط تلفن پیشنهاد گردیده است. در سیستم پیشنهادی قبل از مرحله تصدیق هویت گوینده، ارقام کد شناسائی شخصی گوینده بازشناسی میشوند.

بازشناسی ارقام می پردازیم. بخش ۳ و ۴ نیز به ترتیب به تصدیق هویت گوینده و بیان نتیجه گیری اختصاص دارند.

۱- دادگان صدای FARSDIGITS1

در یک سیستم بازشناسی گفتار و گوینده، داده‌های مختلف بسته به کاربرد سیستم در غالب دادگان گفتار یا عبارتی پایگاههای داده صدا جمع‌آوری و به منظور آموزش و ارزیابی آن سیستم بکار گرفته میشوند. در این کار تحقیقاتی نیز به منظور ارزیابی سیستم تصدیق هویت و بخش بازشناسی کد شناسائی شخص گوینده و برآورد کارائی تکنیکهای مورد استفاده در طراحی آن و بدلیل در اختیار نبودن یک دادگان تلفنی استاندارد در زبان فارسی اقدام به طراحی و ضبط یک پایگاه داده بنام FARSDIGITS1 نمودیم. از آنجا که اهمیت سیستم‌های تصدیق هویت گوینده هنگامی بارز میشود که از آنها از راه دور و از طریق خطوط مخابراتی و تلفنی استفاده گردد، لذا دادگان مورد استفاده در این مقاله از طریق خطوط تلفن ضبط شده است. بدین معنا که گویندگان گفتار خود را در یک طرف ارتباط تلفنی بیان و در طرف دیگر دستگاه ضبط صدا، گفتار آنها را ضبط نموده است. صدای ضبط شده از طریق خطوط تلفن دارای محدودیت پهنای باند و نویز بیشتری بوده و تحت تاثیر خصوصیات کابل و کانالهای مخابراتی، اکو و مانند آن قرار میگیرد و در نتیجه با صدائی که مستقیماً و از طریق یک میکروفون و در محیط عاری از نویز و سرو صدا ضبط میشود کاملاً متفاوت است. خصوصیات صدای تلفنی باعث میشود که بازشناسی گفتار و گوینده به مراتب مشکلتر گردد.

برای تشکیل این دادگان از ۱۰۰ گوینده شامل ۶۱ گوینده مذکر و ۳۹ گوینده مؤنث با متوسط سنی ۳۰ سال خواسته شد تا ارقام صفر الی نه فارسی را طی ۱ الی ۲ جلسه با فاصله زمانی ۷ الی ۳۰ روز بیان نمایند. گفتار این گویندگان با استفاده از یک مدار واسط تلفنی و کارت صدا با فرکانس نمونه‌برداری ۱۱۰۲۵ هرتز و دقت ۱۶ بیت بازاء هر نمونه، نمونه‌برداری شد. گویندگان عمدتاً دارای لهجه تهرانی بودند. هر گوینده حداقل ده بار ارقام ۰ الی ۹ را بیان نموده است و هر تکرار این ارقام در یک فایل ذخیره شده است. ابتدا و انتهای ارقام موجود در هر فایل نیز با استفاده از اسپکتروگرام آن بصورت دستی تعیین و در نتیجه بازاء هر فایل صدا یک فایل برچسب ایجاد گردید که در این

بازشناسی ارقام و در حالت کلی تر، بازشناسی کلمات، یکی دیگر از شاخه‌های پردازش خودکار گفتار می‌باشد که کاربردهای متعددی چون ارائه فرمانهای صوتی برای آن متصور می‌باشد. در بازشناسی ارقام نیز به جای استفاده مستقیم از سیگنال گفتار ویژگی‌هایی از آن استخراج میگردد که حاوی اطلاعات مفهومی گفتار می‌باشد. اصولاً کلیه ویژگی‌هایی که برای بازشناسی گوینده در بالا بر شمردیم در بازشناسی ارقام مورد استفاده قرار میگیرند. لیکن ویژگی که بیش از همه تا کنون برای این منظور بکار گرفته شده است ضرائب کپستروم حاصل از آنالیز فوریه در معیار مل و مشتقات اول و دوم آنها میباشد که حتی در محیط‌های نویزی و تلفنی نیز کارائی خوبی از خود نشان داده‌اند. در بازشناسی ارقام نیز نیاز به استفاده از تکنیک‌هایی برای مدل نمودن هر رقم می‌باشد که کلیه تکنیک‌های مدل سازی که برای تصدیق هویت بر شمردیم برای این منظور نیز بکار گرفته شده‌اند.

اصولاً پیاده سازی یک سیستم بازشناسی ارقام و بازشناسی گوینده و نیز آموزش، تست و ارزیابی آن نیازمند استفاده از یک دادگان گفتار مناسب می‌باشد. از آنجا که ما در سیستم طراحی شده، بازشناسی ارقام فارسی و نیز تصدیق هویت با استفاده از گفتار مربوط به همین ارقام را صورت میدهیم و از طرفی مایل به استفاده از سیستم طراحی شده از طریق تلفن می‌باشیم، بدلیل آنکه یک دادگان تلفنی از ارقام فارسی در اختیار نداشتیم، اقدام به ضبط یک دادگان صدای تلفنی نمودیم.

در این مقاله از شبکه عصبی پیشگو در بازشناسی ارقام و از تلفیق درخت عصبی بر آمدگی و الگوریتم ژنتیکی در تصدیق هویت گوینده استفاده شده است. ویژگی مورد استفاده در بازشناسی ارقام ضرائب کپستروم حاصل از آنالیز فوریه بر مبنای معیار مل و ویژگی مورد استفاده در تصدیق هویت ضرائب کپستروم حاصل از آنالیز پیشگوئی خطی می‌باشند. روشهای مختلف تعیین سطح آستانه و تصمیم گیری مورد ارزیابی قرار گرفته‌اند. روش چندی سازی برداری نیز، هم برای بازشناسی ارقام و هم برای تصدیق هویت گوینده مورد ارزیابی و نتایج حاصل از این روش با نتایج روش‌هایی که قبلاً بیان شد، مقایسه و مورد تحلیل قرار گرفته‌اند.

در بخش ۱ این مقاله ابتدا به توصیف دادگان گفتار ضبط شده و در بخش ۲ به شرح روش پیشنهادی در

هم پیچیدن زمانی می باشد. در ادامه این بخش پس از پیاده سازی شبکه عصبی پیشگو میزان کارایی آنرا با استفاده از دادگان ارقام تلفنی ضبط شده، مورد ارزیابی قرار می دهیم.

۲-۱- مدل شبکه عصبی پیشگو

در سالهای اخیر شبکه عصبی و بخصوص شبکه عصبی چند لایه پرسپترون در زمینه های مختلف پردازش اطلاعات بویژه بازشناسی گفتار مورد استفاده قرار گرفته است [۴] و [۵]. از خصوصیات مهم اینگونه شبکه ها می توان به الگوریتم یادگیری مؤثر بنام انتشار خطا به عقب^{۱۷} [۴] و همچنین قابلیت آنها در نگاشت ورودی به خروجی در فضای مسئله اشاره نمود. شبکه عصبی چند لایه پرسپترون که در اینجا به اختصار شبکه عصبی نامیده می شود در این مقاله به عنوان طبقه بندی کننده کلمات موجود در مجموعه کلمات یا لغتنامه^{۱۸} که شامل ارقام ۰ الی ۹ می باشد بکار رفته است. روشی که در اینجا بررسی می گردد نگرش دیگری به مسئله بازشناسی ارقام با استفاده از شبکه های عصبی چند لایه است. در این روش از مدل پیشگوی عصبی استفاده می شود. در این مدل شبکه های عصبی چند لایه بعنوان پیشگویی کننده نمونه ها بکار گرفته می شوند. این مدل شامل یک دنباله شبکه عصبی چند لایه پرسپترون جهت پیشگویی غیرخطی هر کلاس می باشد. در واقع این مدل از ساختار زمانی گفتار جهت بازشناسی آن بهره می گیرد، چرا که ارتباط زمانی بین بردارهای ویژگی متوالی در گفتار از آنجا که شامل اطلاعات مهمی در بازشناسی گفتار است حائز اهمیت زیادی می باشد. در این روش تغییرات سرعت در بیان گفتار با بکارگیری الگوریتم برنامه ریزی پویا نرمالیزه [۶] می شود. با توجه به شکل زیر، شبکه پیشگوی عصبی با استفاده از بردارهای ویژگی a_{i-T} تا a_{i-1} که قبل از بردار ویژگی a_i قرار دارند، بردار ویژگی a_i را تخمین می زند. در اینجا T معرّب تعداد بردارهای ویژگی قبلی است که در تخمین بردار a_i مورد استفاده قرار گرفته اند. می توان مقدار واقعی a_i را با مقدار تخمینی آن یعنی \hat{a}_i مقایسه نمود و اختلاف آنها که خطای پیشگویی می باشد را بدست آورد یعنی:

$$e = \left\| \hat{a}_i - a_i \right\|^2 \quad (1)$$

فایل علاوه بر مشخصات کلی مربوط به ضبط صدای گوینده، اطلاعات مربوط به شروع و انتهای موقعیت ارقام نیز قرار داده شد. به منظور استفاده از این دادگان در سیستم تصدیق هویت گوینده، به هر گوینده یک رشته ۷ رقمی تصادفی از ارقام به عنوان کد شناسائی شخصی او تخصیص داده شده است که با استفاده از فایل های برچسب و فایل های صدای متناظر آنها و با انتخاب و کنار هم گذاشتن ارقام مربوط به کد ۷ رقمی مورد نظر از فایل های حاوی ارقام ۰ تا ۹، این کد ساخته می شود. یعنی مثل اینکه گوینده این کد را دقیقاً به همان صورت بیان نموده باشد. به این ترتیب این امکان فراهم گردید که برای هر گوینده کد شناسائی او که توسط خود او بیان شده باشد و نیز کد شناسائی او که توسط سایر گوینده ها بیان شده باشد را تولید نمائیم.

