

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 شورای عالی اطلاع رسانی
	عنوان زیرپروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

عنوان زیرپروژه:

تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 گروه کار اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

فهرست مطالب

شماره صفحه	عنوان
۶	۱ مقدمه
۷	۲ مقدمه ای بر سیستم تشخیص گفتار
۷	2-1 مقدمه
۹	۲-۲ پارامترهای بازشناسی گفتار
۹	۱-۲-۲ وابسته یا مستقل از گوینده
۱۰	۲-۲-۲ گفتار مجزا/ متصل/ پیوسته
۱۰	۳-۲-۲ اندازه کتاب لغت
۱۰	۴-۲-۲ محدودیت‌های زبانی
۱۱	۵-۲-۲ گفتار مکالمه‌ای
۱۱	۶-۲-۲ محیط
۱۱	۳-۲ اجزای یک سیستم بازشناسی
۱۲	۱-۳-۲ نمونه برداری از سیگنال صوتی
۱۲	۲-۳-۲ استخراج ویژگی از سیگنال گفتار
۱۳	۳-۳-۲ تطبیق الگو
۱۴	۴-۳-۲ پردازش زبان
۱۴	۴-۲ انواع مدل سازی در ASR
۱۵	۱-۴-۲ مدل مخفی مارکوف
۱۹	۵-۲ جمع بندی
۲۰	۳ پیش پردازش، نمایش و خلاصه سازی اطلاعات صوتی
۲۰	۱-۳ مقدمه
۲۰	۲-۳ تبدیل صوت از آنالوگ به دیجیتال
۲۱	۳-۳ مشخصات فیزیکی سیگنال گفتار
۲۱	۴-۳ استخراج ویژگی
۲۱	۵-۳ پیش پردازش سیگنال
۲۲	۶-۳ پنجره بندی
۲۲	۷-۳ پردازش بانک فیلتر
۲۳	۸-۳ ضرایب کپسترال، دلتا، دلتا دلتا و انرژی
۲۵	۴ مدل مخفی مارکوف و کاربرد آن در بازشناسی گفتار

	عنوان پروژه:		 مؤسسه ملی اطلاع‌رسانی
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		
	عنوان زیر پروژه:		
	تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی	کد زیر پروژه: پیک‌متن‌فارس - ۲ - خ	ویرایش: ۱/۰
	تاریخ: ۱۳۸۸/۰۳/۱۹		

۱-۴ مقدمه ۲۵

۲-۴ فرآیندهای تصادفی ۲۵

4-3 زنجیره مارکوف ۲۵

۱-۳-۴ زنجیره مارکوف زمان گسسته ۲۵

۲-۳-۴ زنجیره‌های مارکوف زمان پیوسته ۲۶

۴-۴ تعریف مدل مخفی مارکوف ۲۷

۱-۴-۴ اجزاء مدل مخفی مارکوف ۲۷

4-4-2 سه مسأله اساسی مدل مخفی مارکوف ۲۹

۳-۴-۴ نحوه ارزیابی HMM - الگوریتم پیشرو ۲۹

۴-۴-۴ نحوه کدگشایی HMM - الگوریتم ویتربی: ۳۱

۵-۴-۴ نحوه تخمین پارامترهای HMM - الگوریتم Baum - Welch: ۳۲

۵-۴ مدل مخفی مارکوف پیوسته (CDHMM) ۳۶

۱-۵-۴ مدل مخفی مارکوف با چگالی مخلوطی پیوسته ۳۶

۲-۵-۴ مقیاس‌گذاری ۳۸

۶-۴ دنباله‌های چندین مشاهده ای ۴۰

۵ ملاحظات عملی در استفاده از HMM ها ۴۳

۱-۵ تخمین‌های اولیه ۴۳

۲-۵ توپولوژی مدل ۴۴

۳-۵ ضوابط آموزش: یکنواخت کردن پارامترها ۴۵

۶ روشهای جستجو ۴۷

6-1 مقدمه ۴۷

۲-۶ استدلال جلو رو در مقابل عقب‌رو ۴۷

6-3 توابع هیوربستیک ۴۸

۱-۳-۶ روش جستجوی اول عمق ۴۸

۲-۳-۶ روش جستجوی اول سطح ۴۹

6-4 ارزیابی روشهای جستجو ۵۰

۷ مؤلفه‌های فضای جستجو در سیستمهای بازشناسی گفتار ۵۲

۱-۷ مقدمه ۵۲

۲-۷ مدل آکوستیکی ۵۲

۳-۷ مدل زبانی ۵۲

۴-۷ درخت واژگان ۵۳

	عنوان پروژه:		 مؤسسه ملی اطلاع‌رسانی
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		
	عنوان زیر پروژه:		
	تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی	کد زیر پروژه: پیک‌متن‌فارس - ۲ - خ	ویرایش: ۱/۰
	تاریخ: ۱۳۸۸/۰۳/۱۹		

۸ الگوریتم‌های جستجو برای بازشناسی گفتار ۵۵

۱-۸ مقدمه ۵۵

8-2 جستجوی ویتربی ۵۶

۳-۸ جستجو در بازشناسی کلمات مجزا ۵۸

۴-۸ جستجو در بازشناسی گفتار پیوسته ۵۸

۵-۸ نقش اشاره‌گر backtracking ۶۰

۹ مقاوم سازی بازشناسی گفتار ۶۱

۱-۹ مقدمه ۶۱

۲-۹ مدلی از تاثیر محیط بر سیگنال گفتار ۶۲

۳-۹ روشهای مبتنی بر داده ۶۳

۱-۳-۹ تفاضل میانگین ضرایب کپسترال ۶۴

۲-۳-۹ نرمال سازی ضرایب کپسترال با استفاده از میانگین و واریانس ۶۴

۳-۳-۹ آنالیز پیشگویی خطی و درکی گفتار و فیلتر RASTA ۶۵

۴-۳-۹ ضرایب کپسترال ریشه ای ۶۸

۵-۳-۹ ضرایب خود همبستگی فاز (PAC) ۶۹

۶-۳-۹ ویژگیهای مبتنی بر تاخیر گروه ۷۰

۴-۹ روشهای مبتنی بر مدل ۷۰

۱-۴-۹ معیار تصویردهی وزن دار ۷۱

۵-۹ شیوه‌هایی برای تطبیق با گوینده جدید ۷۳

۱۰ تطبیق گوینده ۷۵

۱-۱۰ مقدمه ۷۵

۲-۱۰ تطبیق گوینده برای سیستمهای بازشناسی گفتار ۷۵

۳-۱۰ انواع تطبیق ها ۷۶

۴-۱۰ روش نگاشت طیفی ۷۶

۵-۱۰ روش نگاشت مدل ۷۷

۶-۱۰ تطبیق مدل با روش MLLR ۷۹


۷-۱۰ تبدیل اشتراک ۸۱

۱۱ بازشناسی هویت از طریق گفتار ۸۳

۱-۱۱ مقدمه: ۸۳

۲-۱۱ روشهای پیاده سازی سیستم های تصدیق گوینده: ۸۴

12 سنتز گفتار (Speech Synthesis) ۸۷

	عنوان پروژه:		 انستیتو مطالعات زبان و ادبیات
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		
	عنوان زیر پروژه:		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	کد زیر پروژه: پیکرمتن-فارس - ۲ - خ

۸۷..... ۱-۱۲ مقدمه

۸۷..... ۱۲-۲ انواع متدهای تبدیل متن به گفتار

۸۸..... ۱-۲-۱۲ سنتز شمرده به شمرده لغات (Articulatory Synthesis)

۸۸..... ۲-۲-۱۲ سنتز فرمانت (Formant Synthesis)

۸۸..... ۳-۲-۱۲ سنتز اتصالی (Concatinative Synthesis)

۹۱..... ۱۲-۳ ارزیابی سنتز گفتار

۹۲..... ۱۳ دادگان های گفتاری

۹۲..... ۱-۱۳ مقدمه

۹۲..... ۲-۱۳ دادگان گفتاری فارسدات

۹۲..... ۳-۱۳ دادگان گفتاری فارسدات بزرگ

۹۳..... ۴-۱۳ دادگان گفتاری فارسدات تلفنی (مونولوگ)

۹۳..... ۵-۱۳ دادگان گفتاری فارسدات تلفنی بزرگ (محاوره ای)

۹۴..... ۶-۱۳ دادگان اعداد و ارقام منفصل و پیوسته فارسی

۹۴..... ۷-۱۳ دادگان گفتاری TIMIT

۹۴..... ۸-۱۳ دادگان گفتاری TIDIGITS

۹۵..... ۹-۱۳ دادگان گفتاری AURORA2

۹۷..... ۱۴ نتیجه گیری

۱۰۲..... ۱۵ منابع و مآخذ

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 اُورا محالاطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

۱ مقدمه

امروزه کاربردهای پردازش گفتار فراگیرتر شده اند. نیاز به داشتن سیستمهای اتوماتیک که همانند انسان با گفتار ارتباط داشته باشند روز به روز محسوس تر می گردد. از شاخه های پردازش گفتار می توان به کدینگ، بهبود کیفیت گفتار، بازشناسی گفتار (گسسته یا پیوسته)، بازشناسی گوینده (تعیین و تایید گوینده) اشاره کرد. برای آموزش سیستمهای گفتاری نیاز به داشتن دادگان آموزشی می باشد که با توجه به کاربرد بایستی دارای ویژگیهای خاصی باشند. در این ارتباط برای زبان انگلیسی دادگان مناسبی از قبیل TIMIT و Aurora طراحی شده است. برای زبان فارسی در سالهای اخیر گامهای مطلوبی برداشته شده است اما هنوز خلا داشتن یک پایگاه داده مناسب، استاندارد و قابل اطمینان حس می شود. برای طراحی دادگان مناسب جهت هر کاربرد پردازش گفتار بایستی نیاز آن کاربرد را برای دادگان آموزشی به خوبی مطالعه کنیم. در این مجموعه تلاش شده تا با مطالعه اجمالی روی برخی از کاربردهای پردازش گفتار به شاخص های مهمی که باید در طراحی پایگاه داده استاندارد دیده شود، اشاره شود. به همین منظور چارچوب گزارش بفرم زیر تنظیم گردیده است:

در فصل دوم مقدمه ای بر سیستم تشخیص گفتار خواهیم داشت. سپس در ادامه به مباحث سیگنالینگ، پیش پردازش، نمایش و خلاصه سازی اطلاعات صوتی خواهیم پرداخت. فصل چهارم پرکاربردترین الگوریتم یادگیری و مدل کردن در پردازش گفتار را معرفی می کند، در واقع به معرفی مدل مخفی مارکوف و کاربرد آن در بازشناسی گفتار می پردازد. ملاحظات عملی در استفاده از HMM ها نیز در فصل پنجم آورده شده است. روشهای جستجو، مؤلفه های فضای جستجو در سیستمهای بازشناسی گفتار و الگوریتمهای جستجو برای بازشناسی گفتار نیز مباحث دیگری هستند که در فصول هفتم تا نهم بررسی شده اند. مقوله دیگری در پردازش گفتار، بازشناسی گفتار در محیطهای نویزی است که تحت عنوان مقاوم سازی بازشناسی گفتار در فصل دهم نگاه اجمالی شده است. تطبیق گوینده در سیستم های بازشناسی گوینده در فصل یازدهم آورده شده است. تشخیص هویت گوینده نیز بعنوان کاربرد دیگری از پردازش گفتار در فصل یازدهم مورد توجه قرار گرفته است. فصل دوازدهم حاوی اطلاعاتی در مورد سنتز گفتار می باشد. دادگان های گفتاری مهمترین بخش این گزارش می باشد که به شاخصهای دادگان گفتاری می پردازد، برای مشاهده این مقال به بخش سیزدهم مراجعه کنید. در انتها نیز جمع بندی کلی آورده شده است.

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 ژورنال اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

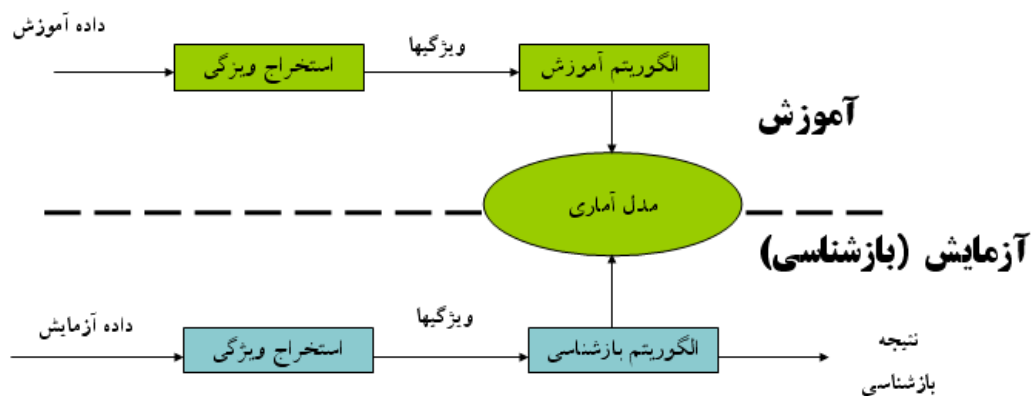
2 مقدمه ای بر سیستم تشخیص گفتار

1-2 مقدمه

سیگنال گفتار را می‌توان به صورت یک پوش طیفی که به آرامی تغییر می‌یابد، در نظر گرفت. این پوش طیفی بوسیله انسان دریافت می‌گردد و به دنباله‌ای از کلمات و معانی آنها ترجمه می‌گردد. سیستم‌های بازشناسی خودکار گفتار نیز بطور مشابه تلاش می‌کنند که این پوش طیفی را به دنباله‌ای از کلمات تبدیل نمایند. مشکلات فراوانی مانند تغییرات پوش طیفی گفتار در این راه وجود دارد. تغییرات پوش طیفی گفتار به دلایلی نظیر جنسیت گوینده، نحوه بیان و تأکید در گفتار گوینده و نیز تغییر محیط آکوستیکی که سیستم بازشناسی در آن عمل می‌کند، بوجود می‌آید. طراحی سیستم بازشناسی خودکاری که بتواند با توانایی انسان در مقابله با این تغییرات برابری نماید، هنوز یک مسأله و چالش اساسی محسوب می‌شود. بسیاری از سیستم‌های بازشناسی گفتار از روش‌های آماری برای مقابله با دسته‌ای از تغییرات پوش طیفی استفاده می‌کنند. کارایی این سیستم‌های آماری تا به امروز به طور قابل توجهی افزایش یافته است، به طوری‌که برای یک کتاب لغت نامحدود مستقل از گوینده صحت بازشناسی بیش از ۹۰٪ بدست آمده است. به علاوه، برای کتاب لغت‌هایی با اندازه محدود، صحت بازشناسی بیش از ۹۵٪ نیز قابل دسترسی است. با این درصد کارایی، این سیستم‌ها قابل بهره‌برداری به نظر می‌رسند، ولی باید توجه کرد که اکثر این سیستم‌ها در محیط‌هایی یکسان و ساکت (بدون حضور نویز) مورد آموزش و آزمایش قرار گرفته‌اند. این در حالی است که در شرایط عملی، سکوت به ندرت وجود دارد و محیط‌های آکوستیکی آزمایش و آموزش نیز عموماً با یکدیگر متفاوتند.

	عنوان پروژه:		 مؤسسه ملی اطلاع‌رسانی
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		
	عنوان زیر پروژه:		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	کد زیر پروژه: پیک‌متن‌فارس - ۲ - خ

در یک سیستم بازشناسی گفتار، پس از تقسیم سیگنال‌ها و نمونه‌های گفتاری به دو دسته دادگان آموزش و آزمایش، سیگنال ورودی به فواصل زمانی ۲۵ تا ۳۰ میلی ثانیه (یا در نظر گرفتن درصدی همپوشانی) قاب‌بندی می‌شود. سپس در مرحله استخراج ویژگی، از هر قاب ویژگی‌های گفتاری (نظیر ضرایب مل کپستروم یا ضرایب پیشگویی خطی ادراکی) استخراج می‌گردند. در مرحله آموزش، از این ویژگی‌ها برای آموزش مدل‌های آماری بازشناسی نظیر مدل مخفی مارکف (HMM) و شبکه عصبی مصنوعی (ANN) استفاده می‌شود. در مرحله آزمایش یا بازشناسی، ویژگی‌ها از طریق یک الگوریتم بازشناسی نظیر ویتربی با این مدل آماری مقایسه می‌شوند. این مراحل و پردازش‌ها در شکل ۱ قابل مشاهده‌اند.



