

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشکده فنی و مهندسی

پایان نامه کارشناسی ارشد مهندسی برق الکترونیک

عنوان :

تأیید هویت با استفاده از پردازش سیگنال گفتار

استاد راهنما :

دکتر ساجدی

نگارش :

سید مصطفی موسوی بورا

آبان ماه ۱۳۸۹

تقدیم به

پدر و مادرم

آن آموزگاران نخستین

و

تمامی آموزگارانم

تقدیم به

همسر مهربانم

که در انجام مراحل مختلف پروژه مرا یاری نمود

تقدیر و سپاس :

در انجام مراحل مختلف این پروژه از راهنماییهای ارزشمند و حمایت بی دریغ استاد ارجمند و گرامی جناب دکتر ساجدی که همواره تعیین کننده بستر حرکت، مسیر پیشرفت کار و مشوق من بوده است، بر خود لازم می دانم که بدین وسیله زحمات ایشان را ارج نهم و صمیمانه تشکر و قدردانی نمایم.

چکیده

یکی از روشهای تأیید و تعیین هویت افراد، استفاده از صدای فرد می باشد که مقوله شناسایی گوینده و تصدیق و تعیین هویت گوینده نام دارد. از جمله روشهایی که در پردازش گفتار و بطور خاص در شناسایی گوینده کارآیی بسیار خوبی از خود نشان داده است، مدل چندی سازی برداری و مدل شبکه عصبی می باشد که وسیله ای بسیار قوی برای مدل کردن سیگنال های تصادفی و ایستا می باشند. در این پروژه ما از مدل چندی سازی برداری و هم از مدل شبکه عصبی از نوع پرسپترون چند لایه برای تصدیق و تعیین هویت گوینده های از روستای بالابورا از توابع منطقه بندپی غربی شهرستان بابل که می توانند به سه زبان انگلیسی، فارسی و زبان مازندرانی (مازنی) صحبت کنند، بکار گرفته شده است. با توجه به تمهیدات بکار گرفته شده در این پروژه، بر روی سه پایگاه داده از ارقام برای هر سه زبان متشکل از ۲۰ نفر (۱۲ مرد و ۸ زن) که در سنین مختلف ۱۲ تا ۶۱ سال بوده اند و ارقام صفر تا نه را بیان کرده اند، کارآیی سیستم شناسایی گوینده طراحی شده در این پروژه برای تعیین هویت گوینده در هر کدام از این سه زبان مشخص شده است. بطوریکه نتایج حاصل نشان می دهد که درصد شناسایی این سیستم توسط این دو روش در هر سه زبان متفاوت بوده، بطوریکه در مرحله اول دقت سیستم توسط روش VQ تأیید گوینده به زبانهای انگلیسی، فارسی و مازنی بترتیب با MFCC۲۶ برابر با 94.7 %، 92.6 % و 91.3 % شده است که در مرحله دوم با استفاده از یک نوع شبکه عصبی پیشنهادی توسط LPC۲۶ بترتیب برابر با 99.6%، 98.7% و 97.9% بهبود یافته است.

در مراحل مختلف این پروژه تأثیر تعداد دفعات ارائه داده های آموزشی به شبکه عصبی، 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 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.

واژه های کلیدی :

SPEAKER VERIFICATION-NEURAL NETWORKM-VECTOR QUANTIZATION

فهرست مطالب

ت	چکیده
ح	فهرست شکل‌ها
خ	فهرست جدول‌ها
۱	فصل ۱- مقدمه
۲	۱-۱- تعریف مسئله
۳	۲-۱- سیستم شنوایی در انسان
۴	۳-۱- واحدهای گفتار برای بازنمایی گوینده‌ها
۵	۴-۱- مسائل مطرح در بازشناسی هویت گوینده
۷	۵-۱- روش‌های کلی بازشناسی گوینده
۸	۶-۱- بازنمایی مدل‌های گوینده‌ها
۸	۷-۱- معیار شباهت ویژگی‌ها
۸	۸-۱- مقاوم‌سازی سیستم در برابر تغییرات
۹	۹-۱- محاسبه آستانه
۱۰	۱-۹-۱- معیار EER
۱۱	۱۰-۱- طریقه تصمیم‌گیری
۱۲	۱۱-۱- تحقیقات قبلی
۱۵	۱-۱۱-۱- پارامترها و ویژگی‌های استخراج شده از سیگنال گفتار
۱۵	۲-۱۱-۱- روش مدل‌کردن و طبقه‌بندی گوینده‌ها
۱۶	۳-۱۱-۱- نحوه تصمیم‌گیری نهایی
۱۶	۱۲-۱- روش‌های متداول پیاده‌سازی یک سیستم شناسایی گوینده
۱۶	۱-۱۲-۱- روش برنامه‌ریزی پویا
۱۶	۲-۱۲-۱- روش‌های احتمالاتی
۱۷	۳-۱۲-۱- روش چندی‌کردن برداری
۱۷	۴-۱۲-۱- استفاده از شبکه‌های عصبی