۲- بازشناسی ارقام مستقل از گوینده

توسط شبکه عصبی پیشگو

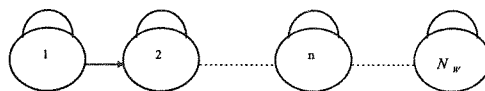
قبل از این گفتیم که هر گوینده به منظور معرفی خود، کد شناسائی شخصی خود را که یک کد ۷ رقمی است بیان می نماید. یکی از وظایف این سیستم آن است که این کد شناسائی شخصی را بازشناسی نموده و ارقام آنرا تعیین نماید. لذا برای این منظور به یک سیستم بازشناسی ارقام فارسی ۰ تا ۹ نیاز خواهیم داشت که در این بخش به بیان الگوریتم های بکار رفته برای این منظور پرداخته و نتایج حاصله را بیان خواهیم نمود. لازم بذکر است که بازشناسی ارقام در اینجا بصورت مستقل از گوینده صورت گرفته و فرض بر آن است که گوینده ارقام را بصورت مجزا و جدا از هم و از طریق تلفن بیان می نماید. بازشناسی ارقام بیان شده بصورت ناپیوسته شامل مراحل متعددی است که عمده ترین آنها تعیین محدوده ارقام، استخراج ویژگی ها و مدل نمودن آنها می باشد. تعیین محدوده کلمات (ارقام) و الگوریتم مورد استفاده برای این روش می بایست به صورت خودکار صورت گیرد. از آنجا که در این بخش هدف ارزیابی الگوریتم ارائه شده در بازشناسی ارقام می باشد. لذا تعیین محدوده کلمات بصورت دستی صورت گرفته تا بدین ترتیب اشتباهات ناشی از تعیین خودکار محدوده ارقام، به حساب الگوریتم بازشناسی گذاشته نشود. در زیر به بیان مدل شبکه عصبی پیشگو در بازشناسی ارقام فارسی می پردازیم. شبکه عصبی پیشگو تلفیقی از شبکه عصبی چند لایه پرسپترون و برنامه ریزی پویا بروش در

$$d_w(t, n) = \|\hat{a}_t(w, n) - a_t\|^2 \quad (3)$$

که e بنام مانده پیشگویی معروف است. از مانده پیشگویی می‌توان بعنوان تابع خطا جهت اصلاح وزنهای عصبی پیشگو استفاده نمود. این موضوع در قسمت الگوریتم آموزش بیشتر توضیح داده خواهد شد.

۲-۲- الگوریتم بازشناسی کلمات

در روش پیشگویی عصبی چند لایه پرسپترون، مدل هر لغت شامل دنباله ای از پیشگوهای عصبی می‌باشد. برای مثال لغت w با این روش بصورت زیر مدل می‌شود:



شکل (۱) مدل هر لغت در روش ارائه شده.

هر گره خود یک پیشگوی عصبی بوده و تعداد این گره‌ها برای این لغت N_w می‌باشد. این شکل در واقع نحوه تخصیص بردارهای ویژگی یک گفتار را به هر پیشگوی عصبی نشان می‌دهد. گفتار ورودی توسط مدل لغت w به N_w بخش تقسیم بندی و هر بخش از آن به یک پیشگوی عصبی داده می‌شود و بعنوان مثال قسمت n ام گفتار توسط پیشگوی n ام تخمین زده می‌شود. چنانچه $n(t)$ بیانگر این باشد که کدام پیشگوی عصبی مدل لغت w برای تخمین بردار t بکار میرود، دنباله $\{n(1), n(2), \dots, n(T)\}$ یک مسیر را روی صفحه محاسبات برنامه ریزی پویا نشان می‌دهد. شرط:

$$n(t) = n(t-1) \text{ or } n(t-1) + 1 \text{ و } n(N) = N_w \text{ و } n(1) = 1$$

در سایر نقاط مسیر می‌بایست ارضاء گردند. مسئله بازشناسی یک کلمه را می‌توان بصورت یافتن دنباله‌ای از $\{n^*(t)\}$ تعریف نمود، بقسمی که مانده پیشگویی کل حداقل شود. با توجه به شرایط و محدودیت‌های مذکور و الگوریتم برنامه ریزی پویا میتوان فرمول بازگشتی زیر را جهت حداقل سازی مانده پیشگویی کل ارائه نمود:

$$g_w(t, n) = d_w(t, n) + \min[g_w(t-1, n), g_w(t-1, n-1)] \quad (2)$$

که بصورت زیر تعریف می‌شود: $d_w(t, n)$ اندازه فاصله محلی (یا مانده پیشگویی) است

۲-۳- الگوریتم آموزشی مدل کلمات

منظور از آموزش مدلها یافتن وزنهای پیشگوهای عصبی است بقسمی که مانده پیشگویی کل بازای مجموعه نمونه‌های آموزشی حداقل شود. تابع هدف بصورت متوسط مانده پیشگویی کل تمامی تکرارهای لغت w بصورت زیر:

$$D(w) = \frac{1}{M_w} \sum_{m=1}^{M_w} D(w, m) \quad (4)$$

تعریف می‌شود به گونه ای که M_w تعداد دفعات تکرار لغت w و $D(w, m)$ مانده پیشگویی کل به ازای تکرار m ام لغت w می‌باشد. الگوریتم آموزش برای بهینه‌سازی تابع هدف، ترکیبی از الگوریتم برنامه ریزی پویا و الگوریتم آموزش شبکه عصبی چند لایه یعنی انتشار خطا به عقب می‌باشد که آنرا بصورت زیر بیان می‌نمائیم:

۱- وزنهای پیشگوهای عصبی را بصورت تصادفی انتخاب می‌کنیم.

۲- مراحل زیر را برای تمامی دفعات تکرار لغت w یعنی M_w به گونه‌ای که $(1 \leq m \leq M_w)$ اجرا می‌کنیم.

۳- مانده پیشگویی کل یعنی $D(w,m)$ را با استفاده از الگوریتم برنامه‌ریزی پویا محاسبه و نیز مسیر بهینه $\{n^*(t)\}$ را با تکنیک بازگشت به عقب معین می‌نمائیم.

۴- وزنه‌های هر پیشگوی عصبی را توسط الگوریتم انتشار خطا به عقب بر روی مسیر بهینه $\{n^*(t)\}$ اصلاح می‌کنیم. در اینجا خروجی مطلوب a_t به خروجی واقعی $a_t(w, n(t))$ در $n(t)$ امین پیشگوی عصبی تخصیص داده می‌شود.

معیار تصحیح وزنها، کاهش مانده پیشگویی کل $D(w,m)$ در امتداد مسیر بهینه $\{n^*(t)\}$ در صفحه محاسبات برنامه‌ریزی پویا میباشد. در نقطه $\{t, n^*(t)\}$ در صفحه مذکور مانده پیشگویی یعنی:

$$e = \|\hat{a}_t(w, n^*(t)) - a_t\|^2 \quad (5)$$

را می‌توان بعنوان تابع خطای انتشار به عقب در نظر گرفت که باید کاهش یابد و بنابراین اصلاح وزنها توسط قانون انتشار خطا به عقب، که در آن خروجی مطلوب a_t متناظر خروجی واقعی $a_t(w, n^*(t))$ در $n^*(t)$ امین پیشگوی عصبی می‌باشد، انجام می‌پذیرد.

۲-۴- استخراج ویژگی و آموزش مدل‌های ارقام

گفتار ۵۸ نفر از گویندگان (۳۶ نفر مذکر و ۲۲ نفر مونث) از دادگان ارقام فارسی FARSDIGITS که در بخش ۱ شرح داده شد، برای انجام بررسی‌های مربوط به بازشناسی ارقام استفاده شده است. این مجموعه ۵۸ گوینده‌ای به دو مجموعه یکی شامل ۵۰ گوینده برای آموزش و دیگری شامل ۸ گوینده برای ارزیابی تقسیم گردید. به منظور ساختن مدل ارقام، گفتار مربوط به بیان هر رقم بد ۲۰ فریم مساوی تقسیم و از هر فریم پس از پیش تاکید^۲ و ضرب در پنجره همینگ^۱، یک بردار ویژگی شامل ۱۲ پارامتر کپستروم بر مبنای معیار مل MFCC استخراج گردید. در هر بردار، ضرائب کپستروم به ۱/۱ برابر بزرگترین آنها نرمالیزه شدند.

همانگونه که قبلا اشاره شد در این مقاله مدل نمودن ارقام توسط مدل پیشگوی عصبی صورت گرفته است. تعداد پیشگوهای عصبی در مدل ارائه شده برای تمامی ارقام برابر ۱۰ انتخاب گردید. ساختار شبکه عصبی مورد استفاده در مدل پیشگوی فوق دارای سد لایه با 2×12 گره در لایه

ورودی، ۴ گره در لایه مخفی و ۱۲ گره در لایه خروجی تعیین گردید. تعداد گره‌های لایه ورودی و خروجی به اندازه‌ای انتخاب شده اند که با قرار دادن بردارهای ویژگی خروجی دو فریم در ورودی شبکه بتوان بردار ویژگی فریم بعدی را پیشگوئی نمود. مدل هر رقم با استفاده از ویژگی‌های بدست آمده از ۲ تکرار از هر گوینده که جمعا ۱۰۰ تکرار را شامل میشود آموزش داده شد. داده‌های آموزشی فوق طی ۱۰۰ تکرار در اختیار شبکه عصبی پیشگو قرار گرفت و بدین ترتیب مدل هر رقم بدست آمد.

۲-۵- ارزیابی روش بازشناسی کد شناسائی شخصی

بعد از آموزش مدل کلیه ارقام ۰ تا ۹، اقدام به ارزیابی میزان کارائی روش ارائه شده در بازشناسی ارقام نمودیم. راندمان بازشناسی برای نمونه‌های آموزشی ۹۵/۶٪ درصد بدست آمد. از طرفی طی آزمایش صورت گرفته بر روی داده‌های آزمایشی شامل ۱۰ تکرار از ارقام ۰ تا ۹ که توسط ۸ گوینده بیان شده بودند و سیستم هیچ آشنائی قبلی با آنها نداشت، راندمان ۸۳/۷٪ بدست آمد. طی آزمایشات دیگری اقدام به افزایش تعداد گره‌های لایه مخفی، تغییر تعداد تکرارها، برش زدن خروجی شبکه و تغییر پارامترهای شبکه عصبی شامل η و ضریب مومنت α نمودیم. آزمایشات صورت گرفته گویای نکات زیر می‌باشد:

الف- با افزایش تعداد گره‌های لایه مخفی اگرچه خطای بازشناسی نمونه‌های آموزشی افزایش می‌یابد اما در عین حال خطای بازشناسی نمونه‌های آزمایشی کاهش می‌یابد. علت این امر تنظیم بیشتر وزنه‌های شبکه متناسب با نمونه‌های آموزشی و در عوض کاهش تعمیم‌پذیری آن می‌باشد. جدول ۱ نتایج آزمایش با تعداد لایه مخفی مختلف را نشان می‌دهد.

جدول (۱) نتایج بازشناسی با توجه به تعداد لایه‌های مخفی.

تعداد لایه مخفی	داده‌های آموزشی	داده‌های آزمایشی
۲	۹۰/۴٪	۷۷/۵٪
۴	۹۵/۶٪	۸۳/۷٪
۶	۹۶/۲٪	۷۱/۲٪
۹	۹۶/۸٪	۷۵/۰٪

ب- افزایش بیش از حد تعداد تکرار آموزش اگرچه باعث دقیقتر شدن پیشگویی‌ها می‌شود اما بدلیل آموزش بیش از

حد شبکه، تعمیم پذیری مدل کاهش می‌یابد. جدول ۲ نتایج بازشناسی را با توجه به تعداد تکرار آموزش نشان می‌دهد.

