

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشکده برق و رباتیک
گروه الکترونیک
پایان نامه کارشناسی ارشد

تشخیص گوینده در محیط شامل چند گوینده با استفاده از ماشین بردار پشتیبان

دانشجو: مرضیه لشکر بلوکی

استاد راهنما:

دکتر حسین مروی

استاد مشاور:

دکتر حسین صامتی

تیر ماه 1390

تقدیم به:

خانواده عزیزم که همواره نقش اساسی در موفقیت هایم داشته اند.

دوستان عزیزی که همواره همراه و مشوقم بوده اند.

تقدیر و تشکر

سپاس و تشکر از زحمات و راهنمایی های علمی استاد گرانقدر جناب آقای دکتر حسین مروی (استاد راهنمای اینجانب) و جناب آقای دکتر حسین صامتی (استاد مشاور اینجانب). برای ایشان همواره آرزوی سلامتی و موفقیت روز افزون دارم.

چکیده:

شناسایی گوینده یکی از مباحث مطرح در بحث پردازش گفتار می باشد. شناسایی گوینده عبارت است از فرآیندی که طی آن با استفاده از سیگنال صحبت تشخیص دهیم چه کسی چه موقع واقعا صحبت می کند. هدف طراحی سیستمی است که بتواند تغییر در گوینده را مشخص نماید و گفتار هر گوینده را برای سیستم برچسب گذاری نماید. یعنی مشخص نماید که کدام گوینده، در چه بازه هایی صحبت کرده است. امروزه این عمل با یک عنوان جدید که هر دو فرآیند جداسازی و برچسب گذاری را در بر می گیرد بنام Speaker Diarization مشهور گشته است. هدف از بخش بندی تقسیم سیگنال گفتاری به بخش هایی است که تنها شامل گفتار یک گوینده هستند و هدف از خوشه بندی نیز شناسایی بخش های گفتاری مربوط به یک گوینده و اختصاص یک برچسب واحد به آنهاست.

هدف از انجام این پایان نامه طراحی و پیاده سازی یک سیستم بخش بندی و خوشه بندی گوینده با استفاده از الگوریتم های جدید و همچنین بهبود نتایج این الگوریتم ها برای این موضوع می باشد. این سیستم باید بطور صحیح نقاط تغییر گوینده را بدون دانستن اطلاعات قبلی از گوینده تشخیص داده و در نهایت تمام قسمت های صوتی مربوط به یک گوینده را در یک خوشه قرار می دهد.

در این پایان نامه، سیستم تشخیص گوینده، از سه مرحله اصلی تشکیل شده است. در مرحله اول قسمت - های غیر گفتاری، از بخش های گفتاری فایل صوتی حذف می شوند، تا دقت و سرعت عملیات سیستم در مراحل بعدی افزایش پیدا کند. سپس فایل گفتاری به بخش هایی همگن که در آن فقط گفتار یک گوینده وجود دارد، تقسیم می شود. در مرحله سوم با استفاده از خوشه بندی مناسب، بخش های گفتاری مرحله قبل، که متعلق به یک گوینده هستند، در یک خوشه جای می گیرند. جهت پیاده سازی سیستم از چهار نوع بردار ویژگی MFCC, root-MFCC, TDC و root-TDC و سه نوع پایگاه داده استفاده شده است و دقت مرحله بخش بندی 80٪ بوده است و دقت مرحله خوشه بندی نیز 59٪ با استفاده از ماشین بردار پشتیبان بدست آمده است.

