

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه صنعتی شریف

دانشکده مهندسی کامپیوتر

پایان نامه کارشناسی ارشد
گرایش هوش مصنوعی

عنوان

شناسایی مستقل از متن گوینده از بین گویندگان متعدد

نگارش

حسین زینلی

استاد راهنما

دکتر حسین صامتی

شهریور ۱۳۹۱

به نام خدا
دانشگاه صنعتی شریف
دانشکده مهندسی کامپیوتر

رساله کارشناسی ارشد

عنوان: شناسایی مستقل از متن گوینده از بین گویندگان متعدد

نگارش: حسین زینلی

کمیته ممتحنین:

امضاء:

استاد راهنما: دکتر حسین صامتی

امضاء:

استاد مشاور: دکتر علی محمد افشین همت یار

امضاء:

استاد مدعو: دکتر محمد مهدی همایون پور

تاریخ: شهریور ۱۳۹۱

تقدیم بہ پدر زحمت کش، مادر فداکار و، مہسر مہربانم

سپاس‌گزاری

گرچه کلام در سپاس از آنان که مشوق و راهنمایم بوده‌اند ناتوان است، در اینجا فرصت را غنیمت می‌شمارم تا از تمامی آنان بخصوص آقای حرم به خاطر گفتگوهای متعدد و راهنمایی‌های ارزشمندشان، آقای دکتر باعلی و دکتر ویسی به خاطر مشورت-های سودمند، شرکت عصرگوش پر داز به خاطر در اختیار قرار دادن امکانات، مهم‌تر از همه استاد گرامی آقای دکتر صامتی به خاطر راهنمایی‌ها و کمک‌های بی‌دینش و سایر دوستانی که بنده را در انجام این پایان‌نامه کمک کرده‌اند کمال تشکر را بنمایم. علاوه بر اینها، به خاطر کمک‌های فراوان همسر مهربانم و همچنین تحمل سختی‌های این مدت از ایشان نیز کمال سپاس‌گزاری را دارم. در پایان از اساتید محترمی که قبول زحمت کردند و در جلسه‌ی دفاع اینجانب شرکت کردند نیز سپاس‌گزارم.

چکیده

گفتار انسان حاوی اطلاعات زیادی از قبیل محتوای معنایی، احساسات و حتی هویت گوینده می‌باشد. هدف ما در این پایان‌نامه شناسایی گوینده مستقل از متن از بین گویندگان متعدد است. در سال‌های اخیر، زمان شناسایی (آزمایش) به یکی از مهمترین چالش‌های سامانه‌های بلادرنگ شناسایی گوینده تبدیل شده است. زمان شناسایی به روش محاسبه‌ی درست‌نمایی گفتار آزمایش در مدل گویندگان ثبت‌نامی وابسته است. برای کاربردهای بلادرنگ، این سامانه‌ها باید خیلی سریع گوینده‌ی هدف را مشخص کنند و به همین دلیل نمی‌توان از روش‌های مرسوم شناسایی گوینده استفاده کرد. هدف اصلی این پایان‌نامه ارائه روش‌هایی است که بتواند زمان شناسایی را کاهش دهد به طوری که روی دقت شناسایی تأثیر منفی نگذارد.

ابتدا، ساختار کلی سامانه‌های شناسایی گوینده و مراحل مختلف آن را معرفی کرده و روش مدل مخلوط گاوسی را به عنوان روش پایه برای این کار معرفی می‌کنیم. سپس روش‌های متداول مورد استفاده برای کاهش زمان شناسایی را مورد بررسی قرار می‌دهیم. در ادامه روشی دومرحله‌ای ارائه می‌کنیم که در مرحله اول از روش نمایش تنک برای شناسایی گوینده استفاده می‌کند. آزمایش‌های ما بر روی دادگان تیمیت نشان می‌دهند که این روش بدون کاهش دقت، سرعت شناسایی را تا ۱۸ برابر افزایش می‌دهد. علاوه بر این روش، روش دومرحله‌ای دیگری پیشنهاد شده است که در مرحله اول از فاصله بین مدل‌های مختلف برای شناسایی استفاده می‌کند. نتایج بر روی همان دادگان نشان می‌دهد که این روش نیز بدون کاهش دقت، سرعت شناسایی را ۲۸ برابر افزایش می‌دهد. در نهایت روشی دو مرحله‌ای بر مبنای فاصله نزدیکترین همسایه ارائه شده که قادر است سرعت را تا ۳/۴ برابر افزایش دهد.

واژه‌های کلیدی:

شناسایی گوینده، مستقل از متن، گویندگان متعدد، سامانه‌های بلادرنگ، زمان شناسایی، مدل مخلوط گاوسی، روش دومرحله‌ای، نمایش تنک، نزدیکترین همسایه.

فهرست مطالب

صفحه	عنوان
۱	فصل ۱ مقدمه
۱	۱-۱ انگیزه
۳	۲-۱ تعریف مسئله
۳	۳-۱ رئوس مطالب
۵	فصل ۲ مروری بر روش‌های شناسایی گوینده
۵	۱-۲ مقدمه
۵	۲-۲ دسته‌بندی شناسایی گوینده
۶	۱-۲-۲ شناسایی گوینده‌ی وابسته به متن
۶	۲-۲-۲ شناسایی گوینده‌ی مستقل از متن
۷	۳-۲-۲ شناسایی گوینده‌ی مجموعه‌ی بسته
۷	۴-۲-۲ شناسایی گوینده‌ی مجموعه‌ی باز
۸	۳-۲ اجزاء و مراحل سامانه‌های شناسایی گوینده
۹	۴-۲ پردازش اولیه و استخراج ویژگی
۹	۱-۴-۲ مقدمه
۱۱	۲-۴-۲ ویژگی‌های کوتاه مدت
۱۱	۳-۴-۲ ویژگی‌های بلند مدت
۱۲	۴-۴-۲ ضرایب کپسترال مقیاس فرکانس مل
۱۵	۵-۴-۲ فیلتر پیش-تاکید
۱۶	۶-۴-۲ ویژگی‌های پویا
۱۷	۵-۲ مدل کردن گویندگان

۱۷	مقدمه	۱-۵-۲
۱۸	مدل مخلوط گاوسی	۲-۵-۲
۲۰	الگوریتم بیشینه‌سازی امید ریاضی	۳-۵-۲
۲۲	مدل پس‌زمینه‌ی جهانی (UBM)	۴-۵-۲
۲۴	مدل وفق یافته با استفاده از تخمین بیشینه‌ی پسین	۵-۵-۲
۲۷	مدل وفق یافته با استفاده از درون‌یابی خطی بیشینه‌ی درست‌نمایی	۶-۵-۲
۲۹	جمع‌بندی	۶-۲
۳۰	فصل ۳ روش‌های متداول برای افزایش سرعت	۳-۲
۳۰	۱-۳ مقدمه	۳-۳
۳۰	۲-۳ پارامترهای موثر در زمان	۳-۳
۳۱	۳-۳ روش‌های متداول کاهش زمان	۳-۳
۳۳	۴-۳ پیش - چندی‌سازی	۳-۳
۳۴	۵-۳ هرس گویندگان	۳-۳
۳۵	۶-۳ استفاده از مدل جهانی	۳-۳
۳۷	۷-۳ روش خوشه‌بندی مدل گویندگان	۳-۳
۴۱	فصل ۴ روش‌های پیشنهادی برای افزایش سرعت	۴-۲
۴۱	۱-۴ مقدمه	۴-۲
۴۲	۲-۴ الگوریتم پیشنهاد شده بر پایه‌ی نمایش تنک	۴-۲
۴۲	۱-۲-۴ شناسایی گوینده با استفاده از نمایش تنک	۴-۲
۴۴	۲-۲-۴ روش ارائه شده برای افزایش سرعت	۴-۲
۴۵	۳-۴ روش ارائه شده بر اساس فاصله‌ی GMMها	۴-۲
۴۸	۴-۴ روش ارائه شده بر اساس فاصله‌ی نزدیکترین همسایه	۴-۲
۵۱	۵-۴ جمع‌بندی	۴-۲