جدول (۲) نتایج بازشناسی با توجه به تعداد تکرار آموزش.

تعداد تکرار ها	داده های آموزشی	داده های آزمایشی
۵۰	۹۰/۴٪	۷۷/۵٪
۱۰۰	۹۵/۶٪	۸۳/۷٪
۲۰۰	۹۶/۲٪	۷۱/۲٪

ج- جهت افزایش یادگیری سیستم و جلوگیری از حالت اشباع، خروجی گرہها را برش زدیم [۷]. میدانیم که وقتی تابع سیگموئید برای تابع فعالیت انتخاب شود داریم:

$$f(x) = 1/(1+\exp(-x)); f'(x) = f(x) [1 - f(x)] \quad (۶)$$

بنابراین وقتی $f(x)$ نزدیک 0 و یا 1 شود، مشتق آن بسمت صفر میل می‌کند در نتیجه یادگیری کاهش می‌یابد. برای جلوگیری از این مسئله مقدار تابع فعالیت گرہ های لایه مخفی به فاصله $[0/1, 0/9]$ محدود شدند. جدول ۳ نتایج بازشناسی را با توجه به برش و عدم برش خروجی گرہها نشان می‌دهد. چنانچه در این جدول دیده می‌شود نتایج تقریباً مانند هم می‌باشد که این امر به این خاطر است که در شرایط آموزش با این داده های آموزشی، $f(x)$ به 0 و یا 1 نزدیک نشده است. این نتایج بازاء ۱۰۰ بار تکرار بدست آمده است.

جدول (۳) نتایج بازشناسی با توجه به برش خروجی گرہها.

برش	داده های آموزشی	داده های آزمایشی
خیر	۹۴/۴٪	۸۳/۳٪
بلی	۹۵/۶٪	۸۳/۷٪

د- برای بررسی نقش تعداد ویژگیهای استفاده شده، آزمایشات دیگری انجام شد. در این آزمایشات، تعداد ویژگیهای برابر ۱۰ و ۱۲ انتخاب شدند. نتایج حاصل در جدول ۴ برای ۱۰۰ تکرار آورده شده است. این جدول نشان می‌دهد که افزایش ویژگیها در افزایش راندمان بازشناسی مؤثر است.

جدول (۴) نتایج بازشناسی با توجه به تعداد ویژگیها.

تعداد ویژگیها	داده های آموزشی	داده های آزمایشی
۱۰	۹۴/۱٪	۸۳/۳٪
۱۲	۹۵/۶٪	۸۳/۷٪

ه- ارزیابی دیگری نیز به منظور بررسی نقش ویژگیها و مشتق اول آنها که گویای اطلاعات گذرا و دینامیک گفتار می باشند صورت گرفت. برای این منظور دو آزمایش انجام شد که در آنها تعداد ویژگیها مساوی اختیار شدند، با این تفاوت که در آزمایش اول از ۱۶ ضریب کپستروم و در آزمایش دوم ۸ ضریب کپستروم بعلاوه ۸ مشتق اول آنها استفاده گردید. این آزمایشات یک بار بازاء ۲۰۰ تکرار و بار دیگر بازاء ۵۰۰ تکرار صورت گرفتند. جدول ۵ نتایج این بررسی را نشان می‌دهد. این نتایج گویای آن است که استفاده از مشتق اول ضرائب کپستروم یا عبارتی استفاده از اطلاعات دینامیک موجود در گفتار موجب بهبود بازشناسی ارقام میگردد.

جدول (۵) نتایج بازشناسی با توجه ویژگیهای ایستا و گذرا.

ویژگی	تعداد تکرار	داده های آموزشی	داده های آزمایشی
8MFCC+8 ΔMFCC	۲۰۰	۹۵/۷٪	۸۲/۹٪
8MFCC+8 ΔMFCC	۵۰۰	۹۶/۱٪	۸۰/۸٪
16 MFCC	۲۰۰	۹۵/۰٪	۷۸/۷٪
16 MFCC	۵۰۰	۹۷/۹٪	۷۵/۸٪

و- مسئله مهم دیگر پارامترهای ثابت یادگیری η و ضریب مومنت α می‌باشد که در یادگیری مؤثر می‌باشند. مقادیر مختلفی را برای η و α بصورت سعی و خطا آزمایش کردیم و نهایتاً مناسب ترین راندمان بازاء $\eta = 0.01$ و $\alpha = 0.6$ بدست آمد.

۳- تصدیق هویت گوینده

قبلاً اشاره شد که در یک سیستم تصدیق هویت، پس از اعلام هویت توسط گوینده، میزان شباهت صدای این گوینده و مدل او تعیین و سپس با یک سطح آستانه مقایسه شده و نهایتاً نسبت به رد یا قبول او اقدام میگردد. این بخش اختصاص دارد به ارائه روشی برای مدل کردن گویندگان که از روش چندی سازی برداری و شبکه های عصبی استفاده می‌نماید. لارم بذکر است که آموزش این مدل توسط الگوریتم های ژنتیکی صورت میگیرد. همچنین در این بخش به انواع خطاها و نیز موضوع تعیین سطح آستانه تصمیم گیری و نیز ارزیابی سیستم های تصدیق

هویت گوینده خواهیم پرداخت و نتایج حاصل از آزمایشات صورت گرفته را ارائه خواهیم کرد.

۳-۱- مدل شبکه عصبی درختی

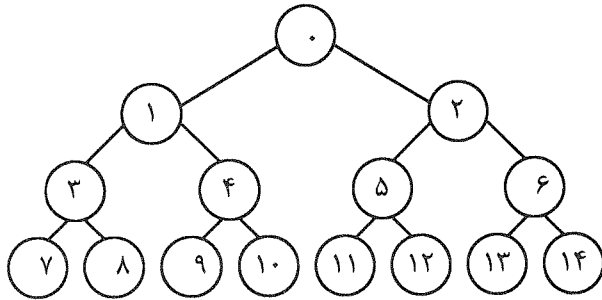
در زمینه بازشناسی الگو، الگوریتمهای مختلفی جهت خوشه‌بندی داده‌ها وجود دارد که عمدتاً در یک سطح داده‌ها را به خوشه‌های مختلف تقسیم بندی می‌نمایند. در این مقاله روشی مورد استفاده قرار می‌گیرد که در آن خوشه‌بندی و طبقه‌بندی داده‌ها بصورت سلسله مراتبی صورت می‌گیرد. تفاوت عمده این روش سلسله مراتبی با سایر الگوریتمهای خوشه‌بندی، سرعت سریع آن جهت یافتن خوشه دربرگیرنده داده ورودی می‌باشد. آموزش این روش بر اساس الگوریتمهای ژنتیکی و شبکه عصبی بنا شده است.

یادگیری بخش مهمی از سیستم‌های هوشمند را تشکیل می‌دهد. یادگیری انواع مختلفی دارد که از آن جمله می‌توان به یادگیری بدون نظارت و یادگیری با نظارت اشاره نمود. در یادگیری بدون نظارت، در فاز آموزش داده‌های آموزشی به خوشه‌هایی تقسیم می‌شوند بنحوی که داده‌های موجود در هر یک از این خوشه‌ها بیشترین شباهت را به هم داشته باشند. البته باید توجه داشت که تعریف شباهت بسته به کاربرد مشخص می‌شود. در فاز آزمایش نیز می‌بایست تعلق داده ورودی را به یکی از این خوشه‌ها مورد بررسی قرار داد. الگوریتمهای مختلفی جهت خوشه‌بندی داده‌ها وجود دارد که عموماً در یک سطح داده‌ها را به خوشه‌های مختلف طبقه‌بندی می‌نمایند. یک سطحی بودن خوشه‌بندی به این معنا می‌باشد که برای یافتن خوشه در برگیرنده داده ورودی لازم است تمامی خوشه‌های موجود در سیستم بررسی شوند. بنابراین جستجو برای یافتن خوشه در این سیستم از مرتبه $O(n)$ می‌باشد. از این الگوریتمها می‌توان به الگوریتم LBG، الگوریتم K-means و الگوریتم حداقل واریانس^{۲۲} اشاره نمود. تک سطحی بودن خوشه‌بندی موجب اتلاف وقت در هنگام بازیابی می‌شود. برای افزایش سرعت دسترسی به خوشه‌های موجود در سیستم می‌توان آنها را بصورت سلسله مراتبی بنحوی سازماندهی نمود که ساختار حاصل شکل یک درخت دودویی تبدیل گردد. در این صورت زمان بازیابی و یافتن خوشه‌ها بصورت لگاریتمی و درجه سیستم از مرتبه $O(n)$ به $O(\log_2 n)$ کاهش خواهد یافت. روش

خوشه‌بندی که در اینجا مورد استفاده قرار می‌گیرد بر اساس ایده فوق و درخت برآمدگی^{۲۳} که در آن از توزیع‌های نرمال جهت خوشه‌بندی سلسله مراتبی استفاده می‌شود بنا شده است. فاز آموزش این روش در واقع یافتن پارامترهای توزیع‌های نرمال این درخت می‌باشد. در این بررسی جهت یافتن پارامترهای بهینه از الگوریتمهای ژنتیکی استفاده شده است.

۳-۱-۱- شبکه درختی برآمدگی

ساختار درخت دودویی زیر را در نظر بگیرید:



شکل (۲) یک درخت باینری کامل. شماره گذاری گره‌های این درخت به این صورت است که ریشه دارای شماره صفر و فرزندان چپ و راست هر گره با شماره n به ترتیب $2n+1$ و $2n+2$ می‌باشند.

این درخت از یک ریشه، تعدادی گره میانی و تعدادی برگ در آخرین سطح تشکیل شده است. متناظر هر گره میانی تابعی به شکل زیر اعمال میشود [۸]:

$$y = \prod_j \left[\frac{1}{\sqrt{2\pi\sigma_j}} e^{-\left(\frac{x_j - \mu_j}{2\sigma_j^2}\right)^2} \right] \quad (7)$$

در این تابع y حاصلضرب توابع نرمال با میانگین μ_j و واریانس σ_j متناظر عضو زام بردار ورودی می‌باشد. فضای تصمیم‌گیری چنین گره‌ای یک ابر بیضی می‌باشد که قسمتی از فضای ورودی گره را محصور می‌نماید. درخت برآمدگی یک ساختار مشابه درخت باینری فوق است که شرط زیر را نیز دارا می‌باشد [۹]:

"مقدار تابع هر گره میانی باید در تمامی نقاط روی فضای ورودی از توابع متناظر شاخه‌های زیرین خود بزرگتر باشد."

