

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

هست کلید درج حکیم (نظامی کنجوی)



دانشکده مهندسی برق و رباتیک

پایان نامه کارشناسی ارشد مهندسی مخابرات سیستم

تشخیص گوینده بر اساس ویژگی‌های استخراجی ترکیبی در محیط با چند

گوینده

نگارنده: میترا جهانیان

استاد راهنما:

دکتر حسین مروی

استاد مشاور:

دکتر سید مسعود میر رضایی

شهریور ۱۳۹۵

پروردگارا:

به پاس تعبیر عظیم و انسانی‌شان از کلمه ایثار و از خودگذشتگان

به پاس عاطفه سرشار و گرمای امیدبخش وجودشان که در این سردترین روزگاران

بهترین پشتیبان است

به پاس قلب‌های بزرگشان که فریادرس است و سرگردانی و ترس در پناهشان به

شجاعت می‌گراید

و به پاس محبت‌های بی‌دریغشان که هرگز فروکش نمی‌کند

این مجموعه را به پدر و مادر عزیزم ، برادران دوست‌داشتنی‌ام و دوستان عزیزم

تقدیم می‌کنم

تقدیم به

معلمان فرهیخته‌ای که صادقانه و عاشقانه تلاش می‌کنند تا نقالی دانسته‌ها را به

نقادی اندیشه‌ها تبدیل سازند

و دانش‌آموزانی که روح خلاق گریزپایشان، تکرار را نمی‌خواهد

تقدیر و تشکر:

با سپاس فراوان از لطف خدای مهربان.

با تشکر از دو استاد بزرگوارم که شایسته‌ی هر نوع سپاس، تجلیل و تکریم‌اند؛

جناب آقای دکتر حسین مروی؛ استاد راهنمای ارجمند که با ایجاد عشق به نوشتن، صبورانه، با ارائه‌ی رهنمودها، انتقادهای و پیشنهادهایشان، در تمامی مراحل اجرای پایان‌نامه مرا حمایت و تشویق نمودند و جناب آقای دکتر سید مسعود میر رضایی؛ استاد مشاور محترم که با نظرهای اصلاحی ارزنده‌ی خود، ضمن دلگرمی بنده، موجب تکمیل این اثر شدند.

همچنین از جناب دکتر معروضی و دکتر سلیمانی به جهت زحمت فراوان در داوری این اثر و از استادان محترمی که در طول دوران تحصیلی‌ام، جهت آموزش و ارتقای علمی بنده، زحمت کشیده‌اند سپاسگزارم.

همچنین از جناب آقای دکتر گرایلو مسئول محترم آزمایشگاه پردازش سیگنال و دوستان عزیز و گرامی آقایان حسینی و غفاریان فر و خانم‌ها آزاده اصغری و یلدا عابدینی و مریم صادقی کمال تشکر رادارم و از مسئولین محترم دانشکده آقایان یونسپان و رضایی و خانم جعفری صمیمانه متشکرم. باشد که خداوند عمر باعزت و روحی متعالی نصیب تک‌تک این عزیزان نماید.

میترا جهانیان

شهریورماه سنه‌ی یک هزار و سیصد و نودوپنج خورشیدی

تعهدنامه

اینجانب **میترا جهانیان** دانشجوی دوره کارشناسی ارشد رشته مهندسی برق مخابرات سیستم دانشکده برق و رباتیک دانشگاه شاهرود نویسنده پایان نامه تشخیص گوینده بر اساس ویژگی‌های استخراجی ترکیبی در محیط با چند گوینده تحت راهنمایی دکتر حسین مروی و مشاوره و راهنمایی دکتر مسعود میر رضایی متعهد می‌شوم.

- تحقیقات در این پایان نامه توسط این جانب انجام شده است و از صحت و اصالت برخوردار است.
- در استفاده از نتایج پژوهش‌های محققان دیگر به مرجع مورداستفاده استناد شده است.
- مطالب مندرج در پایان نامه تاکنون توسط خود یا فرد دیگری برای دریافت هیچ نوع مدرک یا امتیازی در هیچ جا ارائه نشده است.
- کلیه حقوق معنوی این اثر متعلق به دانشگاه شاهرود می‌باشد و مقالات مستخرج بانام « دانشگاه شاهرود » و یا « Shahrood University » به چاپ خواهد رسید.
- حقوق معنوی تمام افرادی که در به دست آمدن نتایج اصلی پایان نامه تأثیرگذار بوده‌اند در مقالات مستخرج از پایان نامه رعایت می‌گردد.
- در کلیه مراحل انجام این پایان نامه ، در مواردی که از موجود زنده (یا بافت‌های آن‌ها) استفاده شده است ضوابط و اصول اخلاقی رعایت شده است.
- در کلیه مراحل انجام این پایان نامه، در مواردی که به حوزه اطلاعات شخصی افراد دسترسی یافته یا استفاده شده است اصل رازداری ، ضوابط و اصول اخلاق انسانی رعایت شده است .

تاریخ

امضای دانشجو

مالکیت نتایج و حق نشر

کلیه حقوق معنوی این اثر و محصولات آن (مقالات مستخرج، کتاب، برنامه‌های رایانه‌ای، نرم‌افزارها و تجهیزات ساخته شده است) متعلق به دانشگاه شاهرود می‌باشد. این مطلب باید به نحو مقتضی در تولیدات علمی مربوطه ذکر شود.

چکیده:

امروزه در دنیای ارتباطات پردازش داده با سرعت و دقت بالا از اهمیت بسزایی برخوردار است. پردازش دادگان گفتاری نیز به دلیل کاربرد وسیع در تمامی جهات زندگی بشر سهم بسزایی را ایفا می‌کند. نمایه سازی گوینده یعنی تشخیص دهیم چه کسی چه زمانی صحبت می‌کند. هدف طراحی سیستم‌های نمایه سازی گوینده این است که تغییر در گوینده در فایل صوتی تشخیص داده شود و گفتار هر گوینده به درستی برچسب گذاری و دسته بندی شود. این فرآیند امروزه بانام Speaker Diarization شناخته شده است. در این پایان نامه، سیستمی طراحی کردیم که با استفاده از ویژگی‌های آکوستیکی¹ MFCC و مشتقات مرتبه اول و دوم آن به همراه ویژگی‌های انرژی و نرخ عبور از صفر ویژگی‌های گفتار و غیر گفتار را استخراج کند سپس با استفاده از فریم‌های قطعاً سکوت و قطعاً آهنگ به مدل سازی این دو می‌پردازد و در یک فرایند دوبخشی شامل حذف سکوت و حذف آهنگ به جداسازی گفتار از غیر گفتار و آهنگ موجود در فایل صوتی می‌پردازد. در ادامه با استفاده از بردارهای i در فضای برداری ویژگی‌ها به کاهش بُعد می‌پردازیم. در ادامه برای تشخیص تغییر در گوینده معیار فاصله را به کار برده و با خوشه بندی توسط برنامه ریزی خطی عدد صحیح (ILP^2) گفتار هر گوینده را جدا و برچسب گذاری و خوشه بندی می‌کنیم. به بهینه سازی پارامترها پرداختیم. پایگاه داده‌ی مورد استفاده AMI corpus می‌باشد. نتایج خوبی در خطای بازشناسی گوینده (DER) گزارش داده شد.

کلید واژگان: نمایه سازی گوینده، پردازش گفتار، برنامه عدد صحیح خطی، بردارهای i

¹ Mel Frequency Cepstral Coefficient

² Integer Linear Programming

فهرست مطالب

۱	فصل ۱ : مقدمه
۲	۱-۱- اهمیت گفتار
۳	۲-۱- پردازش گفتار
۳	۳-۱- تشخیص گوینده
۵	۴-۱- ساختار پایان نامه
۷	فصل ۲ : بیان مسئله و پیشینه تحقیق
۸	۱-۲- پیشینه بازشناسی گفتار
۹	۲-۲- انواع روش‌های بازشناسی گوینده
۱۰	۳-۲- انواع ویژگی‌های مورداستفاده در جداسازی گوینده
۱۳	۴-۲- بخش‌بندی
۱۳	۱-۴-۲- بخش‌بندی متریک
۱۶	۲-۴-۲- بخش‌بندی بر اساس مدل
۱۹	۵-۲- تشخیص فعالیت گفتار
۲۰	۶-۲- خوشه‌بندی کردن
۲۱	۱-۶-۲- مدل‌های گوینده برای خوشه‌بندی کردن
۲۴	۲-۶-۲- الگوریتم خوشه‌بندی

۲۸.....	۷-۲- روش‌های ارزیابی
۲۹.....	۲-۷-۱- معیار DER
۳۲.....	۲-۸-۱- پایگاه داده‌های مورد استفاده
۳۲.....	۲-۸-۱- AMI curpus
۳۳.....	۲-۸-۲- TIMIT
۳۳.....	۲-۹- خلاصه و جمع‌بندی فصل
۳۵.....	فصل ۳: روش پیشنهادی در بازشناسی گوینده
۳۶.....	۳-۱- مقدمه
۳۷.....	۳-۲- استخراج ویژگی
۳۸.....	۳-۳- تشخیص فعالیت گفتار
۴۰.....	۳-۳-۱- حذف سکوت
۴۱.....	۳-۳-۲- حذف موزیک
۴۳.....	۳-۳-۳- اندازه‌گیری‌های تشخیص فعالیت گفتار
۴۳.....	۳-۴-۱- ارزیابی تشخیص فعالیت گفتار
۴۳.....	۳-۴-۱- اندازه‌ی مدل‌های مخلوط گوسی‌ها حین جداسازی سکوت
	۳-۴-۲- ماتریس کوواریانس مدل‌های مخلوط گوسی در دسته‌بندی تکراری حذف سکوت و حذف
۴۴.....	آهنگ

- ۴۵..... اثر متوالی کردن حذف سکوت و حذف آهنگ ۳-۴-۳
- ۴۵..... نرخ عبور از صفر برای آموزش مدل‌های آهنگ ۴-۴-۳
- ۴۶..... بخش‌بندی گوینده ۵-۳
- ۴۷..... دسته‌بندی گوینده ۶-۳
- ۴۷..... انتخاب مدل گوینده ۱-۶-۳
- ۴۸..... استخراج بردار i ۲-۶-۳
- ۵۰..... انتخاب الگوریتم خوشه‌بندی ۳-۶-۳
- ۵۰..... HAC ۴-۶-۳
- ۵۱..... خوشه‌بندی ILP ۵-۶-۳
- ۵۵..... فصل ۴ : نتایج
- ۵۶..... مقدمه ۱-۴
- ۵۶..... پارامترهای بخش‌بندی ۲-۴
- ۵۷..... ارزیابی خوشه‌بندی گوینده‌ها ۳-۴
- ۵۸..... آزمایشات خوشه‌بندی سلسله‌مراتبی متراکم با مدل‌های گوینده‌ی مخلوط گوسی ۱-۳-۴
- ۵۸..... معیار آستانه‌ی فاصله ۲-۳-۴
- ۵۹..... HAC-۳-۳-۴ با مدل‌های گوینده‌ی بردار i
- ۵۹..... ILP-۴-۳-۴ با آزمایشات با مدل‌های گوینده‌ی مدل مخلوط گوسی

۴-۳-۵- خوشه‌بندی برنامه عدد صحیح خطی با مدل‌های گوینده‌ی بردار i ۶۰

فصل ۵: نتیجه‌گیری و پیشنهادات ۶۵

۵-۱- نتیجه‌گیری ۶۶

۵-۲- پیشنهادات ۶۷

فهرست منابع و مراجع ۶۹

فهرست اشکال:

- شکل ۱-۲: سیستم کلی نمایه سازی : الف (نمایه کلی خوشه بندی کلی ب) معماری خوشه بندی نمایه سازی کلی
- شکل ۲-۲: BIC برای تشخیص تغییر در گوینده
- شکل ۳-۲: پنجره‌ی متغیر برای تشخیص تغییر. فاصله‌ی بین نیمه‌های پنجره‌ها محاسبه و برحسب زمان نمایش داده شده است. قله و اوج در فاصله نمایانگر تغییر است.
- شکل ۴-۲: جستجوی افزایشی برای تشخیص تغییر گوینده. هر جستجو برای تنها یک تغییر است. بعد از هر یافتن نقطه تغییری جستجو از سر گرفته می‌شود. [۱۶]
- شکل ۵-۲: روش خوشه‌بندی متراکم سلسله مراتبی
- شکل ۶-۲: مثالی از محاسبه‌ی DER : الف) منبع ب) سیستم ما
- شکل ۱-۳: بلوک دیاگرام سیستم پیشنهادی
- شکل ۲-۳: حذف سکوت با استفاده از خوشه‌بندی تکراری
- شکل ۳-۳: حذف موسیقی با استفاده از یک آشکارساز موسیقی گفتار
- شکل ۴-۳: استخراج بردارهای i
- شکل ۵-۳: خوشه‌بندی ILP روی یک گراف کامل از مدل‌های گوینده [۲۹]
- شکل ۱-۴: HAC با حد آستانه‌ی فاصله برای مدل‌های مخلوط گوسی
- شکل ۲-۴: HAC با آستانه‌ی فاصله برای مدل‌های مخلوط گوسی گوینده
- شکل ۳-۴: ILP با آستانه فاصله برای مدل‌های مخلوط گوسی گوینده
- شکل ۴-۴: کارکرد خوشه یابی برنامه عدد صحیح با مدل‌های گوینده بردار i با ابعاد متفاوت از زیر فضای تغییر.

فهرست جدول‌ها:

جدول ۱-۲: ساده سازی محاسبه ی خطا

جدول ۲-۲: مثالی از بخش بندی برای محاسبه ی DER

جدول ۱-۳: خطای SAD برای اندازه های مختلف مدل های مخلوط گوسی برای مدل های گفتار

و سکوت (برحسب درصد.)

جدول ۲-۳: خطای SAD حالت ماتریس کوواریانس مدل های مخلوط گوسی (برحسب

درصد.)

جدول ۳-۳: متوالی کردن سیستم حذف آهنگ و حذف سکوت برای SAD (برحسب درصد.)

جدول ۴-۳: تعریف مجدد مدل های آهنگ با استفاده از فریم های با ZCR بالا (برحسب درصد.)

جدول ۱-۴ : مقادیر DER با بهترین الگوریتم های خوشه بندی

جدول ۲-۴: نتایج خطای بازشناسی بر اساس فواصل مختلف

جدول ۳-۴: نتایج خطای DER از دو مدل گوینده و دو الگوریتم خوشه بندی

جدول ۴-۴ : مقایسه با روش های دیگر

فهرست اختصارات

MFCC	mel frequency cepstral coefficient	ضرایب کپسترال فرکانسی در مقیاس مل
LPC	linear prediction coefficient	ضریب پیشگویی خطی
LFCC	linear frequency cepstral coefficient	ضریب کپسترال فرکانسی خطی
ZCR	zero crossing rate	نرخ عبور از صفر
CLR	cross likelihood ratio	نرخ احتمال متقابل
NCLR	normalized cross likelihood ratio	نرخ احتمال متقابل نرمال
BIC	bayesian inference criterion	معیار استنتاج بیزین
GMM	gussian mixture model	مدل مخلوط گوسی
ASR	activity speech recognition	تشخیص فعالیت گفتار
UBM	universal background model	مدل پس‌زمینه کلی
HAC	hierarchical agglomerative clustering	خوشه‌بندی سلسله مراتبی سنتی
ILP	integer linear program	برنامه خطی عدد صحیح
HMM	hidden marcov model	مدل مخفی مارکوف
DER	diarization error rate	نرخ خطای بازشناسی
MSR	missed speech rate	نرخ گفتار از دست‌رفته
FASR	false alarm speech rate	نرخ هشدار اشتباه گفتار
WCCN	within class covariance normalization	کوواریانس بین کلاسی نرمال
SAD	Speech activity detection	تشخیص فعالیت گفتار

فصل ١ : مقدمه

۱-۱- اهمیت گفتار

در دنیای مدرن امروز به جرئت می‌توان گفت متداول‌ترین راه ارتباطی جوامع بشری با یکدیگر استفاده از گفتار و صحبت کردن است. اهمیت گفتار در تمامی جنبه‌های زندگی انسان بروز کرده است. در تمام ارتباطها اعم از مکالمه متداول بین افراد، سخنرانی‌های عمومی، کلاس‌های آموزشی، سخنرانی‌های دیپلماتیک، گفتگوهای جنگی و یا حتی مکالمات سری نیز گفتار به کار می‌رود. لذا بررسی این ابعاد از راه ارتباطی ضرورت بسیار زیادی دارد. در میان داده‌های جمع‌آوری شده توسط بشر گفتار سهم بسزایی را به خود اختصاص می‌دهد. مؤسسات و ارگان‌ها و حتی خود ما به جمع‌آوری این داده‌ها می‌پردازیم. کارایی این اصوات ذخیره‌شده بسیار متنوع است و کاربردهای متفاوت و اهمیت‌های مختلفی دارند. از این‌رو بررسی و آنالیز و واکاوی این داده‌ها خود امری جدا در حوزه‌ی علوم را شکل داده است. جستجو و بازیابی و استخراج اطلاعات دلخواه از این حجم اطلاعات نیازمند سیستم‌های کامپیوتری و نرم‌افزارهای کامپیوتری دارد. ساختار بندی و آرشیو این اطلاعات مختلف پیکربندی خاصی را دارد. گزارشات تلویزیونی و رادیویی و مکالمات تلفنی بخش اعظم اطلاعات گفتاری را تشکیل می‌دهند. الگوهای گفتار را می‌توان به دو صورت دسته‌بندی کرد: دسته‌ی اول بر اساس اندام‌های صوتی و اندازه و شکل گلو و دهان و ویژگی‌های تارهای صوتی دسته‌بندی کرد و دسته‌ی دوم بر اساس الگوهای رفتاری همانند تحصیلات و موقعیت اجتماعی و سبک سخت گفتن دسته‌بندی می‌شوند. در اغلب قالب‌های امروزی دسته‌ی اول رایج‌تر است. تحقیقات وسیعی در این زمینه در سال‌های اخیر صورت گرفته است و نتایج بسیار خوبی حاصل شده است. کاربردهای متنوع این پردازش‌ها همچون تشخیص مجرم و جدا کردن صحبت‌های مهم یک شاهد و تشخیص گفتارهای متنوع از هم بر کسی پوشیده نیست.

۱-۲- پردازش گفتار

هدف از سیستم نهایی در این حوزه پاسخ به این سؤال است که چه کسی در چه زمانی صحبت می‌کند؟ این حوزه خود دارای بخش‌های مختلفی از جمله: قطعه‌بندی گوینده^۱، تشخیص گوینده^۲، رونویسی قوی^۳ و اندیس گذاری قوی^۴ است. از چنین سیستم‌هایی برای جابه‌جایی راحت در داده‌های صوتی در فایل‌های صوتی که متعلق به چند گوینده است، استفاده می‌شود. هدف نهایی در این کاربردها، پیاده‌سازی روش‌هایی برای افراز فایل صوتی به نواحی است که در آن‌ها گوینده‌ای خاص صحبت می‌کند.

با افزایش تعداد مدارک متنی موجود در اینترنت، نیاز به فن‌هایی نظیر فهرست‌نگاری متن به‌منظور تسهیل دسترسی و جستجو در این مدارک افزایش پیدا کرد. نظیر همین نیاز نیز با افزایش تعداد مدارک صوتی نظیر سخنرانی‌ها و مصاحبه‌ها و گردهمایی‌ها و ... ایجاد شد. به‌طور مشخص دسترسی به مدارک صوتی بسیار سخت‌تر از دسترسی به متن است و گوش دادن به فایل صوتی ضبط‌شده بیشتر از خواندن متن زمان‌بر است و فهرست‌نگاری دستی مدارک صوتی در مقایسه با فهرست‌نگاری متن، مشکل است. راه‌حل پیشنهادی جهت رفع این مشکل فهرست‌نگاری خودکار مدارک صوتی^۵ است.

۱-۳- تشخیص گوینده

اولین بار سیستم‌های تشخیص گوینده توسط کمپانی NIST در سال ۱۹۹۹ ارائه شد [۱]. در سال ۲۰۰۱، پلکان و سیدهارون به همراه گروهشان با استفاده از کم کردن اثر نویز روی سیگنال،

¹ Speaker Segmentation

² Speaker Recognition

³ Robust Transcription

⁴ Robust indexing

⁵ Automatic Audio Indexing

بهبودهایی در نتایج سیستم ایجاد کردند. در سال ۲۰۰۷، بولیان و کنی با به‌کارگیری بردارهای ویژگی دیگری و ادغام روش‌های قبلی و استفاده از مدل‌های گوسی در سیستم، نتایج متفاوتی به دست آوردند [۲]. در همان سال با استفاده از ویژگی‌های سیگنال‌های صوتی مانند فرکانس پیچ، انرژی، فرکانس‌های ماکزیمم سیگنال و سه ویژگی دیگر توسط یامانا و ماتسونوگا بهبود قابل‌توجهی در بخش-بندی گوینده حاصل شد. در سال‌های بعد با انجام روش‌های دیگر بر روی قسمت‌های مختلف آن تا به امروز سیستم تشخیص گوینده در حال تکمیل بوده‌اند. [۳]

هدف از این پایان‌نامه، طراحی و پیاده‌سازی سیستمی است که بتواند در یک فایل صوتی که شامل گفتار چندین گوینده است، تغییر در گوینده را تشخیص و جدا کند و تاحدامکان، گفتار هر گوینده را دسته‌بندی نماید. این سیستم می‌تواند شامل دو بخش اساسی به‌صورت زیر باشد:

- بخش‌بندی گوینده

- خوشه‌بندی گوینده

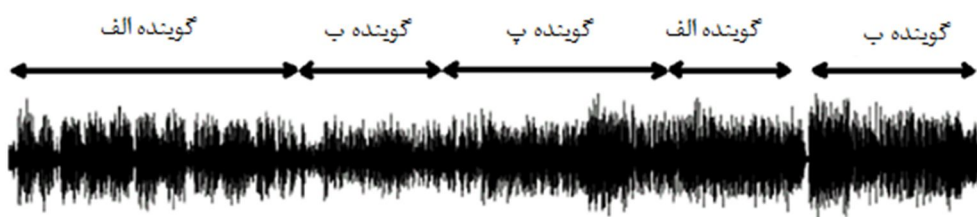
در قسمت بخش‌بندی^۱، تقسیم سیگنال گفتاری به قسمت‌هایی که تنها شامل گفتار یک گوینده هستند، صورت می‌گیرد. در مرحله‌ی خوشه‌بندی^۲، شناسایی و دسته‌بندی بخش‌های گفتاری مرتبط با بازشناسی یا فهرست‌نگاری^۳ گفتار در محیطی که چندین گوینده ممکن است در آن اقدام به سخن گفتن کنند، مانند یک جلسه یا یک کنفرانس و اخبار و نظایر آن، کاربرد دارد. این کار نه‌تنها می‌تواند به سیستم‌های بازشناسی گفتار پیشرفته جهت بهبود نتایج بازشناسی گروهی کمک نماید بلکه در شناسایی و متن‌نگاری مکالمه‌ها نیز به آن‌ها کمک می‌کند.