فصل ۵ پیاده‌سازی و آزمایش‌ها	۵۲
۱-۵ مقدمه	۵۲
۲-۵ دادگان مورد استفاده	۵۲
۳-۵ سامانه‌ی شناسایی گوینده	۵۴
۴-۵ روش ارزیابی سامانه‌ی شناسایی گوینده	۵۴
۵-۵ استفاده از MLLR یا MAP	۵۵
۶-۵ انتخاب بردارهای ویژگی	۵۷
۷-۵ تأثیر تعداد گوینده بر زمان شناسایی	۵۷
۸-۵ نتایج بدست آمده از نمایش تنک	۵۹
۹-۵ نتایج بدست آمده از روش فاصله‌ی GMMها	۶۱
۱-۹-۵ نتایج آزمایش فایل‌ها بصورت جداگانه	۶۲
۱-۹-۵ نتایج آزمایش فایل‌ها با هم	۶۵
۱۰-۵ نتایج بدست آمده از روش نزدیکترین همسایه	۶۷
۱۱-۵ جمع‌بندی	۷۰
فصل ۶ نتیجه‌گیری و پیشنهادها	۷۱
مراجع	۷۴

فهرست اشکال

صفحه	عنوان
۸	شکل ۱-۲ نمودار جعبه‌ای مرحله‌ی ثبت‌نام
۹	شکل ۲-۲ نمودار جعبه‌ای مرحله‌ی شناسایی
۱۳	شکل ۳-۲ نمودار جعبه‌ای مراحل استخراج ضرائب کپسترال
۱۴	شکل ۴-۲ بانک فیلتر مل
	شکل ۵-۲ نمودار چگالی طیف توان، (الف) سیگنال اصلی با فرکانس نمونه‌برداری ۴۴۱۰۰، (ب) پیش- تاکید شده همان سیگنال
۱۵	شکل ۶-۲ نمای کلی دو روش ساخت مدل جهانی، (الف) ادغام داده‌های آموزش دو زیرگروه و ساختن مدل جهانی، (ب) ساخت مدل جداگانه برای هر زیرگروه و سپس ادغام مدل‌های ساخته شده برای ساخت مدل جهانی نهایی
۲۴	شکل ۱-۵ طریقه‌ی کار روش تطبیق درون‌یابی خطی بیشینه‌ی درست‌نمایی
۵۶	شکل ۲-۵ نمودار زمان شناسایی گوینده‌ی روش پایه بر حسب تعداد گویندگان ثبت‌نام شده در سامانه ۵۸
۶۰	شکل ۳-۵ نمودار تغییرات دقت شناسایی روش نمایش تنک بر حسب انتخاب n بهترین گوینده
۶۳	شکل ۴-۵ نمودار تغییرات دقت شناسایی بر حسب تعداد گویندگان انتخاب‌شده برای آزمایش اول
۶۶	شکل ۵-۵ نمودار تغییرات دقت شناسایی بر حسب تعداد گویندگان انتخاب‌شده برای آزمایش دوم
	شکل ۶-۵ نمودار تغییرات دقت شناسایی روش نزدیکترین همسایه بر حسب تعداد گویندگان انتخاب شده
۶۹	

فهرست جداول

صفحه	عنوان
۲	جدول ۱-۱ مقایسه روش‌های مختلف زیست‌سنجی
۳۹	جدول ۱-۳ فاکتور افزایش سرعت و دقت روش خوشه‌بندی
۶۰	جدول ۱-۵ نتایج شناسایی گوینده برای آزمایش اول با انتخاب ۱۵ بهترین گوینده
۶۱	جدول ۲-۵ نتایج شناسایی گوینده برای آزمایش دوم با انتخاب ۱۰ بهترین گوینده
۶۴	جدول ۳-۵ نتایج بدست آمده از روش فاصله‌ی GMMها برای nهای مختلف برای آزمایش اول
۶۶	جدول ۴-۵ نتایج بدست آمده از روش فاصله‌ی GMMها برای nهای مختلف برای آزمایش دوم
۶۹	جدول ۵-۵ نتایج بدست آمده از روش نزدیکترین همسایه

فصل ۱

مقدمه

۱-۱ انگیزه

در سال‌های اخیر با پیشرفت علم تغییرات چشم‌گیری در تکنولوژی‌های مختلف بوجود آمده است. یکی از این تکنولوژی‌ها شناسایی افراد با استفاده از زیست‌سنجی^۱ است. زیست‌سنجی، علم شناختن افراد با استفاده از خصوصیات فیزیکی و ویژگی‌های رفتاری آنها است. نیاز به امنیت باعث شده که زمینه‌ی زیست‌سنجی در سال‌های اخیر پیشرفت سریعی داشته باشد. از روش‌های متداول زیست‌سنجی می‌توان به شناسایی با استفاده از اثر انگشت^۲، صورت^۳، عنبیه^۴، شبکیه^۵ و گفتار^۶ اشاره کرد. در جدول ۱-۱ مقایسه‌ی این روش‌ها نشان داده شده است [de Luis-Garcia_03]. این مقایسه بر مبنای ویژگی‌های دقت یا قابلیت اطمینان، سادگی استفاده، سادگی پیاده‌سازی، کاربر پسندی و هزینه روش، انجام شده است.

¹ Biometric
² Finger-Print
³ Face
⁴ Iris
⁵ Retina
⁶ Voice

جدول ۱-۱ مقایسه روش‌های مختلف زیست‌سنجی

هزینه	سادگی پیاده‌سازی	کاربر پسند بودن	سادگی استفاده	دقت	ویژگی روش
متوسط	زیاد	کم	متوسط	زیاد	اثر انگشت
کم	متوسط	زیاد	زیاد	کم	صورت
بالا	متوسط	متوسط	متوسط	متوسط	عنبریه
متوسط	کم	کم	کم	زیاد	شبکیه
کم	زیاد	زیاد	زیاد	متوسط	گفتار

همان طور که از جدول بالا مشخص است استفاده از گفتار یکی از روش‌های مناسب زیست‌سنجی است که علاوه بر کاربرپسند بودن و سادگی پیاده‌سازی بالا، هزینه پایینی دارد. استفاده از گفتار برای شناسایی افراد در مقایسه با روش‌های دیگر مزایایی دارد که از آنها می‌توان به موارد زیر اشاره کرد:

- گفتار همیشه با انسان است و در زمان نیاز می‌توان از آن استفاده کرد.
- گفتار قابل سرقت نیست و به همین دلیل گزینه مناسبی می‌باشد.
- گفتار نسبت به اثر انگشت کمتر در سوانح آسیب می‌بیند.
- شناسایی با استفاده از گفتار نیاز به تجهیزات و سنسورهای پیچیده ندارد.
- به راحتی در کاربردهای راه دور مثل پشت تلفن قابل استفاده می‌باشد.