شکل ۱ بخش‌ها و مراحل پردازش در یک سیستم بازشناسی گفتار

یکی از مراحل اصلی در روند بازشناسی گفتار فوق، مرحله استخراج ویژگی و تولید پارامترهایی از سیگنال گفتار است که علاوه بر کاهش حجم داده ورودی به بازشناس گفتار، خصوصیات برجسته‌ای را تعیین کند که واحدهای مختلف گفتاری نظیر واج‌ها، هجاها یا کلمات را از یکدیگر متمایز نماید. در این راستا اکثر سیستم‌های بازشناسی از ویژگی‌های مل کپستروم استفاده می‌نمایند که نحوه عملکرد گوش و سیستم شنوایی انسان را شبیه‌سازی می‌کند. با این وجود، بکارگیری خصوصیات شنوایی و سیستم تولید گفتار در انسان، لزوماً راه‌حل بهینه‌ای برای استخراج ویژگی در بازشناسی گفتار نیست. زیرا این ویژگی‌ها برای ایجاد تمایز میان واحدهای گفتاری بهینه نیستند و علاوه بر این، نیازمند روش‌های تکمیلی برای هنجارسازی گویندگان و مقابله با نویز محیط هستند.

برای رفع این نقیصه و ایجاد تمایز میان واحدهای گفتاری و همچنین برای مقابله با نویز راه‌حل‌های متعددی پیشنهاد شده‌اند. برای مقابله با نویز سه روش کلی وجود دارد:

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 ژورنال اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

۱. حذف نویز از سیگنال گفتار با روش‌هایی نظیر تفاضل طیف و بکاربردن فیلترهای حوزه زمان یا فرکانس نظیر فیلتر وینر.
۲. جبران تأثیر نویز بر ویژگی‌ها که با اعمال تبدیل بر ویژگی‌ها سعی در حذف نویز از ویژگی‌ها دارد. از جمله روش‌های شناخته‌شده این دسته می‌توان به تفاضل میانگین ضرایب کپسترال و یا هنجارسازی ضرایب کپسترال بر حسب میانگین و واریانس اشاره نمود. روش‌های SDCN، SPLICE نیز در این گروه جای می‌گیرند.
۳. روش‌های تطبیق مدل تلاش می‌کنند که بجای اصلاح سیگنال یا پارامترها، مدل آکوستیک محیط یا مدل گوینده را در مرحله بازشناسی اصلاح نمایند. مزیت این روش‌ها آن است که در آنها داده‌های مشاهده شده تغییر نمی‌کنند و هیچ نوع فرض یا تصمیم‌گیری قبلی درباره سیگنال گفتار ضروری نیست. برخی از روش‌های مبتنی بر مدل عبارتند از: تصویر وزن دار، ترکیب موازی مدل‌ها (PMC) و بازگشت خطی با بیشترین درست‌نمایی (MLLR). روش MLLR در مسائل مقاوم سازی نسبت به گوینده و به عبارتی تطبیق گوینده بیشتر بکار گرفته شده است.

2-2 پارامترهای بازشناسی گفتار

پارامترهای مختلفی در یک سیستم بازشناسی گفتار موثر هستند. این پارامترها، تعیین کننده درجه پیچیدگی سیستم می‌باشند. این پارامترها عبارتند از: وابسته و مستقل بودن از گوینده، بازشناسی کلمات مجزا و گفتار پیوسته، اندازه کتاب لغت، محدودیت‌های زبانی، گفتار مکالمه‌ای و شرایط محیطی که بازشناسی در آن انجام می‌گیرد. در این زیربخش این پارامترها به اختصار مورد بررسی قرار می‌گیرند.

۲-۲-۱ وابسته یا مستقل از گوینده

یک سیستم وابسته به گوینده در تعریف فقط برای استفاده یک گوینده طراحی می‌شود، درحالی‌که یک سیستم مستقل از گوینده برای استفاده هر گوینده‌ای طراحی می‌گردد. بطور معمول سیستم‌های وابسته به گوینده دقیق‌تر از یک سیستم مستقل از گوینده هستند و نتایج بهتری را ارائه می‌دهند. ایراد عمده سیستم وابسته به گوینده این است که هر بار برای بازشناسی گوینده جدید نیاز به آموزش دارد. سیستم میانه سیستم‌های وابسته و مستقل از گوینده، سیستم چند گوینده است که برای تعداد گوینده‌های ثابت و کم بکار می‌رود.

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 مؤسسه ملی اطلاع‌رسانی
	عنوان زیرپروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

۲-۲-۲ گفتار مجزا/ متصل / پیوسته

در بازشناسی کلمات مجزا^۱، هر کلمه بصورت جداگانه و واضح بیان می‌شود و سیستم بازشناسی با هر کلمه بطور مستقل سرکار دارد. در بازشناسی کلمات متصل^۲، دنباله‌ای از کلمات برای بازشناسی مورد توجه قرار می‌گیرند، ولی کلمات جمله باید بطور مجزا و با فواصل زمانی سکوت از هم جدا شوند. در بازشناسی گفتار پیوسته^۳، کلمات با مکث‌های از پیش تعریف شده از یکدیگر جدا نمی‌شوند و تلفظ لغات نیز تحت تأثیر آثار هم ادایی^۴ قرار می‌گیرد. بنابراین واضح است که نسبت به دو مورد قبل مشکل‌تر است.

۲-۲-۳ اندازه کتاب لغت

تعداد کلمات موجود در کتاب لغت عامل مهمی در تشخیص کارایی یک سیستم بازشناسی گفتار است. در سیستم‌های کتاب لغت کوچک (کمتر از ۱۰۰ لغت) معمولاً می‌توان به دقتی در حدود ۱۰۰٪ (حتی در سیستم‌های مستقل از گوینده) دست یافت. سیستم‌های با دایره لغات کوچک در کاربردهایی نظیر تشخیص کارت اعتباری و تشخیص شماره تلفن کاربرد دارند. به هر حال دقت بازشناسی به لغات موجود در کتاب لغت نیز وابسته می‌باشد. اگر کلمات مشابه و گیج‌کننده باشند، رسیدن به دقت ۱۰۰٪ حتی برای کتاب لغت‌های بسیار کوچک، دشوار است.

۲-۲-۴ محدودیت‌های زبانی^۵

محدودیت‌های زبانی را می‌توان با یک مدل از زبان بیان نمود. مدل زبانی یک زبان طبیعی مرکب از چهار جزء می‌باشد. نمادها، دستور زبان^۶، معنا^۷ و جنبه عملی^۸. نمادهای زبان واحدهای طبیعی هستند که همه پیغام‌ها از آنها تشکیل می‌گردند و بیانگر کلمات یا واحدهای کوچکتر از کلمه، نظیر هجاها و واجها هستند. دستور زبان، مرکب از محدودیت‌های واژگانی^۹ و نحوی^{۱۰} است که بیانگر شکل گرفتن کلمات از

¹ Isolated Word Recognition

² Connected Word Recognition

³ Continues Speech Recognition

⁴ Co-Articulation

⁵ Linguistic constraints

⁶ Grammar

⁷ Semantic

⁸ Pragmatic

⁹ Lexical

	عنوان پروژه:		 ژورنال اطلاع رسانی	
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی			
	عنوان زیر پروژه:			
تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		کد زیر پروژه: پیکرمتن فارسی - ۲ - خ	ویرایش: ۱/۰	تاریخ: ۱۳۸۸/۰۳/۱۹

واحدهای کوچکتر از کلمه و نیز شکل گرفتن جملات از کلمات می‌باشند. جنبه معنایی نیز مرتبط به نحوه ترکیب کلمات برای شکل دادن پیغام‌های با معنی است. به عنوان مثال جمله «اسب صحبت می‌کند» از لحاظ نحوی صحیح ولی از لحاظ معنایی نادرست است. در بالاترین سطح، جنبه عملی یک زبان جای دارد که بیانگر وابستگی ادا کردن و معنی کلمه به گوینده‌ها و محیط است. محدودیت‌های معنایی و جنبه عملی به ندرت در سیستم‌های بازشناسی گفتار استفاده می‌شوند. چرا که این محدودیت‌ها را به دشواری می‌توان به صورت فرمول بیان کرد. ولی محدودیت‌های دستوری تقریباً در تمامی سیستم‌های بازشناسی گفتار پیوسته به صورت محدودیت‌های واژگانی و نحوی مورد استفاده قرار می‌گیرند و تعداد جملات مجاز برای بازشناسی را کاهش می‌دهند.

۲-۲-۵ گفتار مکالمه‌ای^{۱۱}

گفتار بطور طبیعی به شکلی فوری و مکالمه‌ای بیان می‌شود که تشخیص آن بوسیله ماشین‌ها بسیار دشوار است. در چنین گفتاری، جملات ناقص، شروع‌های مجدد، خنده‌های بلند و سرفه کردن وجود دارد که گفتار را از حالت روان بودن خارج نموده و کتاب لغت را عملاً نامحدود می‌سازد.

۲-۲-۶ محیط^{۱۲}

محیطی که سیستم بازشناسی در آن عمل می‌کند، بر کارایی بازشناسی مؤثر است. شرایط نامناسبی همچون نویز محیط، ضعف میکروفن و اثرات کانال انتقال ممکن است کارایی را به میزان قابل توجهی کاهش دهد. در حالت کلی چنانچه یک سیستم بازشناسی برای یک محیط عاری از نویز طراحی شده باشد، بکار بردن آن، در شرایط نامناسب و نویزی، بدون انجام اصلاحات، کارایی را به شدت کاهش خواهد داد.


3-2 اجزای یک سیستم بازشناسی

یک سیستم بازشناسی گفتار شامل اجزای مختلفی می‌شود. در ادامه این فصل به بررسی اجزای یک سیستم بازشناسی گفتار، از مرحله نمونه‌برداری تا مرحله پردازش زبان می‌پردازیم. بسته به نوع

¹⁰ Syntactic

¹¹ Spontaneous speech

¹² Environment

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 مؤسسه ملی اطلاع‌رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

بازشناسی برای گفتار مجزا، متصل یا پیوسته، ممکن است برخی از این اجزاء در سیستم موجود نباشند یا با جزئیات بیشتری در آن سیستم در نظر گرفته شده باشند.

۲-۳-۱ نمونه برداری از سیگنال صوتی


برای استفاده از الگوریتم‌های پردازش سیگنال گسسته، باید موج پیوسته سیگنال گفتار ورودی را به شکل گسسته تبدیل نمود. به این منظور باید از موج گفتار ورودی نمونه‌برداری کرد. بنابر قضیه نایکوئیست، فرکانس نمونه برداری موج پیوسته باید حداقل دو برابر بزرگترین مؤلفه فرکانسی موجود در موج پیوسته باشد. در کاربردهای عملی، موج گفتار پیوسته عموماً از خروجی یک میکروفن و یا از خط تلفن دریافت می‌شود. بنابراین در بسیاری موارد می‌توان پهنای باند خطوط انتقال تلفن را به عنوان معیاری برای تعیین فرکانس نمونه‌برداری در نظر گرفت.

یکی دیگر از مسائل مهم نمونه‌برداری از سیگنال پیوسته که دقت نمونه برداری را تعیین می‌کند، تعداد بیت‌های مورد استفاده در هر نمونه است. در مرحله نمونه برداری، اندازه هر نمونه به یکی از L سطح ممکن نسبت داده می‌شود. این امر سبب اضافه شدن خطا به اطلاعات می‌شود که به آن خطای چندی کردن اطلاق می‌گردد.

۲-۳-۲ استخراج ویژگی از سیگنال گفتار

پس از قطعه‌بندی به اجزای سازنده گفتار، باید ویژگی‌هایی را از سیگنال گفتار استخراج نمود تا از آنها در مرحله تطبیق الگو استفاده شود. سیگنال گفتار ویژگی‌های زیادی دارد که عموماً با طیف لحظه‌ای سیگنال گفتار یا شکل مجرای گفتار و... مرتبط می‌باشند. پردازش این همه ویژگی برای کاربردی بخصوص، همانند بازشناسی گفتار، کاری منطقی و عملی نخواهد بود. بدین منظور تبدیل‌هایی روی سیگنال گفتار انجام می‌شود تا بتوان ویژگی یا ویژگی‌های مفید را استخراج نمود. استخراج ویژگی به دو دلیل انجام می‌گیرد. اول آنکه سبب تمرکز روی اطلاعات موجود در سیگنال می‌شود و این امر منجر به بهبود میزان شباهت و عدم شباهت میان کلاس‌های مختلف می‌شود. ثانیاً داده‌ها را به نحو قابل ملاحظه‌ای کاهش داده، محاسبات به میزان زیادی کم می‌شود.

به منظور استخراج بردارهای ویژگی باید یک سری پردازش‌ها روی سیگنال انجام شود. این پردازش‌ها

	عنوان پروژه:		 مؤسسه ملی اطلاع‌رسانی	
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی			
	عنوان زیر پروژه:			
تحلیل ابعاد دادگان گفتاری و امکان‌سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		کد زیر پروژه: پیک‌متن‌فارس - ۲ - خ	ویرایش: ۱/۰	تاریخ: ۱۳۸۸/۰۳/۱۹

عبارتند از: قاب‌بندی^{۱۳}، پیش‌تأکید کردن^{۱۴}، اعمال پنجره، تبدیل فوریه زمان کوتاه و....

۲-۳-۳ تطبیق الگو

برای ایجاد بانکِ قواره‌های مرجع در سیستمی که برای بازشناسی گفتار یک گوینده تعلیم می‌بیند، از هر الگوی تلفظ شده توسط آن گوینده، یک یا چند قواره بدست آمده و ذخیره می‌شود. الگوی ذخیره شده می‌تواند یک مدل آماری بیانگر مشخصات آن الگو باشد. در هر حال همواره لازم است که میان مجموعه بردارهای ورودی که هر یک حاوی اطلاعات طیفی اخذ شده از بخشی کوتاه از سیگنال گفتار ورودی هستند و الگوی مرجع، یک تطبیق الگو و مسیریابی زمانی^{۱۵} انجام گردد تا بتوان تغییرات سرعت بیان عبارت را نیز در سیستم منظور کرد. به این ترتیب امکان محاسبه یک معیار شباهت مابین کلمه ورودی و قواره‌های ذخیره شده در سیستم فراهم می‌گردد.

در بعضی از سیستم‌های بازشناسی گفتار بر مبنای کلمه، تطبیق ورودی با قواره‌های مرجع با استفاده از یک روش بهینه‌سازی مشهور به برنامه ریزی پویا^{۱۶} انجام می‌گردد. حل مسأله مسیریابی زمانی با روش برنامه‌ریزی پویا به پیش‌ش زمانی پویا^{۱۷} با علامت اختصاری DTW معروف است. در این روش برای هر کلمه یا عبارت، یک قواره در حافظه نگهداری می‌شود و هنگام تشخیص الگوی ورودی، میزان انطباق آن با تمامی قواره‌ها بررسی می‌گردد. این کار برای تعداد لغات زیاد، نیازمند حجم حافظه و محاسبات بسیار زیادی است. از طرف دیگر، این روش در بازشناسی گفتار پیوسته کارایی خود را از دست می‌دهد، چرا که مرز دقیق کلمات مشخص نیست. به علاوه این روش برای گویندگان متعدد نیز دارای کارایی نیست، زیرا نمی‌تواند تغییرات آکوستیک بین گوینده‌های مختلف را بخوبی مدل کند.