شکل (۱-۱) نحوه کار این سیستم را به‌خوبی نشان می‌دهد.

¹ Segmentation

² Clustering

³ Indexing



شکل ۱-۱: نمایش بخش‌بندی گویندگان روی گفتار ورودی

در شکل (۱-۱) فایل صوتی مورد بررسی یک صوت ضبط‌شده تک کاناله است که شامل چندین منبع صوتی می‌باشد. هر منبع توسط یک یا چند منبع صوتی مشابه یا متفاوت ناشی شده است و می‌تواند شامل گفتار گوینده، موسیقی، انواع نویز، مکث و... باشند. ورودی سیستم بازشناسی فایلی مشابه به شکل (۱-۱) را در بردارد. نوع و جزییات منابع صوتی موجود در فایل به ویژگی کاربردی آن فایل بستگی دارد.

به‌طور کلی سیستم‌های جداسازی گفتار در حوزه‌های مختلف پیش روی خود چالش‌های متفاوتی از جمله کیفیت ضبط صوت، پهنای میکروفون، نویز، میزان و نوع منابع غیر گفتاری، تعداد گویندگان، سبک و ساختار گفتار، طول مدت گفتار و ترتیب گویندگان روبرو هستند. بنابراین کاربرد هر حوزه، هر کدام مشکلات و مسائل مربوط به خود را دارا می‌باشند. در همه‌ی این سیستم‌ها تلاش بر بدست‌آوردن نتایج قابل قبول و مناسب می‌باشد. [۴]

۱-۴- ساختار پایان‌نامه

پایان‌نامه پیش رو از ۵ فصل تشکیل شده است که به‌صورت ذیل شکل گرفته است:

فصل اول مقدمه‌ای بر پردازش سیگنال و صوت می‌باشد. در این فصل به بیان آنچه در گفتار می‌گذرد و اهمیت بررسی و تحلیل داده‌های صوتی می‌پردازیم. به نیازمندی این حوزه از پردازش

اطلاعات پرداخته و چرایی آن را مورد بررسی قرار می‌دهیم.

در فصل دوم به ارائه‌ی پیشینه‌ی مبحث پیش رو می‌پردازیم. بیان آغاز ایده‌ی بازشناسی گوینده و آنچه در طی سالیان متمادی تحقیق و پژوهش دانشمندان و محققان و دانشجویان راه علم به وجود آمده است می‌پردازیم. چه ایده‌هایی ترکیب شده‌اند و چه چیزهایی دستخوش تغییر قرار گرفته‌اند. چه کسانی در این زمینه کار کرده‌اند و پیشگامان این عرصه چه کسانی بوده‌اند.

در ادامه به معرفی الگوریتم‌های در دسترس و روش انجام شده می‌پردازیم. به تعریف کار خود و الگوریتم مورد استفاده در این پایان‌نامه پرداخته و آزمایش‌های انجام شده در این متن را معرفی و بررسی می‌کنیم. مزایا و معایب طرح معرفی شده را بیان می‌کنیم .

در فصل پایانی نتیجه‌گیری‌های این آزمایش‌ها و آنچه به دست آمده و مقایسه‌های به دست آمده را ذکر کرده و به ارائه پیشنهادها قابل انجام می‌پردازیم.

فصل ۲ :

بیان مسأله و پیشینه تحقیق

۱-۲- پیشینه بازشناسی گفتار

پردازش گفتار به‌عنوان یکی از زیرشاخه‌های پردازش سیگنال، به‌سرعت در حال گسترش است. تکنیک‌های پیچیده و نوآوری‌های روزافزون این دانش، همگی در راستای دستیابی به این آرزو هستند که امکان بیابیم مفاهیم در قالب ابزارهای ریاضی فراهم گردد. در این متن، به بیان خلاصه‌ای از انواع روش‌های بازشناسی گوینده می‌پردازیم.

هدف بلندمدت سیستم‌های بازشناسی خودکار گفتار، طراحی ماشینی است که سیگنال صوتی مربوط به یک جمله بیان‌شده را به دنباله‌ای از کلمات نوشته‌شده تبدیل نماید. سیستم‌های بازشناسی خودکار گفتار اطلاعات متنوعی از منابع دانش گوناگون را در جهت دستیابی به جمله بیان‌شده از روی سیگنال صوتی دریافت شده، به کار می‌گیرند. اما مشکلات متعددی در بازشناسی گفتار پیوسته بدون قید وجود دارد که عبارت‌اند از:

۱. تحت تأثیر قرار گرفتن کیفیت سیگنال صوتی به‌وسیله نویز محیط و تابع انتقال سیستم انتقال مانند میکروفون، تلفن و...

۲. عدم وضوح مرز مابین کلمات و واج‌ها در سیگنال صوتی.

۳. تنوع وسیع سرعت بیان.

۴. دقت ناکافی در بیان کلمات و به‌خصوص انتهای آن‌ها در گفتار محاوره‌ای نسبت به گفتار مجزا.

۵. تأثیر تنوع متعدد گوینده از جمله جنسیت، شرایط فیزیولوژیک و روانی بر گفتار.

۶. به‌کارگیری محدودیت‌های معنایی- نحوی زبان برای گفتار زبان طبیعی به روشی مشابه ارتباط انسان با انسان در سیستم بازشناسی.

۲-۲- انواع روش‌های بازشناسی گوینده

در جهت غلبه بر مشکلات مذکور تاکنون روش‌های متنوعی پیشنهاد شده است که از جمله آن‌ها روش‌های آماری مبتنی بر قانون تصمیم‌گیری بیز، روش‌های مبتنی بر شبکه عصبی و در برخی موارد ترکیب روش‌های آماری و شبکه عصبی است. با بررسی روش‌های فوق می‌توان دریافت که شناسایی کلمه یا واج بدون خطا بدون استفاده از دانش سطوح بالاتر به‌خصوص در بازشناسی گفتار پیوسته با حجم لغت‌نامه بزرگ، امکان‌پذیر نیست. به‌عنوان یک نتیجه، یک سیستم بازشناسی گفتار که با انبوهی از فرض‌ها درباره واج‌ها، کلمات و جملات مواجه است، در حالت ایده آل بایستی محدودیت‌های سطوح بالا را که به‌وسیله واژگان، نحو، معانی و ادراک مشخص می‌شود، در نظر بگیرد. در سیستم‌های مبتنی بر قانون تصمیم‌گیری بیز برخی از این محدودیت‌ها توسط مدل زبانی به سیستم بازشناسی اعمال می‌شود.

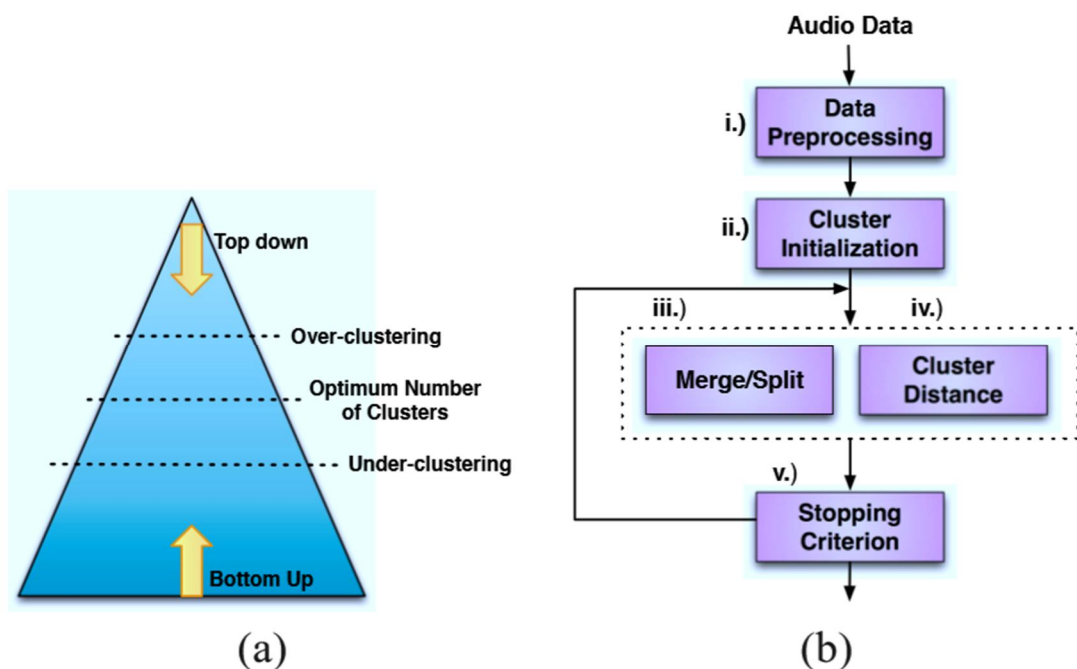
نتایج مطالعات و بررسی‌ها نشان داده است که مدل‌های زبانی‌ای که در حالت کلی توالی واحدهای زبانی را مدل می‌کنند، در کاهش خطای بازشناسی نقش عمده‌ای ایفا می‌کنند. در این میان، استفاده از مدل‌های زبانی مبتنی بر شبکه‌های عصبی با وجود قابلیت این شبکه‌ها در یادگیری زنجیره نمادها و نیز به دلیل قابلیت هموارسازی و خاصیت تعمیم‌دهی آن‌ها بر روش‌های آماری مزیت دارد.

سیستم‌های نمایه زنی گفتار شامل دو سیستم میشوند :

- سیستم بالا به پایین

- سیستم پایین به بالا

این دو سیستم در مدل سازی اولیه با یکدیگر متفاوت هستند که در شکل زیر نمای کلی این دو سیستم را می‌بینید :



شکل ۲-۱: سیستم کلی نمایه سازی (الف) نمایه کلی خوشه بندی کلی (ب) معماری خوشه بندی نمایه سازی کلی

در شکل بالا در بخش ب روند کلی سیستم را مشاهده میکنید. سیستم بالا به پایین همانطور که دیده میشود در ابتدا خوشه بندی و مدل سازی را با تعداد بیشتری نسبت به گویندگان احتمالی آغاز میکنیم و در ادامه تا رسیدن به نتایج ایده ال خوشه ها را حذف و ادغام میکنیم .

در سیستم پایین به بالا معمولاً با یک خوشه بندی کار را آغاز میکنیم و روند خوشه سازی را تا رسیدن به نتایج ایده ال ادامه میدهیم .

۲-۳- انواع ویژگی‌های مورد استفاده در جداسازی گوینده

در سیستم‌های بازشناسی گفتار نیازمند استخراج ویژگی‌هایی برای جداسازی و شناخت گوینده‌ها هستیم. ویژگی‌های آکوستیکی‌ای که اطلاعات گوینده را در اسپکتوگرام متمایز می‌کند و در مکالمات تلفنی ثابت هستند، مدنظر می‌باشد. در استخراج ویژگی این نکته اهمیت دارد که ویژگی استخراجی به گونه‌ای باشد که بتواند تقریباً تمامی اطلاعات را در برگیرد و همزمان در مقایسه با ویژگی‌های دیگر اجزا دارای مشخصه و تفاوت برای جداسازی و متمایز بودن را داشته باشد.

ضرایب کپسترال مل فرکانسی^۱ یا ضرایب پیشگویی خطی^۲ اگرچه برای متمایز ساختن گوینده‌ها طراحی نشده‌اند اما به‌طور وسیعی در حوزه شناسایی گوینده و تشخیص گوینده به کار گرفته می‌شوند. از زمانی که مدل‌سازی گویندگان در حوزه بازشناسی گوینده معمول شد، ضرایب کپسترال مل فرکانسی و دیگر ویژگی‌های کپسترال رایج‌ترین ویژگی‌های مورد استفاده می‌باشند.

به‌طور معمول در طی قسمت‌بندی گویندگان ۱۳ ضرایب کپسترال مل فرکانسی در کنار انرژی زمان کوتاه استفاده می‌شود. در حالی که در طول خوشه‌بندی کردن گویندگان استفاده از مشتقات مرتبه بالاتر ضرایب کپسترال مل فرکانسی رایج‌تر است. تعداد ویژگی‌ها و دقت و حجم محاسباتی بر اساس تعداد و نوع ویژگی مورد استفاده افزایش پیدا می‌کند. [۵]

دیگر ویژگی به نام ضرایب کپسترال خطی فرکانسی^۳ به‌دست‌آمده با فیلتر بانک خطی به‌جای فیلتر بانک در مقیاس مل [۶] و نیز ضرایب کپسترال پیشگویی خطی^۴ نیز مورد استفاده قرار گرفته‌اند [۷]. اما هیچ نتیجه‌ای مبنی بر بهتر بودن کارایی این سیستم‌ها در کارهای انجام‌شده گزارش داده نشده است.

اندازه معمول پنجره تحلیل و آنالیز به‌طور میانگین در اغلب سیستم‌های پردازش ۲۵-۳۰ میلی‌ثانیه با حرکت پنجره‌ی ۱۰ میلی‌ثانیه می‌باشد. حرکت پنجره به این معنا می‌باشد که بین قسمت اول و دوم همپوشانی فاصله زمانی خاصی قرار دارد تا از دست رفتن اطلاعات در مرزهای قسمت‌ها و نیز در خطاهای پنجره‌بندی جلوگیری شود [۷].

برای سیستم شناسایی گفتار ویژگی‌های آکوستیکی که بین گفتار و غیر گفتار را متمایز می‌کنند

¹ Mel Frequency Cepstral Coefficient (MFCC)

² Linear Prediction Coefficient (LPC)

³ Linear Frequency Cepstral Coefficient (LFCC)

⁴ Linear Prediction Cepstral Coefficient (LPCC)

مورد استفاده قرار می‌گیرند. ویژگی‌هایی از قبیل انرژی [۷] و نرخ عبور از صفر^۱، طیف مرکزی^۲، شار طیفی^۳ [۸] در شناسایی گفتار مورد استفاده قرار می‌گرفتند.

استفاده از حداقل یکی از این بردارهای ویژگی که در بالا ذکر شدند همیشه در اتصال با ویژگی‌های کپسترال دیده شده‌اند. این ویژگی‌ها در کنار ویژگی‌های کپسترال در اکثر موارد موجب بهبود نتیجه و بالا رفتن دقت سیستم شناسایی می‌شود اگرچه در برخی کاربردها ترکیب ویژگی‌های مختلف برای اهداف و کاربردهای نادرست دقت را پایین می‌آورد.

جدا از ویژگی‌های آنالیز زمان کوتاه ذکر شده در بالا ویژگی‌های مدولاسیون فرکانسی^۴ هرگز که مشخصه‌های بلندمدت سیگنال آکوستیک را بیان می‌کند نیز مورد بررسی قرار گرفته‌اند [۹]. این ویژگی‌ها با چالش‌های زیادی روبرو شده‌اند که در هر حال بعد بالای ویژگی‌ها و هزینه‌های محاسباتی با آن پیوند خورده‌اند. ویژگی‌های ترکیبی بلندمدت محاسبه شده در طول پنجره ۵۰۰ میلی‌ثانیه مانند میانه‌ی پیچ^۴، طیف میانگین بلندمدت^۵، انحراف فرمتهای ۵ ام و ۴ ام^۶، همسازها به نرخ نویز^۷، پراکندگی فرمت^۸ و غیره برای کلاس‌بندی سریع مورد استفاده قرار می‌گیرند [۱۰]. این ویژگی‌ها اطلاعات آوایی و دستگاه صوتی منبع را فراهم می‌کنند و شناخت گوینده بهتری نسبت به استفاده از ضرایب کپسترال مل فرکانسی رادارند.

اخیراً، اسلینی و همکارانش ویژگی‌های مشتق شده به‌عنوان فعال‌ساز لایه تنگنای یک شبکه عصبی را استفاده کرده است. شبکه عصبی مصنوعی برای مشخص کردن قسمت ۵۰۰ میلی‌ثانیه‌ای بطوریکه به گوینده‌ای متفاوت یا همسان تعلق دارد آموزش داده می‌شود [۱۱]. در کار دیگر پیشرفت

¹ Zero Crossing Rate(ZCR)

² Spectral Centroid

³ Spectral Flux

⁴ Median Pitch

⁵ Long time average spectrum

⁶ Deviation of the 4th and 5th Formants

⁷ Harmonics to Noise Ratio

⁸ Formant Dispersion

نسبی ۵۰ درصدی برای سیستم تشخیص گفتار زمانی که بیشینه فعالیت دوبعدی با یک شبکه عصبی عمیق و ۱۳ ضرایب کپسترال مل فرکانسی متصل شده‌اند گزارش داده شده است [۱۲].

یک فضای ویژگی جالب دیگر که در سال ۲۰۱۰ منتشر شد، خطای بازشناسی را اندکی قربانی کرد تا به یک سرعت ۱۰ برابری با استفاده از ویژگی‌های با مقدار دودویی برای اجرای دسته‌بندی برسد [۱۳]. در این کار ویژگی‌های آکوستیکی ضرایب کپسترال مل فرکانسی با استفاده از احتمال گرفته شده از مدل مخلوط گوسی به یک فضای ویژگی دودویی انتقال داده شد.

۲-۴- بخش بندی

در بخش بندی گفتار هدف ساخت قطعه‌های متصل به هم و همگن از فایل صوتی می‌باشد که با قسمت همسایه تفاوت دارد. به این کار همچنین تشخیص تغییر آکوستیک^۱ نیز گفته می‌شود. در ادامه به دو روش بخش بندی که بیشتر در بازشناسی گفتار مورد استفاده قرار می‌گیرند اشاره می‌کنیم:

۲-۴-۱- بخش بندی متریک^۲

یکی از رایج ترین روش های بخش بندی کردن صوت امروزه بخش بندی متریک است. این روش در بخش بندی آهنگ بسیار مشهور است. در بخش بندی متریک یک فاصله متریک در ابتدا بین دو بخش صوت تعریف می‌شود که شباهت آن‌ها را مشخص می‌کند. سپس یک راهکار تشخیص تغییر با استفاده از این متریک اجرا می‌شود. در مقایسه با روش بعدی (بر پایه مدل) این روش ویژگی‌های بسیاری همچون عدم نیازمندی به اطلاعات قبلی از داده‌ها را داراست.

برای بخش بندی آهنگ، فاصله به صورت مستقیم بین ویژگی‌ها محاسبه می‌شود. اگرچه در پردازش صوت ویژگی‌های استفاده شده (به صورت کلی ویژگی‌های کپسترال) برای محاسبات فاصله

¹ Acoustic Change Detection

² Metric Based Segmentation

برای مقایسه شباهت گوینده به دلیل تغییر کامل آن‌ها با تلفن‌ها مناسب نمی‌باشند.

برای جمع‌آوری اطلاعات گوینده از بخش‌های بلندتر، فرض می‌شود که ویژگی‌های هر بخش از یک توزیع احتمال می‌آیند. مقایسه فاصله بین این توزیع‌های احتمال با استفاده از اندازه‌گیری‌های شباهت آماری همانند دیورژانس^۱ KL، نرخ شباهت متقابل^۲، معیار استنتاج بیزین^۳ و غیره انجام می‌شود. معمول‌ترین روش توزیع احتمال برای مدل‌سازی قطعات بردارهای ویژگی حین بخش‌بندی گوینده توزیع چند متغیره کوواریانس کامل گوسی است.

۲-۴-۱-۱- معیار استنباط بیزین

معیار استنباط بیزین یک معیار انتخاب مدل است. این معیار آماری مدل‌های موجود را برای شرح دادن داده‌ها مقایسه می‌کند. هدف در حین این انتخاب، محاسبه‌ی هر نوع عدم تناسب است. برای یک مجموعه از بردارهای X ، BIC همانند زیر تعریف می‌شود.

$$\Delta BIC = R - \lambda P \quad (1-2)$$

عبارت اول احتمال داده‌های مدل را به دست می‌دهد که شامل توزیع مدل‌های قسمت‌هایی است که می‌خواهیم بین آن دو معیار را بررسی کنیم. این بخش شامل دو عبارت احتمال و تابع توزیع جدا برای هر قسمت به همراه توزیع و احتمالی جدا برای قسمت شامل اتصال دو قطعه به یکدیگر و تشکیل یک قطعه واحد می‌باشد. به‌طور خلاصه عبارت اول را می‌توان به‌صورت زیر تعریف کرد.

$$R = \frac{N_X}{2} \log(|\Sigma_X|) - \frac{N_{X_1}}{2} \log(|\Sigma_{X_1}|) - \frac{N_{X_2}}{2} \log(|\Sigma_{X_2}|) \quad (2-2)$$

¹ Kullback–Leibler divergence

² Cross Likelihood Ratio

³ Bayesian Inference Criterion(BIC)

عبارت دوم شامل پارامتر تجربی λ است که به صورت تجربی به دست می آید و به طور معمول عددی بین ۰.۵ تا ۲ می باشد که باید بهینه شود. تابع P اندازه و بعد داده‌ی آموزش مدل را مشخص می کند. عبارت دوم پیچیدگی مدل نام دارد و به صورت زیر می توان نشان داد.

$$P = \frac{1}{2} \left(p + \frac{1}{2} p(p+1) \right) \quad (3-2)$$

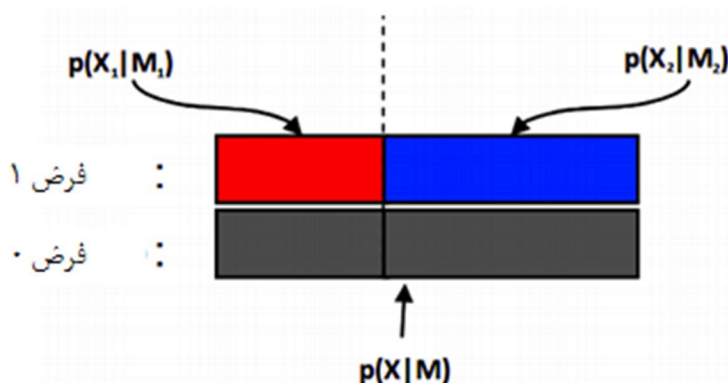
معیار بیزین برای تشخیص تشابه دو مجموعه از بردار ویژگی‌ها با توزیع یکسان یا غیر همسان بکار می رود. برای اندازه گیری میزان شباهت بین بلوک X_1 و X_2 فرض های زیر باید مقایسه شوند:

H_0 : بردار ویژگی از X_1 و X_2 توزیع همسان دارند

H_1 : بردار ویژگی X_1 و X_2 از دو توزیع متفاوت می آیند.