فصل ۲- اجزای یک سیستم تأیید هویت با استفاده از پردازش سیگنال گفتار.....	۱۹
۱-۲- مقدمه	۲۰
۲-۲- احتساب سیگنال صحبت و نمونه برداری از آن.....	۲۲
۳-۲- پیش پردازش	۲۳
۴-۲- تشخیص ابتدا و انتهای کلمات	۲۵
۵-۲- قطعه بندی سیگنال صحبت	۲۸
۶-۲- ویژگی های مناسب برای سیستم تشخیص گوینده	۲۹
۷-۲- معرفی چند ویژگی از سیگنال صحبت	۳۰
۱-۷-۲- آنالیز بانک فیلتر	۳۱
۲-۷-۲- ضرایب خود همبستگی	۳۲
۳-۷-۲- آنالیز پیشگویی خطی	۳۳
۴-۷-۲- آنالیز کپسترال	۳۵
۵-۷-۲- استفاده از مقیاس MEL در آنالیز کپسترال	۴۰
۶-۷-۲- ضرائب MFCC	۴۱
۸-۲- نرمال سازی ضرایب	۴۲
۹-۲- مقاوم سازی پارامترها و ضرایب در برابر نویز و اثرات کانال انتقال	۴۳
۱-۹-۲- اثرات کانال انتقال و نویز	۴۳
۱۰-۲- اندازه گیری فاصله	۴۹
۱-۱۰-۲- معیار فاصله مینکووسکی	۴۹
۲-۱۰-۲- وزن دهی به معیار فاصله اقلیدوسی	۵۱
۳-۱۰-۲- معیار WLR	۵۲
۴-۱۰-۲- معیارهای فاصله LLR و IS	۵۳
۱۱-۲- نتیجه گیری	۵۵
فصل ۳- سیستم تأیید هویت گوینده وابسته به متن با استفاده از روش چندی کردن برداری	۵۶
۱-۳- مقدمه	۵۷
۲-۳- روش مدل سازی گوینده با Vector Quantization	۵۸
۳-۳- بازشناسی گوینده با استفاده از روش چندی سازی برداری	۵۹
۴-۳- هدف از ایجاد کدبوک	۶۰
۵-۳- روش های دسته بندی برای ایجاد کدبوک	۶۰

۶۰Vq	۳-۶- نحوه آموزش بردارهای ویژگی گوینده در روش
۶۴ Vq	۳-۷- نحوه تست(شناسایی) بردارهای ویژگی گوینده در روش
۶۵ Vq	۳-۸- نمای کلی از تأیید هویت گوینده با استفاده از
۶۶	۳-۹- نحوه اجرای بازشناسی گوینده توسط Vq با استفاده از نرم افزار مطلب
۷۳	۳-۱۰- جمع بندی و نتیجه گیری
۷۴	فصل ۴- سیستم تأیید هویت گوینده وابسته به متن با استفاده شبکه های عصبی
۷۵	۴-۱- مقدمه
۷۶	۴-۲- استخراج ویژگی
۷۸	۴-۳- شبکه های عصبی
۷۸	۴-۳-۱- شبکه عصبی پرسپترون چندلایه
۷۹	۴-۳-۲- شبکه عصبی تاخیر زمانی
۸۰	۴-۴- سیستم بازشناسی گوینده مطرح شده
۸۱	۴-۵- ارزیابی
۸۵	۴-۶- نتیجه گیری و پیشنهادات
۸۶	فهرست مراجع