رویکرد دیگر آن است که از هر الگوی لغت‌نامه، یک مدل آماری ساخته شده و هر الگوی لغت‌نامه که مدل آماری آن بیشترین میزان شباهت (احتمال وقوع) را به الگوی ورودی مشاهده داشته باشد، به عنوان الگوی ورودی بازشناخته گردد. معمول‌ترین مدل مورد استفاده در این مورد، مدل مخفی مارکوف^{۱۸} نام دارد. در مدل مخفی مارکوف، هر واحد از لغت‌نامه توسط مجموعه‌ای از حالت‌ها به همراه احتمالات انتقال از حالتی به حالت دیگر نمایش داده می‌شود. در بسیاری از سیستم‌های بازشناسی مبتنی بر مدل مخفی مارکوف، مدل بر پایه واج‌ها تولید شده و سپس در مرحله نهایی بازشناسی، اطلاعات مربوط به

¹³ Framing

¹⁴ Pre-Emphasizing

¹⁵ Time Alignment

¹⁶ Dynamic Programming

¹⁷ Dynamic Time Wrapping

¹⁸ Hidden Markov Model

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 گروه کا اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

واژگان با اطلاعات واجی ترکیب می‌گردد تا کلمات مورد شناسایی قرار گیرند. مرحله آموزش مدل مخفی مارکوف نسبت به DTW پیچیده‌تر است و به داده‌های آموزشی بیشتری نیاز دارد.

یکی دیگر از شیوه‌های تطبیق الگو استفاده از شبکه‌های عصبی مصنوعی است. شبکه‌های عصبی امکانات خوبی را در جهت پردازش‌های موازی و تطبیق‌یابی در اختیار قرار می‌دهند. در بعضی از تحقیقات انجام شده، روش‌های شبکه عصبی با روش‌های دیگری نظیر مدل مخفی مارکوف ترکیب شده‌اند. با توجه به اینکه سیستم تشخیص مورد استفاده در این پایان‌نامه از مدل مخفی مارکوف استفاده می‌کند، این مدل به طور جداگانه، به اختصار شرح داده می‌شود.

۲-۳-۴ پردازش زبان

صرفنظر از آنکه واحد پایه تشخیص گفتار چیست (کلمه، هجا، آوا یا واج)، برای تعیین چگونگی ادغام این واحدها از نظر ترتیب، متن و معنا، از محدودیت‌های زبان استفاده می‌شود. چنانکه قبلاً گفته شد اجزای مدل یک زبان عبارتند از: نمادها، دستور، معنا و جنبه عملی. این اجزاء به عنوان محدودیت‌های زبان استفاده می‌شوند تا در بازشناسی گفتار، تمامی اطلاعات موجود در ارتباط گفتاری در نظر گرفته‌شود.

2-4 انواع مدل سازی در ASR

بیشتر سیستم‌های ASR شامل سه قسمت عمده می‌باشند:

- signal processing front-end

این واحد سیگنال گفتار را به دنباله‌ای از بردارهای ویژگی تبدیل می‌کند. استخراج بردارهای ویژگی به منظور کلاس بندی انجام می‌گیرد. و موجب کاهش نرخ اطلاعات برای پردازش نسبت به سیگنال اصلی می‌شود.

- acoustic modeling

- language modeling

سیستم‌های بازشناسی یا تشخیص خودکار گفتار برای تبدیل گفتار به متن مورد استفاده قرار می‌گیرند. یک سیستم بازشناسی گفتار قابل تجزیه به بلوکهای مجزای عملیاتی است و به هر بلوک می‌توان یک مجموعه ورودی و یک مجموعه خروجی نسبت داد. از جمله می‌توان به موارد زیر اشاره نمود: پردازش سیگنال، ایجاد و مدیریت پایگاه داده گفتار، محاسبات آماری و حسابی، ارزیابی و پیاده‌سازی انواع گرامرها، الگوریتمهای جستجوی کارا، حذف نویز، بهبود گفتار و غیره که باید به نحو مطلوب پیاده‌سازی گردند.

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 مؤسسه ملی اطلاع‌رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

روشهای مختلفی برای بازشناسی گفتار وجود دارد. الگوریتمهای اصلی بازشناسی گفتار بر پایه روشهای زیر می‌باشند:

۱- مدل مخفی مارکوف^{۱۹} - پیچش زمانی پویا^{۲۰} - شبکه عصبی^{۲۱}.

سیستمهای مبتنی بر مدل مخفی مارکوف موفقیت خوبی در پردازش گفتار دارند چرا که روش مدل سازی مارکوف، در عین اینکه دقت بالای بازشناسی را حفظ می‌کند قابلیت مقاوم بودن را برای سیگنالهای گفتار فراهم می‌کند. مدل‌های مخفی مارکوف را بر اساس تابع چگالی آنها به دو دسته مدل پیوسته و مدل گسسته تقسیم می‌کنند که مدل پیوسته دقت بازشناسی بهتری نسبت به مدل گسسته^{۲۲} دارد.

در سیستمهای بازشناسی گفتار که بر اساس مدل مخفی مارکوف می‌باشند، دو قسمت بیشترین زمان بازشناسی را به خود اختصاص می‌دهند: یکی «عملیات جستجو» است که بهترین کلمه منطبق را از روی مقایسه گفتار ورودی با مدل‌های گفتاری مرجع پیدا می‌کند و دیگری قسمت مربوط به «محاسبه احتمال خروجی» در مدل پیوسته مارکوف می‌باشد. معمولاً از الگوریتم «ویتربی» برای انجام عملیات جستجو در سیستمهای بازشناسی گفتار استفاده می‌شود.

در سیستم ASR، بازشناسی یعنی جستجو و برگرداندن دنباله حالتی که با دنباله مشاهده داده شده بیشترین مطابقت را داشته باشد. دقت بازشناسی به تعداد موارد آموزشی مدل و ظرفیت مدل برای مشاهده در مرحله آموزش بستگی دارد.

هر چند روشهایی مثل برنامه‌ریزی پویا و شبکه‌های عصبی نیز در زمینه بازشناسی گفتار فراوان استفاده شده اند ولی تکنولوژی ASR مبتنی بر مدل مخفی مارکوف استفاده بیشتری داشته و مدل مخفی مارکوف یکی از مدل‌های آماری بسیار قوی در توصیف رفتار پیچیده فونهای گفتاری است.

۲-۴-۱ مدل مخفی مارکوف

مدل مخفی مارکوف ابزاری بسیار قدرتمند برای مدل کردن فرآیندهای تصادفی می‌باشد و کاربردهای زیادی در پردازش گفتار دارد. مدل مخفی مارکوف یک فرایند تصادفی دولایه است که در آن فرآیند تصادفی اصلی پنهان است و با استفاده از مجموعه‌ای دیگر از فرآیندهای تصادفی که سری مشاهدات را

¹⁹ hidden Markov model (HHM)

²⁰ dynamic time wrapping (DTW)

²¹ neural network

²² discrete hidden Markov model (DHMM)

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 ژورنال مطالعات زبانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

ایجاد می‌نمایند، آشکار می‌شود.

عناصر لازم برای تعریف مدل مخفی مارکوف، عبارتند از:

N ، تعداد حالات

M ، تعداد نمادهای مشاهده شده در هر حالت

A ، ماتریس احتمالات انتقال بین حالت i به j (a_{ij})

B ، توزیع احتمال مشاهده نماد k در حالت j ($b_j(k)$)

π ، توزیع احتمالات اولیه برای حالت i (π_i)

اگر سیگنال به یک دنباله $\mathbf{X} = \{x_1, x_2, \dots, x_t, \dots, x_T\}$ از بردارهای ویژگی‌ها تبدیل شود، مدل مخفی مارکوف با یک روال به صورت زیر بردارهای ویژگی را تولید می‌کند.

در زمان $t=1$ مدل با احتمال π_{q_1} در حالت q_1 قرار می‌گیرد.

اگر در زمان t مدل در حالت q_t باشد، در این حالت با احتمال $b_{q_t}(x_t)$ بردار ویژگی x_t را تولید خواهد کرد.

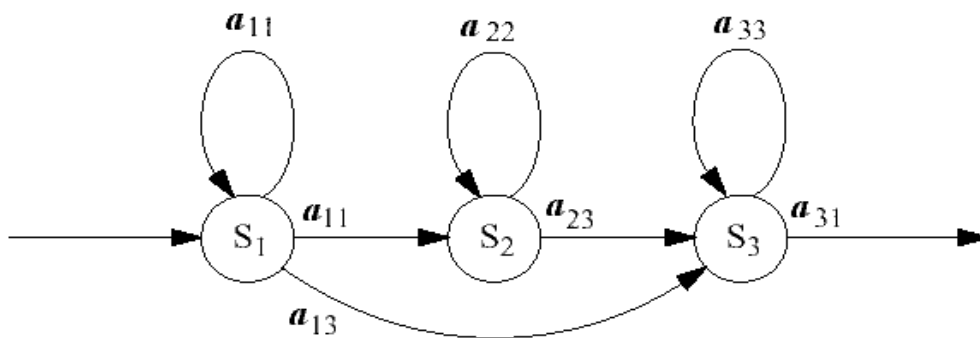
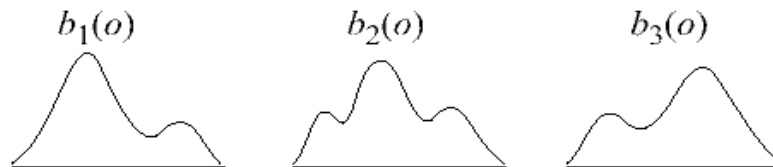
اگر در زمان t در حالت q_t باشد، در این حالت با احتمال $a_{q_t, q_{t+1}}$ به حالت q_{t+1} در زمان $t+1$ خواهد رفت.

با توجه به روال فوق، احتمال تولید رشته بردار \mathbf{X} توسط مدل، با فرض طی کردن رشته حالات $\mathbf{Q} = \{q_1, q_2, \dots, q_T\}$ از رابطه (۱) بدست می‌آید:

$$P(\mathbf{X} | \mathbf{Q}, I) = p_{q_1} b_{q_1}(x_1) a_{q_1 q_2} \dots a_{q_t q_{t+1}} b_{q_{t+1}}(x_{t+1}) \dots b_{q_T}(x_T) \quad (1)$$

رابطه (۱) را می‌توان چنین بیان کرد که مدل مخفی مارکوف λ با پارامترهای π_i و a_{ij} و $b_i(\mathbf{x})$ ، با فرض طی شدن رشته حالات \mathbf{Q} ، رشته بردار \mathbf{X} را با احتمال $P(\mathbf{X} | \mathbf{Q}, \lambda)$ تولید می‌کند.

با توجه به $b_i(\mathbf{x})$ ، دو نوع مختلف از مدل مخفی مارکوف خواهیم داشت. در حالت پیوسته، $b_i(\mathbf{x})$ می‌تواند به صورت مجموع چند تابع چگالی احتمال گوسی تعریف شود. شکل ۲ یک مدل مخفی مارکوف پیوسته سه حالت را نشان می‌دهد.



شکل ۲ مثالی از مدل مخفی مارکوف سه حالت

تابع چگالی احتمال گوسی به صورت زیر تعریف می‌شود:

$$p(x) = N(x, m, \Sigma) = \frac{1}{\sqrt{\det(\Sigma)(2\pi)^d}} \exp\left[-\frac{(x-m)^T \Sigma^{-1}(x-m)}{2}\right] \quad (۲)$$

که در آن m بردار میانگین با طول d و Σ ماتریس کواریانس با ابعاد $d \times d$ است.

مدل مخفی مارکوف دو مسأله اساسی دارد:

مسأله آموزش: مجموعه مشاهده شده $X = \{x_1, x_2, \dots, x_t, \dots, x_T\}$ مفروض است، پارامترهای مدل چگونه تخمین زده شود که $P(X|\lambda)$ بیشینه شود؟ در مدل پیوسته روابط بازتخمین به صورت زیر بدست می‌آید:

$$\bar{c}_{jk} = \frac{\sum_{t=1}^T g_t(j, k)}{\sum_{t=1}^T \sum_{k=1}^M g_t(j, k)} \quad (۳)$$

$$\bar{m}_{ik} = \frac{\sum_{t=1}^T g_t(j, k) \cdot X_t}{\sum_{t=1}^T g_t(j, k)} \quad (۴)$$

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 شورای عالی اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

$$\bar{\Sigma}_{i,k} = \frac{\sum_{t=1}^T g_t(j,k)(X_t - m_{jk})(X_t - m_{jk})^T}{\sum_{t=1}^T g_t(j,k)} \quad (5)$$

در روابط فوق $g_t(j,k)$ احتمال بودن در حالت j در زمان t است با شرط اینکه k امین مخلوط مبین X_t باشد.

مسأله بازشناسی: مجموعه مشاهده شده $\mathbf{X} = \{X_1, X_2, \dots, X_t, \dots, X_T\}$ مفروض است. هدف محاسبه $P(X|I)$ با داشتن پارامترهای مدل است.

$$p(X|I) = \sum_{all Q} p(X|Q, I)P(Q|I) \quad (6)$$

برای حل این احتمال از رویه‌های پیشرو و پسرو استفاده می‌شود. در رویه پیشرو متغیر پیشرو به صورت زیر تعریف می‌شود:

$$a_t(i) = p(X_1, \dots, X_t, q_T = i | I) \quad (7)$$

آنگاه با استفاده از استقرا، احتمال $P(X|\lambda)$ به صورت زیر محاسبه می‌گردد:

$$a_1(j) = p_j b_j(X_1) \quad , \quad 1 \leq j \leq N$$

$$a_{t+1}(j) = \left[\sum_{i=1}^N a_t(i) a_{ij} \right] b_j(X_{t+1}) \quad , \quad 1 \leq t \leq T-1, \quad 1 \leq j \leq N$$

$$p(X|I) = \sum_{i=1}^N a_T(i)$$

در رویه پسرو متغیر اتفاقی پسرو به صورت زیر تعریف می‌شود:

$$b_t(i) = p(X_{t+1}, X_{t+2}, \dots, X_T | q_T = i, I) \quad (8)$$

پس احتمال $P(X|\lambda)$ با استفاده از استقرا به شکل زیر محاسبه می‌گردد:

$$b_T(i) = 1 \quad , \quad 1 \leq i \leq N \quad (9)$$

	عنوان پروژه:		 انستیتو ملی اطلاع رسانی
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		
	عنوان زیر پروژه:		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	کد زیر پروژه: پیک متن فارسی - ۲ - خ


$$b_t(i) = \sum_{j=1}^N a_{ij} b_j(x_{t+1}) b_{t+1}(j) \quad , \quad 1 \leq i \leq N, \quad t = T-1, T-2, \dots, 1 \quad (10)$$

$$p(X | I) = \sum_{j=1}^N p_j b_j(x_1) b_1(j) \quad (11)$$

شرح کامل مدل مخفی مارکوف در مراجع موجود است.

5-2 جمع بندی

در این بخش درباره پارامترهای تعیین کننده پیچیدگی که شامل اندازه کتاب لغت، محدودیت‌های زبانی، گفتار مکالمه‌ای و ... می‌باشند، بحث شد. همچنین اجزای مختلف یک سیستم بازشناسی از مرحله نمونه‌برداری و استخراج ویژگی تا مرحله تطبیق الگو و پردازش زبان مورد بررسی قرار گرفت. و در نهایت یک سیستم بازشناسی گفتار به عنوان سیستم پایه معرفی شد.