مدل برای H_0 فرض اول M نام دارد که با داده‌ی $X_1 - X_2$ متصل شده به هم آموزش دیده می شود. مدل برای H_1 فرض دوم M_1 و M_2 به ترتیب برای X_1 و X_2 نام گذاری می شود.

$$\Delta BIC = BIC(M) - BIC(M_1) - BIC(M_2) \quad (4-2)$$



شکل ۲-۲: BIC برای تشخیص تغییر در گوینده

در شکل (۱-۲) فرض‌های ممکن شرح داده شده است. دو قطعه‌ای که مورد آزمایش هستند تا تصمیم بگیریم به یک گوینده تعلق دارند یا خیر با دورنگ و بخش با اندازه متفاوت نشان داده شده‌اند. تفاوت و تشابه توزیع هر دو بخش با تفاوت و تشابه رنگ در شکل مدل‌سازی شده است. یک مقدار مثبت برای ΔBIC عدم تشابه بین دو بلوک را مشخص می‌کند و این نشان‌دهنده‌ی تغییر در گوینده از بلوک X_1 به بلوک X_2 می‌باشد.

آقای چن و همکارانش یک سیستم بدون نظارت کامل را ساختند. روش آن‌ها بارها و بارها در تعداد زیادی سیستم‌های تشخیص تغییر گوینده و محیط مورد استفاده قرار گرفت [۸، ۱۴].

روش آن‌ها در [۱۵] و [۱۶] با استفاده از روش BIC با کاهش محاسبات ناشی از کاهش دقت اندک به سرعت اجرای بالاتری دست یافتند.

روش‌های بخش‌بندی متریک گوینده در دو استراتژی اجرا می‌شوند: یک روش پنجره‌بندی ثابت و یک پنجره‌بندی گسترشی. در مورد اول، یک پنجره با اندازه ثابت داریم که برای تشخیص تغییر مرکزیت از آن استفاده می‌شود [۱۶]. اگر بردارهای ویژگی با جداسازی توزیع‌ها در یک سمت نقطه‌ی وسط بهتر مدل شود که باعث فاصله‌ی بیشتر بین توزیع‌ها می‌شود، نقطه‌ی وسط به‌عنوان نقطه‌ی تغییر انتخاب می‌شود. اندازه‌ی پنجره‌ی بخش‌بندی به‌طور معمول ۵ ثانیه است و دو بخش ۲.۵ ثانیه‌ای برای شباهت مورد استفاده قرار می‌گیرند. با پنجره با سایز بزرگ‌تر دو بخش بهتر مدل می‌شوند. اگرچه ممکن است با این کار احتمال از دست دادن تغییر در گوینده باشد. چون ممکن است تغییر در گوینده‌ی بیشتر از یکی در پنجره تحت نظر ما وجود داشته باشد و به همین دلیل نتوان آن را تشخیص داد.

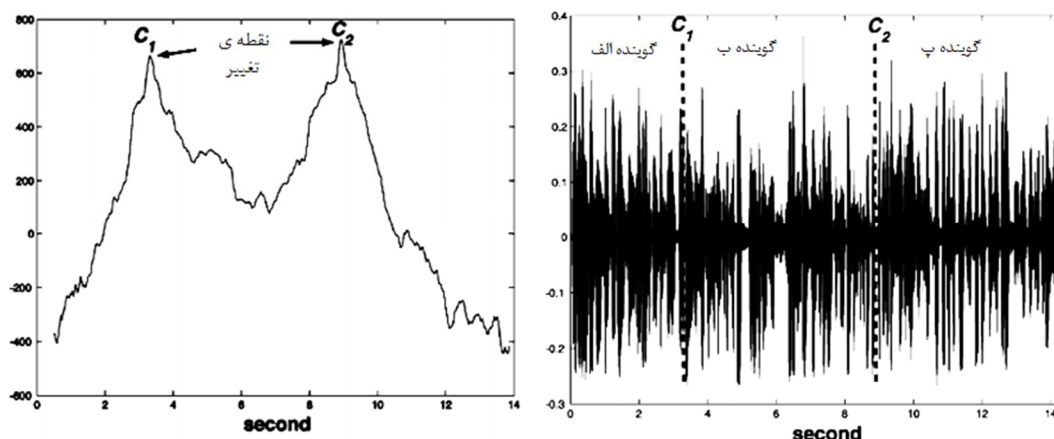
۲-۴-۲- بخش‌بندی بر اساس مدل

روش‌های بخش‌بندی گوینده بر اساس مدل برای هر کلاس بخش‌بندی یک مدل مخلوط

گوسی^۱ آموزش داده می‌شود. این مدل مخلوط گوسی‌ها به‌عنوان یک تابع توزیع چگالی در مدل مخفی مارکوف استفاده می‌شود. در مدل مخفی مارکوف هر حالت به هر حالت دیگر با احتمال انتقال حالت خاصی متصل است. الگوریتم دیکد کردن ویتربی با استفاده از این مدل مخفی مارکوف یک بخش‌بندی صوت ضبط‌شده را اجرا می‌کند.

بزرگ‌ترین معایب این نوع بخش‌بندی این است که مدل مخلوط گوسی‌ها نیازمند این هستند که از قبل شناخته‌شده باشند و این بدین معنی است که باید با برخی داده‌های خارجی آموزش‌دیده شوند.

در شکل (۲-۲) مشاهده می‌کنید که برای یافتن تغییر در گوینده بین بخش‌های پشت سر هم فاصله‌ی میانه‌ی دو پنجره محاسبه می‌شود و برحسب زمان در طول پنجره رسم می‌شود. زمانی که فاصله به بیشترین حد خود می‌رسد محل تغییر در گوینده و محل بیشترین تفاوت در دو پنجره است. این خود نشان‌دهنده‌ی تغییر در ویژگی‌ها و یا همان تغییر در گوینده است. ابديم معنی که این دو قسمت مجاور متعلق به دو گوینده‌ی متفاوت می‌باشند.

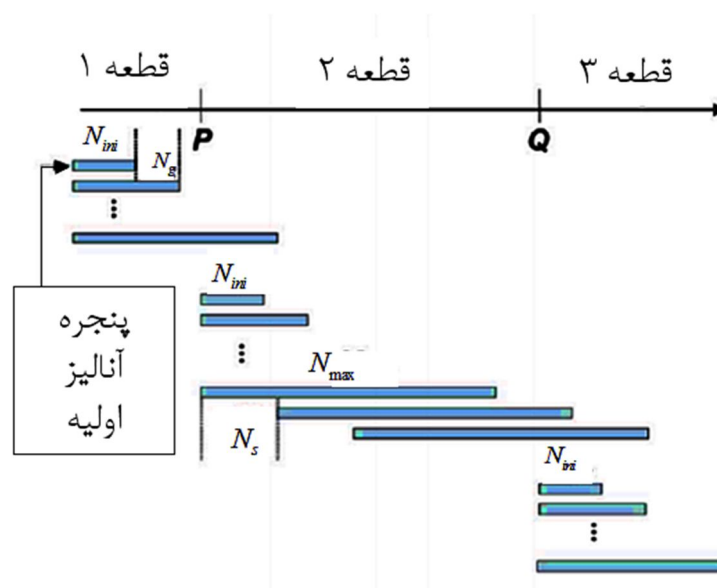


شکل ۲-۳: پنجره‌ی متغیر برای تشخیص تغییر. فاصله‌ی بین نیمه‌های پنجره‌ها محاسبه و برحسب زمان نمایش داده‌شده است. قله و اوج در فاصله نمایانگر تغییر است.

¹ GMM : Gussian Mixture Models

در شکل (۳-۲) روش جستجوی پنجره‌ای و جاروب کل قطعه نشان داده شده است. روش کار به این صورت است که در ابتدا پنجره‌ای با کمترین طول و طول اولیه در نظر گرفته و به جستجو برای تغییر ویژگی‌های آکوستیکی می‌پردازیم. سپس طول پنجره را اندک‌اندک در هر مرحله افزایش می‌دهیم تا زمانی که نقطه‌ی تغییر در گوینده را بیابیم. سپس جستجو را از سر گرفته و از نقطه‌ی تغییر پنجره گذاری را مجدد آغاز می‌کنیم تا به نقطه‌ی تغییر بعدی برسیم. دقت داشته باشید برای طول پنجره ماکسیمم خاصی در نظر می‌گیریم که اگر تا به این حد نقطه‌ی تغییری پیدا نشد از مرحله‌ی بعد آن پنجره‌ی بعدی را به میزان خاصی جلوتر شروع می‌کنیم. این بدین معنی است که به میزان مشخصی از اطلاعات بررسی شده از نقطه‌ی تغییر قبلی را که بررسی کردیم از لیست بررسی مجدد در جستجوی در حال انجام حذف می‌کنیم و به جستجو در بین ادامه‌ی بخش‌ها می‌پردازیم.

زنجیره صوتی



شکل ۳-۲: جستجوی افزایشی برای تشخیص تغییر گوینده. هر جستجو برای تنها یک تغییر است. بعد از هر یافتن نقطه

تغییری جستجو از سر گرفته می‌شود. [۱۶]

در زمان بخش‌بندی و دسته‌بندی کردن پی‌درپی، خروجی مجموعه‌ای از بخش‌های گوینده را به‌عنوان داده‌ی آموزش برای مدل مخلوط گوسی‌های اجرای بعدی بکار می‌برد. این کار باعث تعریف و

نوسازی بخش‌بندی گوینده می‌شود. یک پیش‌دسته‌بندی به‌طور معمول برای به دست آوردن اطلاعات دسته‌های اولیه بخش‌های گوینده [۱۰] اجرا می‌شود. هر گروه اطلاعات مشابه گوینده را نشان می‌دهد.

اغلب روش‌های بر پایه مدل برای بخش‌بندی به‌عنوان یک پس‌پردازش برای دستیابی به بخش‌بندی اصلاح‌شده به کار می‌رود [۱۷]. برخی از این تکنیک‌های بر پایه مدل در بخش‌بندی در تشخیص اتوماتیک گفتار مشهور هستند [۱۸] که در این کاربرد تغییر آکوستیکی در بین گفتار و غیر گفتار جستجو می‌شود.

در برخی سیستم‌های تشخیص گفتار تلفنی [۱۹] یک پیش‌بخش‌بندی صوت ضبط‌شده بر پایه جنسیت و پهنای باند با استفاده از مدل مخلوط گوسی آموزش داده‌شده برای هر ۴ کلاس (۲ پهنای باند* ۲ جنسیت) استفاده می‌شود.

۲-۵- تشخیص فعالیت گفتار^۱

یافتن بخش‌های همسان گفتار در یک فایل صوتی و جدا کردن آن‌ها از دیگر انواع اصوات تشخیص فعالیت گفتار نامیده می‌شود. در سیستم‌های پردازش گفتار سودمند است که بخش‌هایی از سیستم که بخش‌های شامل گفتار هستند بجای تمام فایل صوتی بررسی شود.

این امر در بهبود حجم محاسباتی و زمان پردازش و استفاده‌ی بهینه از منابع کمک شایانی می‌کند. جدا از مزیت‌های محاسباتی، در نبود سیستم تشخیص فعالیت گفتار در سیستم‌های پردازش گفتار خطاهایی وارد می‌شود. سیستم تشخیص فعالیت گفتار یک‌رویه‌ی پایه‌ای در تقریباً همه‌ی

^۱ SAD:Speech Activity Detection

بخش‌های پردازش و تشخیص و بهبود صوت می‌باشد [۱۸].

در بازشناسی گوینده، خطای متریک خود اهمیت وجود سیستم تشخیص فعالیت گفتار را نمایان می‌کند. چون هر دو خطای گفتار ازدست‌رفته و هشدار خطای گفتار در خطای کلی بازشناسی لحاظ شده‌اند. از آن بیشتر با داده‌ای محدود از گوینده از بخش کوچکی از گفتار، حضور داده‌های غیر گفتاری که مدل‌های گوینده را تحت تأثیر و خطا قرار داده‌اند باعث تغییر در کارایی عملکرد سیستم بازشناسی می‌شوند.

سیستم تشخیص فعالیت گفتار اغلب خوشه‌بندی بر اساس فریم را انجام می‌دهد. مدل‌های آماری روی یک فضای ویژگی مناسب برای جداسازی گفتار و غیر گفتار آموزش می‌بینند و تخمین زده می‌شوند. در بسیاری از موارد، مدل‌های مخلوط گوسی مدل‌های آماری هستند که استفاده می‌شوند و فضای ویژگی در اغلب موارد ویژگی‌های کپسترال هستند. در برخی کارها استفاده از ویژگی‌های آکوستیکی مانند انرژی [۷] و نرخ عبور از صفر [۲۰] و شار طیفی [۸] نشان داده شده است.

۲-۶- خوشه‌بندی کردن

خوشه‌بندی کردن مشکل معمول در آنالیز داده‌های آماری است. در بسیاری از شاخه‌های علمی من جمله داده‌های اکتشافی معدن تا تشخیص جوامع در شبکه‌های اجتماعی مورد کاربرد است. خوشه‌بندی کردن فرایند گروه‌بندی کردن مجموعه‌ای از اهداف -مانند اهداف در هر گروه- که دسته نام دارند است که بیشترین شباهت را نسبت به بقیه‌ی اهداف در دیگر گروه‌ها نسبت به یکدیگر دارند. هدف‌ها باید در فضای برداری یا مدل‌های آماری نشان داده شوند. شباهت گفته شده در بالا یک فاصله‌ی تعریف شده بین اهداف است که توسط کاربر تعریف می‌شود.

واژه شباهت به این دلیل استفاده می‌شود که اندازه‌گیری‌های تعریف شده نیازمند هیچ‌گونه از معیارهایی چون غیر منفی بودن و تقارن و نامساوی مثلثی ندارد. واژه شباهت و فاصله در اینجا

به صورت معادل بکار می‌روند بدین منظور که فاصله‌ی کمتر نشان‌دهنده‌ی شباهت بیشتر است و برعکس.

فرایند دسته‌بندی کردن عموماً انتقالی ثابت است و جایگاه هر هدف در فضای خودش نسبت به فضای دسته‌بندی خودشان مرتبط‌تر هستند. در مسئله شناسایی گوینده، هدف دسته‌بندی بخش‌های سیگنال صوت بر اساس گوینده‌ی فعال در هر بخش است. هر دسته باید به صورت ایده آل یک گوینده را نشان دهد.

به جرئت می‌توان گفت ابعاد زمانی و فرکانسی اسپکتوگرام یک تک بخش، بزرگ است و در نتیجه مقایسه‌ی بین اسپکتوگرام دو بخش از لحاظ محاسباتی به دلیل بسیار زیاد بودن حجم داده‌های محاسباتی غیرقابل انجام است و نمی‌توان با تکیه بر مقایسه‌ی دو اسپکتوگرام به نتایج دلخواه دست‌یابیم. در نتیجه نیازمند این هستیم که هر بخش به فضای با بعد کمتر برده شود تا بتوان شباهت اطلاعات گویندگان‌شان مورد بررسی قرار بگیرد. هر بخش باید دارای یک بردار نشانه یا مدل آماری باشد که اطلاعات گوینده را دربرداشته باشد.

۲-۶-۱- مدل‌های گوینده برای خوشه‌بندی کردن

از آنجایی که بازشناسی گوینده نیازمند دربرگیری اطلاعات گوینده از بخش‌های صوتی می‌باشد مدل‌های گوینده معمولاً در شناسایی گوینده و کاربردهای مرتبط بکار گرفته می‌شود.

۲-۶-۱-۱- مدل‌های گوسی مخلوط

مدل‌های گوسی مخلوط معمولاً از ویژگی‌های کپسترال برای مدل‌سازی گوینده‌ها استفاده می‌روند. یک مدل گوسی مخلوط ابزار معروفی برای مدل‌سازی داده‌های چند مدل و پردازش به صورت زیر است:

$$p(x) = \sum_{i=1}^N \omega_i N(\mu_i, \Sigma_i, x) \quad s.t. \sum_{i=1}^N \omega_i = 1 \quad (5-2)$$

از آنجایی که طول مدت بخش می‌تواند کوچک باشد تعداد بردارهای ویژگی قابل دسترس از هر بخش در برخی موارد برای تخمین و ساخت یک مدل گوسی کامل ناکافی است. برای حل این مشکل مدل‌های پس‌زمینه‌ای همگانی^۱ از پیش آموزش یافته برای ساخت مدل گوینده هر بخش تنظیم شده است [۲۱] مدل‌های پس‌زمینه‌ای همگانی یک مدل جامع برای داده‌های ترکیب شده از چند گوینده که اطلاعات آن‌ها را دربردارد می‌باشد. برای مدل مخلوط گوسی از ویژگی‌های کپسترال اندازه‌گیری‌های شباهت آماری مختلفی همانند دیورژانس KL متقارن^۲ و نسبت احتمال متقابل و غیره به کار برده شده‌اند [۲۲]. دیورژانس KL یک اندازه‌گیری تئوری اطلاعات بدین صورت است که دو توزیع احتمال چقدر از یکدیگر متفاوت هستند. نرخ شباهت متقابل دو مقدار $P(X_1|M_2)$ و $P(X_2|M_1)$ که توابع گوسی مدل سازی شده روی ویژگی‌های استخراجی می‌باشند، را مقایسه می‌کند.

$$CLR(X_1, X_2) = \log \frac{P(X_1, M_1)}{P(X_1, M_2)} + \log \frac{P(X_2, M_2)}{P(X_2, M_1)} \quad (6-2)$$

$$NCLR(X_1, X_2) = \frac{1}{|X_1|} \log \frac{P(X_1, M_1)}{P(X_1, M_2)} + \frac{1}{|X_2|} \log \frac{P(X_2, M_2)}{P(X_2, M_1)} \quad (7-2)$$

M_i مدل تخمین زده شده روی X_i است. همان‌طور که می‌بینند اگر بردارهای ویژگی X_1 و X_2 متعلق به یک گوینده باشند، مدل X_1 به خوبی بر مدل X_2 منطبق است و در نتیجه شباهت متقابل افزایش یافته و در نتیجه فاصله کاهش می‌یابد. همان‌طور که در (۷-۲) ملاحظه می‌کنیم عبارت فاصله‌ی CLR به صورت نرمال شده خود به عنوان معیار فاصله‌ای مجزا تعریف و مورد کاربرد است.

¹ Universal Background Model (UBM)

² Kullback–Leibler Divergence

³ Normalized Cross Likelihood Ratio

به دلیل تغییرپذیری بالای ماتریس‌های کوواریانس و وزن‌های ترکیبی نسبت به گفته‌های [۲۱]، شاخص قابل اطمینانی برای مشخص کردن اطلاعات گوینده وجود ندارد. در نتیجه به جای محاسبه امتیازات شباهت بالا، متوسط‌های یک مدل مخلوط گوسی را با به دست آوردن یک بردار در فضای برداری بعد بالا به هم پیوند داده می‌شوند. اندازه‌های فاصله همانند فاصله‌ی کسینوسی و فاصله‌ی مهالانوبیس [۲۳] در این فضا مورد بررسی قرار گرفتند. برای مقایسه‌ی بین دو ابر بردار، مدل‌های مخلوط گوسی نیازمند این هستند که با یک مدل پس‌زمینه‌ای همگانی وفق داده شوند تا اطمینان پیدا کنیم که بردارهای میانگین متناظر مدل مخلوط گوسی بین بخش‌ها مقایسه شده است. الگوریتم تطبیق در منبع [۲۱] شرح داده شده است.

۲-۶-۱-۲- نمایش بردار i

مفهوم بردارهای i برای اولین بار در بحث اثبات گوینده به‌عنوان یک استخراج ویژگی از مدل مخلوط گوسی معرفی شد تا بعد پارامترهای مدل مخلوط گوسی را کاهش دهند. سیستم استخراج بردار i ، سیستمی است که دنباله‌ی بردارهای (به‌طور معمول ضرایب کپسترال) به دست آمده از دنباله صوتی را به یک بردار با طول ثابت نگاشت می‌کند به طوری که مدل‌های مخلوط گوسی با بعد‌های مشخص که مدل‌های کلی پس‌زمینه‌ای خوانده می‌شوند برای جمع‌آوری داده‌های الگوریتم باوم-ولچ از دنباله‌ی صوتی مورد استفاده قرار گیرند. مدل پس‌زمینه‌ای همگانی اندازه‌ی آن‌ها در برخی مدل مخلوط گوسی، از مرتبه‌ی ۵۱۲ و ۱۰۲۴ و یا حتی ۲۰۴۸ می‌باشند و اندازه‌ی ابر بردارها^۱ برای انجام محاسبات خیلی بزرگ هستند. در عوض، آنالیز برای کاهش بعد ابر بردارها به یک بردار معرف با بعد چند صدتایی می‌رسد. فرض بر این شده است که این زیر فضا، -که زیر فضای تغییر^۲ نام دارد- شامل اطلاعات طیفی گوینده و پس‌زمینه می‌باشد.