کلمات کلیدی:

بخش بندی آماری گوینده

بخش بندی گویندگان

تشخیص بخش های صوتی

خوشه بندی گویندگان

فهرست مطالب

فصل اول: معرفی سیستم های تشخیص گوینده

- 1-1-1- مقدمه..... 2
- 2-1-1- مراحل مختلف کاری سیستم های تشخیص گوینده..... 6
- 1-2-1-1- قطعه بند آکوستیکی..... 7
- 2-2-1-1- تشخیص گفتار از غیر گفتار..... 8
- 3-2-1-1- تشخیص جنسیت گوینده..... 9
- 4-2-1-1- تشخیص تغییر گوینده..... 9
- 3-1-3-1- روش های بخش بندی و خوشه بندی گویندگان..... 10
- 1-3-1-1- روش های بر اساس فاصله..... 10
- 2-3-1-1- روش های بر اساس مدل..... 11
- 3-3-1-1- روش های هیبرید یا ترکیبی..... 11
- 4-1-4-1- خوشه بندی نمودن..... 11
- 5-1-5-1- خلاصه..... 12

فصل دوم: تشخیص گفتار از نواحی غیر گفتاری

- 1-2-1-2- مقدمه..... 14
- 2-2-2-2- ساختار قسمت تشخیص گفتار از غیر گفتار..... 16
- 1-2-2-2-1- پیش پردازش..... 16
- 2-2-2-2-2- استخراج ویژگی..... 17
- 1-2-2-2-2- انرژی..... 18
- 2-2-2-2-2- نرخ عبور از صفر..... 19
- 3-2-2-2-2- استخراج ویژگی به کمک ضرایب کپسترال فرکانسی در مقیاس مل..... 19
- 4-2-2-2-2- ضرایب LPC..... 23
- 5-2-2-2-2- آنتروپی..... 24
- 6-2-2-2-2- اندازه متناوب بودن..... 26
- 7-2-2-2-2- اطلاعات زیر باند..... 28
- 8-2-2-2-2- سایر پارامترها..... 28

29	3-2-2- محاسبه آستانه.....
29	4-2-2- تصمیمات VAD.....
30	1-4-2-2- تصمیم گیری مبتنی بر مدل مخفی مارکوف.....
31	2-4-2-2- تصمیم گیری مبتنی بر شبکه های عصبی.....
33	5-2-2- تصحیح نتایج VAD.....
33	3-2- بلوک دیاگرام چند VAD استاندارد.....
33	1-3-2- استاندارد ETSI AMR.....
34	2-3-2- الگوریتم GSM.....
35	4-2- خلاصه.....
فصل سوم: آشکارسازی تغییر گوینده	
37	1-3- مقدمه.....
38	2-3- بخش بندی گوینده.....
38	1-2-3- بخش بندی بر اساس فاصله.....
40	2-2-3- بخش بندی بر اساس مدل.....
40	3-2-3- بخش بندی هیبرید.....
40	3-3- مقایسه روش های بخش بندی.....
41	4-3- روش های متداول آشکارسازی گوینده.....
41	1-4-3- معیار اطلاعات بیزین (BIC).....
42	2-1-4-3- بخش بندی با استفاده از مدل آماری گوینده.....
45	2-4-3- ترکیب آماره T^2 و BIC.....
47	1-2-4-3- سرعت و بهره بیشتر در بخش بندی T^2 -BIC.....
49	3-4-3- فاصله نرخ درستنمایی عمومی (GLR).....
49	4-4-3- فاصله KL2.....
51	5-4-3- آشکارسازی تغییر گوینده با استفاده از DSD.....
52	6-4-3- BIC متقاطع (Cross-BIC (XBIC).....
53	7-4-3- درستنمایی مدل مخلوط گوسی (GMM-L).....
53	5-3- خلاصه.....

فصل چهارم: روش های دسته بندی

- 55.....1-4-مقدمه.....
- 56.....2-4-اجزا سیستم خوشه بندی.....
- 57.....3-4-روش های خوشه بندی.....
- 58.....1-3-4-روش های خوشه بندی سلسله مراتبی.....
- 59.....1-1-3-4-تکنیک های خوشه بندی بالارونده.....
- 60.....2-1-3-4-تکنیک های خوشه بندی پایین رونده.....
- 61.....2-3-4-روش های خوشه بندی افزایی.....
- 61.....4-4-روش های خوشه بندی متداول در سیستم های خوشه بندی گوینده.....
- 63.....5-4-دسته بندی کننده ماشین های بردار پشتیبان.....
- 63.....1-5-4-دسته بندی کننده ماشین بردار پشتیبان خطی.....
- 63.....1-1-5-4-دسته بندی کلاس های جداپذیر.....
- 68.....2-1-5-4-دسته بندی کلاس های جدا ناپذیر.....
- 71.....3-1-5-4-دسته بندی داده های چند کلاسه با ماشین های بردار پشتیبان.....
- 72.....2-5-4-ماشین های بردار پشتیبان غیر خطی.....
- 74.....6-4-خلاصه.....