$Index = 0$ /* this is root index */

While (Not isaLeaf (Index))

Begin

$$y_{left}^{Index} = \prod_j \left[\frac{1}{\sqrt{2\pi\sigma_{left,j}^{Index}}} e^{-\frac{(x_j - \mu_{left,j}^{Index})^2}{2\sigma_{left,j}^{Index}}} \right]$$

$$y_{right}^{Index} = \prod_j \left[\frac{1}{\sqrt{2\pi\sigma_{right,j}^{Index}}} e^{-\frac{(x_j - \mu_{right,j}^{Index})^2}{2\sigma_{right,j}^{Index}}} \right]$$

If ($y_{left}^{Index} > y_{right}^{Index}$) /* Select Left Branch */

Index = Index = 2 * Index + 1

Else /* Select Right Branch */

Index = Index = 2 * Index + 2

EndIf

End While

Leaf = Index

$$z^{Leaf} = \sum_j w_j^{Leaf} x_j$$

If ($z^{Leaf} > Threshold$)

$x \in Leaf$

(8)

Else

$x \notin Leaf$

EndIf

۳-۱-۲- آموزش شبکه درختی

برای آموزش شبکه درختی برآمدگی اولین قدم تعیین ساختار درخت یعنی تعداد لایه‌ها، میانگین و واریانس گره‌های میانی می‌باشد. در مرحله دوم لازم است که شبکه‌های خطی متناظر هر برگ آموزش ببینند. روشهای مختلفی برای تعیین ساختار درخت وجود دارد که در اینجا از الگوریتمهای ژنتیک جهت یافتن ساختار بهینه استفاده شده است. روش دیگر خوشه‌بندی بازگشتی است که در آن داده‌های آموزشی با استفاده از الگوریتمهای خوشه‌بندی معمولی بطور بازگشتی به دو قسمت تقسیم می‌گردند و آنگاه بعد از هر تقسیم بندی عمل خوشه‌بندی با توجه به اینکه باید شامل توزیع‌های نرمال باشند، اصلاح می‌گردند. جهت آموزش شبکه‌های خطی از قانون پرسپترون و تنها در

دو گره فرزند گره ریشه باعث می‌شوند که منطقه محصور توسط گره ریشه به دو بخش تقسیم گردد. این عمل تقسیم بصورت بازگشتی برای هر گره تکرار می‌شود تا منطقه داده‌های ورودی به زیر مناطق کوچکتر تقسیم گردند. در واقع هدف درخت برآمدگی این است که فضای ورودی بطور مکرر به مناطق کوچکتر تقسیم بندی شود تا اینکه نهایتاً در سطح برگها عمل طبقه بندی داده‌ها بصورت جداپذیر خطی قابل انجام باشد. در سطح برگها، متناظر هر برگ یک شبکه عصبی خطی وجود دارد که وظیفه طبقه بندی نقاطی که در منطقه آن برگ قرار می‌گیرند را بر عهده دارد. وقتی داده ورودی X به شبکه درختی اعمال شود، شاخه‌های قویتر و در نهایت برگ فعالتر انتخاب می‌شود و در این گره است که عمل طبقه بندی با توجه به خروجی شبکه خطی صورت می‌گیرد. این خروجی توسط رابطه $z^{leaf} = \sum_j w_j^{leaf} x_j$ بدست می‌آید که در آن

z^{leaf} خروجی نام برگ شماره leaf، w_j^{leaf} نام وزن این برگ و x_j عضو نام بردار ورودی می‌باشد. این روند بصورت شبه کد در روابط ۸ آورده شده است. این الگوریتم در واقع مرحله بازشناسی و بازیابی می‌باشد. نکته مهم در انتخاب تابع گره‌های میانی در این است که با انتخاب توابع گره‌های میانی دیگر، شبکه‌هایی با خصوصیات و کاربردهای دیگر حاصل می‌شود. برای مثال، شکل ۳ توابع دیگر موجود برای گره‌های میانی را نشان می‌دهد. انتخاب هر یک از این توابع، شبکه‌های درختی دیگر را باعث می‌شود. از انواع دیگر این نوع درختها می‌توان به Oct-Tree، Kd-Tree، BoxTree و BallTree اشاره نمود [۹]. تابع قسمت (B) در شکل ۳ در واقع همانی است که در شبکه درختی برآمدگی مورد استفاده قرار می‌گیرد. همانطور که دیده می‌شود این تابع بصورت یک برآمدگی^{۲۴} است و به این دلیل نام این نوع شبکه‌های به نام شبکه درخت برآمدگی مشهور است.



شکل (۳) توابع مختلف برای انتخاب بعنوان توابع گره‌های میانی. انتخاب تابع قسمت (A) باعث بوجود آمدن درختهای Oct_Tree، Kd-Tree و BoxTree می‌شود. توابع قسمت (B) و (C) در درختهای BumpTree و BallTree مورد استفاده قرار می‌گیرد. تابع قسمت (D) در درخت (D) Kd-Tree کارایی بالایی ایجاد می‌کند [۹].

یک مرحله برای حداقل سازی خطای داده‌های آموزشی استفاده می‌کنیم.

۳-۱-۳- تعیین ساختار بینه شبکه درختی

برآمدگی توسط الگوریتم های ژنتیک

گام اساسی در بکارگیری الگوریتمهای ژنتیکی یافتن کدینگ مناسب جهت کد کردن اطلاعات می‌باشد زیرا الگوریتمهای ژنتیکی نه بر پارامترها بلکه بر کد شده آنها اعمال می‌شوند. در واقع می‌توان گفت میزان کارایی الگوریتمهای ژنتیکی بستگی زیادی به نحوه کد کردن دارد. درخت شکل ۲ را مجدداً در نظر بگیرید. این درخت را می‌توان با یک آرایه با طول ثابت ۱۵ بصورت شکل ۴ نشان داد.

۰	۱	۲	۳	۴	۵	۶	...	۱۱	۱۲	۱۳	۱۴
---	---	---	---	---	---	---	-----	----	----	----	----

شکل (۴) نمایش بلاکی درخت باینری شکل ۲. هر بلاک معرف یک گره درخت است. هر بلاک شامل اطلاعات گره متناظر در درخت باینری نیز می‌باشد. ترتیب شماره‌گذاری هر یک از این بلاکها مانند شماره‌گذار بیهای گره‌ها در درخت باینری می‌باشد. در علم ژنتیک به کل بلاکها یک کروموزوم و به هر یک از اطلاعات موجود در هر بلاک یک ژن گفته می‌شود.

هر بلاک (خانه) در شکل فوق معرف یک گره درخت و هر ژن (مقدار) در هر بلاک معرف یک پارامتر آن گره مثلاً میانگین و یا واریانس می‌باشد. بلاک شماره صفر بیانگر ریشه درخت می‌باشد. فرزندان یک گره شماره n در بلاکهای شماره 2n+1 و 2n+2 قرار می‌گیرند. مسئله مهمی که در کدینگ شبکه درختی برآمدگی وجود دارد وجود وابستگی بین پارامترهای گره‌های مختلف می‌باشد. برای مثال با توجه به شرط اصلی این درختها باید پارامترهای گره‌های فرزند در فضای تعریف شده توسط گره پدر قرار گیرند. بنابراین پارامترهای گره‌های فرزند بایستی با توجه به گره‌های پدر انتخاب شوند. استفاده از یک کدینگ معمولی که در آن باید هر پارامتر مستقل از پارامترهای دیگر ارائه‌گردد، موجب می‌شود که با اعمال عملگرهای ژنتیکی (جهش و ترکیب) ساختارهایی بوجود آیند که شرط اصلی شبکه‌های درختی برآمدگی را دارا نباشند. این امر حتی با انتخاب نسل اول از شبکه‌های معتبر با گذشت زمان در نسلهای آتی بوقوع خواهد پیوست. برای رفع این مشکل

روش کدینگ زیر معرفی می‌شود. قبل از معرفی این روش لازم است که نحوه تعریف فضای هر گره مشخص شود. تابع میانی درخت در معادله ۷ را می‌توان بصورت زیر نیز نشان داد:

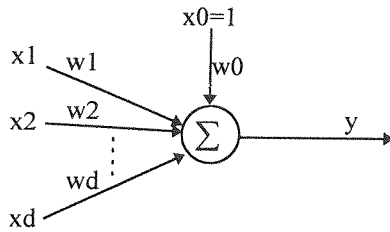
$$y = \frac{1}{\prod_j \sqrt{2\pi\sigma_j}} e^{-\left(\sum_j \frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right)} \quad (9)$$

برای معرفی فضای تعریف شده توسط این گره از توان قسمت نمایی استفاده می‌کنیم یعنی:

$$S = \frac{(x_1 - \mu_1)^2}{2\sigma_1^2} + \frac{(x_2 - \mu_2)^2}{2\sigma_2^2} + \frac{(x_3 - \mu_3)^2}{2\sigma_3^2} + \dots + \frac{(x_d - \mu_d)^2}{2\sigma_d^2} \quad (10)$$

مشاهده می‌شود که فضای S یک ابربیضی به مختصات $(\mu_1, \mu_2, \mu_3, \dots, \mu_d)$ و شعاع $(\sqrt{2}\sigma_1, \sqrt{2}\sigma_2, \sqrt{2}\sigma_3, \dots, \sqrt{2}\sigma_d)$ می‌باشد که ابعاد فضا فرض شده است. بنابراین بین شعاع این ابر بیضی و واریانسهای ورودی رابطه مستقیم وجود دارد. بهمین دلیل شعاعهای ابربیضی و واریانسها گاهاً بجای هم مورد استفاده قرار می‌گیرند. در کدینگ ارائه شده، هر ژن یک مقدار حقیقی در فاصله (+۱ و -۱) می‌باشد. ژن در علم ژنتیک در واقع یک فیلد اطلاعاتی مثلاً میانگین می‌باشد. نحوه تبدیل یک کروموزوم به یک ساختار درخت برآمدگی به شرح زیر است [۸]:

قبل از هر چیز لازم است فضایی که شبکه در آن تعریف می‌شود مشخص گردد. این فضا بصورت یک ابر بیضی به مرکز میانگین داده‌های ورودی و شعاعی که تمامی داده‌ها را در ابعاد مختلف شامل گردد تعریف می‌شود. حال مرکز سیستم مختصات خود را به مرکز این ابر بیضی انتقال داده و در هر بعد، بردار یکه را معادل شعاع در آن بعد انتخاب می‌کنیم. فضا و سیستم مختصات فوق توسط ریشه درخت معین می‌گردند و ژنهایی که مرکز دو گره فرزند ریشه را کد می‌نمایند برحسب این مختصات بیان می‌شوند و همچنین ژنهایی که شعاعهای این دو گره را کد می‌کنند بر روی مرکز مختصات پدر و قطر بیضی فضای تعریف



شکل (۶) نمایش ورودی‌ها و خروجی یک شبکه پرسپترون تک لایه.