¹ Supervector

² Total Variation Subspace

$$m = M + Tx$$

(۸-۲)

که m ابر بردار متوسط وفق داده‌شده‌ی گفتاری است که X بردار i آن بخش می‌باشد. M ابر بردار میانگین مدل‌های پس‌زمینه‌ای همگانی است و T یک ماتریس مرتبه پایین بلند است که زیر فضای تغییر کلی را نشان می‌دهد. این زیر فضا نیازمند این است که روی یک مجموعه داده آموزش داده شود. اگرچه ابر بردارها معمولاً ابعاد ده هزارتایی دارند، این نمایش تمامی ابر بردارها را مجبور می‌کند در یک زیر فضا محدود شوند. بعد این زیر فضا چند صدتا می‌باشد.

استخراج بردار i نیازمند داده‌های آموزش برچسب‌گذاری شده گوینده‌ها با موارد مشابه گوینده می‌باشد. موارد موجود شامل تنوع در نویز زمینه و توازن آوایی می‌باشند. الگوریتم آموزش برای زیر فضای تغییر کلی [۲۳] و استخراج i بردار از الگوریتم باوم-ولچ^۱ در ابزار تشخیص MSR اجرا شده است [۲۴].

۲-۶-۲- الگوریتم خوشه‌بندی

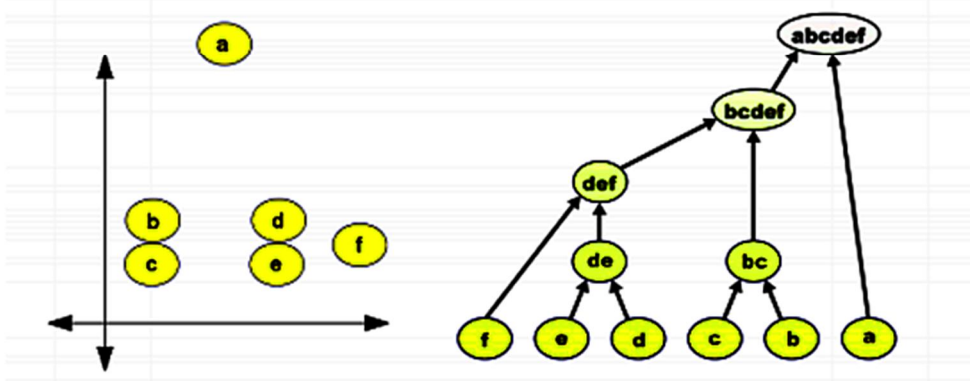
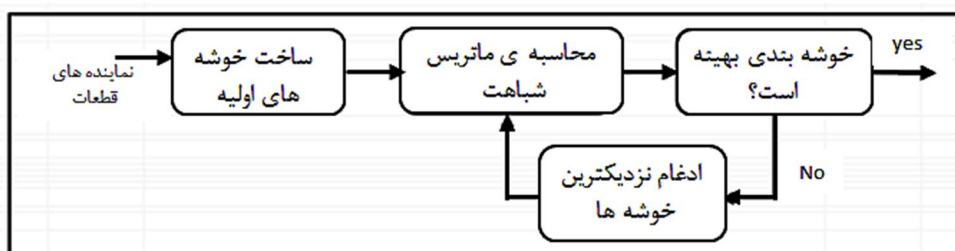
با داشتن ماتریس شباهت بین مدل مخلوط گوسی مدل‌ها و بردارهای i ، باهدف بهترین مجموعه دسته‌ها با حداقل واریانس درون دسته‌ای و حداکثر واریانس برون دسته‌ای به دنبال انواع روش‌های دسته‌بندی هستیم. ما دو نوع الگوریتم دسته‌بندی را مورد بررسی قرار دادیم: ۱- دسته‌بندی متراکم طبقه‌ای^۲ و ۲- دسته‌بندی برنامه خطی صحیح^۳.

هر دو مورد بالا روش‌های حل متفاوتی دارند و نیز معیارهای متفاوتی برای رسیدن به بهترین دسته‌بندی را دارند. در ادامه به معرفی هر کدام به صورت جداگانه می‌پردازیم تا بتوان در مورد هر روش فهم و ادراک متوسطی به دست آورد.

¹ Baum-Welch Statistics

² Hierarchical Agglomerative Clustering (HAC)

³ Integer Linear Program (ILP) clustering (https://en.wikipedia.org/wiki/Integer_programming)



شکل ۲-۵: روش خوشه‌بندی متراکم سلسله مراتبی

۲-۶-۱- خوشه‌بندی متراکم سلسله مراتبی

این الگوریتم یک الگوریتم حریصانه^۱ است بدین معنی که یک انتخاب بهینه‌ی محلی را در هر مرحله با امید یافتن بهینه‌ی کلی اجرا می‌کند. در یک فرایند تکراری، دودسته که بیشترین شباهت را دارند ادغام می‌شوند. تعداد دسته‌ها در هر مرحله یکی کاهش می‌یابد. این فرایند تکراری تا زمانی که تنها یک دسته باقی بماند ادامه پیدا می‌کند. در زمان ادغام دودسته، داده‌ی بخش مرتبط با ۲ دسته باهم یکی مخلوط می‌شوند و یک مدل تک گوینده روی آن‌ها محاسبه می‌شود. فاصله‌ی هر دسته‌ی دیگر با این دسته‌ی جدید مجدداً محاسبه می‌شود و ماتریس شباهت برای هر پله به‌روزرسانی می‌شود. همان‌طور که در شکل (۲-۴) مشاهده می‌کنید در هر مرحله با یافتن بهینه‌ی محلی و ادغام کردن داده‌ها تا جایی که به یک دسته برسند مراحل ادامه پیدا می‌کند و در آخر نتیجه یک دسته‌ی

¹ Greedy Algorithm: (https://en.wikipedia.org/wiki/Greedy_algorithm)

یکتای بهینه است. در زیر مراحل تشریح داده شده است:

مرحله ۰: محاسبه ماتریس شباهت (X_i, X_j)

مرحله ۱: یافتن i^* و j^* به طوری که $(X_{i^*}, X_{j^*}) = \min_{i,j} (X_i, X_j)$ و $i^* \neq j^*$

مرحله ۲: مرحله ادغام: جایگزین کردن X_{i^*} و X_{j^*} را با یک هدف $X_{K^*}, K^* = \min(i^*, j^*)$

مرحله ۳: به روزرسانی کردن ماتریس شباهت (X_i, X_{K^*})

مرحله ۴: رفتن به مرحله اول اگر تعداد دسته‌ها بیشتر از یک

مرحله ۵: محاسبه بهترین مجموعه از دسته‌ها با استفاده از معیار بهینه

مجموعه بهینه از دسته‌ها بر اساس معیار بهینه انتخاب می‌شود. یک معیار بهینه معیاری است که مجموعه‌ای را انتخاب کند که حداقل فاصله برون دسته‌ای بزرگ‌تر از یک حد آستانه است. معیار دیگر توسط [۲۵] معرفی شده است که در آن در طول دسته‌ها در هر مرحله‌ی اجرا، مجموعه دسته‌ای که هیستوگرام فاصله‌های درون دسته‌ای و فاصله‌های برون دسته‌ای در بیشترین فاصله از هم قرار دارند انتخاب می‌شود.

$$\arg \max_k \min_{i \neq j} (X_i^{(k)}, X_j^{(k)}) \geq \theta \quad (9-2)$$

که $(X_i^{(k)}, X_j^{(k)})$ ماتریس شباهت در اجرای k^{th} است.

$$\arg \max \frac{|m_{inter} - m_{intra}|}{\sqrt{\frac{\sigma_{inter}^2}{n_{inter}} + \frac{\sigma_{intra}^2}{n_{intra}}}} \quad (10-2)$$

که m_{inter} و σ_{inter} و n_{inter} میانگین و انحراف استاندارد و تعداد عناصر در فواصل برون دسته‌ای و به‌طور مشابه فواصل درون دسته‌ای هستند.

سیستم^۱ ICSI [۱۰] یک دسته‌بندی اولیه بخش‌های یک‌ثانبه‌ای با استفاده از ویژگی‌های بلندمدت اجرا می‌کند. سپس این بخش‌بندی در یک مدل مخلوط گوسی-مدل پنهان مارکوف^۲ با بردارهای ویژگی ضرایب کپسترال فرکانسی در مقیاس مل (MFCC) دوباره تعریف می‌شوند. هر حالت مدل پنهان مارکوف یک گوینده را نشان می‌دهد و توسط یک مدل مخلوط گوسی مدل شده است. در هر اجرا، تعداد دسته‌ها با ادغام حالت‌های مدل مخلوط گوسی در یک خوشه‌بندی متراکم سلسله مراتبی با استفاده از فاصله‌ی BIC کم می‌شود. سیستم منبع [۷] ضرایب کپسترال پیشگویی خطی^۳ برای ۳۰ دسته از صوت یکنواخت را استفاده می‌کند. سپس به صورت تکراری مدل پنهان مارکوف-مدل مخلوط گوسی را برای خوشه‌بندی متراکم سلسله مراتبی که از فاصله‌ی CLR^۴ استفاده می‌کند اجرا می‌کند. سیستم منبع [۶] همان پیکربندی را استفاده می‌کند ولی به جای استفاده از خوشه‌بندی متراکم سلسله مراتبی یک روش بالا به پایین را بکار می‌گیرد. سیستم‌های زیاد دیگری روش خوشه‌بندی متراکم سلسله مراتبی را برای دسته‌بندی با فواصل مختلف بکار گرفته‌اند. [۱۳، ۲۶]

۲-۶-۲-۲ - خوشه‌بندی بر اساس ILP

میگنر و همکارانش یک آنالیز نسبی روی هردو روش بازشناسی گوینده که شامل روش پله به پله و روش سیستم یکپارچه است انجام داده‌اند [۲۷]. در روش پله‌ای بازشناسی از طریق زیرمجموعه‌های تشخیص فعالیت گفتار، بخش‌بندی، دسته‌بندی انجام می‌شود. در روش همگانی و یکپارچه اطلاعات الگوریتم دسته‌بندی برای بخش‌بندی مجدد گوینده استفاده می‌شود و دسته‌بندی مرتباً تکرار می‌شود. در ادامه‌ی تلاش‌های این تیم و در نتیجه‌ی تلاش‌های روویر و میگنر در سال ۲۰۱۲ [۲۸]، یک روش بهینه‌سازی کلی برای دسته‌بندی گوینده با استفاده از یک برنامه‌ی خطی

¹ International Computer Science Institute

² HMM: hidden markov model

³ linear prediction cepstral coefficient (LPCC)

⁴ Cross Likelihood Ratio

صحیح^۱ اجرا شد. دسته‌بندی به‌عنوان یک مسئله‌ی بهینه‌سازی ترکیبی روی یک گراف کامل (هر گره به گره‌ی دیگر متصل است) مطرح می‌شود. بخش‌های گوینده به‌عنوان گره‌های گراف در نظر گرفته می‌شود و ماتریس تلاقی ماتریس شباهت است. برنامه‌ی خطی صحیح برای یافتن دسته‌ی بهینه یک نوع دیگری از مسئله‌ی K مرکزی است^۲. به‌طور ساده‌تر مسئله‌ی K مرکزی برای انتخاب K شهر برای ساختن مرکزها است. فاصله‌ی بین یک شهر و نزدیک‌ترین مرکز آن حداقل می‌شود. روش ILP برای بازشناسی گوینده‌ی بدون ناظر زمانی که تعداد گویندگان نامعلوم است بکار می‌رود.

با معرفی ILP در شاخه‌ی بازشناسی گوینده در حال حاضر سیستم‌های بازشناسی با استفاده از مدل مخلوط گوسی-مدل پنهان مارکوف یک مرحله پس پردازش دیکد کردن و پتری را بکار می‌برند. در سال‌های اخیر میگر و همکاران [۲۹] روش بالا را برای کم کردن زوائد برای درست کردن سریع دسته‌ها بکار گرفته‌اند. بدین ترتیب در سال‌های اخیر، روش‌های بر پایه ILP برخلاف روش‌های قدیمی بر پایه‌ی خوشه‌بندی متراکم سلسله مراتبی به دلیل کارکرد بهترشان در سرعت و دقت موردتوجه همگانی قرار گرفته‌اند. جعبه‌ابزار LUIM یک خطای DER برابر با ۱۷.۱۹ درصد را با استفاده از مدل مخلوط گوسی بر پایه دسته‌بندی بر اساس فاصله‌ی CLR گزارش داده است و در مقابل یک خطای DER برابر با ۱۵.۴۶ درصد را با استفاده از بردار \vec{a} بر پایه دسته‌بندی ILP نشان داده است. [۲۹]

۲-۷- روش‌های ارزیابی:

در بحث ارزیابی نتایج سیستم بازشناسی گوینده معیارهای ارزیابی متفاوتی وجود دارد. در سیستم‌های داده‌ای مختلف در دنیا بنا به کاربرد، معیار ارزیابی متفاوت است.

¹ Integer Linear Program (ILP)

² https://en.wikipedia.org/wiki/Integer_programming

۲-۷-۱- معیار DER

میزان خطای بازشناسی^۱ درصد زمانی را مشخص می‌کند که گوینده‌ها به اشتباه برچسب‌گذاری شده‌اند [۳۰]. خروجی سیستم با یک قسمت که به صورت دستی برچسب‌گذاری گوینده شده است، مقایسه می‌شود. فرمول زیرسیستم ارزیابی بازشناسی گوینده NIST 2005,2006,2007,2009 را مشخص می‌کند [۳۱]. در این ارزیابی سیستم مورد ارزیابی باید بازشناسی گوینده را بر روی مکالمات ضبط‌شده‌ی گروهی انجام دهد. در این فرمول، برای هر قسمت S با طول زمانی $dur(s)$ تعداد N_{ref} تعداد گوینده‌های مشخص شده در منبع اصلی و N_{hyp} تعداد گویندگان فرض شده توسط سیستم ما را به ترتیب نشان می‌دهد. $N_{correct}$ تعداد گوینده‌ها در قسمت s که تطابق درستی بین گوینده‌های منبع و سیستم ما را دارد مشخص می‌کند.

$$DER = \frac{\sum_{s=1}^S dur(s) \cdot (\max(N_{ref}, N_{hyp}) - N_{correct})}{\sum_{s=1}^S dur(s) \cdot N_{ref}} \quad (11-2)$$

زمانی که همپوشانی نداشته باشیم، فرمول DER ساده‌سازی می‌شود.

جدول ۱-۲: ساده‌سازی محاسبه‌ی خطا

زمان کلی غیر گفتار (NS)	زمان کلی گفتار (S)
غیر گفتار درست برچسب‌گذاری شده (C2) + زمان هشدار اشتباه گفتار (T2)	گفتار به درستی برچسب‌گذاری شده (C1) + زمان گفتار از دست رفته (T1) زمان گفتار نادرست برچسب‌گذاری شده (T3)

^۱DER: Diarization Error Rate

همان طور که در جدول (۲-۲) ملاحظه می‌کنید زمان‌های به دست آمده در جداسازی گفتار و غیر گفتار به صورت آمده در جدول می‌توانند جداسازی شوند. بدین صورت که زمان‌هایی را به عنوان زمان گفتار تلقی می‌کنیم که شامل گفتار درست تشخیص داده شده (C1) و گفتار تشخیص داده نشده یا همان گفتار از دست رفته (T1) و در آخر شامل زمان گفتاری است که به نادرست غیر گفتار برچسب گذاری شده است (T3) می‌باشد. همچنین برای زمان غیر گفتار به صورت جمع زمان‌های غیر گفتار درست تشخیص داده شده (C2) به همراه زمانی است که به غلط گفتار شناسایی شده است (T2) می‌باشد.

DER می‌تواند به صورت سیستماتیک به دو نوع تقسیم شود. فرض کنید صوتی با S ثانیه گفتار و NS ثانیه غیر گفتار به صورت دستی مشخص شده باشد. قسمت‌های غیر گفتار شامل سکوت، مکث گوینده و موزیک و نویز می‌باشد. دودسته به طور کامل و جزئی به صورت شکل زیر می‌توانند دسته بندی شوند.

زمان گفتار از دست رفته^۱ زمانی است که الگوریتم به غلط قسمتی که در اصل گفتار است را نا گفتار تشخیص می‌دهد. از طرف دیگر زمان هشدار اشتباه گفتار^۲ زمانی است که الگوریتم به غلط قسمتی که غیر گفتار می‌باشد را به عنوان گفتار تشخیص دهد. این دو خطای تشخیص فعالیت گفتار که بخش پیشین در تقریباً همه سیستم‌های تشخیص گوینده است رخ می‌دهد. تعداد خطاهای نرخ گفتار از دست رفته^۳ و نرخ تشخیص اشتباه گفتار به ترتیب به صورت $E1=T1*100/S$ و $E2=T2*100/S$ مشخص می‌شوند.

خطای $E3=T3*100/S$ به صورت مجزا توسط هر دو بخش قسمت بندی گوینده و دسته بندی گوینده تولید می‌شود و اغلب بانام خطای گوینده^۴ نام برده می‌شود. تغییر گوینده اگر در حین قسمت بندی از

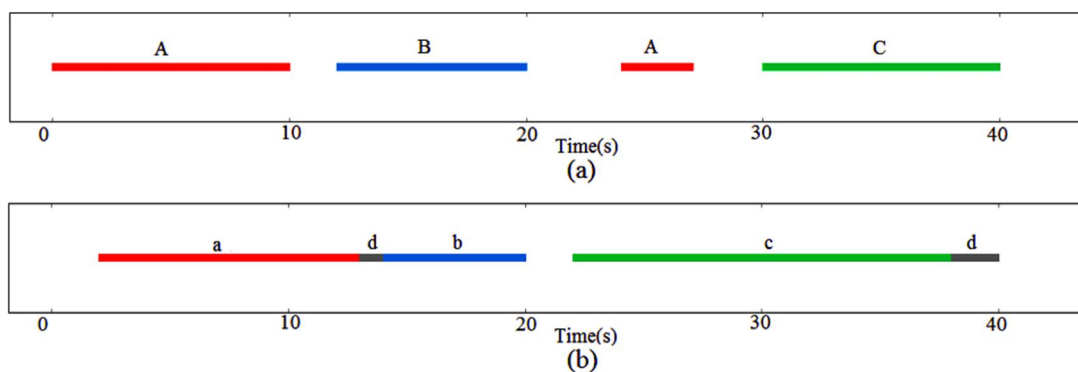
¹ Missed Speech Time

² False Alarm Speech Time

³ MSR: Missed Speech Rate

⁴ Speaker Error (spk_err)

دست برود باعث دسته‌بندی و کلاس‌بندی شدن نادرست قسمت کوتاه‌تر می‌شود. از طرف دیگر اگر سیستم صوت را به تعداد زیادی قسمت که با زیرنویس مشخص شده است تقسیم شود بیش از حد بخش‌بندی¹ نامیده می‌شود. شانس برای درست کردن این خطا حین دسته‌بندی با یکی کردن قسمت‌های همسایه وجود دارد. هرچند اگر قسمت موردنظر کوچک باشد احتمال خطا حین استخراج اطلاعات گوینده از آن بخش بالا می‌رود. قسمت به اشتباه بخش‌بندی شده دلیل اصلی E3 است. در آخر خطای کلی به صورت جمع تمامی خطاهای گفته‌شده در بالا محاسبه می‌شود.



شکل ۲-۶: مثالی از محاسبه‌ی DER (الف) منبع (ب) سیستم ما

محاسبه‌ی `spk_err` برای اجرای یک تطابق بین برچسب‌های الگوریتم و برچسب‌های منبع لازم و ضروری است. برای این منظور بهترین جفت کلاس را بر اساس بیشینه همپوشانی بین دودسته برچسب انجام می‌رود. مثالی در این زمینه در جدول زیر آورده شده است. در جدول ۲-۲ در بخش برچسب‌گذاری منبع همان‌طور که در شکل (۲-۶) مشاهده می‌کنید صوت ورودی ما به بخش‌های مختلف که هرکدام دارای گفتار مربوط به یک گوینده خاص است تقسیم‌شده است. در ستون کناری بخش‌بندی و تخصیص گوینده انجام‌شده است و در کنار داده‌های داشته از منبع برای مقایسه قرار گرفته است. در شرح داده‌های منبع می‌بینیم که صوت در درست دارای سه گوینده می‌باشد که بر اساس قسمت‌های برحسب ثانیه داده‌شده‌اند به‌طور مثال گوینده‌ی اول ثانیه‌ی اول تا دهم صحبت

¹ Oversegmentation

کرده است و از ثانیه‌ی ۱۰ تا ثانیه‌ی ۱۲ شامل غیر گفتار است و از ثانیه‌ی ۱۲ تا ثانیه‌ی ۲۰ گوینده‌ی دوم صحبت کرده است و به همین طریق در ادامه در بخش‌بندی و مشخص‌سازی دستی دیده می‌شود که ثانیه‌ی ۲ تا ۱۳ گوینده‌ی اول تشخیص داده شده است و از ثانیه‌ی ۱۳ تا ۱۴ را غیر گفتار تشخیص داده است که دیده می‌شود تطابق درستی با فرضیات منبع ندارد و این امر باعث به وجود آمدن خطاهای گفته شده می‌باشند. در ادامه‌ی تشریح جدول ۲-۲ می‌توان دید که در انتهای صوت نیز گفتار گوینده سوم را به اشتباه غیر گفتار دسته‌بندی کرده است.