	عنوان پروژه:		 ژورنال اطلاع‌رسانی
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		
	عنوان زیر پروژه:		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	کد زیر پروژه: پیک‌متن‌فارس - ۲ - خ

3 پیش پردازش، نمایش و خلاصه سازی اطلاعات صوتی

3-1 مقدمه

در کاربردهایی که به پردازش صوت نیاز دارند (مانند کد کردن، سنتر و بازشناسی) به نمایش خاصی از اطلاعات گفتار نیاز است. نیاز اصلی برای بازشناسی گفتار استخراج ویژگیهایی از صوت می‌باشد که تشخیص واحدهای پایه مورد نظر در بازشناسی را ممکن می‌سازد و مجموعه اطلاعات گفتار را به گونه ای خلاصه سازی می‌نماید که کار تشخیص به سادگی و با دقت لازم انجام گیرد.


3-2 تبدیل صوت از آنالوگ به دیجیتال

به منظور پردازش دیجیتالی صوت، سیگنال صوتی باید از آنالوگ به دیجیتال تبدیل گردد که این تبدیل نیاز به تجهیزات سخت افزاری دارد: یک میکروفن که یک سیگنال صوتی $p(t)$ را به یک سیگنال الکتریکی $x_c(t)$ تبدیل می‌نماید سپس یک نمونه بردار با فواصل زمانی T_c (با فرکانس $f_c=1/T_c$)، مقادیر و ولتاژهای $x_c(nT_c)=x[n]$ را ارایه می‌نماید. نهایتاً یک مبدل آنالوگ به دیجیتال (A/D)، مقادیر $x[n]$ را به مقادیر مشخص، تبدیل می‌کند (معمولاً ۱۶ بیت).

نمونه بردار و مبدل A/D، اغلب در کارتهای صوتی کامپیوترها قرار داده می‌شوند. با توجه به قانون نایکوئیست، اگر فرکانس نمونه برداری بیشتر از دو برابر پهنای باند فرکانسی گفتار اصلی باشد، سیگنال نمونه برداری شده همان اطلاعات سیگنال پیوسته اصلی را خواهد داشت.

به عبارت دیگر سیگنال گفتار از نظر زمان و دامنه گسسته شده و می‌تواند در کامپیوتر ذخیره و پردازش شود. در این حالت ما یک دنباله N نمونه ای را که به یک نقطه واحد در فضای برداری N بعدی منطبق می‌شود خواهیم داشت. یک مرحله از فرآیند برای بازشناسی فونمها تبدیل این مقادیر N تایی به L مقدار با ارزش می‌باشد که $N>L$ بوده و به پیچیدگی محاسباتی کمتری نیاز دارد. این مرحله با محدودیتهایی انجام می‌شود بطوریکه برای هر قاب یک بردار ویژه‌ای اختصاص داده شود.

بعد از این فرآیند، الگوریتم بازشناسی، یک دسته بند ساده است که مشخص می‌کند که کدام مقادیر فضای L بعدی به هر فونم یا هر واحد پایه بازشناسی اختصاص داده می‌شوند. هر نقطه فضای L بعدی از شکل موج N نمونه‌ای به عنوان یک فونم مطابق با آن مقدار کلاس بندی می‌شود.

	عنوان پروژه:		 ژورنال مطالعات رایان
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		
	عنوان زیر پروژه:		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	کد زیر پروژه: پیک-متن-فارس - ۲ - خ

3-3 مشخصات فیزیکی سیگنال گفتار

صدا با تولید کننده هایی تولید می شود که در طول تولید یک فونم در مدت زمان کوتاهی در حالت ثابت^{۲۳} هستند و با یک حرکت تلفظی دیگر به یک حالت ثابت دیگر می رسند عموماً سیگنالهای صوتی در 80-200ms با وضعیت ثابتی از دستگاه گویایی تولید می شوند. لذا در تحلیل گفتار قابها را حدود 30ms در نظر می گیرند که در اینصورت می توان سیگنال صوتی را ایستان فرض نمود.

3-4 استخراج ویژگی

از قسمتهای ایستان سیگنال گفتار باید ویژگیهایی را استخراج کرد تا از آنها در مرحله بازشناسی استفاده نمود. سیگنال گفتار اطلاعات زیادی دارد که عموماً با طیف لحظه‌ای آن یا شکل مجرای گفتار و ... مرتبط می باشند. پردازش این همه ویژگی ضرورت ندارد بدین منظور تبدیلاتی روی سیگنال گفتار انجام می شود و ضرایبی استخراج می شوند که به محاسبات کمتری نیاز دارند و در عین حال اطلاعات لازم را در بردارند.

یکی از ضرایبی که از سیگنال گفتار استخراج می شود ضرایب MFCC^{۲۴} می باشد هر چند روشهای پردازش دیگری نیز وجود دارد ولی این روش دقت کافی و پیچیدگی محاسبات کمتری در مقایسه با دیگر روشها دارد.

3-5 پیش پرداز سیگنال

اغلب پیک به دست آمده ازجهاز صوتی با فرکانس بالا، دامنه کوچکتتری نسبت به فرمانتهای با فرکانس پایین دارد. لذا یک پیش تاکید برای فرکانسهای بالا به منظور یکنواخت شدن سیگنال لازم است. در حقیقت این فیلتر اثرات طیفی حنجره (دو قطب) و لبها (یک صفر) را حذف می کند. همچنین تغییرات ناگهانی موجود در سیگنال را که بر اثر نویزهای شدید محیط بوجود می آید، حذف می کند و باعث یکنواخت شدن سیگنال می گردد. اینگونه پردازش اغلب با فیلتر کردن سیگنال گفتار با یک فیلتر FIR مرتبه اول بدست می آید.

²³ stable

²⁴ mel frequency cepstral coefficient

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 ژورنال اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

پیش پردازش‌های دیگری نظیر حذف نویز و حذف سکوت نیز در سیستم های ASR مهم هستند. در سیستم ASR مبتنی بر HMM اگر سکوت‌های طولانی از گفتار حذف نشوند کارایی بطور قابل ملاحظه کاهش پیدا می‌کند لذا یک تشخیص دهنده گفتار مناسب لازم است. اگر نسبت سیگنال به نویز تغییرات قابل ملاحظه‌ای نداشته باشد تشخیص دهنده‌های ساده مبتنی بر معیار انرژی، کارایی خوبی نشان می‌دهند.

3-6 پنجره بندی

روشهای معمول برای ارزیابی طیفی در صورتی که سیگنال ایستان (یعنی سیگنالی که خصوصیات آماری آن در زمان متغیر نباشد) باشد کارایی خوبی نشان می‌دهند. برای صوت این فرض فقط در فواصل زمانی کوتاه برقرار است یعنی یک تحلیل زمان کوتاه انجام می‌شود که سیگنال را به دنباله های متوالی پنجره بندی شده تبدیل می‌نماید و به طور جداگانه پردازش می‌شوند.

در سیستم‌های ASR پنجره ای که فراوان مورد استفاده قرار گرفته، پنجره همینگ می‌باشد که پاسخ ضربه آن به صورت زیر است:

$$W(n) = 0.54 - 0.46 \cos \frac{2pn}{N-1} \quad 0 \leq n \leq N-1 \quad (12)$$

3-7 پردازش بانک فیلتر

تحلیل طیف و ویژگی‌هایی از سیگنال گفتار را آشکار می‌کند که اساساً با ترکیب جهاز صوتی مرتبط است. برای این منظور ویژگی‌های طیفی گفتار را عموماً از خروجی بانکهای فیلتر که هر کدام در یک رنج فرکانسی تعریف می‌شوند بدست می‌آورند. معمولاً یک مجموعه فیلتر ۲۴ تایی میان گذر برای شبیه سازی پردازش گوش انسان بکار برده می‌شود. فیلترها اغلب در محور فرکانس به شکل غیر یکسان توزیع می‌شوند. به عنوان یک قانون قسمتی از طیف که زیر ۱ KHz است با بانکهای فیلتر بیشتری پردازش می‌شود، چرا که آن قسمت اطلاعات بیشتری را از جهاز صوتی مانند فرمانت اول دربر دارد. پاسخ فرکانسی بانکهای فیلتر، پردازش ادراکی انجام شده در گوش را شبیه سازی می‌کنند و این فیلترینگ وزن دهی ادراکی نامیده می‌شود.

در سیستم های ASR، مقیاس ادراکی که بیشترین کاربرد را داشته، مقیاس mel می‌باشد. فرکانس مرکزی هر بانک فیلتر Mel قبل از ۱ kHz به طور یکنواخت توزیع می‌شود و بعد از ۱ kHz یک مقیاس لگاریتمی را دنبال می‌کند.

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 ژورنال اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

3-8 ضرایب کپسترال، دلتا، دلتا دلتا و انرژی

ضرایب کپسترال معمولاً با یک ضریب انرژی e_t همراه هستند که لگاریتم انرژی قاب یا پنجره می‌باشد. این پارامتر مفید است چرا که فونم‌های مختلف انرژی‌های مختلفی دارند. پارامترهای کپسترال و انرژی، اطلاعات دینامیکی سیگنال گفتار را به حساب نمی‌آورند از طرفی این اطلاعات گاهاً برای ASR ها مفید هستند. لذا می‌توان دیفرانسیل مرتبه اول و دوم را برای بدست آوردن آن اطلاعات بکار برد که حاوی اطلاعات دینامیک و اطلاعات انتقال میان حالات مختلف گویش هستند.

چنانچه t شماره قاب و i شماره ضریب باشد، ضرایب دلتا کپسترال از روابط زیر بدست می‌آیند:

$$d_t(i) = \sum_{t=1}^N \left[\frac{t(C_{t+t}(i) - C_{t-t}(i))}{2 \sum_{t=1}^N t^2} \right]_{N \leq t \leq T-N} \quad (13)$$

که در این رابطه T تعداد کل قابها و N تعیین کننده طول پنجره‌ای است که روی آن مشتق گرفته می‌شود و C_t ضریب کپسترال در زمان t می‌باشد که از فرمول (۱۴) محاسبه می‌شود.

$$(14)$$

در قابهای ابتدایی و انتهایی برای محاسبه مشتقات ضرایب از روابط زیر استفاده می‌شود

$$d_t(i) = C_{t+1}(i) - C_t(i) \quad t < N \quad (15)$$

$$d_t(i) = C_t(i) - C_{t-1}(i) \quad t \geq T - N \quad (16)$$

مقدار T را می‌توان بطور اختیاری تعیین کرد. ولی بطور معمول مقادیر ۲ یا ۳ مناسب‌ترین مقدار برای آن هستند. گاه درمخرج رابطه (۱۳) از یک مقدار ثابت نیز استفاده می‌شود.

برای محاسبه ضرایب دلتا دلتا کپسترال، ضرایب دلتا کپسترال را باید در روابط ۲، ۴ و ۵ قرار داد. ضریب لگاریتم انرژی برای یک قاب از رابطه (۱۷) بدست می‌آید. این ضریب باید برای همه قابهای موجود در سیگنال گفتار بدست آید.

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 ژورنال اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

$$E = \log \left(\sum_{n=1}^M S^2(n) \right) \quad (17)$$

پس از بدست آمدن ضرایب انرژی، می توان مشتقات آن را بدست آورد. ضرایب دلتا انرژی را می توان از رابطه (۱۸) محاسبه نمود.

$$\Delta E_t = \sum_{t=1}^N \left[\frac{t(E_{t+t} - E_{t-t})}{2 \sum_{t=1}^N t^2} \right] \quad N \leq t < T - N \quad (18)$$

در رابطه فوق E_t ضریب لگاریتم انرژی برای قاب t ، T تعداد کل قابهای موجود در سیگنال گفتار، و N طول پنجره‌ای است که روی آن مشتق گرفته می‌شود.

در قابهای ابتدایی و انتهایی، برای محاسبه دلتا انرژی از روابط (۱۹) و (۲۰) استفاده می‌شود:

$$\Delta E_t = E_{t+1} - E_t \quad t < N \quad (19)$$

$$\Delta E_t = E_t - E_{t-1} \quad t \geq T - N \quad (20)$$

برای محاسبه ضرایب دلتا انرژی باید ضرایب دلتا انرژی را در روابط بالا قرار داد.

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 گنجینه اسناد و اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

4 مدل مخفی مارکوف و کاربرد آن در بازشناسی گفتار

4-1 مقدمه

همان طور که گفته شد یکی از روشهای مدلسازی آماری مورد استفاده در گفتار مدل مخفی مارکوف می باشد. مدل کردن سیگنالهای گفتار، برای تشخیص خودکار گفتار، بهبود گفتار، سنتز گفتار و ترجمه ماشینی به کار رفته و توانایی مناسبی از خود نشان داده است.

4-2 فرآیندهای تصادفی

یک فرآیند تصادفی مجموعه‌ای از متغیرهای تصادفی $\{X(t), t \in T\}$ است که برای هر $t \in T$ ، $X(t)$ یک متغیر تصادفی می‌باشد.

4-3 زنجیره مارکوف

4-3-1 زنجیره مارکوف زمان گسسته

در زنجیره مارکوف زمان گسسته، $X = X_1, X_2, \dots, X_n$ یک دنباله از متغیرهای تصادفی از یک الفبای گسسته محدود $O = \{O_1, O_2, \dots, O_n\}$ می‌باشد که در آن توزیع احتمال شرطی هر وضعیت آینده X_{n+1} بطوریکه وضعیتهای قبلی X_0, X_1, \dots, X_{n-1} و وضعیت فعلی سیستم X_n باشد مستقل از وضعیتهای قبلی است و فقط بستگی به وضعیت فعلی آن دارد یعنی به صورت رابطه‌ای خواهیم داشت:

$$P(X_1, X_2, \dots, X_n) = P(X_1) \prod_{t=2}^n P(X_t | X_1^{t-1}) \quad (21)$$

که $X_1^{t-1} = X_1, X_2, \dots, X_{t-1}$ و

$$P(X_t | X_1^{t-1}) = P(X_t | X_{t-1}) \quad (22)$$

و معادله (21) به صورت زیر در می‌آید:

	عنوان پروژه:		 ژورنال تطبیق ریاضی
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		
	عنوان زیر پروژه:		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	کد زیر پروژه: پیک-متن-فارس - ۲ - خ

$$P(X_1, X_2, \dots, X_n) = P(X_1) \prod_{t=2}^n P(X_t | X_{t-1}) \quad (23)$$

چنین فرایند تصادفی زنجیره مارکوف و معادله (۲۳) فرض مارکوف نامیده می‌شود.

زنجیره مارکوف برای مدل کردن رخدادهای ایستان بکار می‌رود یعنی احتمال انتقال وضعیت s' به وضعیت s مستقل از زمان t می‌باشد. به چنین زنجیره‌های مارکوف، زنجیره مارکوف همگن گفته می‌شود که در آن:

$$P(X_t = s | X_{t-1} = s') = P(s | s') \quad (24)$$

اگر X_t را به یک حالت نسبت دهیم زنجیر مارکوف را با یک فرآیند حالت متناهی می‌توان نشان داد بطوریکه گذار بین حالتها با تابع احتمال $P(s | s')$ مشخص می‌شود. اگر مجموعه حالتها زنجیره مارکوف را به $\{1, \dots, N\}$ و حالت زمان t را با s_t نشان دهیم پارامترهای مدل مارکوف بصورت زیر خواهد بود:

$$a_{ij} = P(s_t = j | s_{t-1} = i) \quad , 1 \leq i, j \leq N \quad (25)$$

$$p_i = P(s_1 = i) \quad , 1 \leq i \leq N \quad (26)$$

a_{ij} احتمال انتقال از حالت i به حالت j است و p_i احتمال اولیه بودن زنجیره مارکوف در حالت i می‌باشد و با داشتن a_{ij} ها یک ماتریس انتقال وضعیت تعریف می‌شود. هر دو احتمال محدودیت‌هایی دارند:

$$\sum_{j=1}^N a_{ij} = 1 \quad , 1 \leq i \leq N \quad (27)$$

$$\sum_{j=1}^N p_j = 1 \quad (28)$$

۴-۳-۲ زنجیره‌های مارکوف زمان پیوسته

در زنجیره‌های مارکوف زمان پیوسته مدت زمانی که سیستم در یک وضعیت باقی می‌ماند خود یک متغیر تصادفی نمایی است در صورتیکه در زنجیره‌های مارکوف زمان گسسته هر انتقال یک مرحله است و زمان ماندن در یک وضعیت مطرح نیست که در حالت کلی بصورت زیر تعریف می‌شود:

فرآیند تصادفی با زمان پیوسته $\{N(t), t \geq 0\}$ که مقادیر غیر منفی و صحیح را اختیار می‌کند را در نظر

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 ژورنال اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

بگیرید فرآیند $\{N(t), t \geq 0\}$ یک زنجیره مارکوف با زمان پیوسته است اگر برای تمام $S, t \geq 0$ و مقادیر غیرمنفی $i, j, X(u)$ و $0 \leq u \leq s$ داشته باشیم:

$$P\{X(t+s) = j | X(s) = i, X(u) = x(u), 0 \leq u \leq s\} = P\{X(t+s) = j | X(s) = i\} \quad (29)$$

4-4 تعریف مدل مخفی مارکوف

در زنجیره‌های مارکوف هر حالت بطور مشخص با یک رخداد قابل مشاهده متناظر است یعنی خروجی تصادفی نیست ولی در بعضی زنجیره‌های مارکوف، مشاهده یک تابع احتمالی از حالت است که مدل مخفی مارکوف نامیده می‌شود و از دو لایه تشکیل می‌شود، یک لایه فرآیند تصادفی قابل مشاهده و لایه زیرین، یک لایه فرآیند تصادفی است که مستقیماً قابل مشاهده نمی‌باشد.