فصل پنجم: پیاده سازی و مشاهدات سیستم ترکیبی پیشنهادی

- 76.....1-5-مقدمه.....
- 77.....2-5-ساختار سیستم پیاده سازی شده.....
- 80.....3-5-پایگاه داده.....
- 82.....4-5-استخراج ویژگی.....
- 84.....5-5-معیار ارزیابی سیستم های تشخیص گوینده.....
- 88.....6-5-نتایج آزمایشات.....
- 88.....1-6-5-اثر اعمال VAD بر روی سیگنال گفتار.....
- 89.....2-6-5-اثر تغییر طول پنجره VAD بر روی دقت سیستم.....
- 89.....3-6-5-اثر تغییر طول پنجره BIC بر روی نتایج بخش بندی.....
- 93.....4-6-5-دقت حاصل از بخش بندی بر دو نوع از دادگان با استفاده از MFCC.....

- 5-6-5- اثر تغییر بردار، ویژگی، بر، روی، دقت، مرحله، بخش بندی..... 93
- 5-6-6- مقایسه، نتایج، مرحله، بخش بندی، با، بکارگیری، بردارهای، ویژگی متفاوت..... 95
- 5-6-7- اثر جنسیت، گویندگان، بر تشخیص، درست، مرزهای، بخش بندی..... 96
- 5-6-8- دقت مرحله خوشه بندی بکارگیری ماشین بردار پشتیبان (SVM) با بردار ویژگی MFCC..... 96
- 5-6-9- دقت مرحله خوشه بندی ماشین بردار پشتیبان با بکارگیری بردار ویژگی root-MFCC..... 97
- 5-6-10- اثر تغییر نوع تابع کرنل ماشین بردار پشتیبان بر روی دقت مرحله خوشه بندی..... 98
- 5-7- خلاصه..... 98

فصل ششم: جمع بندی و پیشنهادات

- 6-1- جمع بندی و خلاصه نتایج..... 100
- 6-2- پیشنهادات..... 101
- منابع..... 103

فهرست شکل ها

- شکل (1-1): نمایش بخش بندی و خوشه بندی گویندگان روی گفتار ورودی..... 4
- شکل (2-1): ساختار کلی سیستم های بخش بندی و خوشه بندی گوینده..... 6
- شکل (1-2): دیاگرام یک VAD ساده..... 16
- شکل (2-2): نمایش پنجره همینگ 512 نقطه ای در حوزه زمان..... 16
- شکل (3-2): شمای کلی سیستم استخراج ویژگی..... 18
- شکل (4-2): مراحل استخراج ویژگی با روش MFCC..... 20
- شکل (5-2): اعمال بانک فیلتر Mel scaled و محاسبه انرژی در هر زیر بانده..... 22
- شکل (6-2): شبکه ای از HMM ها جهت بررسی دنباله احتمالی گفتار و سکوت..... 31
- شکل (7-2): دیاگرام ساده ای از یک VAD مبتنی بر شبکه های عصبی..... 32
- شکل (8-2): دیاگرام ساده ای از الگوریتم AMR2..... 34
- شکل (9-2): دیاگرام الگوریتم GSM..... 35
- شکل (1-3): پنجره های همسایه..... 38
- شکل (2-3): ترکیب گوسین برای یک سیگنال شامل سکوت/گفتار..... 39
- شکل (3-3): منحنی ها با اعمال متریک T^2 -statistic..... 46
- شکل (1-4): انواع دسته بندی..... 55
- شکل (2-4): مراحل خوشه بندی..... 56
- شکل (3-4): روش های خوشه بندی..... 57
- شکل (4-4): روشهای خوشه بندی بالا و پایین رونده..... 58
- شکل (5-4): مثال ساده ای از خوشه بندی سلسله مراتبی..... 60
- شکل (6-4): یک نمونه از مسئله دو کلاسه خطی جداپذیر که نمونه ها توسط دو دسته بندی کننده خطی جدا شده..... 64
- شکل (7-4): حاشیه برای جهت 2 بیشتر از حاشیه در جهت 1 است..... 65
- شکل (8-4): نمونه ای از داده هایی که به صورت خطی به طور کامل از همدیگر جدا نمی شوند..... 68