فهرست شکل‌ها

- شکل (۱-۱) عناصر اصلی سیستم بازشناسی گوینده ۶
- شکل (۲-۱) نحوه تعیین خطای مساوی از روی پارامترهای FA و FR بصورت شماتیک ۱۱
- شکل (۳-۱) منحنی ROC برای سه سیستم مجزای A,B,C (از مرجع [9]) ۱۲
- شکل (۱-۲) : ساختار کلی یک سیستم تأیید هویت گوینده ۲۳
- شکل (۲-۲) - طیف هموار شده واج واکدار ۲۵
- شکل (۳-۲) - پاسخ فرکانسی فیلتر پیش تأکید به ازای $a=0.95$ ۲۶
- شکل (۴-۲) - فریم بدون اثر فیلتر پیش تأکید ۲۷
- شکل (۵-۲) - فریم با اثر فیلتر پیش تأکید ۲۷
- شکل (۶-۲) - شکل موج زمانی دو کلمه ادا شده به زبان فارسی ۲۹
- شکل (۳-۲) - بلوک دیاگرام بانک فیلتر ۳۳
- شکل (۴-۲) : الف) شکل موج زمانی یک قطعه از سیگنال صحبت ۳۴
- ب) ضرایب خود همبستگی مربوط به شکل ۲-۴- الف ۳۴
- شکل (۵-۲) - مدل واقعی از دستگاه تولید صوت در حوزه زمان گسسته ۳۵
- شکل (۶-۲) - یک مدل تخمینی از شکل (۵-۲) با استفاده از ضرایب LP ۳۶
- شکل (۷-۲) : الف) شکل موج زمانی یک قطعه از سیگنال صحبت ۳۷
- ب) ضرایب LP مربوط به شکل ۲-۷- الف ۳۷
- شکل (۸-۲) - ساختار تبدیل همومورفیک ۳۷
- شکل (۹-۲) - بلوک دیاگرام یک سیستم فیلتر همومورفیک برای محاسبه ضرایب کپسترال ۳۹
- شکل (۱۰-۲) - شکل موج های مربوط به محاسبه ضرایب کپسترال ۴۰
- شکل ۱۱-۲ : الف) شکل موج زمانی یک قطعه از سیگنال صحبت ۴۱
- ب) ضرایب کپسترال مربوط به شکل ۲-۱۰- الف ۴۱
- شکل (۱۲-۲) - روند محاسبه ضرایب کپسترال در پیاده‌سازی سیستم تأیید هویت گوینده ۴۱
- شکل (۱۲-۲) - بانک فیلتر با مقیاس Mel ۴۳
- شکل (۱۳-۲) - دیاگرام حذف نویز با روش CMS ۴۸
- شکل (۱۴-۲) - پاسخ فرکانسی فیلتر میان‌گذر در روش RASTA ۵۰
- شکل (۱۵-۲) - مقادیر واریانس و معکوس واریانس ضرایب کپسترال ۵۳

- شکل (۲-۱۶) - بلوک دیاگرام تعیین فاصله با استفاده از روش LLR و IS ۵۵
- شکل (۳-۱) - مدل چندی سازی برداری برای دو گوینده ۶۱
- شکل (۳-۲) - نحوه دسته‌بندی کدبوک‌ها در فرایند آموزش ۶۴
- شکل (۳-۳) - بلوک دیاگرام الگوریتم K-means ۶۵
- شکل (۳-۴) - دسته‌های حاصل از الگوریتم K-means ۶۵
- شکل (۳-۵) - نمودار FR و FA در تعیین ترشولد ۷۱
- شکل (۴-۱) - بلوک دیاگرام الگوریتم MFCC ۸۰
- شکل (۴-۲) - شبکه عصبی پرسپترون چندلایه ۸۱
- شکل (۴-۳) - سیستم بازشناسی گوینده مطرح شده ۸۲
- شکل (۴-۴) - تاخیر زمانی در پرسپترون چندلایه ۸۳

فهرست جداول

- جدول (۳-۱) - مقدار متوسط فواصل مراکز کدبوک گوینده ناشناس به سایر گوینده‌ها به ازای $C = 10$ ۶۹
- جدول (۳-۲) - مقدار متوسط فواصل مراکز کدبوک گوینده ناشناس به سایر گوینده‌ها به ازای $C = 13$ ۷۰
- جدول (۳-۳) - خطای تأیید گوینده (بر حسب درصد) در حالت بیان سه و پنج رقم از اعداد صفر تا نه (به زبان فارسی) ۷۲
- جدول (۳-۴) - خطای تأیید گوینده (بر حسب درصد) در حالت بیان سه و پنج رقم از اعداد صفر تا نه (به زبان انگلیسی) ۷۲
- جدول (۳-۵) - خطای تأیید گوینده (بر حسب درصد) در حالت بیان سه و پنج رقم از اعداد صفر تا نه (به زبان مازنی) ۷۳
- جدول (۴-۱) - نرخ بازشناسی و تأیید هویت برای ۵ الی ۲۰ گوینده با زبان انگلیسی ۸۴
- جدول (۴-۲) - نرخ بازشناسی و تأیید هویت برای ۵ الی ۲۰ گوینده با زبان فارسی ۸۴
- جدول (۴-۳) - نرخ بازشناسی و تأیید هویت برای ۵ الی ۲۰ گوینده با زبان مازنی ۸۵
- جدول (۴-۴) - مقایسه نرخ بازشناسی برای مدل بازشناسی VQ و سیستم پیشنهادی به زبان انگلیسی ۸۶
- جدول (۴-۵) - زمان سپری شده در طول روند تست ۸۶