جدول ۲-۲: مثالی از بخش‌بندی برای محاسبه‌ی DER

فرضیه	برچسب‌گذاری منبع
قسمت (ثانیه ۱۲م الی ثانیه ۱۳م) : (a) (گوینده الف)	قسمت (ثانیه ۱۰م الی ثانیه ۱۱م) = گوینده (الف)
قسمت (ثانیه ۱۳م الی ثانیه ۱۴م) : (d) (غیر گفتار)	قسمت (ثانیه ۱۲م الی ثانیه ۲۰م) = گوینده (ب)
قسمت (ثانیه ۱۴م الی ثانیه ۲۰م) : (b) (گوینده دوم)	قسمت (ثانیه ۲۴م الی ثانیه ۲۷م) = گوینده (الف)
قسمت (ثانیه ۲۲م الی ثانیه ۳۸م) : (c) (گوینده سوم)	قسمت (ثانیه ۳۰م الی ثانیه ۴۰م) = گوینده (ج)
قسمت (ثانیه ۳۸م الی ثانیه ۴۰م) : (d) (غیر گفتار)	

۲-۸- پایگاه داده‌های مورد استفاده

۲-۸-۱- AMI corpus

اولین پایگاه داده مورد استفاده در این پایان‌نامه AMI corpus [۳۲] می‌باشد که شامل ۱۰۰ ساعت گفتگوی ضبط شده بین گویندگان مختلف می‌باشد. مکالمات ضبط شده با استفاده از میکروفون‌های نزدیک و دور و نیز دوربین‌های کنفرانسی و تخته‌های نمایش مجزا است. جلسات در زبان انگلیسی با استفاده از سه اتاق‌های مختلف با خواص مختلف صوتی ثبت شده و شامل سخنرانان عمدتاً غیربومی است. مدت زمان این مکالمات و جنس گویندگان مختلف می‌باشد. مشخصات تک تک

ثانیه‌های مکالمات و ویژگی‌های گویندگان در دسترس می‌باشد. این پایگاه داده توسط گروه تحقیقاتی دانشکده انفورماتیک^۱ دانشگاه ادینبرگ در اسکاتلند^۲ تدوین و منتشر شده است و کارایی‌های مختلفی برای آن ذکر شده است که در این پایان‌نامه برای بازشناسی گوینده مورد استفاده قرار می‌گیرد.

۲-۸-۲- TIMIT

دومین پایگاه داده مورد استفاده در این پایان‌نامه دیتابیس TIMIT [۳۳] است. دادگان TIMIT یک بانک اطلاعاتی از گفتار پیوسته انگلیسی است که با همکاری شرکت TI^۳ و دانشگاه MIT^۴ منتشر شده است. این دادگان حاوی ۶۳۰۰ جمله است که توسط ۶۳۰ گوینده و با ۸ لهجه معمول آمریکای شمالی بیان شده‌اند. ۷۰٪ گویندگان مرد و ۳۰٪ آن‌ها زن هستند. هر گوینده ۱۰ جمله را ادا کرده است که ۲ جمله از این ۱۰ جمله توسط سایر گویندگان نیز ادا شده است. از این پایگاه داده در آموزش داده‌ها و مدل‌ها و در دیگر مواقع مورد نیاز بهره گرفته شده است.

۲-۹- خلاصه و جمع‌بندی فصل

در این فصل برخی مفهومی‌های پس‌زمینه‌ای مورد استفاده در بخش‌بندی و دسته‌بندی و تشخیص فعالیت گفتار مورد بررسی قرار گرفت. مابقی کارهای بازشناسی گوینده همانند تعداد گویندگان نامشخص است. بسیاری از روش‌های نشانه‌گذاری‌های صوت برای مثال تشخیص حضور آهنگ [۳۴] و کمک به یافتن ساختار برنامه‌ی رادیوتلوویزیونی [۳۵] در گذشته مورد بررسی قرار گرفتند. سیستم‌های معمول دیگر شامل موارد مختلفی من جمله ابزارهایی چون SHOUT [۳۶] و INRIA [۳۷] و IDIAO [۳۸] و DIARTK [۳۹] می‌باشند.

¹ <http://www.ed.ac.uk/informatics>

² <http://www.ed.ac.uk/>

³ [Texas Instruments](http://www.ti.com)

⁴ [Massachusetts Institute of Technology](http://www.mit.edu)

فصل ۳ :

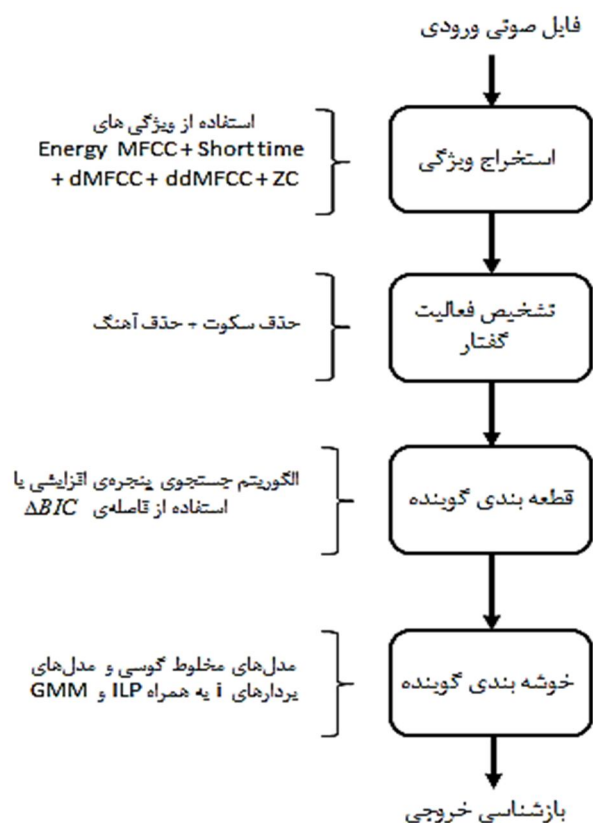
روش پیشنهادی در بازشناسی

گوینده

۳-۱- مقدمه

در این بخش به معرفی سیستم مورد استفاده در این پایان نامه می پردازیم. یک سیستم کامل در نرم افزار متلب برای بازشناسی گوینده فایل های صوتی اجرا شده است. این سیستم تست و ارزیابی شده است. برای ارزیابی معیارهای توضیح داده شده در بخش قبل مورد استفاده قرار گرفته اند. در این بخش توضیح کاملی از این سیستم از استخراج ویژگی تا ارزیابی خطاها را مورد بررسی قرار می دهیم. انتخاب پارامترها و مدل ها در هر بخش توضیح داده شده است.

با پیدایش الگوریتم های دسته بندی بهتر با مدل بردارهای n ، تمرکز روی اجرای جدای از هم بخش بندی گوینده و دسته بندی متمرکز شد. اگرچه استفاده از مدل مخلوط گوسی-مدل مخفی مارکوف همچنان معروف و مشهور است [۱۶]. شکل (۱-۴) بلوک دیاگرام سیستم پیشنهادی که روش پله به پله را نشان می دهد رسم شده است.



شکل ۳-۱: بلوک دیاگرام سیستم پیشنهادی

۳-۲- استخراج ویژگی

ویژگی ضرایب فرکانسی در مقیاس مل^۱ مکرراً در بازشناسی برای تمام زیرمجموعه‌های تشخیص فعالیت گفتار و بخش بندی و دسته بندی مورد استفاده قرار می گیرد. سیستم پیشنهادی ۱۹ ویژگی ضرایب کپسترال فرکانسی در مقیاس مل برای تمام پردازش ها بکار می گیرد. این استفاده به دلیل دربرگیری تمامی اطلاعات فریم های گفتار و غیر گفتار است تا بتوانیم در ادامه در مدل سازی قسمت های مختلف به نتیجه ی بهتری دست پیدا کنیم. علاوه بر این ویژگی ها، انرژی زمان کوتاه و نرخ عبور از صفر و مشتق های مرتبه اول و دوم آن ها در بخش تشخیص فعالیت گفتار بکار می روند. بخش بندی گوینده تنها ۱۹ ویژگی ضرایب کپسترال فرکانسی در مقیاس مل و انرژی زمان کوتاه را

^۱ MFCC : Mel Frequency Cepstral Coefficient

استفاده می‌کند. در بخش دسته‌بندی مشتق‌های مرتبه اول و دوم آن‌ها نیز بکار می‌رود. اندازه فریم‌ها برای آنالیز ۳۰ میلی‌ثانیه با همپوشانی ۲۰ میلی‌ثانیه است. [۲۳]

۳-۳- تشخیص فعالیت گفتار^۱

تشخیص فعالیت عبارت است از گفتار جداسازی قسمت‌های گفتار از غیر گفتار در فایل صوتی ضبط‌شده است. در حین انجام این جداسازی هدف مدنظر ما حداقل کردن خطاهای به وجود آمده در این بخش است. این خطاها شامل نرخ گفتار ازدست‌رفته^۲ و نرخ هشدار خطای گفتار^۳ می‌باشند. درصد گفتاری که به اشتباه در غیر گفتار توسط این بخش دسته‌بندی می‌شود نرخ گفتار ازدست‌رفته نام دارد و درصد غیر گفتاری که اشتبهاً گفتار دسته‌بندی شده است را نرخ هشدار خطای گفتار می‌نامند. عملاً این دو معیارهای ارزیابی بخش تشخیص فعالیت گفتار می‌باشند. معمولاً ۱۱ الی ۳ درصد خطای گفتار ازدست‌رفته و ۲ الی ۴ درصد خطای هشدار خطای گفتار خطاهای معمول امروزی در سیستم‌های تشخیص فعالیت گفتار است. در بازشناسی گفتار، هشدار خطای گفتار باعث خرابی در مدل‌های گوینده طی دسته‌بندی و بخش‌بندی می‌شود و فرایند دسته‌بندی را تحت تأثیر قرار می‌دهد.

یک روش معمول در سیستم‌های تشخیص فعالیت گفتار سعی بر دست‌بندی کردن همه‌ی انواع صداهای موجود در فایل صوتی است. اگر داده‌های بازشناسی شده از قبل معلوم باشند و ویژگی‌های خاصی مثل شاخص‌های صوتی برای مشخص کردن هر بخش داشته باشند، امکان آموزش مدل‌های آماری برای این شاخص‌ها در داده‌ها و دسته‌بندی مستقیم می‌باشد [۳۴]. متأسفانه هنگام بهبود سیستم‌های کلی این چنین ویژگی‌هایی همانند شاخص‌ها و ویژگی‌های صوتی برای جداسازی داده‌های صوتی ورودی نداریم. در این بخش مشکل تخمین مدل‌های غیر گفتار برای تشخیص فعالیت گفتار بدون اطاعات قبلی در مورد داده‌ها بررسی شده است.

¹ Speech Activity Detection

² Missed Speech Rate (MSR)

³ False Alarm Speech Rate (FASR)

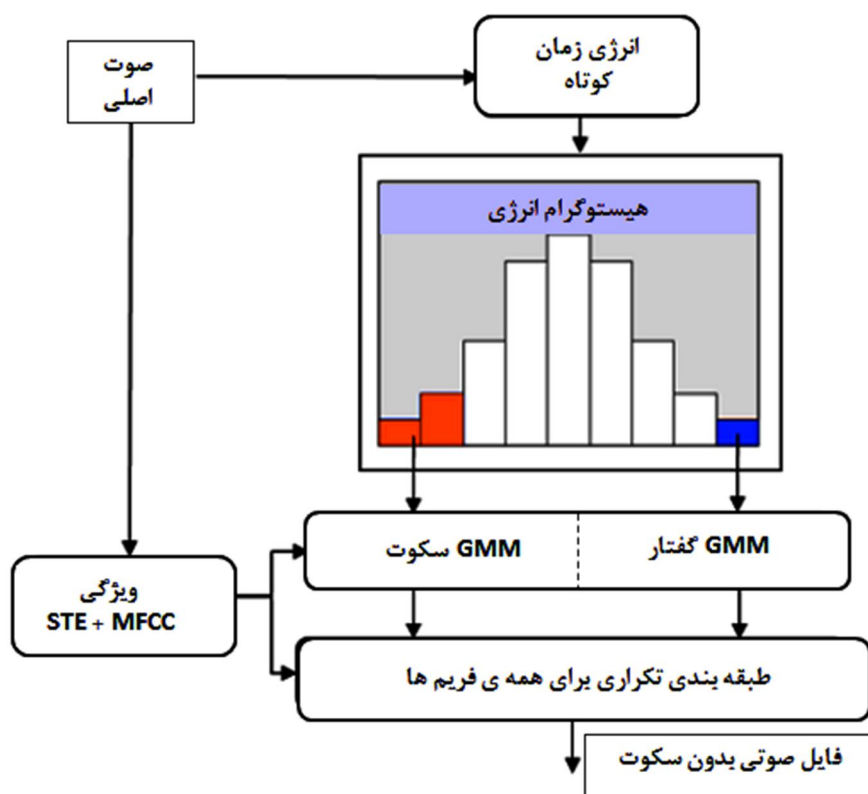
صداها علاوه بر گفتار معمولی که ما در صوت می‌بینیم انواع صداهای تولیدی انسان همانند ملچ‌ملوچ کردن و خندیدن و دست زدن و فوت کردن و یا صداهای ابزاری مانند آهنگ و جیرینگ جیرینگ و غیره نیز وجود دارند. قسمت‌های سکوت و مکث نیز سهم مهمی از بیشتر پایگاه‌های داده صوتی را دربرمیگیرند. این صداها با همدیگر مجموعه‌ی غیر گفتار را تشکیل می‌دهند. روش‌های بر پایه مدل در تشخیص فعالیت گفتار معمول هستند. مدل‌های آماری برای گفتار و غیر گفتار با داده‌های خارجی آموزش دیده شده‌اند. اگرچه ایراد این سیستم‌ها اتکای آن‌ها بر شرایط آکوستیکی خارج از داده‌ی نمونه است. سیستم‌های هیبرید [۶، ۱۰] از یک خوشه بند آموزش دیده شده با داده‌های خارجی استفاده می‌کنند تا یک بخش‌بندی اولیه‌ی متکی بر خود بسازند. سپس مدل‌های گفتار و غیر گفتار مرتبط با بخش‌بندی تعریف می‌شوند تا بتوانند تغییر آکوستیکی در غیر گفتار را مدل‌سازی کنند.

تشخیص‌دهنده‌ی گفتار در این سیستم دسته‌بندی بر اساس مدل است. همچنین مستقل از داده‌های آموزش برای مدل‌سازی گفتار و غیر گفتار است. این نوع سیستم بر پایه مدل تشخیص فعالیت گفتار بر اساس سیستم ارزیابی NIST RT2009 می‌باشد [۷]. در سیستم ما تشخیص فعالیت گفتار در دو مرحله‌ی مجزا صورت می‌گیرد. اولین مرحله شامل جداسازی سکوت از کل فایل صوتی با استفاده از انرژی توسط دسته‌بندی تکراری است. در مرحله دوم آهنگ و دیگر اصوات غیر گفتار تشخیص داده می‌شوند. برای حذف آهنگ، صوت بدون سکوت به سیستم تشخیص گری که در بالا گفته شد داده می‌شود. فریم‌هایی از صوت که دارای آهنگ با اطمینان بالا هستند برای آموزش مدل‌های آهنگ که مکرراً تکرار می‌شود بکار می‌روند. در هر دو مرحله، تنها بخش‌های با طول ۱ ثانیه و یا بیشتر به‌عنوان غیر گفتار برچسب می‌خورند به این دلیل که هر از چند گاهی از غیر گفتار به گفتار منتقل نشویم. این قید و محدودیت در منبع [۳۸] و [۳۶] با استفاده از مدل مخلوط گوسی - مدل مخفی مارکوف بکار گرفته شده است.

۳-۳-۱- حذف سکوت

حذف سکوت در روش اجرایی توسط ۱۹ ضرب کیسترال در مقیاس مل به همراه انرژی زمان کوتاه و مشتقات اول و دوم آن‌ها صورت گرفته است. بخش‌بندی خود راه‌انداز یک مقدار مشخص را در فریم به هر دودسته‌ی سکوت و گفتار اختصاص می‌دهد. مدل سکوت با استفاده از یک گوسی با اندازه‌ی ۴ روی فضای ویژگی ۶۰ بعدی آموزش داده می‌شود. مدل گفتار با همان اندازه از فریم‌های مطمئن گفتار آموزش داده می‌شود.

در یک مرحله‌ی تکراری دسته‌بندی کردن، هر فریم به دودسته سکوت و گفتار دسته‌بندی می‌شود. فریم‌های گفتار با اطمینان بالا و سکوت از این مرحله برای آموزش مدل‌های گفتار و سکوت برای اجرای بعدی بکار می‌روند. همان‌طور که تعداد اجراها بالا می‌رود تعداد گوسی‌های ۶۰ بعدی مورد استفاده برای مدل کردن مخلوط گوسی‌های گفتار و سکوت بالا می‌رود. این روند تا مقدار حداکثری ادامه دارد. همان‌طور که در شکل (۲-۳) مشاهده می‌کنید از فایل صوتی اصلی فریم‌های بانرژی حداکثر و حداقل جداگانه برای تشخیص گفتار و سکوت از هم بکار می‌روند و با استفاده از آن‌ها مدل‌های گوسی برای سکوت و گفتار تولید شده و توسط این مدل‌ها و دسته‌بندی آن‌ها فریم‌های سکوت از فایل اصلی جدا می‌شود و فایل صوتی بدون سکوت در آخر منتج می‌شود.



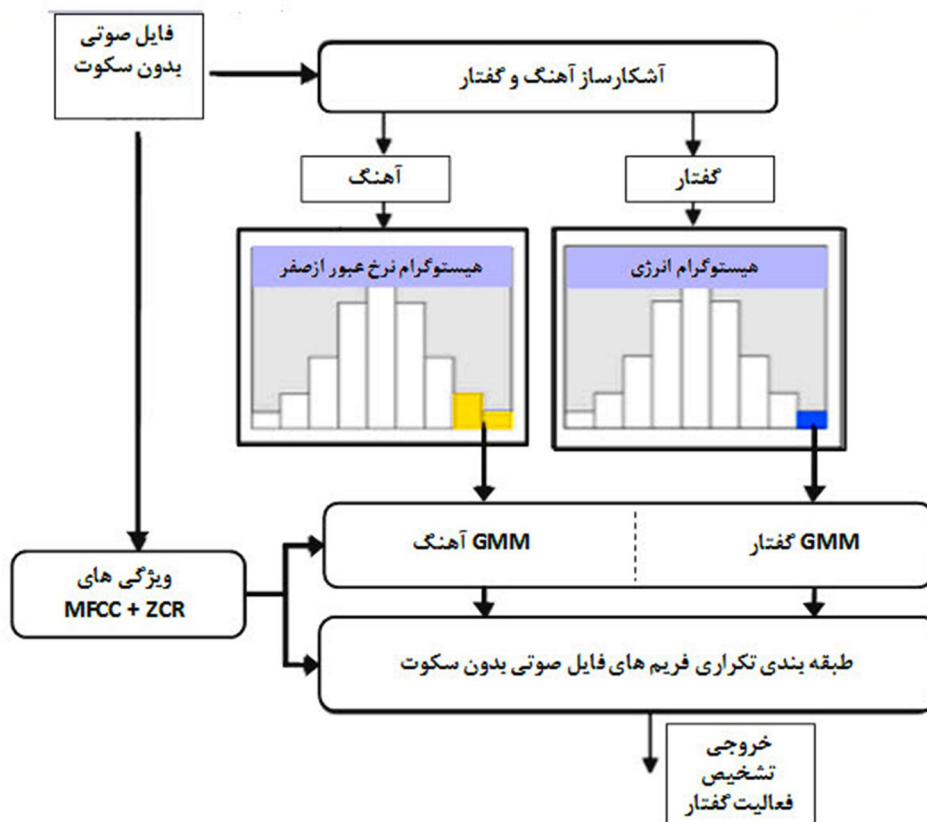
شکل ۳-۲: حذف سکوت با استفاده از خوشه بندی تکراری

۳-۲-۳- حذف موزیک

در سال ۲۰۰۵ در [۲۰] یک مدل سازگار بر پایه دسته بند آهنگ و گفتار که درصد دسته بندی ۹۵ درصد را گزارش داده معرفی شد. نویسنده فایل صوتی را به قطعه های ۱ ثانیه ای بخش بندی کرده است و سپس ۵۰ بردار ویژگی را در طول پنجره های ۲۰ میلی ثانیه استخراج کرده است. این بردارهای ویژگی دوبعدی بودند- بعد اول انرژی کوتاه مدت و بعد دوم نرخ عبور از صفر سیگنال پنجره بندی شده- هیستوگرام انرژی کوتاه مدت و نرخ عبور از صفر برای قطعات ۱ ثانیه ای محاسبه شدند و با هیستوگرام مدل های موزیک و آهنگی که از یک دیتابیس بزرگ موزیک و گفتار به دست آمده اند مقایسه می شوند. این هیستوگرام های ایده آل با توزیع X^2 مدل شدند. قطعه ای مورد بررسی بعد از مقایسه بین هیستوگرام قطعه ای یک ثانیه ای و مدل های X^2 به عنوان آهنگ یا گفتار برچسب گذاری می شود.

این شاخص گفتار و آهنگ [۲۰] زمانی که گفتار و آهنگ باهم هستند شکست می خورد. در

اغلب داده‌های گفتاری برای جداسازی بخش‌های مختلف گفتاری آهنگ‌های پس‌زمینه خاصی به آن اضافه می‌شود که این باعث پیدایش خطا و از دست رفتن اطلاعات اصلی ما می‌شود. از آنجاکه ما خروجی دسته‌بند را به‌عنوان یک بخش‌بندی قوی در نظر می‌گیریم، تخمین‌های اولیه مدل‌ها برای آهنگ و گفتار بر اساس فریم‌های با اطمینان بالا از هر دو کلاس آموزش می‌بینند. یک دسته‌بندی مشابه با سیستم حذف سکوت برای جداسازی دسته‌های موزیک و گفتار و حذف بخش‌های موزیک انجام شد (شکل ۳-۳). ویژگی‌های مورد استفاده ۱۹ ضریب ضرایب کپسترال فرکانسی در مقیاس مل به همراه نرخ عبور از صفر و مشتقات اول و دوم آن‌ها است. انرژی زمان کوتاه در این فرایند استفاده نشده است. بعد از نادیده گرفتن انرژی زمان کوتاه گفتار با آهنگ پس‌زمینه که در دسته‌بند موزیک قرار گرفته بود جداسازی شد.