- شکل (4-9): نمایش ماشین بردار پشتیبان غیر خطی..... 74.....
- شکل (5-1): بلوک دیاگرام سیستم پیاده سازی شده..... 76.....
- شکل (5-2): انتقال اطلاعات گفتار با استفاده از یک VAD..... 77.....
- شکل (5-3): دیاگرام الگوریتم G.729B..... 79.....
- شکل (5-4): بلوک دیاگرام بردار ویژگی TDC..... 83.....
- شکل (5-5): تشخیص خطا در سیستم های تشخیص گوینده..... 87.....
- شکل (5-6): جداسازی قسمت های گفتاری از غیر گفتار..... 88.....
- شکل (5-7): اثر تغییر طول پنجره VAD بر روی دقت سیستم..... 89.....
- شکل (5-8): چگونگی قرار دادن یک آستانه و بعد انتخاب نقاط تغییر گوینده را نمایش میدهد..... 90.....
- شکل (5-9): سیگنال گفتاری گوسی مدل شده در مرحله بخش بندی..... 90.....
- شکل (5-10): اثر افزایش طول پنجره BIC بر روی نتیجه مرحله بخش بندی برای 8 نفر دادگان فارس دات..... 91.....
- شکل (5-11): اثر افزایش طول پنجره BIC بر روی نتیجه مرحله بخش بندی برای 12 نفر دادگان فارس دات..... 92.....
- شکل (5-12): اثر افزایش طول پنجره BIC بر روی نتیجه مرحله بخش بندی برای 18 نفر دادگان فارس دات..... 92.....
- شکل (5-13): مقایسه میزان خطای سیستم با تغییر بردار ویژگی مورد استفاده..... 95.....
- شکل (5-14): تاثیر جنسیت بر روی خروجی مرحله بخش بندی سیستم..... 96.....
- شکل (5-15): مقایسه نتایج خطای حاصل از خوشه بندی با تغییر نوع تابع کرنل بکار گرفته شده..... 98.....

فهرست جداول

- جدول (5-1): مقادیر خطا برای دادگان تهیه شده فارسی آزمایشگاهی..... 93
- جدول (5-2): مقادیر خطا برای دادگان AMI..... 93
- جدول (5-3): مقادیر خطا برای تعداد 3 نفر گوینده در دادگان فارس دات..... 93
- جدول (5-4): مقادیر خطا برای تعداد 5 نفر گوینده در دادگان فارس دات..... 94
- جدول (5-5): مقادیر خطا برای تعداد 8 نفر گوینده در دادگان فارس دات..... 94
- جدول (5-6): مقادیر خطا برای تعداد 11 نفر گوینده در دادگان فارس دات..... 94
- جدول (5-7): مقادیر خطا برای تعداد 14 نفر گوینده در دادگان فارس دات..... 94
- جدول (5-8): مقادیر خطا برای تعداد 17 نفر گوینده در دادگان فارس دات..... 94
- جدول (5-9): مقادیر خطا برای تعداد 20 نفر گوینده در دادگان فارس دات..... 95
- جدول (5-10): خطای حاصل از دسته‌بندی با استفاده از ماشین بردار پشتیبان با بکارگیری MFCC..... 97
- جدول (5-11): خطای حاصل از دسته‌بندی با استفاده از ماشین بردار پشتیبان با بکارگیری root-MFCC..... 97