با توجه به این ویژگی‌ها، در سال‌های اخیر توجه زیادی به این روش شده است. در کاربردهای تلفنی که قابلیت استفاده از روش‌های دیگر بویژه اثر انگشت نیست، این روش بیشتر مورد توجه قرار گرفته است. با وجود مزایای زیادی که استفاده از گفتار دارد، این روش دارای مشکلات و محدودیت‌هایی است. از این محدودیت‌ها می‌توان به کاهش دقت این روش در وجود نویزهای مختلف به خصوص نویز کانال اشاره کرد. علاوه بر این مشکل می‌توان به زمان‌بر بودن فرایند شناسایی وقتی که تعداد گویندگان موجود در سامانه زیاد باشند و در نتیجه عدم استفاده از این روش در سامانه‌های بلادرنگ اشاره کرد. مشکل دیگری که تمام سامانه‌های شناسایی با آن روبه‌رو هستند مشکل جعل و سوء استفاده است. در

سامانه‌های بازشناسی گوینده نیز این مشکل وجود دارد. امکان جعل صدا توسط انسان و ماشین امنیت سامانه‌های بازشناسی گفتار را تحت تأثیر قرار داده است.

۲-۱ تعریف مسئله

شناسایی گوینده از بین گویندگان متعدد را از دو منظر می‌توان مورد بررسی قرار داد. اول اینکه در سامانه‌های شناسایی گوینده با افزایش تعداد گویندگان ثبت‌نامی، احتمال اشتباه شناسایی کردن یک گوینده بیشتر شده و در نتیجه دقت شناسایی کاهش پیدا می‌کند. یکی از دلایل این امر این است که با افزایش تعداد گویندگان، احتمال وجود گویندگانی که ویژگی‌های صوتی شبیه هم داشته باشند بیشتر می‌شود. این شباهت، احتمال اشتباه شناسایی کردن سامانه را بیشتر می‌کند.

دومین منظر این است که در این سامانه‌ها با افزایش تعداد گویندگان ثبت‌نامی، تعداد درست-نمایی‌هایی که در مرحله‌ی شناسایی باید محاسبه شوند بیشتر شده و در نتیجه زمان پاسخ‌گویی افزایش می‌یابد. در روش‌های مرسوم شناسایی گوینده نیاز است که یک جستجوی خطی روی تمام گویندگان ثبت‌نامی انجام شود، که با افزایش تعداد گویندگان ثبت‌نامی این جستجو طولانی‌تر می‌شود. با وجود این مشکل به نظر می‌رسد که در کاربردهای بلادرنگ نمی‌توان از این سامانه‌ها استفاده کرد. در این پایان‌نامه تمرکز ما روی کاهش زمان شناسایی است و همچنین بدلیل اینکه دقت شناسایی مهمترین فاکتور سامانه‌های شناسایی گوینده است، در روش‌های ارائه شده سعی بر این است که افزایش سرعت باعث کاهش دقت نشود.

۳-۱ رئوس مطالب

در فصل ۲، یک بررسی اجمالی بر سامانه‌ی شناسایی گوینده انجام می‌دهیم و دسته‌بندی‌های مختلف شناسایی گوینده را معرفی می‌کنیم. سپس اجزاء و مراحل مختلف سامانه‌های شناسایی گوینده را مورد بررسی قرار می‌دهیم. در ادامه‌ی این فصل، ویژگی‌های مختلف را معرفی کرده و مراحل استخراج ضرایب

کیسترال مقیاس فرکانس مل را به تفصیل شرح می‌دهیم. سپس مدل‌های مختلف برای مدل‌سازی فضای ویژگی را معرفی می‌کنیم. مدل مخلوط گاوسی را به عنوان روش اصلی مدل‌سازی، شرح داده و روش تخمین پارامترهای آن را نیز توضیح می‌دهیم. علاوه بر این روش، دو روش تخمین بیشینه‌ی پسین و درونیابی خطی بیشینه‌ی درست‌نمایی را به عنوان روش‌های موفق برای تطبیق پارامترهای مدل مخلوط گاوسی شرح می‌دهیم. سپس در فصل ۳ روش‌های متداول برای کاهش زمان مرحله‌ی شناسایی را معرفی کرده و تعدادی از آنها را شرح می‌دهیم. در فصل ۴ روش‌های پیشنهادی برای افزایش سرعت را توضیح داده و سپس نتایج این روش‌ها را در فصل ۵ می‌آوریم. در پایان در فصل ۶ نتایج این پایان‌نامه را به صورت خلاصه آورده و زمینه‌هایی برای کار و پژوهش بیشتر در این موضوع را بیان می‌کنیم.

فصل ۲

مروری بر روش‌های شناسایی گوینده

۱-۲ مقدمه

همان‌طور که در فصل قبل شرح داده شد، بازشناسی گوینده یکی از روش‌های زیست‌سنجی است که به دلیل مزایای متعددی که دارد در سال‌های اخیر مورد توجه قرار گرفته است. در این فصل ابتدا دسته‌بندی‌های مختلف شناسایی گوینده را معرفی کرده و اجزاء و مراحل سامانه‌های شناسایی گوینده را شرح می‌دهیم. ویژگی‌های مختلف مورد استفاده برای شناسایی گوینده را بیان کرده و ضرایب کپسترال مقیاس فرکانس مل را به تفصیل شرح می‌دهیم. در ادامه روش‌های مختلف مدل‌سازی گویندگان را معرفی و مدل مخلوط گاوسی را به عنوان مدل اصلی استفاده شده، توضیح می‌دهیم.

۲-۲ دسته‌بندی شناسایی گوینده

شناسایی گوینده را می‌توان از جهت‌های مختلفی دسته‌بندی کرد که یکی از این جهت‌ها اهمیت متنی است که گوینده بیان می‌کند. از این جهت شناسایی گوینده به دو دسته‌ی وابسته به متن^۱ و مستقل از متن^۲ تقسیم می‌شود [Beigi_11].

^۱ Text-Dependent
^۲ Text-Independent

۱-۲-۲ شناسایی گوینده‌ی وابسته به متن

در شناسایی گوینده‌ی وابسته به متن چیزی که گوینده بیان کرده مهم است. در این دسته در مرحله‌ی شناسایی، گوینده باید همان متنی که در مرحله‌ی ثبت‌نام گفته را تکرار کرده باشد. این روش بیشتر برای تصدیق هویت گوینده کاربرد دارد و در عمل برای شناسایی گوینده قابل استفاده نیست. مزیت اصلی سامانه‌هایی که از این روش استفاده می‌کنند، گفتار کم برای آموزش و آزمایش و همچنین دقت بالاتر آنها است.