پس یک مدل مخفی مارکوف یعنی زنجیره مارکوفی که در آن خروجی مشاهده شده از یک متغیر تصادفی X با تابع احتمال خروجی حالتها مرتبط است.

در این مدل برای یک دنباله مشاهده شده، نمی‌توان یک دنباله حالت که تناظر یک به یک با دنباله مشاهدات داشته باشد مشخص کرد یعنی دنباله حالتها، قابل مشاهده نیست و مخفی است برای همین اصطلاحاً مدل مخفی مارکوف نامیده می‌شود.

۴-۴-۱ اجزاء مدل مخفی مارکوف

یک مدل مخفی مارکوف با اجزاء زیر مشخص می‌شود:

یک الفبای خروجی مشاهده شونده یعنی سمبلهای خروجی فیزیکی سیستم مدل شونده

$$O = \{O_1, O_2, \dots, O_m\}$$


$$\Omega = \{1, 2, \dots, N\} \text{ مجموعه حالتها:}$$

بطوریکه S_t حالت زمان t را نشان می‌دهد.

$$A = \{a_{ij}\} \text{ ماتریس احتمال انتقال:}$$

که a_{ij} یعنی احتمال انتقال از حالت i به حالت j :

$$a_{ij} = P(S_t = j | S_{t-1} = i) \quad (30)$$

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 گروه کا اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

ماتریس احتمال خروجی $B = \{b_i(k)\}$

که $b_i(k)$ احتمال تولید سمبل O_k در حالت i می‌باشد.

اگر $X = X_1, X_2, \dots, X_t, \dots$ خروجی مشاهده شده باشد دنباله حالت‌های $S = S_1, S_2, \dots, S_t, \dots$ قابل مشاهده و مشخص نیست و $b_i(k)$ را می‌توان بصورت زیر نوشت:

$$b_i(k) = P(X_t = O_k | s_t = i) \quad (31)$$

توزیع حالت اولیه $p = \{p_i\}$

بطوریکه

$$p_i = P(s_0 = i) \quad 1 \leq i \leq N \quad (32)$$

از آنجائیکه a_{ij} و $p_i, b_{ij}(k)$ مقادیر احتمال هستند باید خصوصیات زیر را داشته باشند:

$$a_{ij} \geq 0, b_i(k) \geq 0, p_i \geq 0 \quad \forall \text{ all } i, j, k \quad (33)$$

$$\sum_{j=1}^N a_{ij} = 1 \quad (34)$$

$$\sum_{k=1}^M b_i(k) = 1 \quad (35)$$

$$\sum_{i=1}^N p_i = 1 \quad (36)$$


HMM دو پارامتر با مقادیر ثابت N, M دارد که N تعداد حالتها و M اندازه الفبای قابل مشاهده را نشان می‌دهند، همچنین مجموعه الفبای قابل مشاهده O ، و سه مجموعه یا ماتریس احتمالات، که در حالت کلی یک HMM را بصورت زیر نمایش می‌دهیم:

$$\Phi = (A, B, p) \quad (37)$$

در مدل مخفی مارکوف مرتبه اول، دو فرض داریم که اولی فرض زنجیره مارکوف می‌باشد:

$$P(s_t | s_1^{t-1}) = P(s_t | s_{t-1}) \quad (38)$$

که s_1^{t-1} دنباله حالت s_1, s_2, \dots, s_{t-1} می‌باشد.

	عنوان پروژه:		 گروه کا اطلاع رسانی	
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی			
	عنوان زیر پروژه:			
	تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی	کد زیر پروژه: پیکرمتن فارسی - ۲ - خ	ویرایش: ۱/۰	تاریخ: ۱۳۸۸/۰۳/۱۹

و فرض دوم استقلال خروجی است:

$$P(X_t | X_1^{t-1}, s_1^t) = P(X_t | s_t) \quad (39)$$

که X_1^{t-1} دنباله خروجی X_1, X_2, \dots, X_{t-1} را نشان می‌دهد.

شرط استقلال خروجی یعنی اینکه احتمال تولید یک سمبل در زمان t فقط به حالت s_t بستگی دارد و از خروجی قبلی، مستقل است.

این فرضها، بدون اینکه تاثیر چندانی در ظرفیت مدل و حافظه آن داشته باشند باعث می‌شوند که ارزیابی، کدگشایی و آموزش مدل ممکن شود و همچنین تعداد پارامترهایی که باید تخمین زده شوند کم شود.

۴-۴-۲ سه مسأله اساسی مدل مخفی مارکوف

سه مسأله اساسی که در رابطه با مدل HMM قابل طرح است به شرح زیر می‌باشد:

مسأله ارزیابی: با یک مدل Φ و دنباله مشاهده $X = (X_1, X_2, \dots, X_T)$ داده شده، احتمال $P(X | \Phi)$ چیست؟ یعنی احتمال اینکه مدل، مشاهده فوق را تولید نماید.

مسأله کدگشایی: با یک مدل Φ مفروض و دنباله مشاهده $X = (X_1, X_2, \dots, X_T)$ ، بهینه‌ترین دنباله حالات $S = (S_0, S_1, \dots, S_T)$ در مدل که مشاهده فوق را تولید نماید کدام است؟

مسأله آموزش: با یک مدل Φ داده شده و یک مجموعه مشاهدات، چطور ما می‌توانیم پارامترهای مدل $\bar{\Phi}$ را برای بیشینه کردن احتمال توأم $\prod_x P(X | \Phi)$ تنظیم نماییم.

۴-۴-۳ نحوه ارزیابی HMM - الگوریتم پیشرو

برای محاسبه احتمال $P(X | \Phi)$ برای دنباله مشاهده، $X = (X_1, X_2, \dots, X_T)$ و مدل مفروض Φ به صورت زیر عمل می‌شود:

$$P(X | \Phi) = \sum_{alls} P(X | S, \Phi) P(S | \Phi) \quad (40)$$

برای یک دنباله حالت $S = (S_0, S_1, \dots, S_T)$ می‌توان نوشت:

$$P(S | \Phi) = P(S_1 | \Phi) \prod_{t=2}^T P(s_t | s_{t-1}, \Phi) = p_{s_1} a_{s_1 s_2} \dots a_{s_{T-1} s_T} \quad (41)$$

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 شورای عالی اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

همچنین داریم:

$$P(X | S, \Phi) = P(X_1^T | S_1^T, \Phi) = \prod_{t=1}^T P(X_t | s_t, \Phi) = b_{s_1}(X_1) b_{s_2}(X_2) \dots b_{s_T}(X_T) \quad (42)$$

با جایگذاری در معادله (۴۰) خواهیم داشت:

$$P(X | \Phi) = \sum_{alls} a_{s_0 s_1} b_{s_1}(X_1) a_{s_1 s_2} b_{s_2}(X_2) \dots a_{s_{T-1} s_T} b_{s_T}(X_T) \quad (43)$$

برای محاسبه معادله (۴۳) پیچیدگی $O(N^T)$ خواهد بود که برای N و T های بزرگ پیچیدگی نمایی بزرگی بدست می‌آید. برای حل این مشکل از الگوریتمی معروف به الگوریتم پیشرو استفاده می‌شود و ایده آن به این صورت است که نتایج میانی را برای استفاده در مراحل بعدی محاسبات دنباله حالات ذخیره می‌نماییم.

احتمال پیشرو را به صورت زیر تعریف می‌کنیم:

$$a_t(i) = P(X_1^t, s_t = i | \Phi) \quad (44)$$

یعنی احتمال بودن در حالت i در زمان t و تولید پاره مشاهده $X_1, X_2, X_3, \dots, X_t$.

$a_t(i)$ را می‌توان بصورت بازگشتی نوشت:

مرحله اول: مقداردهی اولیه

$$a_t(i) = P(X_1^t, s_t = i | \Phi) \quad (45)$$

مرحله دوم: بازگشت

$$a_t(j) = \left[\sum_{i=1}^N a_{t-1}(i) a_{ij} \right] b_j(X_t) \quad 2 \leq t \leq T, 1 \leq j \leq N \quad (46)$$

مرحله سوم: اتمام

$$P(X | \Phi) = \sum_{i=1}^N a_T(i) \quad (47)$$

این الگوریتم دارای پیچیدگی $O(N^2T)$ خواهد بود.

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 ژورنال مطالعات زبانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

۴-۴-۴ نحوه کدگشایی HMM - الگوریتم ویتربی:

از آنجائیکه دنباله حالات مخفی است روش ممکن و کاربردی، پیدا کردن دنباله حالاتی است که بیشترین احتمال را در تولید دنباله مشاهدات دارند.

به عبارت دیگر ما به دنبال دنباله حالات $S = (S_1, S_2, \dots, S_T)$ هستیم که $P(S, X | \Phi)$ را بیشینه نماید و همان چیزی است که در مسأله پیدا کردن مسیر بهینه در برنامه‌ریزی پویا استفاده می‌شود و الگوریتم ویتربی نام دارد. پیچیدگی الگوریتم ویتربی $O(N^2T)$ می‌باشد.

الگوریتم ویتربی را می‌توان تغییر یافته الگوریتم پیشرو قلمداد کرد که بجای جمع کردن احتمالات از مسیرهای متفاوت به یک حالت مقصد، مسیر بهینه را انتخاب و حفظ می‌کند. برای مشخص کردن احتمال مسیر بهینه:

$$V_t(i) = P(X_1^{t-1}, S_1^{t-1}, S_t = i | \Phi) \quad (48)$$

که $V_t(i)$ احتمال شبیه‌ترین دنباله حالت در زمان t می‌باشد که توسط مشاهده X_1^t (تا زمان t) بدست آمده و در حالت i پایان یافته است. رویه استقرار برای الگوریتم ویتربی بصورت زیر می‌باشد:

مرحله اول: مقداردهی اولیه

$$\begin{aligned} V_t(i) &= p_i b_i(X_1) \quad 1 \leq i \leq N \\ B_1(i) &= 0 \end{aligned} \quad (49)$$

مرحله دوم: بازگشت

$$v_t(j) = \max_{1 \leq i \leq N} [V_{t-1}(i) a_{ij}] b_j(X_t) \quad 2 \leq t \leq T, 1 \leq j \leq N \quad (50)$$

$$B_t(j) = \text{Arg Max}_{1 \leq i \leq N} [V_{t-1}(i) a_{ij}] \quad (51)$$

مرحله سوم: اتمام

$$\text{بهترین نمره} = \text{Max}_{1 \leq i \leq N} [V_t(i)] \quad (52)$$

$$S_T^* = \text{Arg Max}_{1 \leq i \leq N} [B_T(i)] \quad (53)$$

مرحله چهارم: مرحله پسرو

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 ژورنال مطالعات رایانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

$$S_T^* = B_{t+1}(S_{t+1}^*) \quad t = T-1, T-2, \dots, 1 \quad (54)$$

$S^* = (S_1^*, S_2^*, \dots, S_T^*)$ بهترین دنباله می باشد.

۴-۵- نحوه تخمین پارامترهای HMM - الگوریتم Baum - Welch:

تخمین پارامترهای مدل $\Phi = (A, B, p)$ برای توصیف دقیق تر دنباله های مشاهده شده خیلی مهم می باشد. راه حل مشخصی برای بیشینه کردن احتمال توأم در یک فرم بسته وجود ندارد ولی با روش تکراری الگوریتم Baum - Welch می توان مسأله را حل کرد که الگوریتم پیشرو- پسرو نامیده شده است.

معرفی احتمال پسرو:

$$b_t(i) = P(X_{t+1}^T | S_t = i, f) \quad (55)$$

که $b_t(i)$ احتمال تولید X_{t+1}^T (از $t+1$ تا انتها) از زمان t و حالت i می باشد.

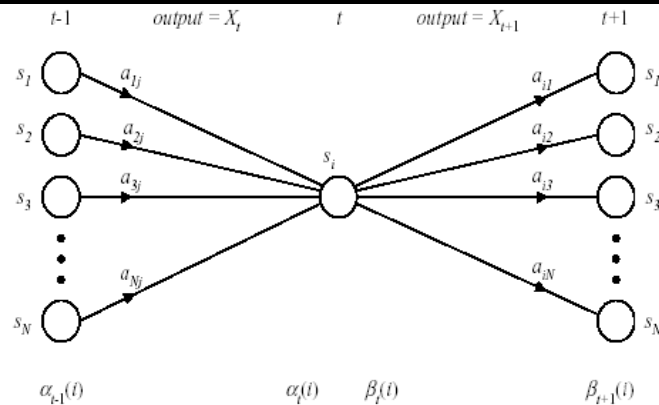
$b_t(i)$ را می توان بروش استقرارء بدست آورد:

مقدار اولیه

$$b_T(i) = \frac{1}{N} \quad 1 \leq i \leq N \quad (56)$$

$$b_t(i) = \left[\sum_j a_{ij} b_j(X_{t+1}) b_{t+1}(j) \right] \quad t = T-1, \dots, 1, 1 \leq i \leq N \quad (57)$$

ارتباط همسایگی در $b, a, (b_t, b_{t+1}, a_{t-1}, a_t)$ در شکل ۳ نشان داده شده است. a بطور بازگشتی از چپ بر راست و b از راست به چپ محاسبه می شوند.



شکل ۳ ارتباط بین $b_t, b_{t-1}, a_t, a_{t-1}$ در الگوریتم پیشرو - پسرو

برای شرح رویه تخمین پارامترهای مدل مخفی مارکوف $x_t(i, j)$ را بصورت زیر تعریف می‌کنیم:

$$x_t(i, j) = P(s_t = i, s_{t+1} = j | X_1^T, \Phi) \quad (58)$$

یعنی احتمال انتقال از حالت i به j برای مدل و دنباله مشاهده مفروض در زمان t .