فصل اول :

معرفی سیستم های

تشخیص گوینده

1-1-مقدمه

امروزه داده های چند رسانه ای بخش قابل توجهی از دانش انسان را در بر می گیرند. حجم پرونده های چند رسانه ای آرشیو شده در موسسه های مختلف در سال های اخیر افزایش چشمگیری داشته است. دسترسی و وضوح بالای این پرونده ها می تواند کمک شایانی به افرادی کند که در جستجوی اطلاعات باشند. بنابراین عملیات جستجو و بازیابی اطلاعات در این حجم بالا کاری است که خود احتیاج به سیستم کامپیوتری دارد. و در نتیجه یکی از حوزه های تحقیقاتی که به تازگی مورد توجه قرار گرفته است، مربوط به ساختار بندی پرونده - های چند رسانه ای است. در میان این داده ها، اطلاعات صوتی اهمیت بالاتری دارد. زیرا بخش اعظم آرشیوها حاوی داده های صوتی از گزارش های تلویزیونی، رادیویی و همچنین مکالمات تلفنی می باشد. در سالهای اخیر تحقیقات وسیعی در این حوزه آغاز شده و نتایج قابل قبولی نیز حاصل شده است. از دیگر کاربردهای این حوزه در تشخیص مجرم، جدا کردن صحبت های مهم یک شاهد یا متهم در دادگاه و ... میتوان اشاره نمود.

در کاربرد صوتی، عمده اطلاعات موجود در پرونده ها، صحبت های تعدادی گوینده است و هدف از سیستم نهایی، پاسخ به این سوال است که چه کسی در چه زمانهایی صحبت کرده است؟ بخش های مختلف این حوزه تحقیقاتی به نامهای مختلفی مانند: قطعه بند گوینده ای¹، تشخیص گوینده²، رونویسی قوی³، و اندیس گذاری گوینده ای⁴ نامیده شده اند. از چنین سیستم هایی برای جابجایی راحت در داده های صوتی، در فایل های صوتی طولانی (مانند: اخبار و ملاقات ها و جلسات یک شرکت و ...) که متعلق به چند گوینده باشند بهره - برداری می شود. مکالمات و محاسبات رادیویی طولانی از محیط هایی هستند که در آنها چند گوینده حضور داشته و با هم صحبت می کنند. هدف نهایی چنین سیستم هایی، پیاده سازی روش هایی مناسب برای افراز پرونده صوتی به نواحی است که در آنها گوینده ای خاص صحبت کرده باشد. دسترسی راحت به

1.Speaker Segmentation
2.Speaker Diarization
3.Rich Transcription
4.Speaker Indexing

بخش هایی از صحبت یک گوینده توسط این سیستم فراهم می گردد. با داشتن حجم بالایی از داده های صوتی اهمیت این سیستم ها بیشتر می گردد.

با افزایش تعداد مدارک متنی موجود در اینترنت، نیاز به تکنیک هایی نظیر فهرست نگاری متن به منظور تسهیل دسترسی و جستجو در این مدارک افزایش پیدا کرد. نظیر همین نیاز نیز با افزایش تعداد مدارک صوتی نظیر سخنرانی ها، مصاحبه ها و گردهمایی ها و ... ایجاد شد. بطور مشخص دسترسی به مدارک صوتی بسیار سخت تر از دسترسی به متن است و گوش دادن به یک فایل صوتی ضبط شده بیشتر از خواندن متن زمان بر است و فهرست نگاری دستی مدارک صوتی در مقایسه با فهرست نگاری متن، مشکل است. راه حل پیشنهادی جهت رفع این مشکل، فهرست نگاری خودکار مدارک صوتی⁵ است.