فصل ۱

مقدمه

۱-۱- تعریف مسئله

یکی از مباحث مطرح در پردازش سیگنال گفتار تعیین و یا تصدیق هویت گوینده بر اساس سیگنال گفتار می‌باشد. در تعیین گوینده هویت فرد مورد نظر بعنوان عضوی از گوینده‌های شناخته شده سیستم شناسایی می‌گردد و در تصدیق گوینده ادعای هویت گوینده بعنوان یکی از گوینده‌های ناشناخته شده سیستم، تصدیق یا رد می‌شود. کاربردهای بسیاری برای بازشناسی گوینده وجود دارد که شامل کنترل دسترسی به اطلاعات محرمانه توسط صدا، ارائه اطلاعات و سرویس‌های شخصی توسط صدا نظیر نامه‌نگاری و خرید تلفنی، برچسپ‌دهی گوینده‌ها در گفتگوهای ضبط شده، تحقیقات دادگاهی و جستجوی مجرمین با نمونه صداها ضبط شده و ... می‌باشد.

بازشناسی گوینده جزء روش‌های تشخیص هویت شخص می‌باشد که بر اساس ویژگی‌های فیزیولوژیکی یا رفتاری شخص عمل می‌کنند و عبارتست از فرآیند خودکار شناسایی شخص گوینده با استفاده از اطلاعات فردی او که از سیگنال گفتار آن شخص استخراج می‌شود. تأیید هویت گوینده ارتباط تنگاتنگی با تشخیص گوینده دارد که در آن یک فرد با هویت مجهول ادعای داشتن هویت یکی از اعضای مجموعه را می‌کند. داده‌های ورودی از گفتار این گوینده تنها با الگوهای گفتار همان عضو از مجموعه گویندگان مقایسه می‌شود و در صورتی که فاصله داده‌های ورودی از الگوی مرجع از یک

آستانه کمتر باشد هویت فرد تأیید و در غیر اینصورت رد می‌شود. همچنین بازشناسی گوینده برای کنترل دسترسی از راه دور مناسبتر است زیرا صدا از طریق خطوط تلفن قابل انتقال است.

۱-۲- سیستم شنوایی در انسان

گفتار تولید شده توسط اندام‌های گویایی بصورت یک موج صوتی از طریق هوا منتقل و به گوش شنونده می‌رسد. این موج وارد گوش شده و باعث تحریک پرده صماخ می‌گردد. این ارتعاشات از طریق استخوانهای لایه گوش به درون حلزونی منتقل شده و باعث تحریک سلسله اعصاب شنوایی که به حلزونی متصل هستند، می‌گردد. در عمل حلزونی به صورت یک بانک فیلتر در محدوده $20000 - 20$ Hz عمل می‌کند [1]. برخی از خصوصیات گوش عبارتند از:

الف - منحنی Mel

فرکانس یک دانگ صدای^۱ خاص که توسط گوش دریافت می‌شود با فرکانس واقعی آن برابر نیست و رابطه خطی میان این دو فرکانس وجود ندارد. منحنی mel رابطه میان این دو فرکانس را نشان می‌دهد. این منحنی در فرکانس‌های کمتر از ۱ کیلوهرتز تقریباً خطی است و بالای آن بصورت لگاریتمی افزایش می‌یابد [1], [2].

ب) فیلترهای باند بحرانی^۲

گوش دارای خاصیت حذف نویز در پهنای باند خاصی حول یک فرکانس می‌باشد. این پهنای باند برای فرکانس‌های مختلف مقدار متفاوتی دارد که برای فرکانس‌های کمتر از ۱ کیلوهرتز تقریباً برابر ۱۰۰ هرتز و برای فرکانس‌های بالاتر بصورت لگاریتمی افزایش می‌یابد. پهنای باند بحرانی بدین معنی است که اگر انرژی نویز در این باند حول فرکانس f_0 بزرگتر از مقدار نباشد آن فرکانس توسط گوش قابل تشخیص خواهد بود و در غیر اینصورت آن فرکانس قابل شناسایی توسط گوش نیست [1], [2].