۲-۲-۲ شناسایی گوینده‌ی مستقل از متن

در این دسته اینکه گوینده چه بیان کرده است و یا اینکه در مرحله‌ی شناسایی متنی متفاوت از متن مرحله‌ی ثبت‌نام بیان کرده باشد، مهم نیست. مستقل از متن بودن چندین سطح دارد. بیشترین استقلال حالتی است که سامانه هم مستقل از متن بیان شده و هم مستقل از زبان^۱ بیان متن باشد. حالت دیگر این است که سامانه، مستقل از متن اما وابسته به زبان باشد. در حالت مستقل از متن و زبان، سامانه باید کار شناسایی را فقط با استفاده از خصوصیات مسیر صوتی گوینده انجام دهد و نباید هیچ فرضی در رابطه با محتوای بیان شده بکند [Beigi_11]. سامانه‌های مستقل از متن نسبت به سامانه‌های وابسته به متن گفتار بیشتری هم برای آموزش و هم برای شناسایی نیاز دارند. به طور کلی این محدودیت باعث شده است که سامانه‌های مستقل از متن نسبت به وابسته به متن دقت کمتری داشته باشند.

علاوه بر این دسته‌بندی، شناسایی گوینده با توجه به نوع مجموعه‌ای که در آن عمل شناسایی

انجام می‌شود، به دو دسته‌ی مجموعه‌ی بسته^۲ و مجموعه‌ی باز^۳ تقسیم می‌شود [Beigi_11]:

^۱ Language-Independent

^۲ Close-Set

^۳ Open-Set

۳-۲-۲ شناسایی گوینده‌ی مجموعه‌ی بسته

در شناسایی گوینده‌ی مجموعه‌ی بسته فرض می‌کنیم که گفتار آزمایش حتماً متعلق به یکی از گویندگانی است که در سامانه ثبت شده‌اند. در این حالت گفتار آزمایش با مدل تمام گویندگان ثبت‌نامی مقایسه می‌شود و گوینده‌ای که مدل آن بیشترین شباهت را با گفتار آزمایش داشته باشد به عنوان گوینده‌ی هدف انتخاب می‌شود. توجه داشته باشید که در این حالت همیشه یک گوینده به عنوان گوینده‌ی هدف انتخاب می‌شود.

شناسایی گوینده مجموعه‌ی بسته بجز کاربردهای خاص، در عمل کمتر مورد استفاده قرار می‌گیرد. به عنوان مثالی از کاربرد این دسته، محیطی را در نظر بگیرید که برای کاربران مختلف شخصی سازی شده است. در این محیط، سامانه‌ای به صورت خودکار با استفاده از صدای فرد تنظیمات شخصی آن را بکار می‌برد. در این مثال اشتباه تشخیص دادن گوینده جریمه‌ی زیادی ندارد. در حقیقت اگر گوینده در مجموعه‌ی گویندگان ثبت‌نام شده وجود نداشته باشد، فرقی نمی‌کند که تنظیمات چه گوینده‌ای انتخاب شده است.

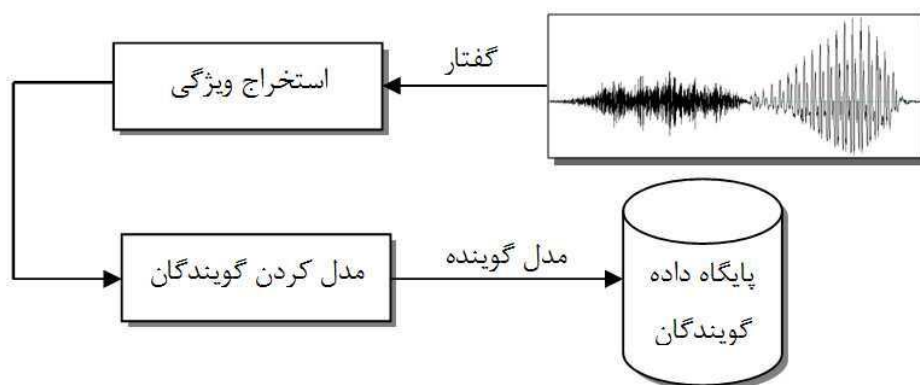
۴-۲-۲ شناسایی گوینده‌ی مجموعه‌ی باز

در این حالت بر خلاف حالت قبلی ممکن است که گفتار آزمایش مربوط به هیچ یک از گویندگان ثبت‌نامی در سامانه نباشد و باید به این سوال که آیا گفتار آزمایش توسط یکی از گویندگان ثبت‌نامی تولید شده است یا نه؟ پاسخ بدهیم. به این دلیل کار شناسایی گوینده در این حالت سخت‌تر از حالت قبل است. شناسایی گوینده‌ی مجموعه‌ی باز ترکیبی از شناسایی گوینده‌ی مجموعه‌ی بسته و تصدیق هویت گوینده است و در حقیقت تصدیق هویت گوینده حالت خاصی از شناسایی گوینده‌ی مجموعه‌ی باز می‌باشد که در هر لحظه فقط یک گوینده در سامانه وجود دارد. به همین دلیل شناسایی گوینده‌ی مجموعه‌ی باز از تصدیق هویت گوینده پیچیده‌تر است.

در این پایان‌نامه به علت اینکه روش وابسته به متن برای شناسایی گوینده کاربرد عملی خاصی ندارد ما روش مستقل از متن را انتخاب کرده‌ایم. علاوه بر این، چون هدف ما در این پایان‌نامه سرعت بخشیدن به مرحله‌ی شناسایی است و از آنجایی که فرآیند زمان‌بر در این مرحله، جستجوی خطی روی تمام گویندگان سامانه است، حالت مجموعه بسته انتخاب شده است. این فرآیند زمان‌بر در هر دو کاربرد مجموعه‌ی باز و مجموعه‌ی بسته مشترک است.

۳-۲ اجزاء و مراحل سامانه‌های شناسایی گوینده

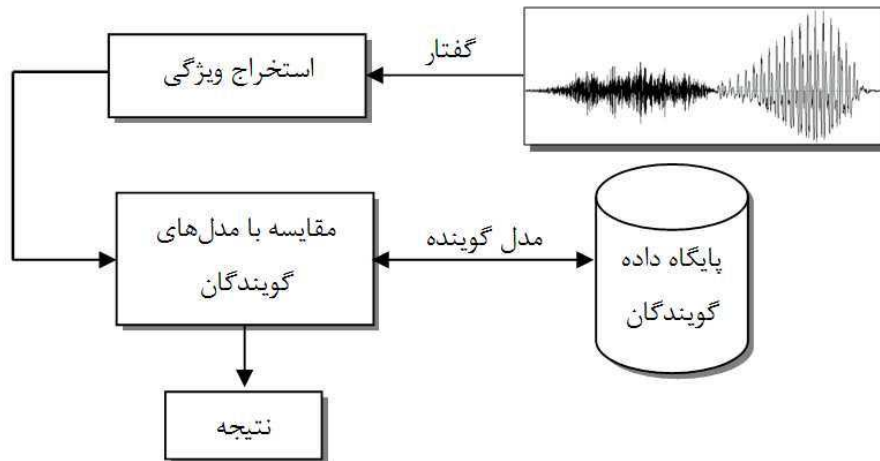
به طور کلی هر سامانه‌ی شناسایی گوینده شامل دو مرحله‌ی مجزا است که مرحله‌ی اول مرحله‌ی ثبت-نام (آموزش) و مرحله‌ی دوم مرحله‌ی شناسایی (آزمایش) است. در مرحله‌ی ثبت‌نام، ابتدا از گفتار آموزشی هر گوینده ویژگی‌هایی استخراج می‌شود. سپس این ویژگی‌ها برای مدل کردن خصوصیات صوتی گوینده استفاده می‌شود. در مرحله‌ی مدل‌سازی، برای هر گوینده مدل واحدی ساخته می‌شود و در نهایت این مدل‌ها در پایگاه داده‌ای از گویندگان ذخیره می‌شوند. شکل ۱-۲ نمودار جعبه‌ای مرحله‌ی ثبت‌نام را نشان می‌دهد.