اولین بار سیستم هایی تشخیص گوینده توسط کمپانی NIST در سال 1999 ارایه شد. در سال 2001، پلکان و سیدهارون به همراه گروهشان با استفاده از کم کردن اثر نویز بر روی سیگنال بهبودهایی در نتایج سیستم دادند و جداسازی بهتر گویندگان را باعث شدند. در سال 2005، بولیان و کنی با بکارگیری بردارهای ویژگی دیگر (یا ادغام روش های قبلی) و استفاده از مدل های گوسی در سیستم نتایج متفاوتی بدست آوردند. در سال 2005 توسط یاماشیتا و ماتسونوگا با استفاده از ویژگی های سیگنال صوتی مانند فرکانس پیچ سیگنال، انرژی، فرکانس های ماکزیمم سیگنال، و سه ویژگی دیگر نتایج در قسمت بخش بندی گوینده این سیستم بهبود داده شد.[1] و در سال های بعدی با انجام روش های مختلف بر روی قسمت های متفاوت آن تا به امروز این سیستم ها در حال تکمیل شدن و بهتر شدن نتایج بوده اند.

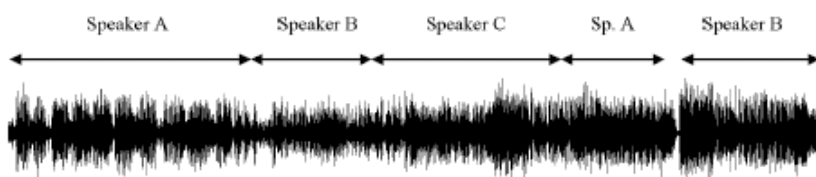
هدف از این پایان نامه، طراحی و پیاده سازی سیستمی است که بتواند در یک فایل صوتی که شامل گفتار چندین گوینده می باشد، تغییر در گوینده را مشخص نماید و تا حد امکان، گفتار هر گوینده را بدون دانستن اطلاعات قبلی از وی، دسته بندی نماید. این سیستم می تواند شامل دو بخش اساسی باشد که عبارتند از:

بخش بندی گوینده

¹.Automatic Audio Indexing

خوشه بندی گوینده

کار قسمت بخش بندی⁶، تقسیم سیگنال گفتاری به سگمنت هایی است که تنها شامل گفتار یک گوینده هستند. در مرحله خوشه بندی⁷، شناسایی و دسته بندی بخش های گفتاری مربوط به یک گوینده و اختصاص یک برچسب واحد به آن انجام می شود. این مطلب در بسیاری از کاربردهای گفتاری که مربوط به بازشناسی یا فهرست نگاری⁸ گفتار در محیطی که چندین گوینده ممکن است در آن اقدام به سخن گفتن بنمایند، مانند یک جلسه، کنفرانس، اخبار و نظایر آن کاربرد دارد. این کار نه تنها می تواند به سیستم های بازشناسی گفتار پیشرفته جهت بهبود نتایج بازشناسی گروهی کمک نماید بلکه در شناسایی و متن نگاری مکالمه ها نیز به آنها کمک می نماید. همانطور که قبلا نیز ذکر شد، امکان استفاده از آن در فهرست نگاری صوتی که امکان جستجو در فایل های صوتی را فراهم می نماید نیز ممکن است. شکل (1-1) نحوه کار این سیستم را بخوبی نشان می - دهد.



شکل (1-1): نمایش بخش بندی گویندگان روی گفتار ورودی

فایل صوتی مورد بررسی یک صوت ضبط شده تک کاناله است که شامل چندین منبع صوتی است. این منابع صوتی متفاوتند و می توانند شامل چند گوینده، موسیقی، انواع نویز و ... باشند. نوع و جزییات منابع صوتی موجود در فایل به ویژگی کاربردی آن فایل بستگی دارد.