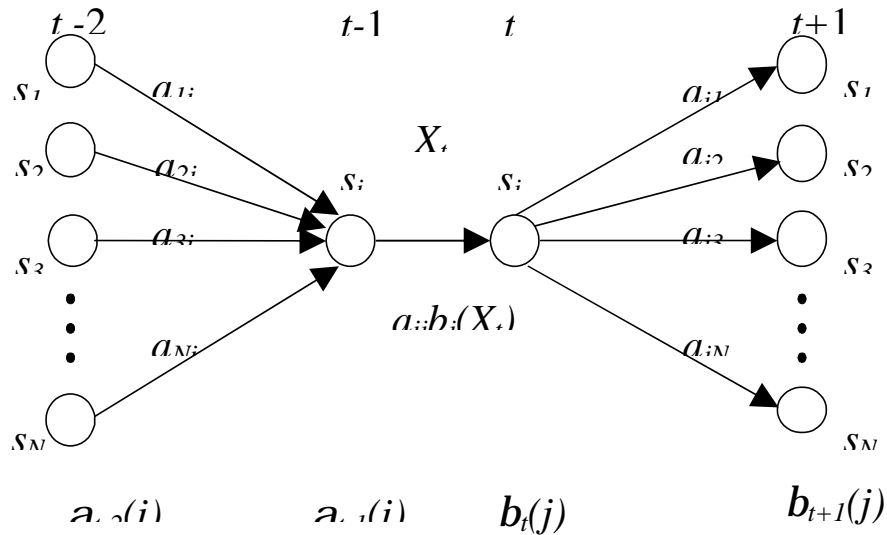
با استفاده از تعریف متغیرهای پیشرو و پسرو داریم:

$$x_t(i, j) = \frac{P(s_t = i, s_{t+1} = j, X_1^T | \mathbf{f})}{P(X_1^T | \mathbf{f})} = \frac{a_t(i) a_{ij} b_j(X_{t+1}) b_{t+1}(j)}{\sum_{k=1}^N a_T(k)} \quad (59)$$

که

$$\sum_{k=1}^N a_T(k) = \sum_{i=1}^N \sum_{j=1}^N a_t(i) a_{ij} b_j(X_{t+1}) b_{t+1}(j) \quad (60)$$

رابطه (۶۰) را می‌توان به فرم شکل ۴ نشان داد:



شکل ۴ نشان دادن عملیات لازم جهت محاسبه $X_t(i, j)$

اکنون متغیر $g_t(i)$ را به صورت زیر تعریف می‌کنیم:

$$g_t(i) = P(S_t = i | X_1^T, \Phi) \quad (۶۱)$$

یعنی احتمال بودن در حالت i در زمان t با دنباله مشاهده مفروض X_1^T وقتی مدل Φ است. با توجه به معادله (۶۱) و متغیرهای پیشرو و پسرو داریم:

$$g_t(i) = \frac{a_t(i)b_t(i)}{P(X_1^T | f)} = \frac{a_t(i)b_t(i)}{\sum_{i=1}^N a_t(i)b_t(i)} \quad (۶۲)$$

خصوصیت زیر با توجه به عامل نرمالیزه (مخرج کسر) برای $g_t(i)$ برقرار است:

$$\sum_{i=1}^N g_t(i) = 1 \quad (۶۳)$$

با توجه به تعریف $g_t(i)$ و $X(i)$ می‌توان آنها را اینگونه به هم مرتبط ساخت:

$$g_t(i) = \sum_{j=1}^N X_t(i, j) \quad (۶۴)$$

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 ژورنال مطالعات زبانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

با جمع بستن $g_t(i)$ بروی t ، کمیتی حاصل می‌شود که می‌توان آن را به عنوان تعداد دفعات مشاهده حالت i یا تعداد گذرهای انجام شده از حالت i در نظر گرفت و جمع $X_t(i, j)$ بر روی t یعنی تعداد انتقال‌ها از حالت i به حالت j . با توجه به روابط بالا، روابط تخمین برای پارامترهای مدل بصورت زیر خواهد بود:

$$\bar{p}_i = g_1(i) \quad (65)$$

یعنی تعداد دفعات قرار داشتن در حالت i در زمان $t=1$.

$$\bar{a}_{ij} = \frac{\sum_{t=1}^{T-1} X_t(i, j)}{\sum_{t=1}^{T-1} g_t(i)} \quad (66)$$

یعنی تعداد گذرها از حالت i به حالت j نسبت به تعداد گذرها از حالت i


$$\bar{b}_j(k) = \frac{\sum_{t=1, s=j, X_k=O_t} g_t(j)}{\sum_{t=1} g_t(j)} \quad (67)$$

یعنی تعداد دفعات بودن در حالت j و مشاهده سمبل X_k نسبت به دفعات بودن در حالت j . مدل تخمین زده شده بصورت $\bar{\Phi} = (\bar{A}, \bar{B}, \bar{p})$ تعریف می‌شود که پارامترهای آن از طرف چپ معادلات (۶۵) تا (۶۷) بدست می‌آید.

به‌وسیله Baum ثابت شده است که یا مدل $\bar{\Phi}$ محتمل‌تر از مدل Φ می‌باشد، $P(X_1^T | \bar{\Phi}) > P(X_1^T | \Phi)$ ، که در نتیجه مدلی پیدا شده که با احتمال بیشتری دنباله مشاهده را تولید می‌نماید و یا اینکه مدل $\bar{\Phi} = \Phi$ می‌باشد. در حالت اول، محاسبات تخمین مجدد تا جایی تکرار می‌شود که یک شرط محدودیت برآورده شود که در نتیجه احتمال مشاهده رشته X_1^T در مدل بهبود می‌یابد و این رویه تخمین، یک تخمین احتمال بهینه یا ML از مدل مخفی مارکوف می‌باشد. محدودیت‌های آماری که در تکرارهای تخمین باید مد نظر قرار گیرند بصورت زیر می‌باشند:

$$\sum_{i=1}^N \bar{p}_i = 1 \quad (68)$$

$$\prod_{i=1}^N \bar{a}_{ij} = 1 \quad (69)$$

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 ژورنال اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

$$\prod_{k=1}^M \bar{b}_j(k) = 1 \quad (70)$$

الگوریتم پیشرو - پسرو که باوم استفاده کرده است یک پیشرفت یکنواخت احتمال را در هر تکرار تضمین می‌کند و نهایتاً به یک احتمال محلی همگرا می‌شود.

4-5 مدل مخفی مارکوف پیوسته (CDHMM)

اگر مشاهدات در یک مجموعه متناهی نباشد بلکه از یک فضای پیوسته باشد توزیع خروجی گسسته قبلی باید تغییر نماید. تفاوت بین HMM گسسته و پیوسته در توابع احتمال خروجی می‌باشد. در بازشناسی گفتار در HMM پیوسته چندی‌سازی یعنی نگاشت بردارهای مشاهده در فضای پیوسته به فضای گسسته انجام نمی‌گیرد و در نتیجه خطای ناشی از چندی‌سازی در نتایج اعوجاج ایجاد نمی‌شود.

4-5-1 مدل مخفی مارکوف با چگالی مخلوطی پیوسته


برای انتخاب یک تابع چگالی مشاهدات پیوسته، برخی محدودیتها باید بر روی فرم تابع چگالی احتمال اعمال شود تا اطمینان حاصل شود که پارامترهای تابع چگالی احتمال می‌توانند به شکل مناسبی تخمین زده شوند. شکل عمومی تابع چگالی احتمال که می‌توان برای آن رویه تخمین مجدد را ارائه نمود مخلوطی متناهی بصورت زیر است:

$$b_j(X) = \sum_{K=1}^M C_{jk} \Pi(X, m_{jk}, \Sigma_{jk}) = \sum_{K=1}^M C_{jk} b_{jk}(X) \quad 1 \leq j \leq N \quad (71)$$

که در آن N تعداد حالات، X بردار مدل شونده، M تعداد مولفه‌های مخلوط، C_{jk} ضریب k امین مخلوط در حالت j ام، Π یک چگالی احتمال بصورت لگاریتمی مقعر یا بیضوی متقارن با بردار میانگین m_{jk} و ماتریس کوواریانس Σ_{jk} برای k امین مخلوط در حالت j ام می‌باشد.

معمولاً Π از نوع توابع چگالی مخلوط گوسی چند متغیره انتخاب می‌شود زیرا این توابع می‌توانند هر تابع چگالی پیوسته را تخمین بزنند. در بعضی کاربردها از توابع چگالی احتمال دیگری نظیر چگالی مخلوط لاپلاسی نیز استفاده می‌شود.

تابع چگالی احتمال گوسی بصورت زیر تعریف می‌شود:

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 ژورنال مطالعات زبانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

$$P(X) = N(X, m, \Sigma) = \frac{1}{\sqrt{\det \Sigma (2\pi)^d}} \exp \left[\frac{-(X - m)^T \Sigma^{-1} (X - m)}{2} \right] \quad (72)$$

این تابع چگالی احتمال d بعدی است که در آن بعد بردار ورودی (ویژگی) X است بردار میانگین d بعدی و ماتریس کوواریانس، ماتریسی متقارن با بعد $d \times d$ می باشد.

$$X = \begin{bmatrix} x_1 \\ \vdots \\ \vdots \\ x_d \end{bmatrix} \quad m = E(X) = \begin{bmatrix} m_1 \\ \vdots \\ \vdots \\ m_d \end{bmatrix} \quad (73)$$

$$\Sigma = \begin{bmatrix} s_{11} & s_{21} & \dots & s_{d1} \\ s_{12} & s_{22} & \dots & s_{d2} \\ \vdots & \vdots & \dots & \vdots \\ \vdots & \vdots & \dots & \vdots \\ s_{1d} & s_{2d} & \dots & s_{dd} \end{bmatrix} = \begin{bmatrix} s_1^2 & s_{21} & \dots & s_{d1} \\ s_{12} & s_2^2 & \dots & \vdots \\ \vdots & \vdots & \dots & \vdots \\ s_{1d} & \dots & \dots & s_d^2 \end{bmatrix} \quad (74)$$

ضرایب مخلوط C_{jk} در محدودیت‌های آماری زیر باید قرار داشته باشد:

$$\sum_{k=1}^M C_{jk} = 1 \quad , \quad 1 \leq j \leq N \quad , \quad C_{jk} \geq 0. \quad (75)$$

بنابراین تابع چگالی احتمال با توجه به این رابطه بطور کامل نرمالیزه می شود یعنی:

$$\int_{-\infty}^{+\infty} b_j(x) dx = 1 \quad (76)$$

رابط تخمین برای C_{jk} و m_{jk} و Σ_{jk} بصورت زیر بدست می آید:

$$\bar{C}_{jk} = \frac{\sum_{t=1}^T g_t(j, k)}{\sum_{t=1}^T \sum_{k=1}^M g_t(j, k)} \quad (77)$$

یعنی تعداد دفعات قرار داشتن در حالت j با استفاده از مخلوط k ام نسبت به تعداد دفعاتی که سیستم در حالت j می باشد.

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 مؤسسه ملی اطلاع‌رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان‌سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

$$\bar{m}_{jk} = \frac{\sum_{t=1}^T g_t(j,k) o_t}{\sum_{t=1}^T g_t(j,k)} \quad (78)$$

هر ترم موجود در صورت را با مشاهده وزن دهی می‌کند و مقدار بخشی از بردار مشاهدات محسوب شده برای k امین جزء مخلوط را بدست می‌دهد.

$$\bar{\Sigma}_{jk} = \frac{\sum_{t=1}^T g_t(j,k) (o_t - m_{jk})(o_t - m_{jk})^T}{\sum_{t=1}^T g_t(j,k)} \quad (79)$$

که در رابطه فوق $g_t(j,k)$ احتمال بودن در حالت j در زمان t است بطوریکه k امین مخلوط مبین O_t باشد عبارتی تعمیم $g_t(j)$ در حالت استفاده از یک مخلوط یا چگالی گسته است. رابطه $g_t(j,k)$ بصورت زیر است:

$$g_t(j,k) = \frac{a_t(j) b_t(j)}{\sum_{j=1}^N a_t(j) b_t(j)} \left[\frac{C_{jk} \Pi(o_t, m_{jk}, \Sigma_{jk})}{\sum_{k=1}^M C_{jk} \Pi(o_t, m_{jk}, \Sigma_{jk})} \right] \quad (80)$$

رابطه تخمین احتمالات انتقال a_{ij} همانند روابط مورد استفاده آن در حالت چگالی گسسته می‌باشد.

۴-۵-۲ مقیاس‌گذاری

زمانیکه احتمالات پیشرو و پسرو را محاسبه می‌کنیم (در الگوریتم پیشرو و پسرو)، اگر طول دنباله مشاهده T به حد کافی بزرگ باشد مقادیر در الگوریتم به صفر میل می‌کنند یعنی برای T های خیلی بزرگ، مقادیر به حدی کوچک می‌شوند که از دقت هر ماشینی کمتر می‌شوند. برای حل این مسأله، می‌توان از یک سری ضرایب برای مقیاس‌گذاری استفاده کرد که در این صورت مقادیر در محدوده دقت ماشین قرار می‌گیرند و ضرایب را در انتهای محاسبات حذف می‌نماییم بدون اینکه در نتیجه بدست آمده تأثیری داشته باشد.

فرض کنید $a_t(i)$ در ضریب S_t زیر ضرب شود:

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 گروه کار اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

$$S_t = \frac{1}{\sum_i a_t(i)} \quad (۸۱)$$

به طوری که $\sum_t S_t a_t(i) = 1$ برای $1 \leq t \leq T$.

$b_t(i)$ را هم می توان در ضریب S_t برای $1 \leq t \leq T$ ضرب نمود. حالت بازگشتی را در محاسبه متغیرهای پیشرو و پسرو می توان در هر مرحله از زمان t با S_t مقیاس گذاری کرد. پس در زمان t ضریب کل مقیاس گذاری که به متغیر $a_t(i)$ پیشرو اعمال می شود به صورت زیر خواهد بود:

$$Scale_a(t) = \prod_{k=1}^t S_k \quad (۸۲)$$

$$Scale_b(t) = \prod_{k=t}^T S_k \quad (۸۳)$$

و ضریب های مقیاس گذاری مجزا در الگوریتم بازگشتی پیشرو و پسرو در هم ضرب می شوند.

با فرض $a'_t(i)$ و $b'_t(i)$ و $x'_t(i, j)$ به عنوان متغیرهای مقیاس گذاری متناظر آنها داریم:

$$\sum_i a'_t(i) = Scale_a(T) \sum_i a_T(i) = Scale_a(T) P(X1\Phi) \quad (۸۴)$$

و احتمال متوسط مقیاس گذاری شده $x'_t(i, j)$ را می توان به صورت زیر نشان داد:

$$x'_t(i, j) = (Scale_a(t-1) a_{t-1}(i) a_{ij} b_j(X_t) b_t(j) Scale_b(t)) / (Scale_a(T) \sum_{i=1}^N a_T(i)) = x_t(i, j) \quad (۸۵)$$

پس احتمالات متوسط را می توان همانند احتمالات مقیاس گذاری نشده به کار برد.

چرا که ضریب مقیاس گذاری حذف خواهد شد:

$$Scale_a(t) Scale_b(t) = \prod_{k=1}^t S_k \prod_{k=t+1}^T S_k = \prod_{k=1}^T S_k = Scale_a(T) \quad (۸۶)$$

که از صورت و مخرج حذف خواهند شد.

بنابراین فرمول تخمین مجدد به همان صورت خواهد بود. بجز اینکه $P(X|\Phi)$ باید با فرمول زیر محاسبه شود:

$$P(X|\Phi) = \sum_i a'_T(i) / Scale_a(T) \quad (۸۷)$$

در مواردیکه احتمال پاریز وجود ندارد مقیاس گذاری لازم نیست و $Scale_a$ را یک در نظر می گیریم.

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 گروه کار اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

یک روش دیگر معمول برای جلوگیری از پاریز استفاده از نمایش لگاریتمی برای همه احتمالات است. این عمل نه تنها باعث می شود که مقیاس گذاری غیرضروری باشد و پاریز رخ ندهد بلکه مقادیر صحیح بدست می آید و نیازی به عملگرهای ممیز شناور نیست که در عمل برای محاسباتی مانند ویتربی مناسب است.

برای ضرب دو عدد می توان لگاریتم آنها را با هم جمع کرد و تقسیم دو عدد نیز معادل تفریق لگاریتم آن دو عدد خواهد بود. فرض کنید در فرمول زیر $P_1 \geq P_2$ باشد:

$$\log_b (P_1 + P_2) = \log_b (b^{\log_b P_1} + b^{\log_b P_2}) = \log_b P_1 + \log_b (1 + b^{\log_b P_2 - \log_b P_1}) \quad (76)$$

اگر ترم دوم معادله بالا کوچک باشد جمع بسادگی برابر $\log_b P_1$ خواهد بود. برای بکارگیری این روش در تخمین مجدد، کلیه روابط باید بصورت لگاریتمی نوشته شود.