شکل ۱-۲ نمودار جعبه‌ای مرحله‌ی ثبت‌نام

بعد از اینکه مدل گویندگان ساخته شد، سامانه آماده‌ی انجام عمل شناسایی است. در مرحله‌ی شناسایی ابتدا از گفتار ورودی ویژگی‌هایی استخراج می‌کنیم. استخراج ویژگی در این مرحله نیز شبیه مرحله‌ی ثبت‌نام است. بعد از استخراج ویژگی نوبت به انجام عمل شناسایی می‌رسد. در این مرحله،

ویژگی‌های استخراج شده با مدل تمام گویندگان مقایسه می‌شود و گوینده‌ای که مدل آن بیشترین شباهت را با این ویژگی‌ها داشته باشد به عنوان گوینده‌ی هدف انتخاب می‌شود. عمل مقایسه بسته به نوع مدلی که برای گویندگان انتخاب می‌شود متفاوت است. شکل ۲-۲ نمودار جعبه‌ای مرحله‌ی شناسایی را نشان می‌دهد.



شکل ۲-۲ نمودار جعبه‌ای مرحله‌ی شناسایی

بیشترین محاسبات در مرحله‌ی شناسایی را عمل مقایسه به خود اختصاص می‌دهد. برای افزایش سرعت سامانه باید تا می‌توان محاسبات این قسمت را کاهش داد.

۴-۲ پردازش اولیه و استخراج ویژگی

۱-۴-۲ مقدمه

هدف از بلوک پردازش اولیه و استخراج ویژگی، این است که یک سری ویژگی‌هایی که توانایی نشان دادن اطلاعات و ویژگی‌های فردی هر گوینده را دارد از سیگنال گفتار استخراج کند. در این بلوک بسته به کاری که قرار است انجام شود ممکن است عملیات دیگری مثل حذف نویز و سایر پیش‌پردازش‌ها^۱ انجام گیرد.

^۱ Pre-Processing

ویژگی‌هایی که در پردازش گفتار و بویژه شناسایی گوینده استفاده می‌شوند به دو دسته‌ی کوتاه مدت^۱ و بلند مدت^۲ تقسیم می‌شوند. ویژگی‌های کوتاه مدت از قاب‌های کوتاهی از گفتار استخراج می‌شوند و نشان‌دهنده‌ی ساختار فیزیکی مسیر صوتی گوینده می‌باشند. از طرف دیگر، ویژگی‌های بلند مدت از محدوده‌ی طولانی‌تری از گفتار استخراج می‌شوند و نشان‌دهنده‌ی اطلاعات آوایی و لغوی گوینده می‌باشند. برای انتخاب و استخراج هرگونه ویژگی، باید نیازمندی‌هایی مورد توجه قرار گیرد. این نیازمندی‌ها ممکن است استفاده یا عدم استفاده از ویژگی‌های مختلف را تحت تأثیر قرار دهد. از این نیازمندی‌ها می‌توان به موارد زیر اشاره کرد [Naik_90]، [Atal_76]:

- قابلیت جداسازی گویندگان را داشته باشد در حالی که نسبت به تغییرات ویژگی‌های صوتی یک گوینده حساس نباشد.
- به راحتی و در زمانی معقول قابل استخراج باشد.
- نسبت به گذشت زمان و تغییرات سنی و همچنین حالات گوینده (مثل سرماخوردگی) مقاوم باشد.
- به صورت طبیعی و متداوم در گفتار تکرار شود.
- تغییرات کمی در محیط‌های مختلف و نویزی داشته باشد.
- قابل تقلید توسط افراد دیگر نباشد.

در کاربردهای عملی امکان اینکه یک ویژگی تمام نیازهای بالا را برآورده کند، وجود ندارد و همیشه یک رابطه‌ی سبک و سنگینی بین آنها وجود دارد و بسته به کاربردهای مختلف اهمیت این نیازها تغییر می‌کند.

¹ Short-term

² Long-term

۲-۴-۲ ویژگی‌های کوتاه مدت

اکثر سامانه‌های بازشناسایی گوینده از ویژگی‌های کوتاه مدت که از بخش کوتاهی از گفتار استخراج می‌شوند، استفاده می‌کنند. از ویژگی‌های کوتاه مدت که بیشتر مورد استفاده قرار می‌گیرند می‌توان به ضرایب کپسترال مقیاس فرکانس مل^۱ (MFCCs) [Davis_80]، ضرایب کپسترال پیشگویی خطی^۲ (LPCCs) [Makhoul_75] و ضرایب پیشگویانه خطی ادراکی^۳ (PLP) [Hermansky_90] اشاره کرد. از آنجایی که ضرایب کپسترال مقیاس فرکانس مل در بازشناسایی گوینده زیاد استفاده شده است و دقت بهتری نسبت به سایر ویژگی‌های کوتاه مدت دارد ما نیز در این پایان‌نامه از آن استفاده کرده‌ایم. در ادامه این ویژگی به تفصیل شرح داده می‌شود.

۳-۴-۲ ویژگی‌های بلند مدت

در سال‌های اخیر تلاش‌های زیادی برای استفاده از ویژگی‌های بلند مدت به منظور بالا بردن دقت و همچنین مقاوم‌سازی روش‌های بازشناسایی گوینده در سامانه‌های مستقل از متن شده است. طرفداران این ویژگی‌ها معتقد هستند که با استفاده از ویژگی‌های کوتاه مدت به تنهایی نمی‌توان دقت سامانه‌های بازشناسایی را بالاتر برد [Campbell_03] و به همین دلیل استفاده از ویژگی‌های بلند مدت را امری ضروری می‌دانند.

ویژگی‌های بلند مدت به سه طبقه‌ی آوایی^۴، لغوی^۵ و نوایی^۶ تقسیم می‌شوند [Buyuk_11]. دنباله‌ی زمانی آواها و تلفظ‌های منحصر به فرد یک گوینده در دسته آوایی قرار می‌گیرند. این ویژگی‌ها اغلب با استفاده از خروجی یک سامانه‌ی بازشناسایی گفتار استخراج می‌شوند. ویژگی‌های لغوی شامل تکرار کلمات خاص یا تکه کلام‌های یک گوینده می‌باشد که اینها هم از خروجی سامانه‌ی بازشناسایی گفتار

¹ Mel-frequency cepstral coefficients

² Linear prediction cepstral coefficients

³ Perceptual linear predictive

⁴ Phonetic

⁵ Lexical

⁶ Prosodic

استخراج می‌شوند. ویژگی‌های نوایی شامل تغییرات آهنگ گفتار، بلندی و طول زمانی قسمت‌های کلمه و حتی جمله می‌باشد. در این دسته فرکانس گام^۱، دیرش^۲ و انرژی از مهمترین ویژگی‌ها می‌باشند.

به طور معمول استخراج ویژگی‌های بلند مدت نیاز به داده‌های آموزشی بیشتری دارند. علاوه بر این، برای استخراج بیشتر این ویژگی‌ها نیاز به یک سامانه‌ی بازشناسی گفتار با دقت بالا است. این محدودیت‌ها باعث شده که استفاده از این ویژگی‌ها رونق چندانی پیدا نکند. با توجه به این محدودیت‌ها ما در این پایان‌نامه فقط از ویژگی‌های کوتاه مدت استفاده کرده‌ایم.