بطور کلی سیستم های بخش بندی و خوشه بندی گوینده در سه حوزه زیر دارای کاربرد می باشند:

- دادگان اخباری

- جلسات ضبط شده

- مکالمات تلفنی

⁶.Segmentation

⁷.Clustering

⁸.Indexing

همانطور که قبلا نیز اشاره شد این سه حوزه تفاوت هایی مانند کیفیت ضبط صوت (پهنای باند، میکروفون ها و نویز) و میزان و نوع منابع غیرگفتاری، تعداد گویندگان، سبک و ساختار گفتار (طول مدت گفتار، ترتیب گویندگان) دارند و هر حوزه جهت کار بخش بندی و خوشه بندی گوینده، مسائل و مشکلات خاص خود را دارد. البته در سیستم های تشخیص گوینده سعی بر آن است تا برای هر سه حوزه کاری، نتایج قابل قبول و مناسبی حاصل شود.[1]

در سطح پایین تر کار چنین سیستمی دسته بندی داده های صوتی در خوشه هایی است که هر یک متعلق به یک گوینده باشد. در همین جا به راحتی میتوان دید که دو دیدگاه ناظرانه⁹ (با سرپرست) و غیر ناظرانه¹⁰ (بدون سرپرست) در این بخش مشاهده می شود. در دیدگاه اول از پیش اطلاعاتی از اینکه چه کسانی در فایل صوتی صحبت می کنند، وجود دارد. ولی در دیدگاه دوم کار سیستم دسته بندی فایل به بازه های زمانی است که در آنها تنها یک گوینده که هویت آن بر ما پوشیده است، صحبت می کند. توجه شود که میتوان از خروجی یک دسته بند غیرناظرانه به عنوان ورودی سیستم های شناسایی¹¹، استفاده کرد و به این ترتیب یک سیستم دسته بندی ناظرانه خواهیم داشت. بنابراین کارایی و همچنین زمان اجرای سیستم ناظرانه بدست آمده بهتر است. از سوی دیگر، عملکرد این سیستم ها، به میزان اطلاعات قبلی مجاز نیز بستگی دارد. این اطلاعات قبلی می تواند نمونه گفتار از گویندگان، تعداد گویندگان موجود در فایل صوتی، یا اطلاعاتی از ساختار فایل ضبط شده باشد. ولی در اکثر سیستم های بخش بندی و خوشه بندی گوینده فرض بر نبود هیچگونه اطلاعات قبلی راجع به گویندگان و تعداد آنهاست. در این پروژه نیز با روش های بکار گرفته شده، فرض بر اینست که هیچگونه اطلاعات قبلی از گویندگان، مانند تعداد آنها، هویت آنها و داده آموزشی موجود نمی باشد و بنابراین مدل های گویندگان را نمیتوان از قبل آماده کرد. شکل (1-2) ساختار کلی سیستم های بخش بندی و خوشه بندی گوینده را نشان می دهد.

-
- 1.Supervised
 - 2.Unsupervised
 - 3.Identification



شکل (1-2): ساختار کلی سیستم های بخش بندی و خوشه بندی گوینده

چنین سیستمی شامل مراحل کاری مختلفی است و میتوان بخش های ذکر شده در قسمت های بعدی را برای آنها در نظر گرفت. [5-6]

1-2- مراحل مختلف کاری سیستم های بازشناسی گوینده

بطور کلی مراحل مختلف یک سیستم بازشناسی گوینده، بصورت زیر خلاصه می گردد:

1-قطعه بندی آکوستیکی¹²

2-تشخیص گفتار از غیر گفتار¹³

3-تشخیص جنسیت گوینده

4-تشخیص تغییر گوینده

5-جمع زدن گوینده های مشابه

این سیستم دارای بلوک های کاری مستقل از هم می باشد که هر بلوک ورودی خود را از خروجی بلوک قبلی دریافت می کند و ورودی لازم برای بلوک کاری پس از خود را تهیه می کند. در برخی سیستم ها، از بلوک سوم کاری صرف نظر می شود. در ادامه شرح مختصری از بخش های مختلف داده شده است. [2-4]

1-2-1- قطعه بند آکوستیکی

در اولین مرحله، باید جریان داده های صوتی به قطعات همگن آکوستیکی تقسیم شود. برای این امر باید نقاطی که تغییر در خواص آکوستیکی داده های صوتی روی میدهد را، بدست آورد. در واقع این نقاط شکست¹⁴ بعنوان ورودی به بلوک کاری بعدی داده می شود. در بسیاری از کاربردهای چند رسانه ای که داده ها علاوه بر صدا دارای تصویر نیز می باشند، عمل تشخیص نقاط تغییر، هم از روی صدا و هم از روی تصویر امکان پذیر است. [2] بنابراین کارایی چنین سیستم هایی نسبت به داده هایی که تنها شامل صوت یا تصویر هستند، بالاتر خواهد بود.