¹ Tone

² Critical band filters

ج) منحنی‌های بلندی هم تراز^۱

میزان حساسیت گوش به فرکانس‌های مختلف صدا یکسان نیست. بیشترین حساسیت در فرکانس ۱ کیلوهرتز است [2]، که از این خاصیت در استخراج بعضی از پارامترها استفاده می‌شود.

د) قانون توان شنوایی^۲

شدت صوتی که توسط گوش شنیده می‌شود با توان واقعی سیگنال برابر نیست بلکه یک رابطه غیر خطی میان این دو وجود دارد. این رابطه به صورت ریشه سوم دامنه سیگنال مدل می‌شود. از خصوصیات فوق در استخراج پارامترهای ادراکی^۳ از سیگنال گفتار استفاده می‌شود که مهمترین آنها پارامترهای Mel-Cepstral و PLP هستند [2],[3].

۱-۳- واحدهای گفتار برای بازنمایی گوینده‌ها

یکی از مهمترین ملاحظات در طراحی سیستم بازشناسی گوینده‌ها انتخاب واحدهای گفتاری برای مدل کردن گوینده‌ها است. انتخاب واحدها شامل واحدهای آوایی یا زبانشناختی به عنوان جملات کامل، کلمات و هجاها و واحدهای شبه آوا و نیز واحدهای صوتی نظیر واحدهای جدا شده از گفتار بر اساس مدل‌های صوتی به جای معیارهای آوایی می‌باشد.

برای بازشناسی گوینده ضرورتی برای بازنمایی گفتار بصورت واحدهای واجی نیست و چون سیستم‌های که از واحدهای واجی استفاده می‌کنند پیچیده تر از سیستم‌هایی هستند که تنها از مدل‌های صوتی استفاده می‌کنند. لذا بیشتر سیستم‌های بازشناسی گوینده از نوع دوم هستند [4].

¹ Equal Loudness curves

² Power loudness curves

³ perceptual

۱-۴- مسائل مطرح در بازشناسی هویت گوینده

هر عبارتی که توسط شخصی بیان می‌شود حداقل دو نوع اطلاعات را در بر دارد که عبارتند از اطلاعات خود پیام و اطلاعاتی در مورد گوینده پیام می‌باشد. عناصر اصلی یک سیستم بازشناسی هویت گوینده در شکل (۱-۲) نشان داده شده است. در تعیین هویت گوینده^۱ سیگنال گوینده ناشناس آنالیز شده و با مدل گوینده‌های شناخته شده مقایسه می‌گردد. گوینده ناشناس بعنوان گوینده‌ای که مدلس بیشترین تطابق را با سیگنال گفتار ورودی دارد معرفی می‌گردد. در مورد تشخیص مجموعه بسته^۲ تصمیمات ممکن برابر تعداد گوینده‌های شناخته شده است ولی در مورد تشخیص مجموعه باز^۳ یک انتخاب دیگر تحت عنوان عدم تطابق سیگنال گفتار با مدل-های شناخته شده نیز وجود خواهد داشت.