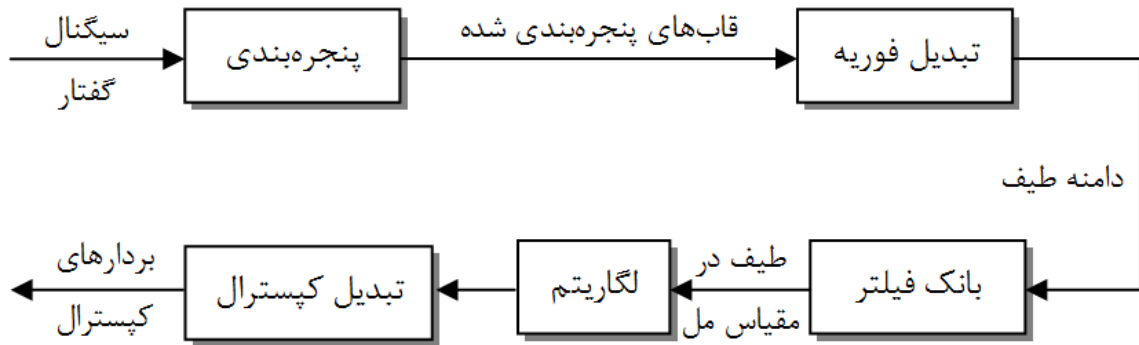
۴-۴-۲ ضرایب کیسترال مقیاس فرکانس مل

شکل ۲-۳ نمودار جعبه‌ای مراحل استخراج ضرایب MFCC را نشان می‌دهد. برای استخراج این ضرایب، ابتدا گفتار را به قاب‌های هم‌پوشان تقسیم کرده و سپس از هر قاب یک بردار ضرایب استخراج می‌کنند. طول قاب‌ها در حدود ۲۰ تا ۳۰ میلی‌ثانیه در نظر گرفته می‌شود. این طول را طوری در نظر می‌گیرند که ویژگی ایستا بودن سیگنال گفتار در هر قاب حفظ شود. جابه‌جایی قاب‌ها برای حفظ هم‌پوشانی قاب‌های مجاور، در حدود ۱۰ الی ۱۵ میلی‌ثانیه است. بعد از اینکه سیگنال گفتار را به قاب‌های هم‌پوشان تقسیم کردند، برای اینکه تأثیر ناپیوستگی در لبه قاب‌ها را کم کنند هر قاب را در پنجره‌ای ضرب می‌کنند. معمولاً از پنجره‌های همینگ^۳ یا هنینگ^۴ برای این کار استفاده می‌شود. ما در اینجا از پنجره‌ی همینگ استفاده کرده‌ایم. فرمول (۱-۲) معادله‌ی حوزه‌ی زمان این پنجره را نشان می‌دهد:

$$w(n) = 0.54 - 0.46 \cos\left(\frac{2\pi n}{N-1}\right) \quad (1-2)$$

در معادله‌ی بالا N نشان دهنده‌ی طول پنجره و n نشان دهنده‌ی عنصر n -ام از پنجره‌ی همینگ است.

¹ Pitch
² Duration
³ Hamming
⁴ Hanning



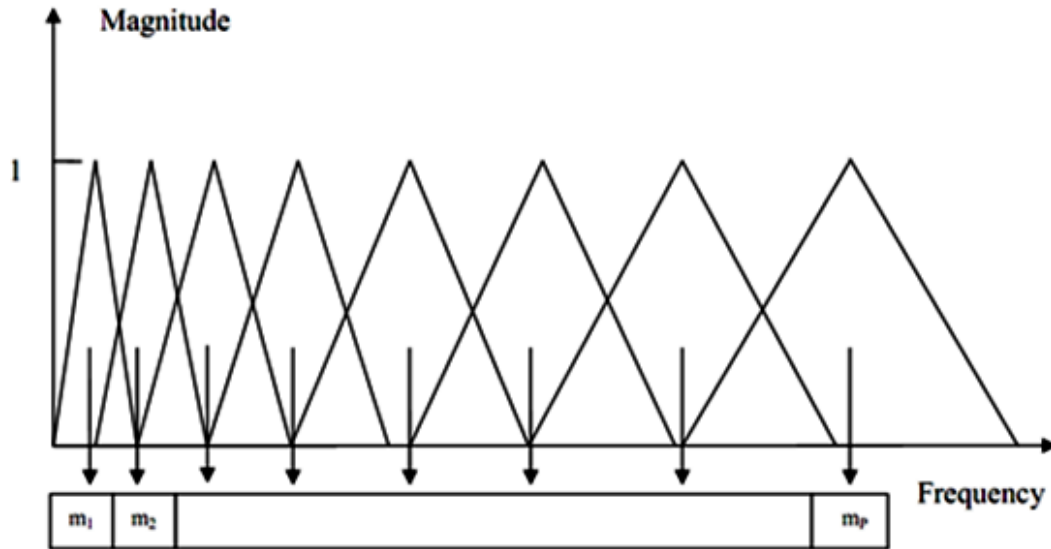
شکل ۳-۲ نمودار جعبه‌ای مراحل استخراج ضرائب کپسترال

بعد از اینکه قاب را در پنجره ضرب کردند از آن تبدیل فوریه گسسته (DFT^1) گرفته و دامنه‌ی تبدیل را در حوزه‌ی فرکانس حساب می‌کنند. تا اینجا دامنه‌ی DFT قاب گفتار بدست آمده است که آن را با X_k نشان می‌دهند، (k اندیس DFT است). حال X_k را در بانک فیلتر مل متناظر ضرب کرده و نتایج را با هم جمع می‌کنند. نتیجه اینکه هر ضریب بانک فیلتر (m_i) یک جمع وزن‌دار از دامنه‌ی طیف در کانال i -ام از بانک فیلتر می‌باشد که مطابق فرمول زیر محاسبه می‌شود:

$$m_i = \sum_k w_{k,i} X_k \quad (2-2)$$

جایی که $w_{k,i}$ وزن فیلتر در کانال i -ام را نشان می‌دهد. این وزن‌ها در خارج از باندهای کانسی فیلتر صفر در نظر گرفته می‌شوند. بانک فیلتر مل استفاده شده در شکل ۴-۲ نشان داده شده است:

¹ Discrete Fourier transform



شکل ۴-۲ بانک فیلتر مل

در محاسبات ضرایب MFCC، بانک فیلترها دارای توزیع یکسانی در حوزه‌ی فرکانس نیستند ولی آنها به صورت متساوی الفاصله در طول مقیاس مل که توسط فرمول (۳-۲) تعریف می‌شود پخش شده‌اند.

$$Mel(f) = 2595 * \log\left(1 + \frac{f}{700}\right) \quad (3-2)$$

علت عدم توزیع یکنواخت این است که پاسخ فرکانسی گوش انسان به صورت غیریکنواخت در حوزه‌ی فرکانس توزیع شده است. همچنین، علت استفاده از توزیع غیرخطی این است که اطلاعاتی که توسط مؤلفه‌های فرکانس پایین منتقل می‌شوند برای انسان‌ها اهمیت بیشتری نسبت به اطلاعات مؤلفه‌های فرکانس بالا دارند. توزیع غیرخطی مل نسبت به سایر توزیع‌های خطی در حوزه‌ی فرکانس باعث افزایش دقت سامانه‌های بازشناسی می‌شود. در پایان با استفاده از فرمول (۴-۲) از لگاریتم ضرایب بانک فیلتر، تبدیل کسینوسی گسسته (DCT¹) گرفته می‌شود. با این کار ضرایب نهایی بدست می‌آیند.