بنابراین $y = \sum_{i=1} w_i x_i + w_0$ که با فرض $x_0 = 1$ بطور خلاصه خواهیم داشت:

$$y = \sum_{i=0} w_i x_i \quad (11)$$

فرض می‌کنیم نمونه‌های آموزشی $\{x^p\}$ و خروجی‌های مطلوب $\{Y^p\}$ متناظر این نمونه‌ها باشد. خطای موجود در شبکه برابر است با:

$$E = 0.5 \sum_p (y^p - Y^p)^2 \quad (12)$$

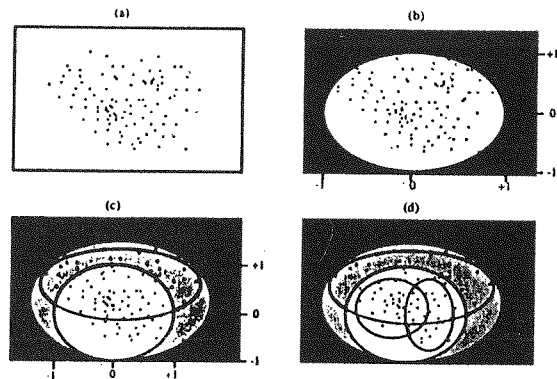
که y^p خروجی واقعی نمونه آموزشی x^p می‌باشد.

$$E = 0.5 \sum_p \left(\sum_i w_i x_i^p - Y^p \right)^2 \quad (13)$$

برای حداقل نمودن خطا باید داشته باشیم:

$$\begin{aligned} \frac{\partial E}{\partial w_j} &= 0 \quad \forall j, j=0,1,\dots,d \\ \frac{\partial E}{\partial w_j} &= \frac{\partial}{\partial w_j} \left[E = 0.5 \sum_p \left(\sum_i w_i x_i^p - Y^p \right)^2 \right] \\ &= \sum_p \left[\left(\sum_i w_i x_i^p - Y^p \right) x_j^p \right] \\ \frac{\partial E}{\partial w_j} &= \sum_i \left[w_i \left(\sum_p x_i^p x_j^p \right) \right] - \sum_p x_j^p Y^p \\ \frac{\partial E}{\partial w_j} &= 0 \Rightarrow \sum_i \left[w_i \left(\sum_p x_i^p x_j^p \right) \right] = \sum_p x_j^p Y^p \end{aligned} \quad (14)$$

توسط آن نگاشت می‌شوند و این روند بصورت بازگشتی برای تمامی گره‌ها انجام می‌پذیرد. این روند بصورت گرافیکی در زیر ترسیم شده است:



شکل (۵) نحوه تبدیل یک کروموزوم به ساختار درخت برآمدگی در فضای دو بعدی. (a) - فضای اولیه داده‌های آموزشی. (b) - فضای نرمال شده و همچنین سیستم مختصات تعریف شده توسط گره ریشه. (c) و (d) - نحوه تقسیم فضای تعریف شده هر گره توسط گره‌های فرزند آن.

با توجه به این کدینگ عملگرهای الگوریتمهای ژنتیکی، دیگر ساختارهای نامعتبر درست نمی‌کنند و همچنین عملگرهای جهش و ترکیب براحتی قابل تعریف می‌باشند. عملگر ترکیب به این صورت تعریف می‌شود که برای ترکیب دو ساختار یک گره بصورت تصادفی اختیار شده آنگاه زیرشاخه‌های متناظر آن گره در دو ساختار درختی با هم جابجا می‌شوند. در عمل جهش یک ساختار درختی یک گره بصورت تصادفی انتخاب شده و پارامترهای آن با احتمال کم با عدد تصادفی در فاصله $(-0.2, +0.2)$ جمع می‌شود. بعد از اعمال عملگر جهش مقادیر پارامترها بین $(-1, +1)$ برش می‌شوند تا کدینگ دچار مشکل نشود. نشان داده شده است که این روش کدینگ دارای خواص پیوستگی، هم ریختی، کامل بودن، بسته بودن و در برداشتن حداقل اطلاعات تکراری می‌باشد.

۳-۱-۴- آموزش شبکه‌های خطی متناظر برگها

در هر برگ یک پرسپترون ساده وجود دارد که بصورت زیر نمایش داده می‌شود:

برای حل معادلات خطی ۱۴ متغیرهای زیر تعریف می‌شوند:

$$M_{i,j} = \sum_p x_i^p x_j^p \quad (15)$$

با این متغیرها داریم:

$$B_j = \sum_p x_i^p Y^p \quad (16)$$

و یا بصورت ماتریسی داریم:

$$\sum_i w_i M_{i,j} = B_j \quad (17)$$

$$WM = B \quad (18)$$

که W^* بصورت زیر محاسبه می‌شود:

$$W^* = BM^+ \quad (19)$$

M^+ شبه معکوس M می‌باشد زیرا ممکن است M معکوس پذیر نباشد. روش ارائه شده در اینجا اولاً قادر به آموزش یک مرحله‌ای شبکه خطی می‌باشد و ثانیاً به حافظه کمی نیاز دارد [۱۰].

۲-۳- انواع خطاها، تصمیم‌گیری و ارزیابی سیستم‌های تصدیق هویت گوینده

قبل از آنکه به بیان توضیحات بیشتر در خصوص اقدامات صورت گرفته برای تصدیق هویت با استفاده از شبکه عصبی و الگوریتم‌های ژنتیکی بپردازیم لازم است ابتدا به شرح مختصری از چگونگی ارزیابی سیستم‌های تصدیق هویت و نحوه تعیین راندمان و کارایی آنها بپردازیم. برای ارزیابی هر سیستم لازم است ماهیت خطاها و اشتباهات موجود در آن سیستم بررسی گردد. در فناوری تصدیق هویت گوینده نیز این ضرورت بدیهی بنظر می‌رسد. همچنانکه در تعریف مسئله تصدیق هویت گوینده لازم است که میزان نزدیکی گفتار ورودی به مدل گوینده ادعا شده معین گردد. اصولاً در تصدیق هویت گوینده یک سطح آستانه جهت اخذ تصمیم در رابطه با پذیرش یا رد

ادعای گوینده لازم است که می‌بایست آنرا در مرحله آموزش سیستم تعیین نمائیم. در این بخش ضمن توضیح انواع خطاهای ممکن در سیستم‌های تصدیق هویت گوینده، به روشهای بدست آوردن سطح آستانه تصمیم‌گیری و چگونگی تصمیم‌گیری جهت رد یا قبول گوینده و تعیین میزان کارایی سیستم در تصدیق هویت گوینده‌ها اشاره خواهیم کرد.

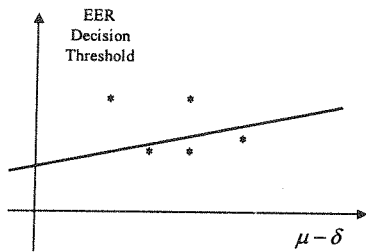
۳-۲-۱- انواع خطاها

در فناوری تصدیق هویت گوینده دو نوع خطا ممکن است رخ دهد که عبارتند از:
- رد اشتباه^{۲۵} یا باختصار FR: رد یک گوینده مجاز یا رد نادرست
- پذیرش اشتباه^{۲۶} یا باختصار FA: قبول یک گوینده غیر مجاز یا تصدیق نادرست

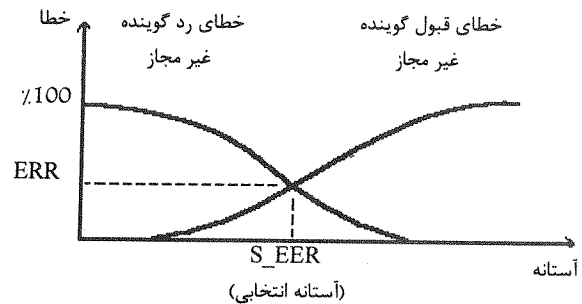
۳-۲-۲- تعیین سطح آستانه

یک مرحله اساسی در پیاده‌سازی سیستم‌های تصدیق هویت گوینده تعیین مقدار آستانه جهت تصمیم‌گیری می‌باشد. یک مقدار آستانه پائین باعث خطای رد اشتباه بالاتر و درعین حال خطای پذیرش اشتباه پائینتر می‌شود. مسلماً بسته به نوع کاربرد مقدار این سطح آستانه متفاوت می‌باشد. مثلاً در کاربردهای نظامی بسیار مطلوب است که نرخ خطای رد بالا باشد یعنی گرچه تعدادی گویندگان مجاز اشتباهاً رد می‌شوند اما درعین حال گویندگان غیر مجاز به سختی ممکن است پذیرفته شوند که این خود باعث افزایش امنیت سیستم می‌شود. اما در کاربردهای دیگر مثلاً یک سیستم بانکی مطلوب است که درصد تصدیق هویت گویندگان مجاز بالا باشد، هر چند ممکن است که درصد قبولی گویندگان غیر مجاز نیز زیاد شود.

روشهای مختلفی جهت پیدا نمودن سطح آستانه‌ای ارائه شده است. یک روش معمول برای انجام اینکار، روش نرخ خطای برابر^{۲۷} و یا بطور مختصر EER می‌باشد. در این روش مقدار آستانه‌ای طوری انتخاب می‌شود که درصد پذیرش اشتباه برابر با درصد رد اشتباه گردد. شکل زیر نحوه یافتن این سطح آستانه را نشان می‌دهد:



شکل (۸) مقادیر آستانه تصمیم‌گیری EER برحسب میانگین منهای واریانس فواصل برون‌گوینده ای نظیر آنها.



شکل (۷) یافتن آستانه تصمیم‌گیری به روش نرخ خطای برابر.

حال با استفاده از روشهای عددی برازش خط، خط زیر را بگونه‌ای از بین نقاط فوق عبور می‌دهیم که خطای مربوط به نمایش این نقاط توسط این خط حداقل شود. به همین دلیل منبع این روش را روش برازش خط می‌نامیم:

$$thresh = c1 * (\mu - \delta) + c2 \quad (20)$$

در واقع با استفاده از روشهای عددی ضرایب $c1$ و $c2$ بدست می‌آیند. با این روش برای بدست آوردن سطح آستانه تصمیم‌گیری یک گوینده نامشخص کافی است که ابتدا فواصل بین‌گویندگی او محاسبه شده و μ ، δ و در نتیجه $\mu - \delta$ این فواصل بدست آید. آنگاه با استفاده از رابطه ۲۰، مقدار آستانه تصمیم‌گیری برای آن گوینده بدست می‌آید. نکته قابل ذکر این است که می‌توان اصلاحیه‌ای که بر روش EER اعمال شد، برای این روش نیز بکار برد یعنی اینکه مقدار آستانه حاصل را در ضریبی ضرب نمود تا کارایی سیستم افزایش یابد.