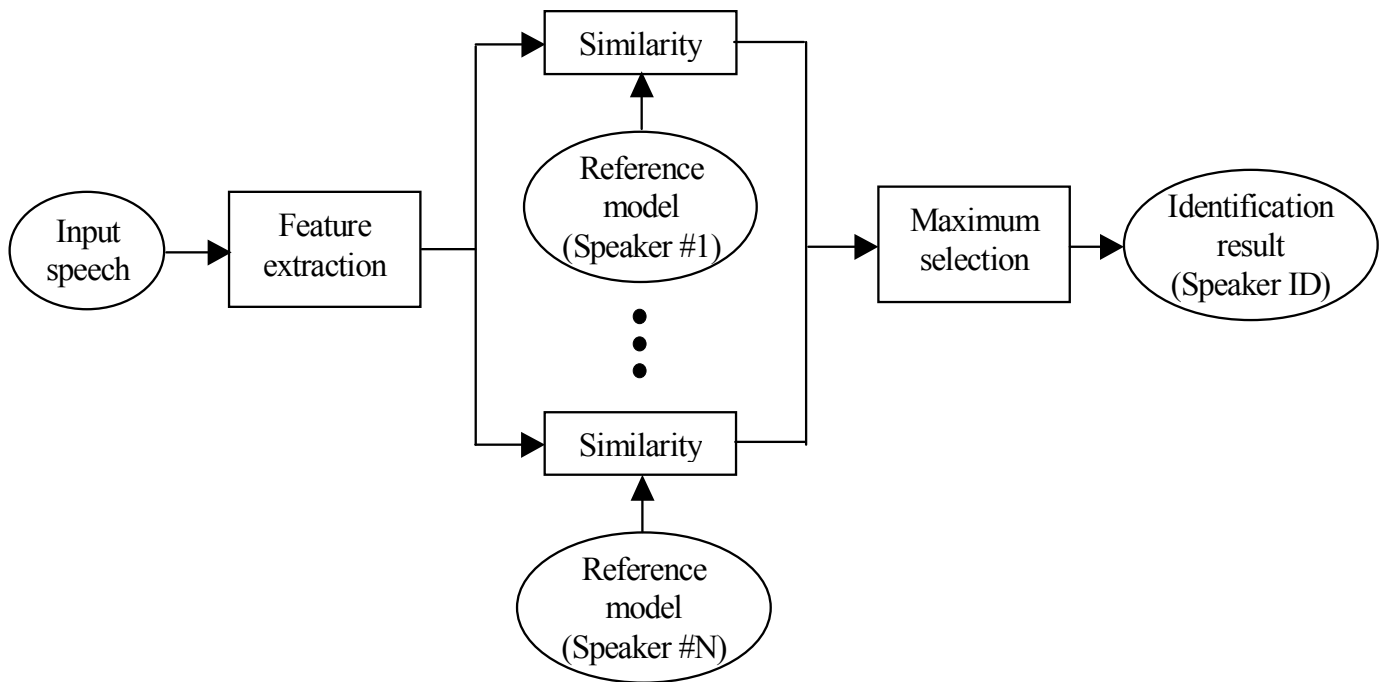
در تصدیق گوینده^۴، ادعای هویت برای یک سیگنال گفتار بعنوان یکی از مدل‌های شناخته شده، تصدیق و یاد رد می‌شود. در این روش گفتار شخص ناشناس یا مدل گوینده ادعا شده مقایسه می‌گردد و اگر تطابق به اندازه کافی باشد ادعای مزبور تصدیق می‌گردد [5]. در شکل صفحه بعد بلوک دیاگرام مربوط به سیستم تعیین و تأیید گوینده نمایش داده شده است.

¹ Speaker Identification

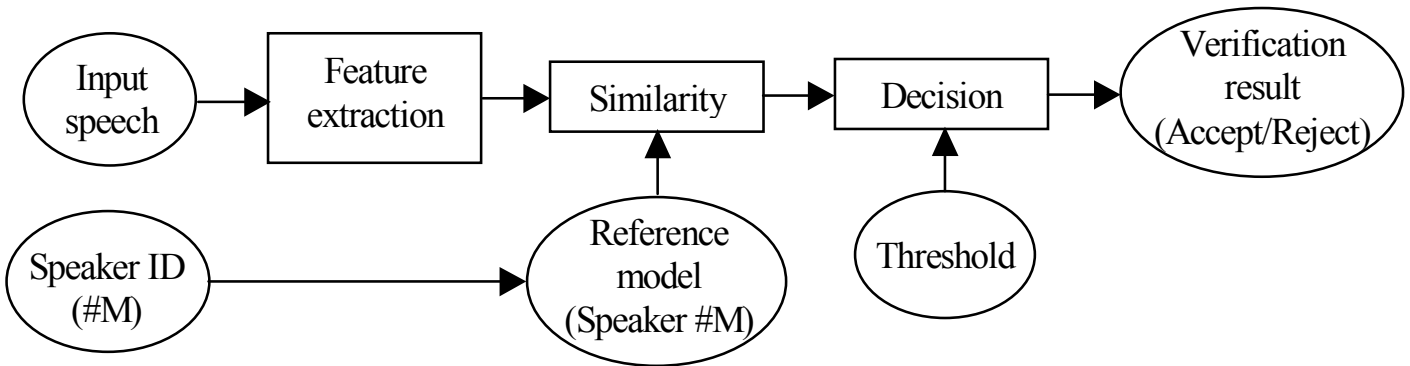
² Close Set

³ Open Set

⁴ Speaker Verification



الف) بلوک دیاگرام تعیین گوینده



ب) بلوک دیاگرام تأیید گوینده

شکل (۱-۱) عناصر اصلی سیستم بازشناسی گوینده [5]: الف) تعیین گوینده
 ب) تأیید گوینده