4-6 دنباله های چندین مشاهده ای

برای آموزش دنباله های مشاهده چندتایی با فرض استقلال دنباله ها از همان الگوریتم پیشرو - پسر می توان استفاده کرد. مسأله اصلی در مدل های چپ به راست آن است که نمی توان مدل را با یک دنباله مشاهده آموزش داد و پارامترهای آن را تخمین زد. بنابراین برای داشتن داده های کافی جهت تخمین مطمئن همه پارامترهای مدل، باید از دنباله های مشاهده چندتایی استفاده کرد و بنابراین باید روند تخمین مجدد اصلاح گردد.

مجموعه k مشاهده به صورت زیر مشخص است:

$$X = [X^{(1)}, X^{(2)}, \dots, X^{(k)}] \quad (77)$$

بطوریکه $X^{(k)} = (X_1^{(k)}, X_2^{(k)}, \dots, X_T^{(k)})$ ، k امین دنباله مشاهده است و هدف تنظیم پارامترهای Φ برای بیشینه کردن احتمال زیر است:

$$P(X|f) = \prod_{k=1}^k P(X^{(k)} | \Phi) = \prod_{k=1}^k P_k \quad (78)$$

از آنجائیکه معادلات تخمین مجدد براساس فرکانس های وقوع حوادث مختلف بنا نهاده شده است فرمولهای تخمین مجدد برای دنباله های چند مشاهده ای با جمع کردن فرکانسهای منفرد وقوع هر دنباله با یکدیگر بدست می آیند. بنابراین فرمولهای تخمین مجدد اصلاح شده برای $\bar{b}_j(l), \bar{a}_{ij}$ در حالت گسسته عبارتند از:

	عنوان پروژه:		 ژورنال مطالعات زبانی
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		
	عنوان زیر پروژه:		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	کد زیر پروژه: پیکرمتن-فارس - ۲ - خ

$$\bar{a}_{ij} = \frac{\sum_{k=1}^k \frac{1}{p_k} \sum_{t=1}^{T_k-1} a_t^k(i) a_{ij} b_j(X_{t+1}^{(k)}) b_{t+1}^k(j)}{\sum_{k=1}^k \frac{1}{p_k} \sum_{t=1}^{T_k-1} a_t^k(i) b_t^k(i)} \quad (79)$$

$$\bar{b}_j(l) = \frac{\sum_{k=1}^k \frac{1}{p_k} \sum_{t=1, S_t=O_t}^{T_k-1} a_t^k(i) b_t^k(j)}{\sum_{k=1}^k \frac{1}{p_k} \sum_{t=1}^{T_k-1} a_t^k(i) b_t^k(i)} \quad (80)$$

در مدل چپ به راست p_i هم تخمین زده نمی‌شود زیرا $p_1 = 1$ و بقیه p_i ها صفر هستند. برای مقیاس کردن روابط، هر سری مشاهده دارای ضریب مقیاس خود است و ایده اصلی حذف عامل مقیاس‌بندی از هر ترم قبل از جمع کردن آنها می‌باشد. روابط تخمین مجدد بر حسب متغیرهای مقیاس شده بصورت زیر نوشته می‌شود:

$$\bar{a}_{ij} = \frac{\sum_{k=1}^k \frac{1}{p_k} \sum_{t=1}^{T_k-1} Scale_a(t) a_t^k(i) a_{ij} b_j(O_t^k + 1) Scale_b(t) b_{t+1}^k(j)}{\sum_{k=1}^k \frac{1}{p_k} \sum_{t=1}^{T_k-1} Scale_a(t) a_t^k(i) Scale_b(t) b_t^k(i)} \quad (81)$$

و رابطه $\bar{b}_j(k)$ هم به همان صورت به دست می‌آید. با در نظر گرفتن چگالی مشاهدات پیوسته، روابط تخمین مجدد برای مخلوط‌های گوسی با در دست داشتن چندین رشته مشاهده به صورت زیر خواهد بود:

$$\bar{C}_{ij} = \frac{\sum_{k=1}^k \sum_{t=1}^{T_k} x_t^{(k)}(j, l)}{\sum_{k=1}^k \sum_{t=1}^{T_k} \sum_{l=1}^M x_t^{(k)}(j, l)} \quad (82)$$

$$\bar{m}_{jl} = \frac{\sum_{k=1}^k \sum_{t=1}^{T_k} x_t^{(k)}(j, l) \cdot O_t^{(k)}}{\sum_{k=1}^k \sum_{t=1}^{T_k} x_t^{(k)}(j, l)} \quad (83)$$

$$\bar{\Sigma}_{jl} = \frac{\sum_{k=1}^k \sum_{t=1}^{T_k} x_t^{(k)}(j, l) (O_t^{(k)} - m_{j,l})(O_t^{(k)} - m_{j,l})^T}{\sum_{k=1}^k \sum_{t=1}^{T_k} x_t^{(k)}(j, l)} \quad (84)$$

	عنوان پروژه:		 ژورنال مطالعات زبانی
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		
	عنوان زیر پروژه:		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	کد زیر پروژه: پیکرمتن-فارس - ۲ - خ
تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی			

که در روابط فوق $x_t^{(k)}(j, l)$ به صورت زیر تعریف می‌شود:

$$x_t^{(k)}(j, l) = \left[\frac{Scale_a(t) a_t^{(k)}(j) Scale_b(t) b_t^{(k)}(j)}{\sum_{j=1}^N Scale_a(t) a_t^{(k)}(j) Scale_b(t) b_t^{(k)}(j)} \right] \left[\frac{C_{jl} \Pi(O_t^{(k)}, m_{jl}, \Sigma_{jl})}{\sum_{m=1}^M C_{jm} \Pi(O_t^{(k)}, m_{jm}, \Sigma_{jm})} \right] \quad (85)$$

در روابط بالا k تعداد دنباله‌های آموزش، T_k طول دنباله و به عبارتی تعداد بردارهای مشاهدات دنباله k ام می‌باشد.

	عنوان پروژه:		 ژورنال اطلاع رسانی
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		
	عنوان زیر پروژه:		
تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		کد زیر پروژه: پیکرمتن فارسی - ۲ - خ	ویرایش: ۱/۰
تاریخ: ۱۳۸۸/۰۳/۱۹			

5 ملاحظات عملی در استفاده از HMM ها

5-1 تخمین‌های اولیه

می‌دانیم که الگوریتم تخمین مجدد در HMM ها باید برای تابع چگالی احتمال به یک بیشینه محلی برسد. با چه تخمینی ماکزیمم محلی برابر ماکزیمم عمومی است؟

در HMM اگر احتمال، صفر فرض شود در حالت گسسته تا آخر صفر خواهد ماند. پس باید یک تخمین مستدل انتخاب گردد. تجربه نشان داده است که برای HMM گسسته می‌توان توزیع یکنواخت را بعنوان تخمین اولیه انتخاب کرد که برای کاربردهای گفتار مناسب است.


اگر HMM های با چگالی مخلوط پیوسته استفاده شود تخمین اولیه مناسب ضروری است. چندین روش برای بدست آوردن اینگونه انتخابهای اولیه وجود دارد:

روش قطعه‌بندی یکنواخت²⁵: بر اساس توزیع بردارهای مشاهده در میان حالات مدل مخفی مارکوف استوار است. در مشاهدات گسسته هر کدام از بردارهای مشاهده در یک حالت، یکی از M سمبل موجود در کتاب کد است. تعداد بردارهای موجود در یک رشته مشاهده بر تعداد حالات تقسیم می‌شود و به این ترتیب به هر حالت تعدادی از بردارهای مشاهده رشته مربوطه منتسب می‌شود. همین امر در مورد سایر رشته‌های مشاهده مربوط به همان نمونه نیز تکرار می‌شود. سپس تعداد در بردارهای مشاهده منتسب به حالت j ام از طریق جمع تعداد بردارهای منتسب به حالت j ام در تمامی رشته‌های مشاهده بدست می‌آید و مقادیر $b_j(k)$ با تقسیم تعداد بردارهای موجود با اندیس k در بردارهای منتسب به حالت j ام بر این مجموع بدست می‌آید. یعنی:

$$b_j(k) = \frac{\text{تعداد بردارهای با اندیس } k \text{ در حالت } j}{\text{تعداد بردارهای موجود در حالت } j}$$

در مشاهدات پیوسته نیز، ابتدا بردارهای مشاهده منتسب به هر حالت j بدست می‌آید و سپس تمامی این بردارها با استفاده از یک الگوریتم خوشه‌بندی k means به M خوشه در حالت j ام تقسیم می‌شوند که M تعداد مخلوطها در یک حالت است و بنابراین هر خوشه بیانگر یکی از M مخلوط چگالی $b_j(O_t)$ است. سپس پارامترهای مدل در حالت j ام مطابق روابط زیر تخمین زده می‌شود.

²⁵ uniform segmentation

	عنوان پروژه:		 گروه کارشناسی زبان فارسی
	عنوان زیر پروژه:		
	تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	کد زیر پروژه: بیک-متن-فارس - ۲ - خ	

نسبت تعداد بردارهای موجود در خوشه m حالت j به تعداد بردارها در حالت j \hat{C}_{jm}

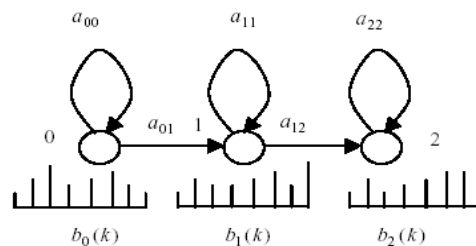
بردار میانگین بردارهای موجود در خوشه m ام حالت j ام \hat{m}_{jm}

ماتریس کوواریانس بردارهای موجود در خوشه m ام حالت j ام $\hat{\Sigma}_{jm}$

به این ترتیب از این مقادیر بعنوان تخمین اولیه پارامترهای مدل مخفی مارکوف با چگالی گسسته یا پیوسته استفاده می‌شود تا بتوان تخمین نهایی مناسبی از این پارامترها به دست آورد.

5-2 توپولوژی مدل

گفتار یک سیگنال غیرایستاد است و هر حالت HMM قابلیت نگهداری بعضی قطعات شبه ایستاد از سیگنال گفتار را دارد. یک توپولوژی چپ به راست مانند شکل ۵ زیر یک انتخاب خوب برای مدل کردن




سیگنال گفتار می‌باشد که انتقال از هر حالت به خود همان حالت نیز ممکن است.

شکل ۵ یک مدل مخفی مارکوف معمول برای مدل کردن واج، شامل سه حالت که هر حالت مرتبط با یک توزیع احتمال خروجی می‌باشد.

زمانیکه یک قطعه گفتار شبه ایستاد تولید می‌شود انتقال از چپ به راست نیز همان تولید گفتار را دنبال می‌کند یعنی برای مدل کردن سیگنالهایی که ویژگیهایشان در طول زمان بطور متوالی تغییر می‌کند مناسب است. در صورتیکه مدل‌های کاملاً متصل برای این کاربردها خیلی مناسب نیستند. توزیع احتمال خروجی وابسته به حالت می‌تواند توزیع گسسته یا تابع چگالی مخلوط پیوسته باشد. در مدل چپ به راست دنباله‌ها از حالت ۱ شروع می‌شود یعنی

$$p_i = \begin{cases} 0 & i \neq 1 \\ 1 & i = 1 \end{cases} \quad (88)$$

و به حالت N ختم می‌شود و همچنین در ماتریس انتقال وضعیت $a_{ij} = 0$ $j < i$ و تغییرات بزرگی در

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 گروه کار اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

اندیسها هم رخ نمی‌دهد یعنی اینکه پرشهای بیش از چند حالت (معمولاً ۲ حالت) مجاز نمی‌باشد یعنی:

$$a_{ij} = 0, \quad j > i + 2 \quad (۸۹)$$

و برای وضعیت نهائی داریم:

$$a_{NN} = 1, \quad a_{Ni} = i, \quad i < N \quad (۹۰)$$

قابل ذکر است که تحمیل محدودیتهایی بر مدل چپ به راست تأثیری در روند تخمین مجدد ندارد. در این توپولوژی مهمترین چیز تعیین تعداد حالتها است. انتخاب توپولوژی مدل به داده‌ها آموزش در دسترس و نوع کاربرد بستگی دارد. اگر هر HMM برای نمایش یک واج بکار رود به سه تا پنج توزیع خروجی نیاز داریم. و برای کلمات، حالت‌های بیشتری نیاز است. (با توجه به تلفظ و طول کلمه) و همچنین برای مدل کردن سکوت نیز 1 یا 2 حالت کفایت می‌کند. در عمل یک انتقال تهی (null) هم تعریف می‌شود که در مواقعی کاربرد دارد که بخواهیم HMM را بدون دیدن هیچ نمادی پیمایش کنیم.

3-5 ضوابط آموزش: یکنواخت کردن²⁶ پارامترها

اگر داده‌های آموزش محدود باشد، باعث می‌شود بعضی پارامترها بصورت ناقص آموزش داده شوند و یا آموزش داده نشوند و دسته‌بندی بر اساس مدل‌های آموزش داده شده ضعیف، باعث افزایش میزان خطا خواهد شد. بعنوان مثال میدانیم محاسبه $b_j(k)$ نیازمند شمارش دفعات بودن در حالت j و مشاهده سمبل O_k بطور همزمان است. اگر دنباله آموزش آنقدر کوچک باشد که در آن چنین حادثه‌ای (مثلاً $X_t = O_k$ و $S_t = j$) رخ ندهد $b_j(k)$ برابر صفر خواهد بود و پس از تخمین مجدد هم صفر می‌ماند و مدل حاصل برای هر دنباله مشاهده شامل $(X_t = O_k$ و $S_t = j)$ احتمال صفر تولید خواهد کرد.

اولین راه حلی که به نظر می‌رسد افزایش مجموعه مشاهدات آموزشی است که همیشه ممکن نیست. دومین راه حل کاهش اندازه مدل مانند تعداد حالتها و تعداد سمبلها در هر حالت می‌باشد که همیشه ممکن نیست و اغلب دلایل فیزیکی مشخصی برای استفاده از یک مدل مفروض وجود دارد. راه حل سوم این است که می‌توان تخمین‌های یک مجموعه از پارامترها را با تخمین‌های مجموعه دیگر از پارامترها درون یابی کرد. راه حل چهارم این است که می‌توان برای کم کردن تعداد پارامترها، آنها را به یکدیگر مرتبط نمود که در HMM شبه پیوسته از این تکنیک استفاده می‌شود. در اکثر این روشها باید از حد آستانه برای پارامترها استفاده کرد تا این اطمینان حاصل شود که تخمین هیچ یک از پارامترهای مدل از

²⁶ .smoothing

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 ژورنال مطالعات زبانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

یک سطح کمتر نمی‌شود. مثلاً برای یک مدل با مشاهدات گسسته داریم:

$$b_j(k) = \begin{cases} b_j(k) & , \quad b_j(k) \geq d_b \\ d_b & , \quad b_j(k) < d_b \end{cases} \quad (91)$$

و برای مدلی با توزیع پیوسته علاوه بر آنکه ماتریس کوواریانس هر مخلوط از نوع قطری انتخاب می‌شود تا تأثیر داده های آموزش ناکافی را خنثی نماید از شرط زیر نیز استفاده می‌شود: (مخصوصاً اگر همبستگی ضرایب ضعیف باشد مثل روش mfcc که ویژگیهای ناهمبسته دارد).

$$\Sigma_{jk}(r, r) = \begin{cases} \Sigma_{jk}(r, r) & , \quad \Sigma_{jk}(r, r) \geq d_l \\ d_l & \quad \Sigma_{jk}(r, r) < d_u \end{cases} \quad (92)$$

پس از اعمال شرایط فوق در روابط تخمین مجدد، باید تمامی پارامترهای باقی مانده طوری تغییر مقیاس داده شوند که چگالیها از محدودیتهای آماری لازم پیروی نمایند.