۳-۲-۳- نرمالسازی گروهی فواصل

همانگونه که قبلاً گفته شد روش کلاسیک برای تصدیق هویت یک گوینده آن است که فاصله بین گفتار وی با مدل گفتار گوینده ادعا شده اندازه‌گیری و سپس با مقایسه این فاصله با یک سطح آستانه تصمیم‌گیری گوینده مورد نظر پذیرفته و یا رد شود. مزیت استفاده از چنین اندازه‌گیری جهت تنظیم سطح آستانه برای تصمیم‌گیری راجع به رد یا قبول هویت گوینده این است که این امتیازدهی انحراف و اختلاف واقعی بین مدل گوینده مجاز و گفتار ورودی را نشان می‌دهد. اما این امتیازدهی خام برای هم گوینده مجاز و هم گویندگان غیر مجاز بخاطر تغییرات درون‌گویندگی، تغییرات محیط ضبط و یا تغییرات فونتیکی ممکن است تغییرات بسیار شدیدی نماید و بنابراین بصورت توزیع‌هایی

بنابراین برای یافتن سطح آستانه، به ازای هر گوینده مجاز نمودار شکل ۷ تشکیل و از روی آن سطح آستانه تصمیم‌گیری برای آن گوینده معین می‌گردد. یک اصلاحیه بر این روش این است که برای پیشی‌بینی وضعیتهای آینده، بجای استفاده از EER، آن را در یک ضریب که معمولاً بزرگتر از ۱ است ضرب می‌کنند [۱۱]. اگر این ضریب ۱ باشد، مقدار نهایی همان EER خواهد بود. اگر این مقدار از ۱ بزرگتر باشد، سطح آستانه بطرف راست انتقال پیدا می‌کند که در اینصورت خطای قبول گوینده غیر مجاز زیاد ولی خطای رد گوینده مجاز کم می‌شود. این عمل برای مواردی که شرایط محیط در هنگام بکارگیری سیستم تغییرات زیادی می‌کند، مثلاً روی خطوط تلفن بسیار مفید می‌باشد.

در روش ارائه شده لازم است که نمودار شکل ۷ برای بدست آوردن سطح آستانه تک تک گویندگان رسم شود و اگر به سیستم گوینده جدیدی اضافه گردد این روند برای بدست آوردن سطح آستانه تصمیم‌گیری وی نیز باید تکرار شود. در ادامه روشی ارائه می‌شود که در آن با توجه به اطلاعاتی که از تعداد محدودی گوینده آموزشی بدست می‌آوریم بتوانیم مقادیر سطوح آستانه را برای گویندگان جدید بدست آوریم [۱۲]. در این روش بازاء هر گوینده آموزشی، مدل وی با تکرارهایی از گفتار او و نیز تکرار هائی از گفتار سایر گوینده‌های آموزشی مقایسه و فواصل درون گوینده ای^{۲۸} و بین گوینده ای^{۲۹} وی بدست می‌آیند. به کمک این فواصل EER و نیز تفاضل میانگین و انحراف استاندارد فواصل بین گوینده ای یعنی $\mu - \delta$ را بدست می‌آوریم. بازاء مقدار EER و $\mu - \delta$ هر گوینده، یک نقطه که توسط * مشخص شده در نمودار شکل ۸ بدست می‌آید. این کار را برای کلیه گوینده‌های آموزشی صورت می‌دهیم.

آماری مثلا میانگین و یا حداقل مقادیر در نظر گرفته شده است [۱۱].

$$\begin{aligned} & \text{if } S_c < kS_{\max} \\ & \text{Score}_N = S_c - \text{stat}[S_i, i = 1, \dots, n] \\ & \text{else} \\ & \text{Score}_N = \infty \end{aligned} \quad (21)$$

که S_c میزان پراکندگی گفتار ورودی با مدل گوینده مجاز، k ضریب سطح آستانه مطلق، S_{\max} حداکثر میزان پراکندگی حاصل از مقایسه گفتار آموزشی گوینده مجاز با مدل خودش، Score_N امتیاز نرمال شده گروهی، Stat تابع نرمالسازی گروهی مثلا میانگین یا حداقل، S_i میزان پراکندگی حاصل از مقایسه گفتار ورودی با مدل (کتابچه کد) نامین عضو گروه این گوینده و n تعداد اعضای گروه گوینده می باشد.

۳-۲-۴- ارزیابی سیستم

برای ارزیابی سیستمهای تصدیق هویت گوینده لازم است کارایی سیستم بصورتی بیان شود تا بتوان سیستمهای مختلف را تحت یک معیار واحد با هم مقایسه نمود. آقایان Haton و Gong نرخ متوسط خطای تصدیق و نرخ خطای تصدیق کل را بصورت زیر پیشنهاد نمودند [۱۲].

$$E_M = \frac{FA + FR}{2} \quad (22)$$

که M تعداد گویندگان موجود در سیستم می باشد [۱۲]. E_M که در رابطه فوق تعریف شده متوسط حسابی می باشد. تعریف E_G بصورت متوسط هندسی نیز بصورت زیر پیشنهاد گردیده است:

$$E_G = \sqrt{FA * FR} \quad (23)$$

۳-۳- تصدیق هویت گوینده توسط تلفیق

درخت برآمدگی و الگوریتم ژنتیکی

در این بخش به تشریح روش تلفیقی حاصل از ترکیب درخت برآمدگی و الگوریتم ژنتیکی برای انجام

با هم پوشانی زیاد مثل آنچه در شکل ۷ روی داده است در آیند. در این شرایط، نرخ خطای یکسان و یا EER زیاد شده و باعث می شود که کارایی تصدیق هویت گوینده کاهش یابد.

یک روش دیگر نگرش به مشکل فوق استفاده از سطح آستانه پویا می باشد که در آن از شباهت نرمال شده گروهی 30 و یا امتیازدهی پراکندگی استفاده می شود. در این روش یک گروه از گویندگان که به گوینده مجاز بسیار نزدیک می باشند، بعنوان اعضای گروه وابسته به آن گوینده در نظر گرفته می شوند. تشکیل گروه را می توان به معنی ایجاد یک محیط محلی در فضای تمامی گویندگان دانست که به کمک آن اندازه های فاصله یا شباهت را تا اندازه ای نرمالیزه می نمائیم. ساده ترین شکل این نرمال سازی محاسبه تفاضل زیر می باشد:

$A = \text{فاصله بین گفتار ورودی و مدل گفتار گوینده مورد نظر}$

$B = \text{تابعی از فاصله بین گفتار ورودی و مدل های}$

$\text{گویندگان موجود در محیط محلی گوینده مورد نظر}$

$A - B = \text{فاصله نرمالیزه شده}$

این روش هم در سیستم مستقل از متن و هم وابسته به متن بکار گرفته شده است. برای باز هم بهتر شدن روند تصمیم گیری میتوان ترکیبی از دو روش سطح آستانه پویا (روش نرمال سازی گروهی) و روش معمول استفاده نمود [۱۱]. بدین صورت که قبل از عمل نرمال سازی گروهی ابتدا از یک سطح آستانه مطلق استفاده می شود که سطح آستانه با استفاده از روش EER محاسبه می گردد. ابتدا فاصله واقعی بین گفتار ورودی و گفتار مدل مرجع بدست آمده و اگر این مقدار از سطح آستانه کوچکتر بود، آنگاه از نرمالسازی گروهی جهت بدست آوردن فاصله جدید استفاده می شود و در غیر اینصورت گوینده گفتار ورودی سریعاً رد هویت می شود. بعد از نرمالسازی گروهی، فاصله جدید با سطح آستانه دیگر جهت تصدیق و یا رد نهایی مقایسه می شود.

در اینجا مثالی برای استفاده از الگوریتم فوق را برای حالتی که هر گوینده توسط یک کتابچه کد 21 حاصل از تکنیک چندی کردن برداری مدل شده است را ارائه می کنیم. لازم به ذکر است که از اصلاحیه روش EER یعنی ضرب یک مقدار ثابت در مقدار نرخ خطای برابر نیز استفاده شده است و همچنین تابع نرمالسازی گروهی یک تابع

تصدیق هویت گوینده و توصیف نتایج حاصل از این روش می‌پردازیم. همانطور که در ابتدای مقاله آوردیم پس از آنکه یک گوینده کد شناسائی شخصی خود که یک کد هفت رقمی است را بیان نمود ابتدا با استفاده از تکنیک گفته شده در قسمت بازشناسی ارقام، کد شناسائی شخصی وی بازشناسی و مشخص می‌شود. قدم بعدی آن است که هویت این گوینده که خود را به گونه ای که بیان شد معرفی نموده است، بررسی گردد. این عمل بدون آنکه گوینده صحبت اضافه دیگری را بیان نماید، با استفاده از همان گفتار اولیه وی در بیان کد شناسائی شخصی، صورت می‌گیرد.

انجام تصدیق هویت گوینده به مراحل متعددی نیاز دارد که شامل استخراج ویژگیهای مناسب از گفتار گوینده، مدل نمودن گوینده، تعیین سطح آستانه تصمیم‌گیری جهت پذیرش یا رد گوینده و نهایتاً مقایسه گفتار گوینده با مدل فرد ادعا شده و قبول یا رد وی می‌باشد. به منظور آموزش و ارزیابی روش ارائه شده از دادگان گفتاری که طی این تحقیق ضبط گردیده و توضیحات مربوط به آن قبل از این ارائه شده است استفاده نمودیم. از گفتار ۵۸ گوینده (۳۶ نفر مرد و ۲۲ نفر زن) این دادگان استفاده گردید. گویندگان ارقام فارسی صفرالی نه را ده بار تکرار و از طریق تلفن ضبط شده اند. با استفاده از این ارقام به هر گوینده یک کد ۷ رقمی منحصر به فرد اختصاص داده شد و بدین ترتیب کد شناسائی شخصی وی ایجاد گردید. مجموعه گفتارهای هر گوینده به دویخش ۶ تایی برای آموزش و ۴ تایی برای تست تقسیم شد. برای آموزش مدل هر گوینده علاوه بر گفتار وی از اطلاعات گفتار گویندگان همجنس او نیز استفاده شد. بدین منظور دو- سوم گویندگان همجنس وی بصورت اتفاقی انتخاب و همراه با اطلاعات وی جهت آموزش سیستم بکار گرفته شدند. در مرحله آزمایش، ۴ تکرار باقیمانده هر گوینده و همچنین یک- سوم گوینده های همجنس باقیمانده و همین تعداد گوینده غیر همجنس مورد استفاده قرار گرفتند. برای آموزش مدل یک گوینده و نیز انجام تستهای تصدیق هویت به کمک این مدل، از گفتار مربوط به بیان کد هفت رقمی توسط این گوینده و سایر گویندگان استفاده گردید.