۱-۵- روش‌های کلی بازشناسی گوینده

سیستم‌های بازشناسی گوینده به دو روش وابسته به متن^۱ و مستقل از متن^۲ عمل می‌کنند. در حالت وابسته به متن، گوینده باید همان جملاتی که برای یادگیری سیستم استفاده شده است را ادا کند ولی در حالت مستقل از متن، گوینده می‌تواند هر جمله‌ای را در زمان تشخیص ادا نماید. از آنجائیکه در حالت وابسته به متن می‌توان ویژگی‌های مربوط به هر واج یا هجا را بررسی کرد لذا دقت‌های بالاتری نسبت به حالت مستقل از متن بدست می‌آید [6].

الف) وابسته به متن

از گوینده خواسته می‌شود که کلمه عبور را ادا نماید و این سیگنال با مدل‌های شناخته شده گوینده‌ها مقایسه و تصمیم‌گیری گرفته می‌شود.

ب) مستقل از متن

در بسیاری از کاربردها کلمه عبور از پیش تعیین شده، قابل استفاده نیست بلکه سیستم باید توانایی شناسایی گوینده با هر گفتاری را داشته باشد. بعلاوه کاربردهای فراوان این نوع تشخیص بیشتر مورد توجه است.

ج) وابسته به متن با کلمات عبوری اتفاقی

هر دو روش قبل به آسانی قابل شکست هستند چون اگر صدای گوینده ضبط و دوباره پخش شود بعنوان یکی از گوینده‌ها پذیرفته خواهد شد. برای جلوگیری از این کار یک مجموعه از کلمات نظیر اعداد به عنوان رمز عبور استفاده می‌شوند. و هر بار از کاربر خواسته می‌شود که یک توالی خاصی از این کلمات را ادا نماید. اخیراً یک چنین سیستمی ساخته شده است که کلمات عبور هر بار بطور کامل عوض می‌شود [7].

¹ Text Dependent

² Text Independent

۱-۶- بازنمایی مدل‌های گوینده‌ها

معمولترین روش برای بازشناسی خودکار گوینده در حالت وابسته به متن ذخیره و ویژگی‌های مقطعی است. در این روش هر گوینده بصورت متوالی از بردارهای ویژگی جمله ادا شده مدل می‌شود و تصمیم بر اساس تطابق الگوها^۱ صورت می‌پذیرد.

۱-۷- معیار شباهت ویژگی‌ها

مهمترین وجه تمایز میان سیستم‌های بازشناسی گوینده نحوه ترکیب پارامترها و مقایسه آنها با مدل‌های مرجع می‌باشد. مقایسه به کمک یک معیار شباهت^۲ یا معیار فاصله^۳ انجام می‌گیرد و با مقایسه ویژگی‌های گوینده ناشناس با مدل‌های شناخته شده نزدیکترین مدل به عنوان گوینده برنده انتخاب می‌شود.

فاصله میان دو بردار V' و V بطول L بصورت: $d = \sum_{i=1}^L |V_i - V'_i|$

یا فاصله اقلیدوسی: $d^2 = \sum_{i=1}^L (|V_i - V'_i|)^2$

یا فاصله وزن دار: $d^2 = \sum_{i=1}^L (W_i(|V_i - V'_i|))^2$

که در آن $W_i = 1/\sigma_i$ که σ_i تخمینی از واریانس ضریب σ_i بردار ویژگی می‌باشد.

۱-۸- مقاوم سازی سیستم در برابر تغییرات

تغییرات در تکرارهای مختلف عامل مهمی در کاهش کارایی سیستم‌های بازشناسی گوینده می‌باشد. این تغییرات ناشی از تغییرات در خود شخص، تغییرات شرایط ضبط و انتقال و نویزی می‌باشد. مطالعه نشان داده است که جملات ضبط شده در یک جلسه همبستگی بیشتری نسبت به جملات ضبط شده در جلسات دیگر دارند لذا مقابله با این تفاوت ضروری است [8],[9]. عملاً دو روش برای حل این معضل وجود دارد که عبارتند از:

¹ Template

² Liklyhod Measure

³ Distance Measure

⁴ Robust

الف) به روز کردن مدل‌ها در هر بار استفاده

در این روش از جملاتی که در مرحلهٔ بازشناسی توسط گوینده ادا می‌شود نیز برای یادگیری سیستم استفاده می‌شود. شرایط یادگیری شامل تعداد جلسات یادگیری، تعداد جملات و شرایط ضبط و انتقال گفتار می‌شود. تغییرات در این شرایط در زمان بازشناسی کارایی سیستم را به شدت تحت تأثیر قرار می‌دهد [9].