شکل ۳-۳: حذف موسیقی با استفاده از یک آشکارساز موسیقی گفتار

۳-۳-۳- اندازه‌گیری‌های تشخیص فعالیت گفتار

در طول حذف سکوت، هیستوگرام انرژی فریم‌ها برای مرتب‌سازی فریم‌ها بر اساس انرژی آن‌ها استفاده می‌شود. فریم‌های در ۲۰ درصد پایین هیستوگرام به‌عنوان فریم‌های قطعاً سکوت و فریم‌های در ۱۰ درصد بالایی نمودار به‌عنوان فریم‌های گفتار با درصد قاطعیت بالا در نظر گرفته می‌شوند. در هر اجرا و مرحله، این فریم‌ها برای آموزش مدل مخلوط گوسی^۱ بکار می‌روند. برای حذف آهنگ، هدف ما جمع‌آوری فریم‌هایی است که گفتار همراه با آهنگ پس‌زمینه هستند اما به‌عنوان غیر گفتار دسته‌بندی شده‌اند. بدین منظور تنها ۴۰ درصد بالای فریم‌های نرخ عبور از صفر از هیستوگرام نرخ عبور از صفر را به‌عنوان فریم‌های خالص آهنگ برای آموزش مدل‌های آهنگ بکار گرفتیم.

۳-۴- ارزیابی تشخیص فعالیت گفتار

در این بخش نتایج ۳ آزمایش مجزا نشان داده شده است. اثر اندازه‌ی مدل‌های مخلوط گوسی در مراحل اجزای SAD و وضعیت ماتریس کوواریانس مدل‌های مخلوط گوسی و اثر اتصال این دو سیستم به هم دیگر مورد بررسی قرار گرفت. نرخ گفتار از دست‌رفته^۲ درصد خطای زمانی است که گفتار به‌صورت غیر گفتار دسته‌بندی شود. نرخ هشدار اشتباه گفتار^۳ خطای زمانی است که صوت غیر گفتاری به اشتباه گفتار تلقی شود. خطای SAD شامل جمع آن دو خطاست.

۳-۴-۱- اندازه‌ی مدل‌های مخلوط گوسی‌ها حین جداسازی سکوت

پس از بخش‌بندی، مدل‌ها برای فریم‌های محدود گفتار و غیر گفتار آموزش می‌بینند که

^۱ GMM: Gaussian Mixture Model

^۲ Missed Speech rate (MSR)

^۳ False Alarm Speech Rate (FASR)

مدل‌های مخلوط گوسی با اندازه‌ی ۴ هستند. به میزانی که داده‌های موجود در دسته‌ها افزایش می‌یابد، اندازه‌ی مدل‌های مخلوط گوسی استفاده‌شده برای مدل‌سازی گفتار و سکوت نیز افزایش می‌یابد. در هر مرحله‌ی اجرا، اندازه‌ها دو برابر می‌شوند تا به یک حد بیشینه‌ای برسند. در این آزمایش‌ها اندازه‌ی مناسبی را که برای هردو بخش بهتر است انتخاب کرده‌ایم. در جدول (۱-۳) درصد خطای SAD را برای اندازه‌های مختلف مدل‌های مخلوط گوسی نشان داده‌شده است. این اندازه‌های بهینه برای گفتار ۳۲ مدل مخلوط گوسی و برای غیر گفتار ۱۶ مدل مخلوط گوسی به دست آمدند.

جدول ۱-۳: خطای SAD برای اندازه‌های مختلف مدل‌های مخلوط گوسی برای مدل‌های گفتار و سکوت (برحسب درصد).

گفتار					
سکوت	اندازه مدل‌های مخلوط گوسی	۸	۱۶	۳۲	۶۴
	۴	۱۳.۴۲	۹.۸۲	۵.۹۴	-
	۸	۱۰.۶۲	۵.۷۹	۵.۵۴	۵.۸۹
	۱۶	۶.۴۴	۵.۷۴	۵.۲	۵.۴۹
	۳۲	۱۳.۵	۶.۵۲	۵.۹۶	-

۳-۴-۲- ماتریس کوواریانس مدل‌های مخلوط گوسی در دسته‌بندی تکراری حذف سکوت و

حذف آهنگ

در این بخش، مدل‌های مخلوط گوسی گفتار، سکوت و آهنگ با استفاده از ۱۶ و ۳۲ و ۳۲ گوسی در آخرین مرحله از دسته‌بندی مدل شدند. یک مدل مخلوط گوسی کوواریانس کامل برای آموزش بهتر کلاس گفتار می‌تواند مفید باشد. این باعث MSR پایین می‌شود هرچند برای دسته‌های سکوت و آهنگ تعداد فریم‌ها محدود شدند اما یک مدل مخلوط گوسی کوواریانس کامل مناسب است. هرچند خطای FASR به‌طور مشخص افزایش می‌یابد.

جدول ۳-۲: خطای برای SAD حالت ماتریس کوواریانس مدل‌های مخلوط گوسی (برحسب درصد).

حالت ماتریس	کامل	قطری
مرحله انجام		
حذف سکوت	۵.۴۹	۵.۲
حذف آهنگ	۸.۱۳	۷.۶۹

۳-۴-۳- اثر متوالی کردن حذف سکوت و حذف آهنگ

استفاده از سیستم حذف موزیک یک نرخ خطای غلط زیادی را می‌دهد. متوالی کردن دو سیستم بهبود خوبی را در مقایسه با استفاده مجزای آن‌ها می‌دهد. همان‌طور که انتظار می‌رفت، MSR افزایش می‌باید ولی FASR به دلیل دسته‌بندی فریم‌های موزیک در غیر گفتار کاهش می‌یابد.

جدول ۳-۳: متوالی کردن سیستم حذف آهنگ و حذف سکوت برای SAD (برحسب درصد).

	MSR	FASR	جمع دو خطا
حذف سکوت	۱.۴۹	۳.۷۱	۵.۲۰
حذف آهنگ	۱.۵۹	۶.۱۱	۷.۶۹
سیستم متوالی	۲.۳	۲.۷۱	۵.۰۱

۳-۴-۴- نرخ عبور از صفر برای آموزش مدل‌های آهنگ

استفاده از فریم‌هایی که به‌طور یقین آهنگ هستند که از مقادیر نرخ عبور از صفر به دست آمده‌اند، برای آموزش مدل‌های آهنگ و نادیده گرفتن فریم‌های با نرخ عبور از صفر پایین باعث کاهش نرخ خطای از دست رفتن گفتار می‌شود. این افزایش به علت برگرداندن فریم‌های گفتار

از فریم‌های گفتار مخلوط با آهنگ پس‌زمینه است. افزایش FASR به دلیل گفتاری است که از فریم‌های با صدای جیرینگ و ... به دست آمده است.

جدول ۳-۴: تعریف مجدد مدل‌های آهنگ با استفاده از فریم‌های با ZCR بالا (برحسب درصد)

	MSR	FASR	جمع هر دو
فریم‌های آهنگ با ZCR بالا	۲.۳	۲.۷۱	۵.۰۱
همه فریم‌های آهنگ	۳.۰۸	۲.۳۲	۵.۴۱

۳-۵ - بخش‌بندی گوینده

الگوریتم بخش‌بندی گوینده مورد استفاده در این پایان‌نامه یک الگوریتم جستجوی پنجره‌ی افزایشی [۱۴] با استفاده از فاصله‌ی ΔBIC که در ادامه معرفی می‌شود، است. در ابتدا، جستجو برای یافتن یک تغییر گوینده صورت می‌گیرد. در هر تغییر پیداشده، جستجو از فریم بعدی شروع می‌شود. پنجره‌ی اولیه با مقدار ۵ ثانیه مقداردهی اولیه می‌شود و یک مقدار ΔBIC برای هر فریم از پنجره محاسبه می‌شود. اگر حداکثر این آرایه از حد آستانه‌ی θ ، فراتر رود، تغییری در حد آستانه تشخیص داده می‌شود. اگر حداکثری در طول پنجره‌ی جستجو یافت نشود، اندازه پنجره به ۲ ثانیه تغییر می‌کند و روند دوباره تکرار می‌شود. دقت شود تنها فریم‌های گفتار به دست آمده از بخش تشخیص فعالیت گفتار مورد پردازش قرار می‌گیرد. بعد از یافتن نقاط تغییر در فریم‌ها، مکان آن‌ها در فایل اولیه و اصلی مشخص می‌شوند و به عنوان نقاط تغییر شناخته می‌شوند.

در دو سیستم قبل [۳۸،۳۶] بخش‌بندی در دو بخش صورت می‌گیرد: مرحله اول ΔBIC تشخیص نقاط تغییر همان‌طور که گفته شد با آستانه‌ی صفر اجرا می‌شود سپس ادغام بخش‌های متوالی برای آن‌هایی که امتیاز ΔBIC شان مثبت است صورت می‌گیرد. دو مرحله‌ای بودن به علت

بیش از حد بخش بندی کردن بر اساس حد آستانه‌ی ΔBIC صفر می‌باشد. برای جلوگیری از این امر، تنها بیشینه‌ای که از یک آستانه‌ی (θ) صفر بزرگ‌تر است انتخاب می‌شود. این امر به مراتب بیش از حد دسته بندی کردن را کاهش می‌دهد:

$$\Delta BIC(x_i) = N \log |\Sigma| - N_1 \log |\Sigma_1| - N_2 \log |\Sigma_2| - \frac{\lambda}{2} (d + \frac{1}{2} d(d+1)) \log N \quad (1-3)$$

در (۲-۳) پارامتر Σ ماتریس کوواریانس داده‌ها و d بعد فضای بررسی ما و N تعداد داده‌های بردارهای ویژگی را در هر مدل گوسی نشان می‌دهد. برای بخش بندی گوینده، ۱۹ ضریب کپسترال فرکانسی در مقیاس مل با انرژی زمان کوتاه استفاده شدند. اغلب سیستم‌های بخش بندی گوینده [۳۸،۳۶،۶] مشتقات ویژگی‌های کپسترال را هنگام بخش بندی استفاده نمی‌کنند.

۳-۶- دسته بندی گوینده

بعد از تشخیص نقاط تغییر در گوینده با استفاده از بخش بندی گوینده زیر بخش دسته بندی گوینده باهدف جمع آوری بخش‌هایی که گوینده یکسان دارند اجرا می‌شود. بدین منظور بخش‌های فراهم شده با مدل‌های گوینده نشان داده می‌شوند. شباهت دوجه دویی بین همه‌ی مدل‌های گوینده به دست می‌آید و الگوریتم دسته بندی برای گروه بندی کردن آن‌ها بکار می‌رود.

۳-۶-۱- انتخاب مدل گوینده

سیستم پیشنهادی دو مدل گوینده را بررسی می‌کند که به طور وسیعی در حیطه‌ی تشخیص گوینده و بازشناسی گوینده مورد مطالعه قرار می‌گیرد: ۱- مدل‌های مخلوط گوسی ۲- مدل‌های بردارهای i . مدل مخلوط گوسی (۳-۵) یک مدل احتمالی روی فضای ویژگی است. ویژگی‌های استفاده شده انرژی زمان کوتاه و ۱۹ ویژگی ضریب کپسترال فرکانسی در مقیاس مل و مشتقات اول و دوم آن‌ها در یک فضای ویژگی ۶۰ بعدی است. شباهت بین مدل مخلوط گوسی، بر پایه‌ی مدل بیشترین شباهت متقابل هر بخش از داده با داده‌ی دیگر است. مدل مخلوط گوسی برای

هر بخش روی بردارهای ویژگی بخش با استفاده از الگوریتم بیشینه امید^۱ آموزش داده می‌شود تا یک مدل مخلوط گوسی با ماتریس کوواریانس قطری با اندازه ۳۲ به دست آید. در حین محاسبه‌ی سیستم با استفاده از مدل‌های گوینده مدل مخلوط گوسی، فاصله‌های CLR^۲ و NCLR^۳((۶-۳)،(۷-۳)) با استفاده از الگوریتم‌های دسته‌بندی HAC^۴ و ILP^۵ محاسبه می‌شوند.

۳-۶-۲- استخراج بردار i

برای به دست آوردن بردار i ، ابتدا یک مدل زمینه کلی روی داده‌های آموزش، آموزش داده می‌شوند. مدل پس‌زمینه کلی^۶ یک مدل مخلوط گوسی با تعداد بسیار زیادی از گوسی‌ها است بنابراین می‌تواند تمامی تنوع‌های ممکن در گفتار را در فضای ویژگی جمع‌آوری کند. در سیستم پیشنهادی با استفاده از دیتابیس TIMIT مدل پس‌زمینه کلی آموزش داده می‌شوند. مدل پس‌زمینه کلی یک مدل مخلوط گوسی با ماتریس کوواریانس قطری با اندازه ۵۱۲ است. آموزش مدل پس‌زمینه کلی تنها یک‌بار انجام می‌شود. متوسط مدل پس‌زمینه کلی و مدل مخلوط گوسی بخش وفق یافته برای به دست آوردن یک ابر بردار با اندازه ۳۰۷۲۰ به هم متصل می‌شوند.

فضای تغییر کلی یک زیر فضا از ابر فضای مدل مخلوط گوسی است که تمام گوینده‌ها و اطلاعات مرتبط با آن‌ها را دربردارد. ماتریس T ماتریس مرتبه پایینی است که ستون‌هایش کل زیرمجموعه‌ی تغییر را در برمی‌گیرد. برای سیستم پیشنهادی، ماتریس T ، با استفاده از همان گویندگانی که برای آموزش مدل پس‌زمینه کلی استفاده شدند آموزش می‌بینند. فرایند آموزش ماتریس T نیز یک فرایند یک مرحله‌ای است. بردار i هر بخش تصویر ابر بردار مدل مخلوط گوسی بر روی زیر فضای تغییر کلی است.

¹ Expectation-Maximization Algorithm

² Cross Likelihood Ratio

³ Normalized Cross Likelihood Ratio

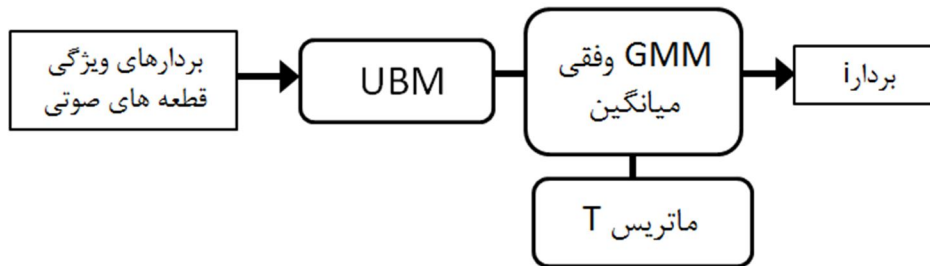
⁴ Hierarchical Agglomerative Clustering

⁵ Integer Linear Program

⁶ Universal Background Model (UBM)

$$m = M + Tx \quad (2-3)$$

M ابر بردار مدل پس‌زمینه کلی است و m ابر بردار مدل مخلوط گوسی متوسط و فقی هر بخش است. پس برای هر بخش استخراج بردار i که x است شامل دو مرحله است: تطبیق مدل پس‌زمینه کلی برای به دست آوردن ابر بردار مدل مخلوط گوسی^۱ مرتبط با آن و استخراج فاکتورهای بردارهای ویژه‌ی تغییر کلی است. الگوریتم آموزش ماتریس T از گویندگان برچسب خورده داده‌های آموزش در [۲۳] شرح داده شده است. سیستم پیشنهادی ابزار تشخیص MSR را برای آموزش مدل پس‌زمینه کلی و آموزش استخراج بردار i استفاده می‌کند. [۲۴]



شکل ۳-۴: استخراج بردارهای i

حین ارزیابی سیستم با بردارهای i، بعد زیر فضاها و انتخاب فواصل بین بردارهای i امتحان می‌شوند. دو فاصله‌ی متریک برای اندازه‌گیری شباهت بین بردارهای i بکار می‌رود: شباهت متریک کسینوسی^۲ (۳-۳) و فاصله متریک مهالانوبیس^۳ (۴-۳) که W ماتریس کوواریانس بین دسته‌های اندازه‌گیری شده از n بردار آموزش i از S گوینده شرح داده شده در بخش‌های قبلی است. فاصله‌ی مهالانوبیس از این پس کوواریانس نرمال بین دسته‌های^۴ نامیده می‌شود. در (۵-۳) محاسبه‌ی WCCN برای فاصله‌ی مهالانوبیس، بردارهای \bar{w}^s متوسط n_s بردار i از S گوینده است.

^۱ ابر بردار برداری است که بعد از اتصال همه‌ی بردارهای میانگین مدل‌های مخلوط گوسی به دست می‌آید

^۲ Cosine Similarity Metric

^۳ Mahalanobis Distance Metric

^۴ within class covariance normalization

$$D(x, y) = 1 - \frac{x^T y}{\|x\| \cdot \|y\|} \quad (3-3)$$

$$D(x, y) = (x - y)^T W^{-1} (x - y) \quad (4-3)$$

$$W = \frac{1}{n} \sum_{s=1}^S \sum_{i=1}^{n_s} (w_i^s - \bar{w}^s)(w_i^s - \bar{w}^s)^T \quad (5-3)$$

۳-۶-۳- انتخاب الگوریتم خوشه‌بندی

سیستم پیشنهادی با دو روش خوشه‌بندی کار می‌کند: الگوریتم خوشه‌بندی متراکم سلسله مراتبی سنتی^۱ و یک الگوریتم خوشه‌بندی گرافیکی به نام الگوریتم ILP^۲ که اخیراً در بحث بازشناسی گوینده در سال ۲۰۱۳ معرفی شد [۳۸].

۳-۶-۴- HAC^۳

در سیستم HAC، در هر مرحله اجرا مدل‌های گوینده‌ی خیلی مشابه برای ادغام انتخاب می‌شوند. در مرحله‌ی ادغام، یک مدل گوینده جدید با استفاده از داده‌های تمام بخش‌های هردو گوینده مشابه ادغام شده تولید می‌شود. ماتریس شباهت با شباهت مدل گوینده‌ی جدید از مدل‌های دیگر به‌روزرسانی می‌شود. در هر مرحله تعداد خوشه‌ها یکی کاهش می‌یابد. فرایند تا زمانی که تنها یک خوشه باقی بماند ادامه پیدا می‌کند. مجموعه بهینه‌ی خوشه‌ها از میان خروجی‌های هر مرحله بر اساس یک معیار بهینگی انتخاب می‌شوند. خوشه‌بندی سلسله مراتبی متراکم می‌تواند با استفاده از هریک از مدل‌های گوینده کار کند. در هر مرحله از خوشه‌بندی سلسله مراتبی متراکم، ادغام نیازمند

¹ Traditional Hierarchical Agglomerative Clustering Algorithm

² Integer Linear Program

³ Hierarchical Agglomerative Clustering

محاسبات اضافی است چون مدل‌ها دوباره باید آموزش ببینند ماتریس فاصله باید طبق مدل گوینده‌ی ادغام‌شده به‌روز شود.



شکل ۳-۵: خوشه‌بندی ILP روی یک گراف کامل از مدل‌های گوینده [۲۹]

همان‌طور که در شکل ۳-۵ دیده می‌شود خوشه‌بندی ILP بر روی گراف مدل‌های گوینده کار می‌کند و با استفاده از معیارهای بهینه‌ی خوشه‌بندی به تصمیم‌گیری برای انتخاب خوشه‌ی بهینه می‌پردازد.

دو معیار خوشه‌بندی بهینه وجود دارد که برای این خوشه‌بندی اجرا شده‌اند. اولین معیار آستانه‌ی فاصله است که در آن، مرحله‌ای از اجرا که در آن حداقل شباهت دوجه‌دویی بزرگ‌تر از یک آستانه باشد، به‌عنوان دسته‌ی بهینه از خوشه‌ها شناخته می‌شود. (۲-۹) معیار بعدی معیار Ts خوشه‌ی بهینه‌ی معرفی شده در [۲۵] است، مجموعه خوشه‌ها با حداقل شباهت برون دسته‌ای و بیشترین شباهت درون دسته‌ای با استفاده از (۲-۱۰) به‌عنوان دسته‌ی بهینه انتخاب می‌شود.

۳-۶-۵- خوشه‌بندی ILP

در خوشه‌بندی ILP، مسئله‌ی K مرکزی برای به دست آوردن مجموعه دسته‌ها اصلاح می‌شود. مسئله‌ی اصلی K مرکزی اصلی تشخیص K شهر از N شهر برای تشکیل مراکز به‌طوری است که

بیشترین فاصله‌ی بین یک شهر و نزدیک‌ترین مرکز به آن حداقل شود. در مورد فرمول سازی ILP، N قسمت مشابه N شهر و K تا از قسمت‌ها به‌عنوان بهترین قسمت‌ها برای داده‌های K گوینده انتخاب شوند. باین حال در مورد بازشناسی، تعداد گوینده‌ها (K) نامعلوم است؛ بنابراین مسئله خارج از مسئله‌ی بهینه‌سازی می‌شود.

مجموعه‌ی دودویی از متغیرهای تصمیم X_{ii} را در نظر بگیرید.

$X_{ii}=1$ مشخص می‌کند که خوشه‌ی i ام یک خوشه نماینده است.

$X_{ij}=1$ مشخص می‌کند که دسته‌ی i ام به خوشه نماینده‌ی j ام تعلق دارد و تخصیص داده‌شده

است. (حتماً لازم است که $X_{jj}=1$ باشد)

توجه کنید که $X_{ij}=1$ و $X_{ji}=1$ معانی متفاوتی دارند اگرچه هر دو ی آن‌ها نشان می‌دهند که

بردارهای i برای بخش‌های i و j به یک خوشه تعلق دارند. اکنون، مسئله‌ی بهینه‌سازی زیر را در نظر

بگیرید.

$$\begin{aligned} \min \quad & \sum_{i=1}^N X_{ii} + \frac{1}{\delta} \sum_{i=1}^N \sum_{j=1}^N d_{ij} X_{ij} \\ \text{s.t.} \quad & \sum X_{ij} = 1 \quad \forall j \\ & X_{ij} \leq X_{ii} \quad \forall j \\ & d_{ij} X_{ij} \leq \delta \quad \forall i, j \\ & X_{ij} \leq \{0,1\} \quad \forall i, j \end{aligned} \quad (6-3)$$

تابع هدف مدنظر برای بهینه‌سازی شامل دو قسمت است، اول N تعداد نمایندگان خوشه‌ها

(تعداد گویندگان) و دوم δ پراکندگی کلی همه‌ی خوشه‌هاست. قید اول تأکید می‌کند که یک قسمت

فقط به یک خوشه تعلق می‌گیرد. قید دوم بر این نکته اشاره دارد که مرکز خوشه به همان خوشه

تعلق داشته باشد. قید سوم جلوی تخصیص یک بردار به خوشه را وقتی که از یک آستانه‌ی خاصی

دورتر از مرکز خوشه است، می‌گیرد.