سیگنال گفتار ورودی مربوط به بیان کد شناسائی گوینده به فریم های ۳۰ میلی‌ثانیه‌ای که ۲۰ میلی ثانیه همپوشانی دارند تقسیم شده و هر فریم تحت آنالیز LPC

مرتب ۱۲ قرار گرفت و از آن ۱۲ ویژگی LPCC استخراج گردید. هر گوینده توسط یک شبکه درختی ۵ لایه کامل شامل ۳۱ گره مدل گردید. تعداد لایه های فوق بدین دلیل انتخاب شدند که اولاً تعداد لایه های بیشتر موجب می‌گردد که بدلیل محدودیت داده های آموزشی، مدل درختی هر گوینده شامل برگهائی تهی از داده آموزشی گردیده و در نتیجه درختهای ناقص ایجاد شود، ثانیاً تعداد لایه های کمتر پوشش ضعیفتری از اطلاعات یادگیری در فضای داده های آموزشی ایجاد می‌نمود.

آموزش درخت برآمدگی به روشهای مختلفی میسر است که یک روش آن توسط Bostock شرح داده شده است [۱۴]. در این روش درخت بصورت بالا به پایین و بازگشتی با استفاده از الگوریتمهای خوشه‌بندی ساخته و آموزش داده می‌شود. روشهای دیگری نیز وجود دارند که بتوان با استفاده از آنها درخت برآمدگی را بنحو بهتری آموزش داد که یکی از آنها استفاده از الگوریتم ژنتیکی می‌باشد که در این تحقیق بمنظور بهینه سازی فرآیند یادگیری درخت برآمدگی مورد استفاده قرار گرفته است. همچنانکه در بخش الگوریتم های ژنتیک گفته شد این الگوریتم ها بعنوان تکنیکی برای بهینه سازی مورد استفاده قرار میگیرند و سعی در یافتن حداقل مطلق در فضای مسئله می‌نمایند. و این در حالی است که در روش ارائه شده توسط بوستاک یافتن یک حداقل محلی مورد نظر است، گرچه هیچ تضمینی در یافتن حداقل مطلق توسط این روش وجود ندارد. مشخصات الگوریتم ژنتیکی مورد استفاده در این تحقیق عبارت است از:

- تعداد تکرار (نسل): ۵۰ تکرار

- جمعیت هر نسل: ده درخت

- نرخ ترکیب: ۷۰٪

- نرخ جهش: ۱۰٪

- تابع معیار: معکوس میانگین مربعات خطاهای شبکه‌های

خطی

بدین ترتیب با استفاده از درخت برآمدگی و الگوریتم ژنتیکی با خصوصیات گفته شده مدلهای گویندگان ساخته شد. در مرحله تست هویت یک گوینده، بردار ویژگی حاصل از هر فریم گفتار ورودی به درخت اعمال شده و پس از پیمایش از ریشه به برگ در بر گیرنده آن، خروجی متناظر با آن برگ محاسبه گردید. این خروجی برای کلیه فریم‌های گفتار کد شناسائی ۷ رقمی گوینده محاسبه و

جدول (۶) مقایسه راندمان روش شبکه عصبی درختی و روش چندی سازی برداری با تعیین آستانه برونش نرخ خطای برابر. (الف) - راندمان شبکه عصبی.

شبکه عصبی درختی		مقدار ضریب
آموزشی	آزمایشی	
۹۵/۹٪	۸۰/۷٪	۰/۹
۹۹/۵٪	۸۹/۲٪	۱
۹۸/۷٪	۹۳/۹٪	۱/۱
۹۶/۰٪	۹۲/۱٪	۱/۲

(ب) - راندمان چندی سازی برداری با کتاب کد ۳۲ عضوی.

چندی سازی برداری (۳۲)		مقدار ضریب
آموزشی	آزمایشی	
۹۰/۷٪	۷۲/۵٪	۰/۹
۹۹/۷٪	۹۰/۴٪	۱
۹۸/۲٪	۹۶/۶٪	۱/۱
۹۳/۸٪	۹۵/۱٪	۱/۲

(ج) - راندمان چندی سازی برداری با کتاب کد ۶۴ عضوی.

چندی سازی برداری (۶۴)		مقدار ضریب
آموزشی	آزمایشی	
۹۷/۲٪	۷۲/۳٪	۰/۹
۹۹/۶٪	۸۹/۰٪	۱
۹۹/۱٪	۹۷/۸٪	۱/۱
۹۶/۰٪	۹۶/۳٪	۱/۲

در آزمایش دیگر بجای استفاده از EER از روش یافتن خط برازش $y = c1 * (\mu - \sigma) + c2$ جهت تعیین آستانه تصمیم‌گیری استفاده شد. نتایج حاصل از تصدیق هویت گوینده با استفاده از درخت برآمدگی و چندی سازی برداری در جدول ۷ آورده شده است. در این جدول مشاهده می‌شود که راندمان روش شبکه عصبی درختی در بهترین حالت بازنه ضریب ۱/۲ برابر ۹۱/۰٪ و راندمان روش چندی سازی برداری بازنه ضریب ۱/۲ و ۹۶/۹٪ و ۹۶/۰٪ بترتیب بازنه ۳۲ و ۶۴ خوشه بدست می‌آید. در آزمایش دیگری سعی شد تصمیم‌گیری نسبت به رد یا قبول گوینده با استفاده از روش نرمالیزه نمودن گروهی صورت گیرد. برای این منظور ۵ گوینده که در فاز آموزش نزدیکترین فاصله را با هر گوی و مومنده داشتند بعنوان گروه آن گوینده انتخاب شدند. برای تعیین فاصله نهایی یک گفتار با مدل هر گوینده، از فواصل گفتار ورودی با تک

میانگین‌گیری شد. این میانگین در واقع میانگینمیزان اختلاف بین مدل و گفتار ورودی می‌باشد. حال با در اختیار داشتن میانگین فوق و با توجه به آنچه که در بخش تصمیم‌گیری تصدیق هویت گوینده بیان نمودیم، با مقایسه این میانگین با سطح آستانه تصمیم‌گیری گوینده مورد نظر میتوان نسبت به رد یا قبول وی اقدام نمود.

برای مقایسه میزان کارائی روش فوق با کارائی روشهای متداول در تصدیق هویت گوینده، آزمایش دیگری صورت دادیم. در این آزمایش، از تکنیک چندی سازی برداری-k means برای مدل نمودن گویندگان استفاده کردیم. هر گوینده با استفاده از ۳۲ و ۶۴ خوشه مدل گردید. در این روش به منظور تصدیق هویت یک گوینده گفتار وی با مدل گوینده ادعا شده مقایسه و فاصله بین بردار ویژگی هر فریم با مرکز ثقل نزدیکترین خوشه در مدل فوق بدست آمد و نهایتاً میانگین فواصل بدست آمده بازنه کلیه فریمها در گفتار تست محاسبه گردید. این میانگین در واقع میزان اختلاف بین مدل گوینده ادعا شده و گفتار ورودی می‌باشد. با مقایسه این میانگین با سطح آستانه تصمیم‌گیری، میتوان نسبت به رد یا قبول گوینده تصمیم‌گیری نمود. نکته مهمی که در تصدیق هویت گوینده وجود دارد تعیین مناسب سطح آستانه تصمیم‌گیری می‌باشد. در آزمایشات فوق از روش نرخ خطای برابر و یا EER استفاده شده که این خطا نظیر محل تلاقی منحنی تغییرات خطاهای FR و FA می‌باشد. همانطور که در بخش بررسی خطاهای تصدیق هویت بیان نمودیم به منظور در نظر گرفتن تغییرات در محیط ضبط صدای گوینده در هنگام تست و همچنین تغییرات تدریجی صدای گوینده با گذشت زمان، در عمل بجای استفاده مستقیم از EER آنرا در یک ضریب بزرگتر ولی نزدیک به یک ضرب می‌نمایند. در جدول ۶ بازنه مقادیر مختلفی برای این ضریب نتایج حاصل از تصدیق هویت با استفاده از درخت برآمدگی و نیز روش چندی کردن برداری با تعداد ۳۲ و ۶۴ برای حجم کتابچه کد، آورده شده است. در این جدول مشاهده می‌شود که راندمان روش شبکه عصبی درختی در بهترین حالت بازنه ضریب ۱/۱ برابر ۹۳/۹٪ و راندمان روش چندی سازی برداری بازنه ضریب ۱/۱ و ۹۶/۶٪ و ۹۷/۸٪ بترتیب بازنه ۳۲ و ۶۴ خوشه بدست می‌آید. افزایش تعداد خوشه‌ها از ۳۲ خوشه به ۶۴ خوشه تغییر قابل توجهی در نتایج حاصل ایجاد نمی‌نماید.

جدول (۸) نتایج تصدیق هویت گوینده با استفاده از نرمال سازی گروهی و نرخ خطای برابر.

مقدار ضریب	داده های آموزشی	داده های آزمایشی
۰/۹	۶۴/۱٪	۶۵/۰٪
۱	۷۳/۴٪	۶۸/۳٪
۱/۱	۷۳/۲٪	۷۱/۶٪
۱/۲	۶۸/۹٪	۶۹/۰٪

جدول (۹) نتایج تصدیق هویت گوینده با استفاده از نرمال سازی گروهی و برازش خط.

مقدار ضریب	داده های آموزشی	داده های آزمایشی
۰/۹	۷۰/۷٪	۶۷/۳٪
۱	۷۲/۸٪	۶۷/۲٪
۱/۱	۶۷/۵٪	۶۷/۰٪
۱/۲	۶۰/۴٪	۶۴/۰٪

ع- نتیجه گیری

آزمایشات صورت گرفته جهت بازشناسی ارقام با تلفیق شبکه عصبی پیشگو و برنامه ریزی پویا نشان داد که افزایش بیش از حد تعداد گره ها در لایه مخفی و نیز ارائه بیش از حد داده های آموزشی به مدل در هنگام آموزش موجب گرایش بیش از حد مدل به داده های آموزشی و کاهش قدرت تعمیم پذیری آن و در نتیجه کاهش راندمان بازشناسی میگردد. در خصوص نوع ویژگیها و تعداد آنها نیز معلوم گردید که که افزایش تعداد ویژگیها تا حد معقول در بهبود راندمان مؤثر میباشد، ضمن آنکه اطلاعات دینامیک گفتار موجود در مشتق ضرائب کپسترال میتواند در بهبود کارائی تاثیر مثبت داشته باشد.