ب) هنجارسازی^۱ ویژگی‌های صوتی

یکی از روش‌های هنجارسازی پارامترها استفاده از روش CMN^2 است که اثرات خطی کانال و تغییرات بلند مدت طیفی را حذف می‌کند. CMS^3 یکی از روش‌های فوق است که در آن ضرایب کپسترال در طول یک توالی از گفتار میانگین‌گیری شده و مقدار میانگین از ضرایب کپسترال در طول یک توالی از گفتار میانگین‌گیری شده و مقدار میانگین از ضرایب کپسترال هر فریم کم می‌شود. این روش اثرات خطی کانال را به خوبی حذف می‌کند ولی بعضی از اطلاعات مفید را هم از بین می‌برد لذا برای گفتار کوتاه مدت مناسب نیست [9].

۱-۹- محاسبه آستانه

پس از اینکه در مرحله آموزش الگوهای مرجع را تولید نمودیم می‌بایست معیاری برای تصمیم‌گیری در مرحله تشخیص (تست) در دست داشته باشیم که به آن آستانه می‌گوییم که نحوه تعیین آن به شرح زیر است :

چون تطبیق الگوها بر اساس کلمات مجزا صورت می‌گیرد، بنابراین می‌بایست مقدار آستانه برای هر کلمه بصورت جداگانه تعیین گردد. مسأله دیگر، نحوه تعیین مقدار آستانه مربوط به هر کلمه می‌باشد که با استفاده از یکی از دو روش زیر صورت می‌گیرد :

الف) برای یک گوینده (Intraspeaker)

ب) برای مجموعه‌ای از گوینده‌گان (Interspeaker)

در حالت اول برای هر کلمه‌ای یک گوینده خاص ادا می‌کند یک آستانه بدست می‌آید ولی در حالت دوم برای هر

¹ Normalization

² Cepstral Mean Normalization

³ Cepstral Mean Subtraction

کلمه‌ای که تمام گوینده‌گان ادا می‌کنند یک آستانه محاسبه می‌گردد. در یک پایگاه داده که برای n نفر تهیه شده باشد که در آن هر گوینده m کلمه را بیان کرده باشد در حالت اول لازمست $m.n$ آستانه، ولی در حالت دوم کافی است که فقط m آستانه محاسبه گردد.

از آنجاییکه در این پروژه از حالت دوم برای تعیین آستانه استفاده شده است، حال به شرح روش‌های مختلفی برای تعیین آستانه می‌پردازیم. معیارهای مختلفی از جمله معیارهای ZFR^1 ، ZFA^2 و EER^3 برای تعیین آستانه بکار گرفته شده است که ما در این پروژه از معیار EER استفاده نموده‌ایم و به شرح آن می‌پردازیم.

۱-۹-۱- معیار EER

دو معیار قبلی نگرش بسیار سخت گیرانه‌ای نسبت به پذیرش^۴ یا رد^۵ یک گوینده ایجاد می‌کنند، بنابراین دیدگاه معتدلتری که در این پروژه مورد استفاده قرار گرفته است معیار EER می‌باشد که در آن هدف این است که آستانه طوری تعیین شود که پذیرش نادرست^۶ و رد نادرست^۷ با یکدیگر برابر شوند.

بعلت اینکه تعداد نمونه‌های مورد بررسی محدود است عملاً امکان رسیدن به احتمال مساوی برای خطاهای فوق کم است به همین دلیل از تخمین رابطه^(۳-۱) استفاده می‌کنیم:

$$EE^{\wedge} = \frac{FA+FR}{2} \quad (3-1)$$

با استفاده از مقادیر T_{min} و T_{max} و فاصله بدست آمده از مقایسه دو کلمه مقدار آستانه T طوری تعیین می‌گردد که احتمال‌های پذیرش نادرست و رد نادرست تا حد امکان به یکدیگر نزدیک شوند. شکل (۲-۱) نحوه تغییر پارامترهای FA و FR را نسبت به یکدیگر نشان می‌دهد.

¹ Zero False Rejection
² Zero False Acceptance
³ Equal Error Rate
⁴ Acceptance
⁵ Rejection
⁶ False Acceptance
⁷ False Rejection
⁸ Equal Error