در نظر داشته باشید که الگوریتم خوشه‌بندی برنامه عدد صحیح خطی^۱ نیازمند هیچ اطلاعات درباره‌ی اهدافی که قصد خوشه‌بندی داریم ندارد و تنها وابسته به ماتریس شباهت است. برنامه‌ی عددی نیازمند این است که به یک ILP تک‌بعدی تبدیل شود.

تنها معایب استفاده از الگوریتم خوشه‌بندی برنامه عدد صحیح خطی این است که مدل‌های گوینده در فرایند به دست آوردن مدل‌های گوینده در خوشه‌بندی سلسله مراتبی متراکم^۲ مرتباً تعریف نمی‌شوند؛ بنابراین اگر اندازه‌ی بخش‌ها کوچک باشند بردارهای i ای که به‌عنوان رهبر بردارها انتخاب می‌شوند ممکن است نتوانند تمامی اطلاعات گوینده را در برگیرند. این الگوریتم خوشه‌بندی می‌تواند همراه با دو مدل گوینده‌ها بکار گرفته شود- مدل‌های مدل مخلوط گوسی و بردارهای i - این به این دلیل است که تنها ماتریس شباهت برای استخراج مجموعه‌ی بهینه‌ی خوشه‌ها لازم است.

¹ Integer Linear Program (ILP)

² HAC

فصل ٤ : نتایج

۱-۴- مقدمه

در این بخش به بررسی و تحلیل و دلایل نتایج به دست آمده می پردازیم و در انتها پیشنهادات و کارهایی پیش رو را بیان می کنیم .

۲-۴- پارامترهای بخش بندی

منظور از پارامترهای بخش بندی دو پارامتر θ و λ در (۱-۴) هستند که برای کم کردن خطاها هر کدام را طوری انتخاب می کنیم که بهینه شوند.

اثر هشدار اشتباه تغییر گوینده این است که تعداد بخش بندی ها را بالا می برد و در نتیجه اندازه ی بخش ها پایین می آید. برای بررسی اثرات ناشی از هشدار اشتباه بهترین الگوریتم برای خوشه بندی (خوشه بندی ILP با بردارهای i) مورد استفاده قرار گرفت. مشاهده شده است که مقادیر پایین θ و λ بخش بندی بیش از حد را به وجود می آورند. به میزانی که θ را افزایش دهیم میانگین طول بخش ها نیز افزایش می یابد. در نتیجه ی این افزایش مدل سازی بهتری از گوینده ها را نتیجه می دهد و باعث کم کردن خطای DER می شود.

جدول ۱-۴ : مقادیر DER با بهترین الگوریتم های خوشه بندی (برحسب %)

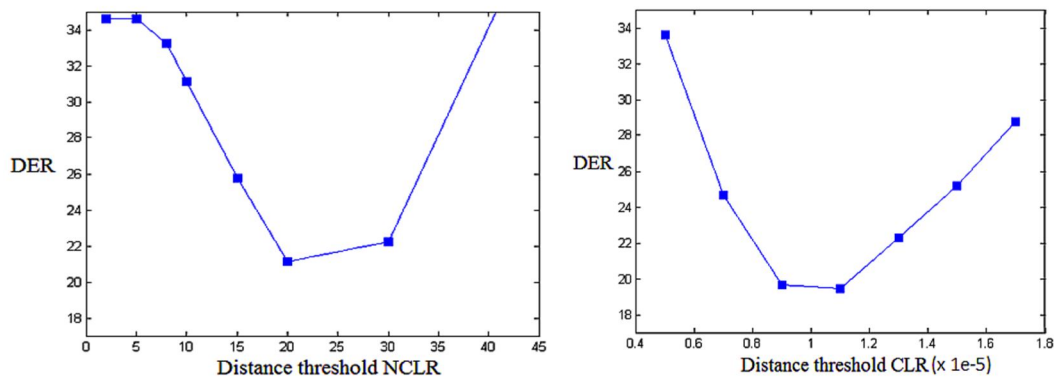
مقادیر پارامترها		λ	
		۱	۱۰
θ	۰	۳۱.۵۲	۳۳.۴۱
	۱۰۰۰	۲۳.۶۷	۱۶.۵۴
	۲۰۰۰	۱۲.۳۵	۱۶.۵۹

همان‌طور که در جدول (۴-۱) مشاهده می‌شود با ثابت نگه‌داشتن λ و افزایش پارامتر دیگر مقدار خطا در هر دو مقدار کاهش پیدا می‌کند. در مورد خود پارامتر λ با ثابت نگه‌داشتن θ و افزایش پارامتر دیگر مقادیر خطا افزایش می‌یابد پس در نتیجه‌ی جدول بالا ترکیب مقادیر $\theta = 2000$ و $\lambda = 1$ اتخاذ می‌شوند.

۴-۳- ارزیابی خوشه‌بندی گوینده‌ها

در این بخش ارزیابی خوشه‌بندی بر اساس خروجی‌های بخش‌های بخش‌بندی و تشخیص فعالیت گفتار را بیان می‌کنیم. در بخش‌های قبل خوشه‌بندی سلسله مراتبی سنتی را با استفاده از مدل‌های مخلوط گوسی و مدل‌های بردار i نشان دادیم. در ادامه الگوریتم خوشه‌بندی برنامه عدد صحیح خطی و مدل‌های مدل مخلوط گوسی و مدل‌های بردار i را فراهم می‌کنیم.

در این قسمت مقادیر DER نرخ خطای بازشناسی کلی هستند که میانگین DER های مجزای هر بخش که بر اساس مدت هر بخش وزن بندی شده است را نشان می‌دهد. در شکل (۴-۱) مشاهده می‌کنید با افزایش حد آستانه‌ی فاصله میزان خطای به وجود آمده در سیستم افزایش پیدا می‌کند و این به دلیل دسته‌بندی داده‌ها به صورت اشتباه در خوشه‌ی نادرست است. چون حد آستانه‌ی فاصله را زیاد کردیم پس داده‌های درون یک خوشه فاصله‌های زیادتری از هم پیدا می‌کنند و پراکندگی داده‌ها در یک خوشه افزایش پیدا می‌کند. این باعث می‌شود که این امکان که داده‌های خوشه‌های دیگر در این خوشه نمایان شوند افزایش پیدا می‌کند. لذا یک حد وسط برای آستانه‌ی فاصله لازم و ضروری است.



شکل ۱-۴ HAC با حد آستانه‌ی فاصله برای مدل‌های مخلوط گوسی

۴-۳-۱- آزمایشات خوشه‌بندی سلسله مراتبی متراکم با مدل‌های گوینده‌ی مخلوط گوسی

خوشه‌بندی سلسله مراتبی روی یک ماتریس شباهت بین دسته‌ها اجرا شد و نیازمندیم معیار توقف که بتواند مجموعه بهینه‌ی خوشه‌ها را مشخص کند دارد. دو معیار توقف بررسی شدند- معیار آستانه‌ی فاصله^۱ و- معیار توقف بهینه^۲ T_s

۴-۳-۲- معیار آستانه‌ی فاصله

در طول خوشه‌بندی سلسله مراتبی همان‌طور که تعداد اجراها زیاد می‌شود، داده‌های بدون خوشه کم می‌شوند تا زمانی که حداقل DER خوشه‌بندی به دست آید. تعداد اجراها مجموعه خوشه‌های بهینه را دنبال می‌کنند. این امکان نیز وجود دارد که بیش از حد خوشه‌بندی هم رخ دهد.

۴-۳-۲-۱- معیار بهینگی T_s

با استفاده از معیار بهینگی [۲۵] که توسط (۲-۱۰) به دست می‌آید، خوشه‌ای با دورترین هیستوگرام برای فواصل درون خوشه‌ای و برون خوشه‌ای انتخاب می‌شود. در جدول ۴-۲ نتایج نشان داده شده است. با استفاده از فاصله‌ی NCLR یک DER ۲۲.۱۵ درصدی به دست آمد در صورتی که با

^۱ distance threshold criterion
^۲ T_s optimal stopping criterion

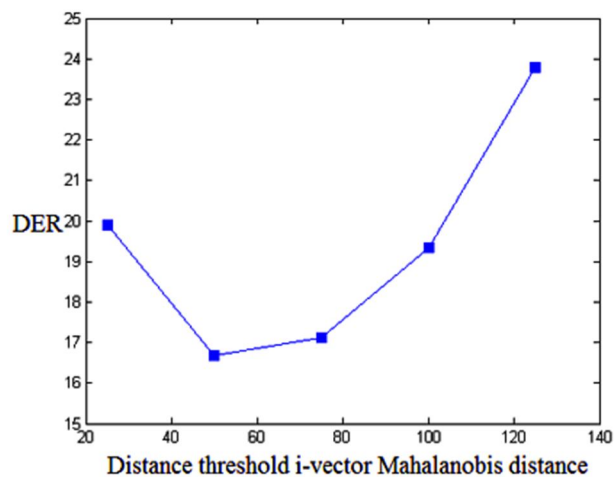
CLR یک DER ۱۹.۸۳ درصدی به دست آمد.

جدول ۲-۴: نتایج خطای بازشناسی بر اساس فواصل مختلف

فاصله مورد استفاده	DER (%)
CLR	۱۹.۸۳
NCLR	۲۲.۱۵

۳-۳-۴ HAC با مدل‌های گوینده‌ی بردار i

با استفاده از خوشه‌بندی سلسله‌مراتبی با بردارهای i ، بردارهای i جدیدی برای هر بخش به دست آمده در مرحله‌ی ادغام خوشه استخراج می‌شود. همان‌طور که در شکل (۲-۴) مشخص است کمترین میزان خطا در کار با فاصله‌ی مهالانوبیس به دست آمد که مقدار ۱۶.۶۹ درصد را به خود اختصاص می‌دهد.



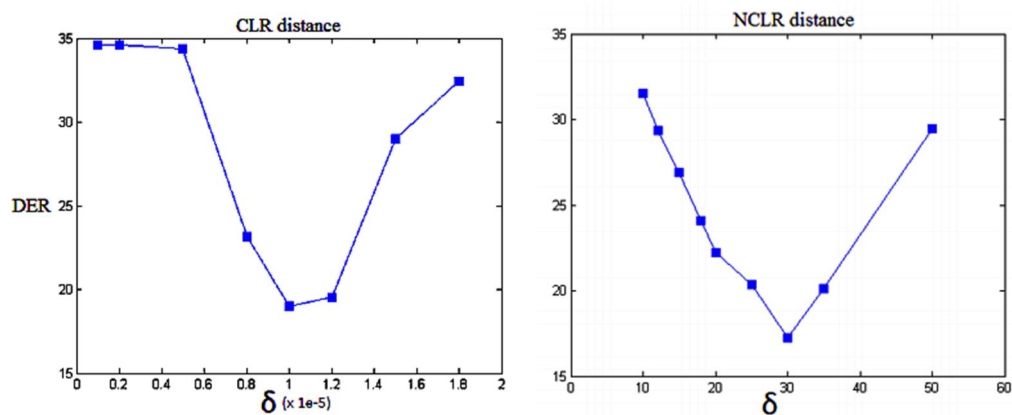
شکل ۲-۴: HAC با آستانه‌ی فاصله برای مدل‌های مخلوط گوسی گوینده

۴-۳-۴ ILP با آزمایشات با مدل‌های گوینده‌ی مدل مخلوط گوسی

خوشه‌بندی برنامه عدد صحیح خطی با استفاده از فواصل CLR و NCLR برای ساخت ماتریس

فاصله اجرا شد. با توجه به شکل (۳-۴) بهترین نتیجه به صورت ۱۹.۰۳ درصد بود در صورتی که با استفاده از فاصله‌ی NCLR بهترین نتیجه ۱۷.۲۷ درصد بود. محور x آستانه‌ی حاضر در قیود مسئله‌ی بهینه‌سازی برنامه عدد صحیح خطی را نشان می‌دهد (۴-۶). NCLR نمایش بهتری از فاصله نسبت به CLR است هرچند برای استفاده در خوشه‌بندی سلسله مراتبی متراکم مناسب نیست. این بدین دلیل است که هر قدر اندازه‌ی خوشه‌ی ادغام‌شده افزایش یابد اندازه‌ی قسمت‌ها نقش بسزایی در کاهش فاصله‌ی NCLR دارد.

از طرف دیگر فرمول سازی برنامه‌ی خطی عددی یک مسیر همبسته برای دستیابی به خوشه‌بندی بهینه را پیشنهاد می‌دهد. برای اثبات این، فرمول سازی برنامه عدد صحیح خطی برای ماتریس‌های شباهت NCLR و CLR تولیدشده توسط مدل‌های گوینده‌ی مدل مخلوط گوسی اجرا شد. در این پایان‌نامه برنامه عدد صحیح خطی تنها با استفاده از بردارهای i محاسبه شد.

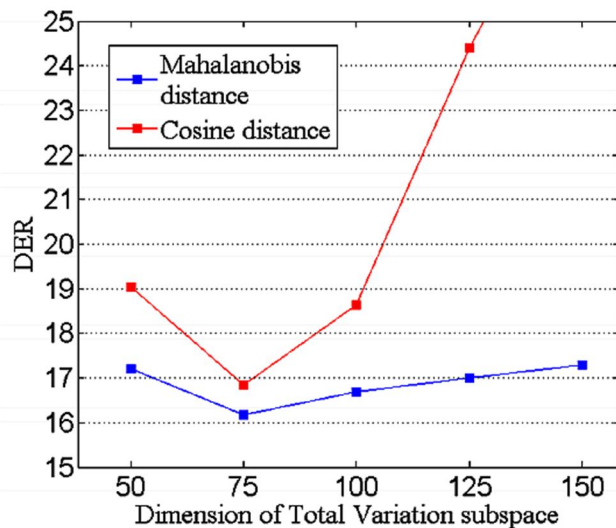


شکل ۳-۴: ILP با آستانه فاصله برای مدل‌های مخلوط گوسی گوینده

۴-۳-۵- خوشه‌بندی برنامه عدد صحیح خطی با مدل‌های گوینده‌ی بردار i

خوشه‌بندی برنامه عدد صحیح خطی با استفاده از بردارهای i آموزش دیده اجرا شده‌اند. با بررسی آزمایش‌های انجام‌شده به بهترین بُعد فضای تغییر کلی و بهترین انتخاب فاصله دست می‌یابیم. مقادیر به دست آمده در فصول قبلی مورد بحث و بررسی قرار گرفت. فاصله‌ی مهالانوبیس یک روش جبران

پس‌زمینه‌ای را پیشنهاد می‌دهد که شباهت بین بخش‌های با گوینده مشابه و پس‌زمینه متفاوت را بهبود می‌دهد.



شکل ۴-۴: کارکرد خوشه‌یابی برنامه عدد صحیح با مدل‌های گوینده بردار i با ابعاد متفاوت از زیر فضای تغییر.

۴-۳-۵-۱- مدل مخلوط گوسی در مقایسه با بردارهای i

مدل‌سازی گوینده‌ها بر اساس مدل مخلوط گوسی یک نمایش بعد بالا از قسمت‌ها را می‌دهد و نیز اطلاعات پس‌زمینه را به خوبی در برمی‌گیرد. شباهت در پس‌زمینه می‌تواند سبب شباهت بین بخش‌های گویندگان مختلف شود. طرح جبران پس‌زمینه نیازمند به کارگیری در فضای ویژگی است. از طرف دیگر بردارهای i به جبران پس‌زمینه از طریق $WCCN^1$ امکان می‌دهد. موضوع دیگر در مورد مدل‌های گوینده با مدل مخلوط گوسی، زمان جبران بالای آن‌ها برای شباهت قسمت‌ها به دلیل عبارت‌های احتمال متقابل در (۲-۶) و (۲-۷) است. در جدول (۴-۳) خطای کلی در انواع ترکیب‌های روش‌های خوشه‌بندی و روش‌های مدل‌سازی را نشان می‌دهد. مشاهده می‌شود که در ترکیب خوشه‌بندی ILP با روش مدل‌سازی بر اساس بردارهای i به کمترین میزان خطا دست‌یافته‌ایم.

¹ Within Class Covariance Normalization

جدول ۴-۳: نتایج خطای DER از دو مدل گوینده و دو الگوریتم خوشه‌بندی

	HAC خوشه‌بندی	خوشه‌بندی ILP
مدل مخلوط گوسی	٪۱۹.۴۵	٪۱۷.۲۷
بردار i	٪۱۷.۱۱	٪۱۶.۱۸

۴-۳-۵-۲ - HAC در مقابل ILP

خوشه‌بندی سلسله‌مراتبی با وجود آنکه یک الگوریتم افزایشی برای خوشه‌بندی است به‌عنوان یک تقریب خوب کار می‌کند. هرچند اگر طی یک مرحله در خوشه‌بندی، یک ادغام اشتباه رخ دهد، به‌طور مشخص کارایی مراحل بعدی را تحت تأثیر قرار می‌دهد. بنابراین یک تخمین مجدد از خوشه برای هر مرحله لازم است. بنابراین خوشه‌بندی سلسله‌مراتبی بیشتر از خوشه‌بندی برنامه عدد صحیح خطی هزینه دارد. خوشه‌بندی برنامه عدد صحیح خطی جستجوی دقیق‌تری نسبت به خوشه‌بندی سلسله‌مراتبی توسط بررسی کردن تمامی ترکیبات خوشه‌های ممکن انجام می‌دهد.

در این فصل به بیان نتایج آزمایشات و بررسی‌های به‌دست‌آمده بر اساس الگوریتم‌ها و روش‌های ارائه‌شده در فصل‌های قبل پرداختیم و به تحلیل و ذکر معایب و کمبودها و برتری‌های این نتایج پرداختیم. در ادامه آنچه اهمیت دارد تکمیل بحث با نتیجه‌گیری کلی از تمامی بحث‌های گفته‌شده است. در انتها ذکر راهکارهای پیش رو می‌تواند در ادامه‌ی روند این راه مؤثر باشد.

۴-۴ - مقایسه نتایج به‌دست‌آمده

با توجه به گستردگی روش‌های خوشه‌بندی و الگوریتم‌های دسته‌بندی و نیز انواع مختلف

ویژگی‌های مورد استفاده در بحث‌های پردازش گفتار و بازشناسی گوینده و نیز با در دست داشتن انواع پایگاه داده‌های مختلف به زبان‌های متنوع بشری این حوزه از علوم گسترده‌تری برخوردار است. این گسترده‌تری حیطة مقایسه و بررسی نتایج با دیگر کارهای انجام شده را محصور کرده و این مقایسات را بسته می‌کند. در منبع [۶۸] نویسنده با ترکیب ویژگی‌های فضایی با ویژگی‌های طیفی بر روی پایگاه داده‌ی AMI به همراه دیگر پایگاه‌های داده با بررسی گفتگوهای مختلف نتایج خوبی روی این پایگاه‌های داده به دست آوردند. آن‌ها با استفاده از این ویژگی‌ها و روش خوشه‌بندی HAC و با مدل‌سازی توسط HMM و GMM بر اساس فاصله‌ی BIC نتایج مطلوبی را در بازشناسی گوینده‌ها به دست آوردند. در منبع [۶۹] نویسنده با استفاده از ویژگی‌های زمانی بلندمدت برای تخمین حالات بعدی سیستم و با استفاده از ویژگی‌های آکوستیکی استخراج شده از فایل ورودی همانند MFCC و LP با مدل‌های GMM در زمینه‌ی بازشناسی گوینده بر روی پایگاه داده‌هایی چون AMI و NIST و ICSI نتایج درخور توجهی به دست آوردند. در جدول (۴-۴) در کنار روش پیشنهادی ما این دو روش نشان داده شده است. ذکر این نکته مهم است که در روش ما فرآیند و مدت زمان به دست آوردن نتیجه به دلیل برون خطی بودن تحلیل داده‌های در دست ما اهمیت چندانی نداشته است و ذکر نشده است هرچند در کاربردهای مختلف اهمیت زمانی مسئله و قابل انجام بودن آن ممکن است اهمیت خاصی داشته باشد.

جدول ۴-۴ : مقایسه با روش‌های دیگر

	روش پیشنهادی	منبع [۶۸]	منبع [۶۹]
درصد خطای DER	٪۱۶.۱۸	٪۱۸.۷	٪۱۲.۳

همان‌طور که در جدول (۴-۴) مشاهده می‌کنید در منبع [۶۸] به دلیل استفاده از ویژگی‌های زمانی بلندمدت ممکن است برخی تغییر در گوینده‌ها در این فواصل زمانی از دست برود و این خود باعث بالا رفتن میزان خطا می‌شود. در منبع [۶۹] به دلیل دربرداشتن ویژگی‌های آکوستیکی زمان کوتاه و

نیز ویژگی‌های بلندمدت کلی توانسته در مدل‌های گوینده به‌خوبی تمامی ویژگی‌ها را در برگیرد و مدل‌سازی قابل‌توجهی را اراده دهد که نتایج خوبی را به دست داده است. در مدل پیشنهادی ما با به‌کارگیری ویژگی‌هایی نظیر MFCC ویژگی‌های آکوستیکی پوشش داده‌شده است و نیز با استفاده از بردار i برای مدل‌سازی توانسته‌ایم کلیت صوت در دست را دربرگیریم.

در آخر قابل‌ذکر است که تحلیل‌های ما فهم شخصی و به‌تنهایی از دانسته‌های خودمان است و نمی‌توان گفت تحلیلی کلی و همه‌جانبه درزمینه‌ی بازشناسی گوینده است زیرا در هر کاربردی و در هر پایگاه داده‌ای ویژگی‌ها و نتایج با دید و رویکرد خاصی باید بررسی شوند.

فصل ۵

نتیجه‌گیری و پیشنهادات

۵-۱- نتیجه‌گیری

هدف از این پایان‌نامه مطالعه‌ی روش‌های به‌روز در بحث بازشناسی گوینده و بهبود سیستم مربوطه بود. سیستم پیشنهادی با استفاده از نرخ خطای بازشناسی متریک ارزیابی شد. سیستم سه بخش اساسی دارد: تشخیص فعالیت گفتار، تشخیص تغییر گوینده بر اساس ΔBIC و بلوک خوشه‌بندی به‌روز شده.