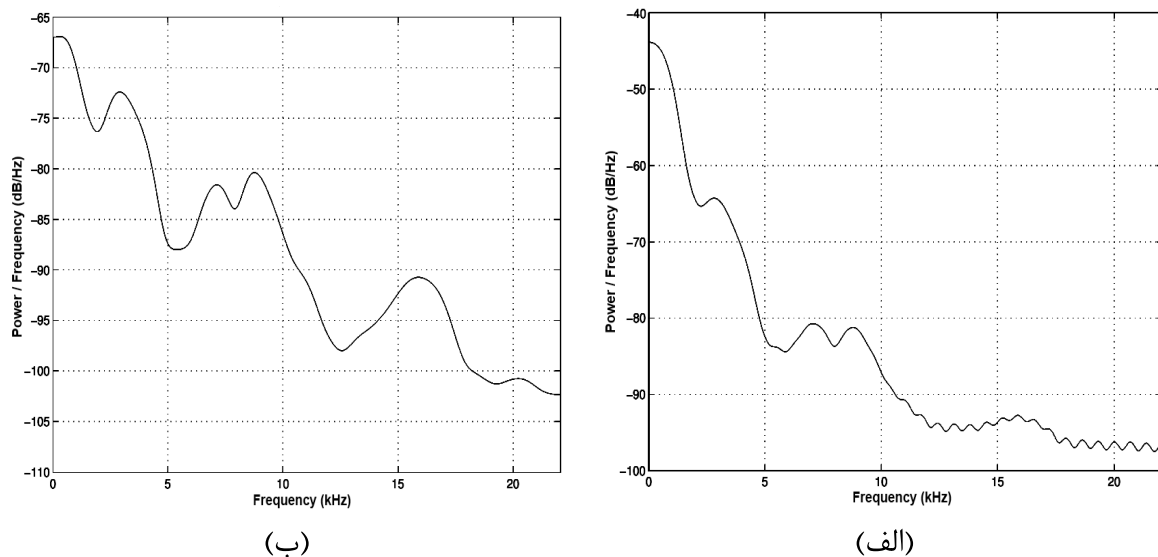
$$c_i = \sqrt{\frac{2}{N}} \sum_{j=1}^N \log(m_j) \cos\left(\frac{\pi i}{N}(j-0.5)\right) \quad (4-2)$$

در فرمول بالا N تعداد کانال‌های بانک فیلتر می‌باشد. تنها تعداد کمی از ضرایب DCT نگه داشته می‌شوند تا یک بردار ویژگی را نشان دهند. این ضرایب را ضرایب ایستا می‌گویند.

¹ Discrete cosine transform

۵-۴-۲ فیلتر پیش-تاکید^۱

خواص سیستم صوتی انسان، جریان هوای حنجره و تشعشع^۲ لب باعث تعدیل رفتن مؤلفه‌های فرکانس بالای گفتارهای صدادار می‌شود [Harrington_99]. برای گفتارهای صدادار، طیف سیگنال خروجی از دهانه‌ی حنجره تقریباً دارای شیب -12db/octave است. زمانی که این سیگنال‌ها از لب‌ها انتشار پیدا می‌کنند، طیف آنها در فرکانس‌های بالا تقویت شده و تقریباً $+6\text{db/octave}$ از شیب قبلی جبران می‌شود و در نتیجه طیف سیگنال خروجی در مجموع با شیب -6db/octave از فرکانس‌های پایین به فرکانس‌های بالا کاهش می‌یابد. در شکل ۵-۲ (الف) کاهش شدید توان سیگنال در فرکانس‌های بالا قابل مشاهده است.



شکل ۵-۲ نمودار چگالی طیف توان، (الف) سیگنال اصلی با فرکانس نمونه‌برداری ۴۴۱۰۰، (ب) پیش-تاکید شده همان سیگنال [Beigi_11]

در سیستم شنوایی انسان، حلوزن گوش^۳ بر اساس بازخورد مغز پردازش‌هایی روی گفتار انجام می‌دهد و بعضی از فرکانس‌ها را تقویت می‌کند. این کار باعث می‌شود که انسان‌ها به راحتی قسمت کم انرژی سیگنال در فرکانس‌های بالا را دریافت کنند. برای شبیه‌سازی سیستم شنوایی انسان و جبران این اثر و جلوگیری از افزایش اثر مؤلفه‌های فرکانس پایین گفتار، قبل از استخراج ویژگی باید فیلتر پیش-تاکید اعمال شود. در شکل ۵-۲ (ب) همان سیگنال در حوزه‌ی فرکانس نشان داده شده است با این

¹ Preemphasis

² Radiations

³ Cochlea

تفاوت که سیگنال از فیلتر پیش-تاکید عبور کرده است. توجه کنید که مقدار توان برای تمام فرکانس‌ها کاهش پیدا کرده است، اما توان نسبی در فرکانس‌های مختلف دارای توزیع بهتری است.

معمولاً فیلتر پیش-تاکید بوسیله‌ی فیلتر کردن سیگنال گفتار توسط یک فیلتر با پاسخ ضربه‌ی محدود (FIR) مرتبه‌ی اول اعمال می‌شود. معادله‌ی این فیلتر به فرم زیر است:

$$F(z) = 1 - kz^{-1} \quad (0 < k < 1) \quad (5-2)$$

در فرمول بالا، k فاکتور پیش‌تاکید است که مقدار متداول برای آن 0.97 می‌باشد. همان طور که مشخص است سیگنال خروجی در حوزه‌ی زمان به فرم زیر بدست می‌آید:

$$y(n) = s(n) - k * s(n-1) \quad (6-2)$$

جایی که $s(n)$ سیگنال ورودی و $y(n)$ سیگنال خروجی فیلتر می‌باشد.

۶-۴-۲ ویژگی‌های پویا

معمولاً دقت سامانه‌های بازشناسی با اضافه کردن اطلاعات پویا (مشتق زمانی) به پارامترهای ایستا بطور چشم‌گیری افزایش می‌یابد. بعد از اینکه پیش‌پردازش‌های اولیه روی ویژگی‌های ایستا اعمال شد، اطلاعات پویا که نمایانگر تغییرات زمانی بردارهای کپسترال هستند به ویژگی‌ها اضافه می‌شوند. بطور معمول از مشتق اول و دوم استفاده می‌شود و مشتقات بالاتر به دلیل اینکه تأثیر اندکی بر دقت دارند مورد استفاده قرار نمی‌گیرند. برای بدست آوردن ضرایب پویای مرتبه‌ی اول (مشتق اول) از فرمول درون-یابی زیر استفاده می‌شود:

$$d_t = \frac{\sum_{\theta=1}^{\ominus} \theta (c_{t+\theta} - c_{t-\theta})}{2 \sum_{\theta=1}^{\ominus} \theta^2} \quad (7-2)$$

جایی که dt ضریب دلتا در زمان t می‌باشد که با استفاده از ضرایب ایستای قبلی و بعدی محاسبه می‌شود. پارامتر \ominus اندازه‌ی پنجره برای بدست آوردن ضرایب دلتا را مشخص می‌کند. برای بدست آوردن

ضرایب مراتب بالاتر از فرمول مشابه‌ای بر روی ضرائب پویای مرتبه‌ی پایین‌تر استفاده می‌شود [Young_06].