روش سلسله مراتبی جهت خوشه بندی داده ها در تصدیق هویت گوینده از مزیت سرعت بازیابی نسبتاً زیاد در مقایسه با دیگر روشهای خوشه بندی برخوردار است. مزیت دیگر آن سازگاری جهت ادغام با روشهای ژنتیکی می باشد. اگرچه شبکه های عصبی چند لایه نیز این قابلیت را دارند اما کندی سرعت آنها باعث کاهش توجه به این قابلیت آنها می شود. لازم به ذکر است که شبکه های عصبی چند لایه برای آموزش از قاعده انتشار به عقب استفاده می کنند که در این الگوریتم لازم است تکرارهای بسیاری جهت یافتن پارامترهای بهینه صورت گیرد، در صورتیکه همچنان که گفته شد در روش شبکه های عصبی درختی تنها در یک مرحله و با سرعتی بالا این امر میسر می باشد.

تک گویندگان هم گروه گوینده مرجع میانگین گیری شد و این میانگین بعنوان فاصله نهایی در نظر گرفته شد. با توجه به این فرضها، نتایج تصدیق هویت با استفاده از روش نرمالیزه نمودن گروهی برای روش درخت برآمدگی و روشهای تعیین سطح آستانه نرخ خطای برابر و برازش خط در جداول ۸ و ۹ آورده شده است. در جدول ۸ مشاهده می شود که با استفاده از نرمال سازی گروهی و روش تعیین سطح خطای برابر در بهترین حالت بازاء ضریب ۱/۱ کارایی برابر ۷۱/۶٪ بدست می آید. همچنین با توجه به جدول ۹ نیز مشاهده می شود که با استفاده از نرمال سازی نتایج گروهی و روش تعیین سطح آستانه بکمک برازش خط، در بهترین حالت بازاء ضریب ۰/۹ راندمان بدست می آید. این نتایج نشان میدهد که استفاده از روش نرخ خطای برابر یا EER در مقایسه با روش برازش خط راندمان بالاتری را نتیجه می دهد. این نتیجه میتواند بدین صورت قابل توضیح باشد که سطح آستانه بدست آمده در روش برازش خط با توجه به شکل ۸ مناسبترین سطح آستانه تصمیم گیری را برای آنکه خطای تصدیق هویت حداقل باشد، نتیجه نمی دهد.

جدول (۷) مقایسه روش شبکه عصبی درختی و روش چندی سازی برداری با تعیین آستانه بروش برازش خط.

الف- راندمان شبکه عصبی

شبکه عصبی درختی		مقدار ضریب
آموزشی	آزمایشی	
۹۳/۰٪	۷۹/۲٪	۰/۹
۹۷/۲٪	۸۴/۱٪	۱
۹۸/۹٪	۸۹/۷٪	۱/۱
۹۸/۲٪	۹۱/۰٪	۱/۲

ب- راندمان چندی سازی برداری با کتاب کد ۳۲ عضوی

چندی سازی برداری (۳۲)		مقدار ضریب
آموزشی	آزمایشی	
۸۸/۶٪	۷۵/۴٪	۰/۹
۹۸/۸٪	۸۷/۵٪	۱
۹۸/۸٪	۹۵/۱٪	۱/۱
۹۵/۲٪	۹۶/۰٪	۱/۲

ج- راندمان چندی سازی برداری با کتاب کد ۶۴ عضوی

چندی سازی برداری (۶۴)		مقدار ضریب
آموزشی	آزمایشی	
۹۴/۳٪	۷۵/۱٪	۰/۹
۹۹/۴٪	۸۸/۵٪	۱
۹۹/۲٪	۹۴/۶٪	۱/۱
۹۶/۹٪	۹۶/۹٪	۱/۲

نتایج حاصل از این بررسی بنوعی روشنگر نکات دیگری نیز می باشد که مهمترین آنها عبارتند از:

- سرعت بازیابی خوشه ها در این روش سریعتر از دیگر روشهای خوشه بندی می باشد.
- زمان آموزش این روش بدلیل بکارگیری الگوریتمهای ژنتیکی که عموماً زمانبر می باشند نسبت به سایر روشها طولانی تر می باشد.
- از مقایسه آزمایش این روش با آزمایش چندی سازی برداری معمولی چنین برمی آید که کارایی این روش کمتر از روش چندی سازی برداری بروش k-means می باشد اما مزیت آن همچنانکه پیشتر شرح داده شد سرعت بسیار بالاتر این روش در مقایسه با روشهای معمول مثل روش k-means می باشد. با فرض آنکه ۳۲ خوشه داشته باشیم، در خوشه بندیهای معمولی لازم است که ۳۲ خوشه بازای هر بردار ورودی مقایسه شوند. اما در درخت برآمدگی با توجه به شکل کامل آن لازم است بازاء هر سطح ۲ گره میانی و و چون ۳ سطح داریم در مجموع ۳×۲ گره برای یافتن برگ فعال مورد بررسی قرار گیرد و در نهایت تنها یک برگ با بردار ورودی مقایسه شود. در حالات کلی رتبه محاسبات لازم برای بازیابی خوشه در برگزینده یک بردار ورودی در الگوریتمهای خوشه بندی معمولی $O(n)$ و در درخت برآمدگی $O(\log_2 n)$ می باشد.

در خصوص تعیین سطح آستانه نیز مشاهده شد که استفاده از روش نرخ خطای برابر یا EER در مقایسه با روش برآزش خط راندمان بالاتری را نتیجه می دهد. لیکن روش نرخ خطای برابر نیاز به تعداد کافی و زیاد تستهای درون گویندگی با استفاده از داده های آموزشی دارد که این امر موجب می گردد که لازم شود گویندگان در تعداد دفعات بیشتری کد شناسایی شخصی خود را جهت آموزش سیستم بیان نمایند. از طرف دیگر روش تعیین سطح آستانه تصمیم گیری به روش برآزش خط گرچه راندمان پایین تری را سبب می شوند، لیکن چون تعیین سطح آستانه به کمک آن تنها به آزمایشات برون گویندگی نیاز دارد، بنابراین این امکان فراهم می شود که گویندگان در تعداد دفعات کمتری کد شناسایی شخصی خود را بیان نمایند که این موضوع بدلیل سهولتی که برای گویندگان در مرحله آموزش فراهم می نماید در سیستمهای واقعی تصدیق هویت گوینده از اهمیت بالایی برخوردار می باشد.

در خاتمه ابراز میدارد که آنچه در این مقاله ارائه گردید ماحصل فعالیتهای تحقیقاتی است که برای ایجاد یک سیستم تصدیق هویت از طریق تلفن صورت گرفته است. روشهای استفاده شده با هدف بهبود کارائی، سرعت و سهولت در انجام کار انتخاب شده و آزمایشاتی برای تعیین توانمندی این روشها صورت گرفته است. لزوماً این روشها، بهترین روشها نبوده و البته این افق برای کلیه محققین در این زمینه باز است که روشها و تکنیکهای دیگر را مورد بررسی قرار داده و با تلاش خود کارائی سیستمهای تصدیق هویت گوینده را بهبود بخشند و از این طریق به امکان استفاده کارآمد تر از این سیستم ها در کاربردهائی که برای آنها متصور است یاری نمایند.

زیر نویسها

- 1-Speaker Verification
- 2-Features
- 3-Pitch
- 4-Formants
- 5-Linear Predictive coefficients
- 6-Reflection coefficients
- 7-Log Area Ratio Coefficients
- 8-Line Spectrum Pair Frequencies
- 9-MFCC: Mel Frequency Cepstrum Coefficients
- 10-LPCC: Linear Frequency Cepstrum coefficients
- 11-Liftering
- 12-Dynamic Time Warping
- 13-Vector Quantization
- 14-HMM: Hidden Markov Model
- 15-Second Order Statistical Measures
- 16-Neural Networks
- 17-Back Propagation
- 18-Vocabulary
- 19-Back Tracking
- 20-Preemphasis
- 21-Hamming Window
- 22-Minimum Variance
- 23-BumpTree
- 24-Bump
- 25-False Rejection
- 26-False Acceptance
- 27-Equal Error Rate
- 28-Intra Speaker Distance
- 29-Inter Speaker Distance
- 30-Cohort Normalization
- 31-Code Book

قدردانی

این پروژه تحقیقاتی در راستای طرح ملی تحقیقات به شماره NRC1357 انجام و از طرف شورای پژوهشهای علمی کشور حمایت گردیده است.

- [1] L. F. Lamel, I. L. Gauvain, "Speaker Verification Over Telephone", *Speech communication*, 31 (2-3), pp. 141-154, 2000.
- [2] C. P. Lim, et al., "Speech Recognition using Artificial Neural Networks", *Proceedings of First Conf. on Web Information Engineering*, Vol. 1, pp. 419-423, 2000.
- [3] W. Siew Chan, L. Chee Peng, R. Osman, "Text-Dependent Speaker Recognition using the Fuzzy ARTMAP Neural Network"; *Proceeding of Intelligent System Technology* Vol. 1, pp. 33-38, 2000.
- [4] J. R. Deller, J. G. Proakis, and J. H. L. Hanson, *Discrete-Time Processing of Speech Signals*, Macmilan, NewYork, 1993.
- [5] M. G. Rahim, *Artificial Neural Networks for Speech Analysis/Synthesis*, Chapman & Hall, University Press, Cambridge, 1994.
- [6] H. Sakoe, R. Isotani, K. Yoshida, K. Iso, and T. Watanabe, "Speaker-Independent Word Recognition Using Dynamic Programming Neural Predication Model", *ICASSP-89*, pp. 29-32, 1989.
- [7] C. G. Looney, *Pattern Recognition Using Neural Networks*, Oxford University Press, 1997.
- [8] B. V. Williams, R. T. J. Bostock, D. Bounds, and A. Harget, "Improving Classification Performance in the Bumptree Network by Optimizing Topology with a Genetic Algorithm", *First Conference on Genetic Algorithms*, pp. 490-495, 1994.
- [9] S. M. Omahandro, "Bumptree for Efficient Function, Constraint, and Classification Learning", *Advances in Neural Information Processing System*, 3, Morgan Kaufmann, 1991.
- [10] S. Renals, and R. Rohwer, "Phoneme Classification Experiments Using Radial Basis Functions", *International Joint Conference on Neural Network*, I, pp. 461-467, Washington D. C., 1989.
- [11] F. Chen, B. Miller, and M. Wagner, "Hybrid Threshold Approach in Text-Independent Speaker Verification", *ICSLP-94*, Yokohoma, pp. 1855-1858, 1994.
- [12] S. Furui, "Cepstral Analysis Technique for Automatic Speaker Verification", *IEEE Trans. on ASSP*, ASSP-29, pp. 254-272, 1981.
- [13] Y. Gong, and J. P. Haton, "Non-Linear Interpolation Methods for Speaker Recognition and Verification", *ESCA Workshop on Automatic Speaker Recognition, Identification, and Verification*, pp. 23-26, 1994.
- [14] R. T. J. Bostock, *Doctoral Dissertation*, Dept. of Computer Science and Applied Mathematics, Aston University, Uk, 1993.