هدف کلی سیستم تشخیص فعالیت گفتار توانایی حذف سکوت‌های غیر گفتاری همانند موزیک از فایل صوتی است. بخش خوشه‌بندی گوینده به مدل‌های جدید گوینده اجازه می‌دهد تا قسمت‌ها را با بردارهای i نمایش دهند که این راه‌گشای کارهای بعدی در زمینه بازشناسی سریع می‌تواند باشد.

سیستم قادر به اجرای تشخیص فعالیت گفتار بدون نیاز به داده‌های آموزش خارجی برای مدل‌های گفتار و غیر گفتار است. انرژی فریم و نرخ عبور از صفر به‌عنوان ویژگی‌های قوی برای ساخت مدل‌های آهنگ و سکوت بکار گرفته شدند. یک تشخیص گر رقابتی گفتار با یک سیستم SAD دومرحله‌ای به دست آمد: شاخص سکوت و به دنبال آن شاخص آهنگ. نتایج با سیستم جدید بر پایه‌ی مدل مخلوط گوسی که با استفاده از داده‌ی خارجی آموزش داده‌شده بود مقایسه شد.

مدل‌های گوینده‌ی بردار i که اکنون جدیدترین روش در بازشناسی گوینده هستند یک نمایش بُعد پایین از اطلاعات گوینده در مقایسه با مدل‌های مدل مخلوط گوسی سنتی ارائه داده است. همچنین آن‌ها یک مزیت محاسباتی ارائه داده‌اند زیرا محاسبه‌ی فاصله‌ی بین بردارهای i سریع‌تر از محاسبه‌ی شباهت بر اساس احتمال متقابل در مدل‌های مدل مخلوط گوسی است؛ بنابراین برای سیستم‌های بازشناسی زمان واقعی بردارهای i به نظر جذاب می‌آیند.

در این پایان‌نامه همان‌طور که در [۲۸] به دست آمد، خوشه‌بندی گوینده با استفاده از یک‌رویه‌ی بهینه‌سازی کلی برای به دست آوردن مجموعه خوشه‌های بهینه، نسبت به روش بهینه‌سازی افزایشی

الگوریتم خوشه‌بندی متراکم سلسله مراتبی عملکرد بهتری را داراست. روش خوشه‌بندی سلسله مراتبی متراکم از نظر محاسباتی روش هزینه بری است و یک ادغام اشتباه طی دسته‌بندی باعث اثراتی در مراحل بعدی می‌شود. از طرف دیگر فرمولاسیون خوشه‌بندی برنامه‌ریزی خطی عددی یک مسیر همبسته برای خوشه‌بندی بهینه پیشنهاد می‌دهد. یک روش گرافیکی خوشه‌بندی برای مسئله‌ی رایج K مرکزی در بهینه‌سازی ترکیبی وفق داده شد. برای دستیابی به کارایی بهتر ILP در مقایسه با خوشه‌بندی سلسله مراتبی متراکم، فرمول‌سازی خوشه‌بندی برنامه عدد صحیح خطی برای ماتریس‌های شباهت NCLR و CLR تولیدشده توسط مدل‌های گوینده‌ی مدل مخلوط گوسی به دست آمد و اجرا شد و ۳.۲۷ درصد بهبود در خطاها را در مقایسه با بهترین خطای الگوریتم خوشه‌بندی مدل مخلوط گوسی -HAC حاصل شد. در متن، خوشه‌بندی برنامه عدد صحیح خطی تنها با استفاده از بردارهای i بررسی شده‌اند.

۵-۲- پیشنهادات

مسیر پیشرفت علم در تمامی حوزه‌ها همیشه بدون وقفه ادامه داشته است و محققان و پژوهشگران همواره در این راه نقش‌های اساسی و بسزایی را ایفا می‌کنند. در بحث بازشناسی گوینده در هر بخش و در هر الگوریتم و روش طرح‌های مبتکرانه و روش‌های جدید تا به امروز تحول و کارکردهای بهینه‌ی زیادی وارد کرده‌اند. در ادامه‌ی کار انجام‌شده در این پایان‌نامه می‌توان به ذکر اندکی از هزار راه‌کار موجود به صورت زیر پرداخت :

- اصلاح خروجی بازشناسی توسط عبور از یک رمزگشا ویتربی می‌تواند موثر باشد.

- بازشناسی متقابل عمل خوشه‌بندی گوینده‌ها در بین اسناد متفاوت است که برای تشخیص قسمت‌های با گوینده‌های یکسان مؤثر است. تلاش‌های حال حاضر در زمینه‌ی تحقیقات بازشناسی در راستای حل این مشکل است.

- بهبود در برنامه عدد صحیح خطی به طور مؤثر با کاهش بخش‌های زائد برنامه عدد صحیح خطی به اصلی اجرای سریع‌تری را دارد هرچند برنامه‌ی متلب^۱ حل چنین مسائل بهینه‌سازی پیچیده‌ای در بعد بالا را پشتیبانی نمی‌کند. حل‌کننده‌هایی چون GUROBI^۲ حل برنامه‌های خطی عددی را پشتیبانی می‌کنند.

-مشاهده شد که طی خوشه‌بندی گوینده، بخش‌هایی که دارای آهنگ پس‌زمینه بودند قادر به نمایش شباهت با بخش‌هایی که پس‌زمینه‌ی مشخص و تمیزی داشتند با استفاده از مدل‌های مدل مخلوط گوسی نبودند. حتی پس از استفاده از روش‌های جبران تغییر پس‌زمینه روی مدل‌های بردار i مسئله و مشکل همچنان پابرجاست. جداسازی گفتار و آهنگ پیش از پارامتر سازی باید مورد توجه قرار بگیرد تا آهنگ پس‌زمینه‌ی گفتار گوینده خنثی شود.

¹ Matlab

² <http://www.gurobi.com/>

فهرست منابع و مراجع

[۱] Martin, A., & Przybocki, M. (2000). The NIST 1999 speaker recognition evaluation—an overview. *Digital signal processing*, 10(1), 1-18.

[۲] Kenny, P., Boulianne, G., Ouellet, P., & Dumouchel, P. (2007). Joint factor analysis versus eigenchannels in speaker recognition. *IEEE Transactions on Audio, Speech, and Language Processing*, 15(4), 1435-1447

[۳] لشکر بلوکی م، (۱۳۹۰)، پایان‌نامه کارشناسی ارشد: "تشخیص گوینده در محیط شامل چند گوینده با استفاده از ماشین بردار پشتیبان"، دانشکده برق و رباتیک، دانشگاه صنعتی شاهرود.

[۴] L.Docio, C.Garcia, "Speaker Segmentation, detection and tracking in multi-speaker long audio recordings", Third COST275 Workshop Bimetrics on the internet. 2005

[۵] Meignier, S., & Merlin, T. (2010, March). LIUM SpkDiarization: an open source toolkit for diarization. In *CMU SPUD Workshop* (Vol. 2010).

[۶] Bozonnet, S., Evans, N. W., & Fredouille, C. (2010, March). The LIA-EURECOM RT'09 speaker diarization system: enhancements in speaker modelling and cluster purification. In *2010 IEEE International Conference on Acoustics, Speech and Signal Processing* (pp. 4958-4961). IEEE.

[۷] Nguyen, T., Sun, H., Zhao, S., Khine, S. Z. K., Tran, H. D., Ma, T. L. N ... & Li, H. (2009, May). The IIR-NTU speaker diarization systems for RT 2009. In *RT'09, NIST Rich Transcription Workshop, May 28-29, 2009, Melbourne, Florida, USA* (Vol. 14, pp. 17-40).

[۸] Arlindo Veiga, Carla Lopes, and Fernando Perdigão. Speaker diarization using gaussian mixture turns and segment matching. Proc. FALA, 2010

[۹] Maganti, H. K., Motlicek, P., & Gatica-Perez, D. (2007, April). Unsupervised speech/non-speech detection for automatic speech recognition in meeting rooms. In *2007 IEEE International Conference on Acoustics, Speech and Signal Processing-*

ICASSP'07 (Vol. 4, pp. IV-1037). IEEE.

- [10] Friedland, G., Janin, A., Imseng, D., Anguera, X., Gottlieb, L., Huijbregts, M., ... & Vinyals, O. (2012). The ICSI RT-09 speaker diarization system. *IEEE Transactions on Audio, Speech, and Language Processing*, 20(2), 371-381
- [11] Yella, S. H., Stolcke, A., & Slaney, M. (2014, December). Artificial neural network features for speaker diarization. In *Spoken Language Technology Workshop (SLT), 2014 IEEE* (pp. 402-406). IEEE
- [12] Ryant, N., Liberman, M., & Yuan, J. (2013, August). Speech activity detection on youtube using deep neural networks. In *INTERSPEECH* (pp. 728-731)
- [13] Anguera, X., & Bonastre, J. F. (2011, May). Fast speaker diarization based on binary keys. In *2011 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (pp. 4428-4431). IEEE.
- [14] Chen, S., & Gopalakrishnan, P. (1998, February). Speaker, environment and channel change detection and clustering via the bayesian information criterion. In *Proc. DARPA Broadcast News Transcription and Understanding Workshop* (Vol. 8, pp. 127-132)
- [15] Anguera, X., & Hernando, J. (2005). Xbic: Real-time cross probabilities measure for speaker segmentation. *Univ. California Berkeley, ICSIBerkeley Tech. Rep*
- [16] Cheng, S. S., Wang, H. M., & Fu, H. C. (2010). BIC-based speaker segmentation using divide-and-conquer strategies with application to speaker diarization. *IEEE transactions on audio, speech, and language processing*, 18(1), 141-157
- [17] Moraru, D., Meignier, S., Fredouille, C., Besacier, L., & Bonastre, J. F. (2004, May). The ELISA consortium approaches in broadcast news speaker segmentation during the NIST 2003 rich transcription evaluation. In *Acoustics, Speech, and Signal Processing, 2004. Proceedings. (ICASSP'04). IEEE International Conference on* (Vol. 1, pp. I-373). IEEE

- [18] Anguera, X., Bozonnet, S., Evans, N., Fredouille, C., Friedland, G., & Vinyals, O. (2012). Speaker diarization: A review of recent research. *IEEE Transactions on Audio, Speech, and Language Processing*, 20(2), 356-370.
- [19] Reynolds, D. A., & Torres-Carrasquillo, P. (2004). *The MIT Lincoln Laboratory RT-04F diarization systems: Applications to broadcast audio and telephone conversations*. MASSACHUSETTS INST OF TECH LEXINGTON LINCOLN LAB
- [20] Panagiotakis, C., & Tziritas, G. (2005). A speech/music discriminator based on RMS and zero-crossings. *IEEE Transactions on multimedia*, 7(1), 155-166.
- [21] Reynolds, D. A., Quatieri, T. F., & Dunn, R. B. (2000). Speaker verification using adapted Gaussian mixture models. *Digital signal processing*, 10(1), 19-41
- [22] Jensen, J. H., Ellis, D. P., Christensen, M. G., & Jensen, S. H. (2007, October). Evaluation of distance measures between Gaussian mixture models of MFCCs. In *International Conference on Music Information Retrieval* (pp. 107-108)
- [23] Dehak, N., Kenny, P. J., Dehak, R., Dumouchel, P., & Ouellet, P. (2011). Front-end factor analysis for speaker verification. *IEEE Transactions on Audio, Speech, and Language Processing*, 19(4), 788-798.
- [24] Seyed Omid Sadjadi, Malcolm Slaney, and Larry Heck, "MSR Identity Toolbox v1.0: A MATLAB Toolbox for Speaker Recognition Research," *Speech and Language Processing Technical Committee Newsletter*, IEEE, November 2013
- [25] Nguyen, Trung Hieu, Engsiong Chng, and Haizhou Li. "T-test distance and clustering criterion for speaker diarization." *INTERSPEECH*. 2008
- [26] Vijayasenan, D., Valente, F., & Boulard, H. (2007, December). Agglomerative information bottleneck for speaker diarization of meetings data. In *Automatic Speech Recognition & Understanding, 2007. ASRU. IEEE Workshop on* (pp. 250-255). IEEE

- [27] Meignier, S., Moraru, D., Fredouille, C., Bonastre, J. F., & Besacier, L. (2006). Step-by-step and integrated approaches in broadcast news speaker diarization. *Computer Speech & Language*, 20(2), 303-330.
- [28] Rouvier, M., & Meignier, S. (2012, June). A global optimization framework for speaker diarization. In *Odyssey* (pp. 146-150)
- [29] Dupuy, G., Meignier, S., Deléglise, P., & Esteve, Y. (2014). Recent improvements on ilp-based clustering for broadcast news speaker diarization. In *Proceedings of Odyssey*.
- [30] Anguera Miró, X. (2006). ROBUST SPEAKER DIARIZATION FOR MEETINGS
- [31] Sun, H., Ma, B., Khine, S. Z. K., & Li, H. (2010, March). Speaker diarization system for RT07 and RT09 meeting room audio. In *ICASSP* (pp. 4982-4985).
- [32] <http://groups.inf.ed.ac.uk/ami/corpus>
- [33] <https://catalog.ldc.upenn.edu/LDC93S1>
- [34] Saunders, J. (1996, May). Real-time discrimination of broadcast speech/music. In *icassp* (Vol. 96, pp. 993-996)
- [35] Liu, Z., Wang, Y., & Chen, T. (1998). Audio feature extraction and analysis for scene segmentation and classification. *Journal of VLSI signal processing systems for signal, image and video technology*, 20(1-2), 61-79
- [36] Huijbregts, M. A. H. (2008). Segmentation, diarization and speech transcription: surprise data unraveled.
- [37] Dalal, N. INRIA Object Detection and Localization Toolkit (2008). *Software available at <http://pascal.inrialpes.fr/soft/olt>*
- [38] Rouvier, M., Dupuy, G., Gay, P., Houry, E., Merlin, T., & Meignier, S. (2013). *An*

open-source state-of-the-art toolbox for broadcast news diarization (No. EPFL-REPORT-192605). Idiap

- [३१] Vijayasenan, D., & Valente, F. (2012). DiarTk: An Open Source Toolkit for Research in Multistream Speaker Diarization and its Application to Meetings Recordings. In *INTERSPEECH* (pp. 2170-2173)
- [ॣ०] Xavier, A. M. (2006). *Robust speaker diarization for meetings* (Doctoral dissertation, PhD Thesis, Speech Processing Group Department of Signal Theory and Communications Universitat Politecnica de Catalunya Barcelona (Espagnol)).
- [ॣ१] Luque, J., Anguera, X., Temko, A., & Hernando, J. (2008). Speaker diarization for conference room: The UPC RT07s evaluation system. In *Multimodal Technologies for Perception of Humans* (pp. 543-553). Springer Berlin Heidelberg.
- [ॣॢ] Friedland, G., Gottlieb, L., & Janin, A. (2009, October). Joke-o-mat: browsing sitcoms punchline by punchline. In *Proceedings of the 17th ACM international conference on Multimedia* (pp. 1115-1116). ACM.
- [ॣॣ] Delgado, H., Fredouille, C., & Serrano, J. (2014). Towards a complete binary key system for the speaker diarization task. In *INTERSPEECH* (pp. 572-576).
- [ॣ।] Tranter, S. E., & Reynolds, D. A. (2006). An overview of automatic speaker diarization systems. *IEEE Transactions on Audio, Speech, and Language Processing*, 14(5), 1557-1565.
- [ॣΔ] Kahn, J., Galibert, O., Quintard, L., Carré, M., Giraudel, A., & Joly, P. (2012, June). A presentation of the REPERE challenge. In *Content-Based Multimedia Indexing (CBMI), 2012 10th International Workshop on* (pp. 1-6). IEEE.
- [ॣϢ] Chan, W. N., Lee, T., Zheng, N., & Ouyang, H. (2006, May). Use of vocal source features in speaker segmentation. In *2006 IEEE International Conference on Acoustics Speech and Signal Processing Proceedings* (Vol. 1, pp. I-I). IEEE.

- [٤٧] Siegler, M. A., Jain, U., Raj, B., & Stern, R. M. (1997, February). Automatic segmentation, classification and clustering of broadcast news audio. In *Proc. DARPA speech recognition workshop* (Vol. 1997).
- [٤٨] Bredin, H., & Poignant, J. (2013). Integer linear programming for speaker diarization and cross-modal identification in tv broadcast. In *the 14rd Annual Conference of the International Speech Communication Association, INTERSPEECH*.
- [٤٩] Zhu, X., Barras, C., Meignier, S., & Gauvain, J. L. (2005, September). Combining speaker identification and BIC for speaker diarization. In *INTERSPEECH* (Vol. 5, pp. 2441-2444).
- [٥٠] Favre, B., Damnati, G., Bechet, F., Bendris, M., Charlet, D., Auguste, R., ... & Fredouille, C. (2013, August). PERCOLI: A Person Identification System for the 2013 REPERE Challenge. In *SLAM@ INTERSPEECH* (pp. 55-60).
- [٥١] Shum, S., Dehak, N., Chuangsuwanich, E., Reynolds, D. A., & Glass, J. R. (2011, August). Exploiting Intra-Conversation Variability for Speaker Diarization. In *INTERSPEECH* (Vol. 11, pp. 945-948).
- [٥٢] Kenny, P., Reynolds, D., & Castaldo, F. (2010). Diarization of telephone conversations using factor analysis. *IEEE Journal of Selected Topics in Signal Processing*, 4(6), 1059-1070.
- [٥٣] Silovsky, J., & Prazak, J. (2012, March). Speaker diarization of broadcast streams using two-stage clustering based on i-vectors and cosine distance scoring. In *2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (pp. 4193-4196). IEEE.
- [٥٤] Johnson, S. E., & Woodland, P. C. (2000). A method for direct audio search with applications to indexing and retrieval. In *Acoustics, Speech, and Signal Processing, 2000. ICASSP'00. Proceedings. 2000 IEEE International Conference on* (Vol. 3, pp.

1427-1430). IEEE.

- [ΔΔ] Kubala, D. L. F. (1999). Fast speaker change detection for broadcast news transcription and indexing.
- [Δϵ] Wu, T. Y., Lu, L., Chen, K., & Zhang, H. (2003, January). Universal Background Models for Real-time Speaker Change Detection. In *MMM* (pp. 135-149).
- [ΔΥ] Povey, D., Chu, S. M., & Varadarajan, B. (2008, March). Universal background model based speech recognition. In *2008 IEEE International Conference on Acoustics, Speech and Signal Processing* (pp. 4561-4564). IEEE.
- [ΔΛ] Reynolds, D. A., Quatieri, T. F., & Dunn, R. B. (2000). Speaker verification using adapted Gaussian mixture models. *Digital signal processing, 10*(1), 19-41.
- [ΔϠ] Majetniak, A., & Tan, Z. H. (2011). Speaker recognition using universal background model on YOHO database.
- [ϵ•] Shum, S. (2011). *Unsupervised methods for speaker diarization* (Doctoral dissertation, Massachusetts Institute of Technology).
- [ϵ∧] Senoussaoui, M., Kenny, P., Dehak, N., & Dumouchel, P. (2010, June). An i-vector Extractor Suitable for Speaker Recognition with both Microphone and Telephone Speech. In *Odyssey* (p. 6).
- [ϵΥ] Povey, D., Burget, L., Agarwal, M., Akyazi, P., Feng, K., Ghoshal, A., ... & Schwarz, P. (2010, March). Subspace Gaussian mixture models for speech recognition. In *2010 IEEE International Conference on Acoustics, Speech and Signal Processing* (pp. 4330-4333). IEEE.
- [ϵϣ] Shum, S. H., Dehak, N., Dehak, R., & Glass, J. R. (2013). Unsupervised methods for speaker diarization: An integrated and iterative approach. *IEEE Transactions on Audio, Speech, and Language Processing, 21*(10), 2015-2028.

- [६ॢ] Bachu, R. G., Kopparthi, S., Adapa, B., & Barkana, B. D. (2008). Separation of voiced and unvoiced using zero crossing rate and energy of the speech signal. In *American Society for Engineering Education (ASEE) Zone Conference Proceedings* (pp. 1-7).
- [६Δ] Geiger, J. T., Wallhoff, F., & Rigoll, G. (2010, September). GMM-UBM based open-set online speaker diarization. In *INTERSPEECH* (pp. 2330-2333).
- [६Ϣ] Siohan, O., & Bacchiani, M. (2013). Ivector-based acoustic data selection. In *INTERSPEECH* (pp. 657-661).
- [६Υ] Lei, H. (2011). Joint Factor Analysis (JFA) and i-vector Tutorial. *ICSI. Web. 02 Oct.*
- [६Ϡ] Zelenak, M., Segura, C., Luque, J., & Hernando, J. (2012). Simultaneous speech detection with spatial features for speaker diarization. *IEEE Transactions on Audio, Speech, and Language Processing*, 20(2), 436-446.
- [६ϡ] Yella, S. H., & Boulard, H. (2014). Overlapping speech detection using long-term conversational features for speaker diarization in meeting room conversations. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 22(12), 1688-1700.

Abstract:

Nowadays, in the world of data processing with high speed and accuracy is very important. Processing speech Data because of the widespread use will contribute in all aspects of human life. Speaker diarization is to recognize who speaks when. The goal is to design a speaker recognition system that identify the speaker change in the audio files and correctly labeled and cluster each speaker's speech. This process is known speaker diarization named today. In this context, the aim is design a system that uses acoustic features MFCC and its first and second order along with energy and zero-crossing rate. Then model silence and music using frame that are absolutely non-speech or music and in tow-steps -silence removal and music removal- seprate speech from audio file. Therefore by using i-vectors in Feature space, reduce dimentional of system. We use integer linear programming (ILP) to label and cluster speaker and modify parameters in this models. Designed system, was examined on AMI corpus. We achieved good result and show in case of DER¹ error.

Keywords: speaker diarization, integer linear programming, i-vector, speech processing

¹ Diarization Error Rate



Shahrood University of Technology

Faculty of Electrical and Robotics Engineering

M.Sc. Thesis in Communication Systems Engineering

**Speaker diarization based on mixed feature extraction in
multi -speaker environment**

By: Mitra Jahanian

Supervisor:

Dr. Hossein Marvi

Advisor:

Dr. Seyyed Masoud Mirrezaei

September 2016