۵-۲ مدل کردن گویندگان

۱-۵-۲ مقدمه

در بازشناسی گوینده، روش‌های مختلفی برای مدل کردن گویندگان وجود دارد. این روش‌ها به دو دسته-ی کلی مولد^۱ و تمایزی^۲ تقسیم می‌شوند. در روش مدل‌سازی مولد، مدل گویندگان به سادگی توسط تابع چگالی احتمال نشان داده می‌شود. از روش‌های مختلف این دسته می‌توان به چندی‌ساز برداری (VQ^3) [Soong_85]، مدل مخلوط گاوسی، مدل مخفی مارکوف [Deshpande_08] و پیچش زمانی پویا (DTW^4) [Wutiw WATCHAI_99] اشاره کرد. در طرف دیگر، در روش‌های مختلف دسته‌ی تمایزی از جمله ماشین بردار پشتیبان ($SVMs^5$) [Xing_12]، [Ruiling_11] و شبکه‌های عصبی مصنوعی ($ANNs^6$) [Rao_10]، [Hossain_07] یک مرز بین گویندگان مختلف تخمین زده می‌شود.

اگرچه در سال‌های اخیر روش‌های متفاوتی که از مدل‌سازی تمایزی استفاده کرده و در حالت مستقل از متن کار بازشناسی گوینده را انجام می‌دهند موفقیت‌های قابل قبولی داشته‌اند، ما در این پایان‌نامه از روش‌های مولد استفاده کرده‌ایم. علت این انتخاب این است که در روش‌های مولد از جمله مدل مخلوط گاوسی حجم محاسبات در مرحله‌ی شناسایی بسیار بالا است و این در حالی است که محاسبات مورد نیاز برای روش‌های تمایزی کمتر است. از آنجا که هدف اصلی ما در این پایان‌نامه کم کردن محاسبات مرحله‌ی شناسایی و در نتیجه افزایش سرعت این مرحله است، از مدل‌های مولد استفاده شده است.

¹ Generative
² Discriminative
³ Vector quantization
⁴ Dynamic time warping
⁵ Support vector machine
⁶ Artificial neural networks

۲-۵-۲ مدل مخلوط گاوسی

مدل مخلوط گاوسی پرکاربردترین روش مدل‌سازی بویژه در کاربردهای مستقل از متن است. این مدل به ویژگی‌های زمانی و دنباله‌ای گفتار حساس نیست و به همین خاطر در کاربردهای وابسته به متن کمتر استفاده می‌شود. خصوصیات زمانی و چگونگی تغییرات گفتار در طول زمان حاوی اطلاعاتی برای شناسایی گوینده است که در کاربردهای وابسته به متن مورد توجه قرار می‌گیرد. تابع چگالی احتمال در مدل مخلوط گاوسی، مجموع وزن‌داری از M تابع گاوسی (مخلوط) به صورت زیر است:

$$\Pr(\mathbf{x} | \lambda) = \sum_{i=1}^M w_i g(\mathbf{x} | \mu_i, \Sigma_i) \quad (۸-۲)$$

جایی که \mathbf{x} یک بردار D بعدی با مقدار پیوسته و w_i وزن مخلوط i -ام است. وزن مخلوط‌ها باید به گونه‌ای بدست بیایند که شرط زیر را ارضاء کنند:

$$\sum_{i=1}^M w_i = 1 \quad (۹-۲)$$

در رابطه (۸-۲)، هر مخلوط یک توزیع گاوسی می‌باشد، $g(\mathbf{x} | \mu_i, \Sigma_i)$ ، که با بردار میانگین μ_i و

ماتریس کواریانس^۱ Σ_i مشخص می‌شود. فرمول توزیع گاوسی در رابطه‌ی زیر نشان داده شده است:

$$g(\mathbf{x} | \mu_i, \Sigma_i) = \frac{1}{(2\pi)^{D/2} |\Sigma_i|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \mu_i)' \Sigma_i^{-1} (\mathbf{x} - \mu_i)\right] \quad (۱۰-۲)$$

مجموعه‌ی تمام پارامترهای مدل مخلوط گاوسی توسط سه‌تایی زیر نشان داده می‌شود:

$$\lambda = \{w_i, \mu_i, \Sigma_i\} \quad i = 1, \dots, M$$

^۱ Covariance matrix

فرم کلی مدل‌سازی در این روش از ماتریس کواریانس کامل استفاده می‌کند، اما در عمل از ماتریس کواریانس قطری استفاده می‌شود. دلیل استفاده از ماتریس قطری این است که ماتریس کواریانس قطری از نظر محاسباتی نسبت به ماتریس کواریانس کامل بهینه‌تر است. علاوه بر این مزیت، دقت شناسایی در این حالت نسبت به ماتریس کامل بیشتر است [Reynolds_00]. زمانی که از ماتریس کواریانس کامل استفاده شود، تعداد پارامترهایی که باید تخمین زده شوند بیشتر می‌شود و این افزایش تعداد پارامترها نیاز به داده‌های آموزشی بیشتری نسبت به حالت قطری دارد. در کاربردهای عملی معمولاً میزان داده‌های آموزشی کم می‌باشد که در این حالات بهتر است از ماتریس کواریانس قطری استفاده کنیم. علاوه بر این، از آنجایی که مخلوط‌های گاوسی با یکدیگر تمام فضای ویژگی‌ها را مدل می‌کنند، استفاده از ماتریس کواریانس کامل حتی اگر ویژگی‌هایی که استفاده می‌شوند به صورت آماری مستقل نباشند ضروری نیست. ترکیب خطی توابع گاوسی با کواریانس قطری امکان مدل کردن رابطه بین ابعاد بردارهای ویژگی را دارد [Reynolds_95a]. تاثیر استفاده از مجموعه‌ای از توابع گاوسی با ماتریس کواریانس کامل را می‌توان با مجموعه‌ای بزرگتر از توابع گاوسی با ماتریس کواریانس قطری ایجاد کرد.

همان طور که مشاهده نمودید هر مدل مخلوط گاوسی دارای پارامترهای نامعلوم زیادی است. این پارامترها باید برای هر گوینده به روشی تخمین زده شوند. برای تخمین پارامترهای مدل، دو رویکرد کلی وجود دارد. در رویکرد اول پارامترهای مدل هر گوینده به صورت مجزا تخمین زده می‌شوند. رویکرد دوم به این صورت می‌باشد که ابتدا یک مدل جهانی از داده‌های آموزشی تمام گویندگان می‌سازیم. سپس با استفاده از پارامترهای این مدل برای هر گوینده، یک مدل منحصر به فرد بدست می‌آوریم. برای هر کدام از این دو رویکرد روش‌های مختلفی وجود دارد. موفقترین روش برای بدست آوردن مستقیم پارامترهای مدل هر گوینده الگوریتم بیشینه‌سازی امید ریاضی (EM^1) است. از روش‌های رویکرد دوم می‌توان به

¹ Expectation-Maximization (EM) algorithm