امروزه روش های کاربردی تعیین نقاط تغییر آکوستیکی، همگی بر پایه ی محاسبه فاصله آماری بین دو قطعه مجاور استوار هستند. تفاوت عمده ی میان آنها معیار فاصله ای است که در آنها بکار می رود. از روش های غیر آماری مورد استفاده میتوان به شبکه عصبی¹⁵ و ماشین بردار پشتیبان¹⁶ اشاره نمود، که در بخش های بعدی توضیح داده خواهند شد.

از دیدگاهی قطعه بندی، یک مساله بهینه سازی¹⁷ است. زیرا هدف نهایی یافتن نقاطی است که در آنها معیار فاصله به ماکزیمم محلی¹⁸ برسد. یکی از پرکاربردترین معیارهایی که امروزه برای تعیین نقاط شکست آکوستیکی بکار می رود، معیار بیزین¹⁹ است. پیش از این، روش های آماری دیگری از سال 1997 ابداع شده بود، که همگی آنها در مقایسه با معیار بیز جواب مناسبی نمی داده اند. [1] آرایه این روش اعتبار روش های دیگر را تا حدودی کمتر نمود.

-
1. Break Point
 2. Artificial Neural Network
 3. Support Vector Machine
 4. Optimization
 5. Local Maximum
 1. Bayesian Information Criterion

1-2-2- تشخیص گفتار از غیر گفتار (دسته بندی²⁰ صوتی)

برای پیاده سازی این سیستم ها، قبل از هر کار دیگری بخش های گفتاری صوت ضبط شده را از بخش های غیر گفتاری آن مانند (سکوت، موسیقی، نویز خیابان، صدای سرفه ، صدای ورق زدن و ...) جدا می نمایند. با حذف بخش های غیرگفتاری میزان بار محاسباتی سیستم کاهش پیدا می کند و سرعت سیستم بیشتر می شود و سپس مراحل بخش بندی و خوشه بندی اجرا می شود. بعد از یافتن نقاط تغییر آکوستیکی، میتوان جریان داده های صوتی را مانند قطعات همگن در نظر گرفت. به عبارت دیگر یک قطعه نباید هم شامل گفتار، هم موسیقی و سکوت با هم باشد. اگر یک قطعه شامل گفتار دو گوینده باشد، باز هم همگن نخواهد بود. بنابراین این بلوک کاری خروجی قطعه بند صوتی را دریافت کرده و از آن قطعاتی را که حاوی داده های صوتی غیرگفتاری اند را حذف می کند. در یک سیستم تشخیص گفتار، معمولا داده های صوتی به 5 کلاس [2] زیر تقسیم می شوند:

1-موسیقی خالص

2-گفتار خالص

3-گفتار همراه با نویز

4-سکوت

5-سکوت همراه با نویز

البته در یک سیستم تشخیص گوینده، تنها احتیاج به تشخیص موارد 2و3 وجود دارد. زیرا هدف سیستم کار با گفتار بوده و هر چیزی غیر از گفتار از جریان داده ی صوتی حذف می شود تا بلوک های کاری پس از این بلوک با تمرکز بر روی گفتار عمل نمایند. روشی که برای رسیدن به هدف این سیستم وجود دارد، بیشترین میزان شباهت (ML)²¹ مبتنی بر مدل مخلوط گوسی (GMM)²² می باشد.

2. Classification

¹.Maximum Likelihood.

². Gaussian Mixture Model