در عمل، در صورتیکه داده های آموزش ناکافی داشته باشیم درصد خطای بازشناسی گفتار بین ۵ تا ۲۰ درصد با روشهای مختلف یکنواخت کردن پارامترها، کم می‌شود.

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 مؤسسه ملی اطلاع‌رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان‌سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

6 روشهای جستجو

6-1 مقدمه

به هر روند جستجو می‌توان به صورت عبور از یک گراف جهت دار نگریند که در آن هر گره معرف یک وضعیت مسأله است و هر شاخه معرف رابطه بین وضعیت‌هایی است که توسط گره‌های دو سر شاخه ارائه می‌شود.

در طی جستجو باید مسیری در درون گراف یافت شود به طوری که این مسیر از وضعیت آغازی شروع و در وضعیت‌های نهایی خاتمه یابد. گراف جستجو را می‌توان به طور کامل به کمک قواعدی که حرکت‌های مجاز را تعریف می‌کنند، ساخت اما در عمل قسمت زیادی از گراف اصلاً ساخته نمی‌شود. اکثر برنامه‌های جستجو کار جالبی می‌کنند، آنها به جای آنکه اول به طور صریح گراف را بسازند و سپس جستجو کنند، گراف را به صورت غیر صریح توسط قواعد ارائه می‌کنند و به طور صریح فقط قسمت‌هایی از گراف را که مورد جستجو قرار می‌گیرند می‌سازند.

پیش از بحث در مورد تک تک روشها به دو نکته زیر که در همه تکنیکها بروز می‌کند توجه می‌کنیم:

- جهت اجرای جستجو
- استفاده از یک تابع هیوریستیک جهت پیش بردن و راهنمایی جستجو

6-2 استدلال جلو رو در مقابل عقب‌رو

هدف روال جستجو، کشف یک مسیر از میان فضای مسأله از یک وضعیت آغازی به وضعیت هدف است. چنین جستجویی می‌تواند در دو جهت حرکت کند:

- به طرف جلو، از وضعیت‌های آغازی
- به طرف عقب، از وضعیت‌های هدف

این دو روش قرینه یکدیگر هستند. فاکتورهای زیر ما را برای انتخاب جهت جستجو راهنمایی می‌کنند: تعداد وضعیت‌های آغازی بیشتر است یا هدف؟ ما همیشه مایلیم که از تعداد کمتر وضعیت‌ها به طرف تعداد بیشتر برویم.

	عنوان پروژه:			
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی			
	عنوان زیر پروژه:			
تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		کد زیر پروژه: پیک متن فارسی - ۲ - خ	ویرایش: ۱/۰	تاریخ: ۱۳۸۸/۰۳/۱۹

در کدام جهت فاکتور شاخه شاخه شدن بزرگتر است ؟ (این فاکتور تعداد متوسط گره‌هایی است که مستقیماً می‌توان از یک گره تک به آنها رسید) ما حرکت در جهتی که این فاکتور کوچکتر است را ترجیح می‌دهیم.

3-6 توابع هیوریستیک

یک هیوریستیک عبارت است از تکنیکی که با کشف راه حل‌هایی به مسائل کمک می‌کند اما تضمین نمی‌کند که به بیراهه منتهی نگردد. هیوریستیک‌هایی وجود دارند که موارد کاربرد عمومی و زیادی دارند ولی برخی از آنها فقط کاربردهای خاصی دارند. استراتژیهای زیر همگی همه منظوره هستند اما برای اینکه اینها در یک قلمرو به خصوص به درستی و خوب عمل کنند باید آنها را با هیوریستیک‌هایی که ویژه محسوب می‌شوند همراه کرد. یک روش برای این کار استفاده از تابع هیوریستیک است که هر وضعیت مسأله را بررسی می‌کند و میزان خوبی و مناسب بودن آن را تعیین می‌کند.

یک تابع هیوریستیک تابعی است که توصیف وضعیت مسأله را بر میزانهای مناسب بودن نگاشت می‌کند. جستجوی هیوریستیک ابزاری قوی برای حل مسائل مشکل است. استراتژیی که جهت کنترل اینگونه جستجو به کار می‌رود اهمیت به سزایی در تعیین میزان موفقیت این روش دارد. در زیر به برخی از این استراتژیهای کنترل اشاره خواهد شد.

- جستجوی اول عمق²⁷
- جستجوی اول پهنا²⁸
- جستجوی اول بهترین²⁹

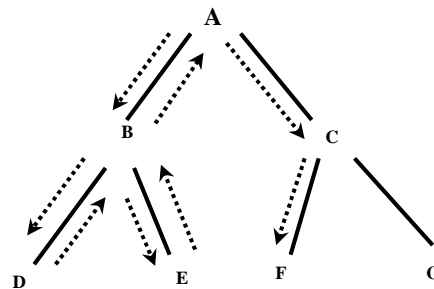
۱-۳-۶ روش جستجوی اول عمق

در جستجوی اول عمق ابتدا هر مسیر احتمالی تا رسیدن به نتیجه (یا هدف) بررسی شده و سپس به مسیر دیگر پرداخته می‌شود. به منظور درک دقیق این شیوه جستجو، درخت زیر که در آن هدف است را در نظر بگیرید:

²⁷ Depth first

²⁸ Breadth first

²⁹ Best first

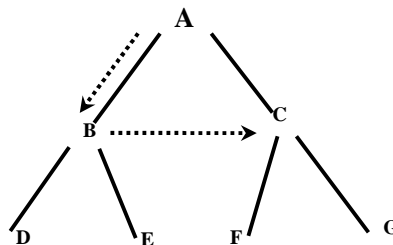


شکل ۶ روند جستجو در شیوه اول عمق

در یک جستجو به شیوه اول عمق این نمودار با ترتیب $ABDBEBACF$ بررسی می‌شود. در این بررسی آنقدر به سمت چپ می‌رویم تا اینکه به یک گره انتهایی یا هدف برسیم. اگر به گره انتهایی رسیده باشیم، آنگاه یک گره به سمت عقب برگشته، به سمت راست می‌رویم و سپس تا هنگامی که یا به هدف و یا گره انتهایی دیگر برسیم به سمت چپ می‌رویم و این عمل تا هنگامیکه یا به هدف برسیم و یا اینکه آخرین گره موجود در فضای جستجو نیز بررسی شود ادامه می‌یابد. یک جستجو انجام شده به شیوه اول عمق قطعاً به جواب می‌رسد زیرا در بدترین حالت، این جستجو به یک جستجوی همگانی تبدیل می‌شود. در نمودار فوق اگر هدف G باشد یک جستجوی همگانی خواهیم داشت.

۶-۳-۲ روش جستجوی اول سطح

شیوه جستجوی اول سطح درست عکس شیوه جستجوی اول عمق است. در این روش تمام گره‌هایی که در یک سطح درخت هستند، پیش از بررسی گرهی در دیگر سطح مورد امتحان قرار می‌گیرند. در زیر این شیوه بین گره‌ها ارائه شده است (در اینجا هدف C است):



شکل ۷ روند جستجو در شیوه اول عمق

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 گروه مخابرات رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

این نمایش نشان می‌دهد که جستجو ابتدا با گره‌های ABC برخورد می‌کند. همانند شیوه جستجوی اول عمق، روال جستجوی اول پهنا نیز در صورت وجود جواب تضمین می‌کند که آن را نهایتاً بیابد اما مشکلاتی هم دارد:

- احتیاج به حافظه زیادی دارد.
 - احتیاج به کار زیادی دارد، به خصوص اگر مسیر حل طولانی باشد.
- در مواردی که چندین راه حل طولانی وجود دارد بهتر است از همان روش اول عمق استفاده شود و نه از اول پهنا تا بتوان سرعت بیشتری به دست آورد.


4-6 ارزیابی روشهای جستجو

ارزیابی قابلیت اجرایی یک سیستم جستجو می‌تواند عملی بسیار پیچیده و دشوار باشد در واقع این ارزیابی بخش عمده‌ای از تحقیقات هوش مصنوعی را به خود اختصاص داده است. با این وجود در توصیف ما از این ارزیابی، دو اندازه‌گیری اولیه از اهمیت خاصی برخوردارند:


- یک جستجو با چه سرعتی یک راه حل را می‌یابد؟
- صحت این راه حل چقدر است؟

انواع مختلفی از مسائل وجود دارند که در مورد آنها همه موضوعات در این نکته خلاصه می‌شود که یک راه حل (هر چه که می‌خواهد باشد) را با حداقل کوشش به دست آورید. در این نوع از مسائل، اندازه‌گیری مورد اول اهمیت دارد با این وجود در انواع دیگری از مسائل آنچه که مهم است آن است که راه حل یافت شده به راه حل بهینه تا حد ممکن نزدیکتر باشد. درک تفاوت بین یافتن یک راه حل «بهینه» و یافتن یک راه حل «خوب» مهم است. این تفاوت بر اساس این واقعیت به وجود می‌آید که یافتن یک راه حل بهینه اغلب اوقات مستلزم انجام یک «جستجوی همگانی» است زیرا این تنها شیوه‌ای است که به وسیله آن می‌توان تعیین کرد که آیا بهترین راه حل یافت شده است یا نه. با این وجود یافتن یک راه حل خوب به معنی یافتن راه حلی است که در داخل مجموعه راه‌حل‌ها قرار می‌گیرد و ممکن است یک راه حل بهتر از آن نیز وجود داشته باشد.

کلیه تکنیکهای جستجویی که ارائه شد در برخی مسائل مشکل، نسبت به یکدیگر دارای برتری‌هایی هستند. از این رو از نظر کلی به دشواری می‌توان گفت که یک شیوه جستجو «همواره» از شیوه‌ای دیگر بهتر است. با این وجود برخی از این تکنیکها به طور متوسط نسبت به ما بقی احتمال بهتر بودن بیشتری

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 شورای عالی اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

دارند. همچنین خود شیوه‌ای که یک مسأله توسط آن تعریف می‌شود نیز در برخی موارد یک شیوه جستجوی مناسب را پیش روی ما قرار می‌دهد.

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 ژورنال اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

7 مؤلفه‌های فضای جستجو در سیستم‌های بازشناسی گفتار

7-1 مقدمه

مدل آکوستیکی، مدل زبانی و درخت واژگان از اصلی‌ترین مؤلفه‌های تشکیل دهنده‌ی فضای جستجو در سیستم‌های بازشناسی گفتار هستند که می‌توانند در محدود کردن جستجو و افزایش سرعت آن، نقش مهمی را ایفاء کنند.

7-2 مدل آکوستیکی


یک فرض کلیدی در پردازش تصادفی گفتار این است که سیگنال گفتار در محدوده‌های زمانی کوتاه، ایستان است. با این فرض بخش استخراج ویژگی، سیگنال گفتار پیوسته را به یک دنباله از بردارهای ویژگی تبدیل می‌کند. هدف این بخش از سیستم بازشناسی گفتار، ارائه تغییرات زمانی و طیفی سیگنال گفتار به صورت دنباله‌ای از بردارها می‌باشد.

بعد از استخراج دنباله‌ی بردارهای ویژگی X از سیگنال گفتار، لازم است که به کمک مدل‌های آکوستیکی مقدار احتمال یا امتیاز $P(X | W)$ محاسبه شود. در کاربردهایی که تعداد کلمات واژگان زیاد می‌باشد، محاسبه $P(X | W)$ برای تک تک کلمات موجود در بانک، غیرعملی می‌باشد. بنابراین دنباله کلمات به واحدهای صوتی پایه‌ای بنام واج شکسته می‌شوند.

برای هر واج یک HMM در نظر گرفته شده و به کمک داده‌های موجود در دادگان آموزشی، آموزش داده می‌شود. HMM یک ماشین حالت متناهی است که انتقال بین حالت‌های آن بر اساس یک توزیع مارکوف است و یک تابع چگالی احتمال، خروجی هر حالت آن را مدل می‌کند، این تابع نقش مهمی را در مدل کردن تغییرات طیفی سیگنال گفتار بازی می‌کند. بنا به پیچیدگی مساله‌ی بازشناسی، تابع چگالی احتمال به صورت پیوسته یا گسسته مدل می‌شود.

7-3 مدل زبانی

مدل زبانی، محدودیت‌هایی روی دنباله‌های کلمات مورد بازشناسی اعمال می‌کند و در واقع مکانیزمی

	عنوان پروژه:		 شورای عالی اطلاع رسانی
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		
	عنوان زیر پروژه:		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	کد زیر پروژه: پیک متن فارسی - ۲ - خ

برای تخمین احتمال رخداد کلمات در یک دنباله از کلمات با فرض مشخص بودن دنباله کلمات پیرامون آن است. مدل زبانی به طور ضمنی دانش زبانی، دانش دامنه و هر گونه اطلاعات دیگری را در جهت کاهش فضای جستجو، جمع می‌کند. میزان محدودیتی که توسط مدل زبانی در یک سیستم بازشناسی ایجاد می‌شود، پیچیدگی آن مدل زبانی نامیده می‌شود و در واقع بیانگر میانگین فاکتور انشعاب است.

از آنجایی که احتمال بیان یک کلمه در هر مرحله از گفتار، در اکثر مواقع به کلمات بیان شده‌ی قبلی وابسته است، یک راه ساده و مؤثر برای به کارگیری مدل زبانی این است که هر دنباله‌ی M کلمه‌ای را توسط یک زنجیره‌ی مارکوف مرتبه n ام مدل کنیم که به اصطلاح N -gram نامیده می‌شود.

4-7- درخت واژگان

درخت واژگان، یکی از اساسی‌ترین اجزاء یک سیستم بازشناسی گفتار با واژگان زیاد می‌باشد. دو راه برای توصیف واجی واژگان در بازشناسی گفتار پیوسته وجود دارد: یکی ارائه خطی است که چندان مؤثر و کارا نیست چون از شباهت واجی بین کلمات سودی نمی‌برد و توصیف واجی هر واژه به صورت مستقل از بقیه ارائه می‌شود. پس فضای جستجو با افزایش تعداد واژگان به صورت خطی افزایش پیدا می‌کند. مثلاً اگر دو کلمه «برادر» و «برادران» که اولی پیشوند دومی می‌باشد، جزء واژگان باشد، توصیف واجی هر کدام به صورت مستقل ارائه می‌شود که در صورت زیاد بودن اینچنین کلماتی مقدار زیادی حافظه به هدر می‌رود.

روش دوم ارائه توصیف واجی واژگان به صورت درختی می‌باشد که به اختصار درخت واژگان³⁰ یا درخت پیشوندی³¹ نامیده می‌شود و در مواردی که تعداد واژگان زیاد باشد، متداول‌ترین روش محسوب می‌شود. جستجو بر مبنای درخت واژگان یک فاکتور اساسی برای ایجاد یک سیستم بازشناسی گفتار بلادرنگ با تعداد واژگان زیاد می‌باشد.

شکل ۸ نمونه‌ای از یک درخت واژگان را نشان می‌دهد. همان‌طور که دیده می‌شود واح‌های مشابه ابتدایی به اشتراک گذاشته شده‌اند. هر برگ این درخت متناظر با یک کلمه در واژگان می‌باشد. به دو دلیل زیر یک یال تهی، که به صورت خط‌چین نشان داده شده است، برای ارائه گره برگی هر کلمه به کار گرفته می‌شود:

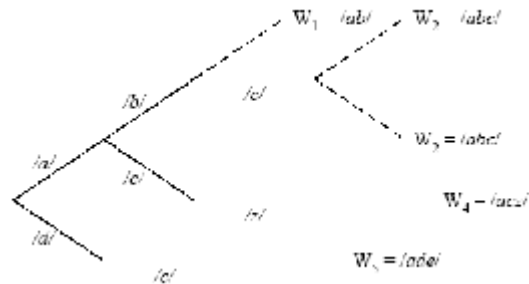
وقتی توصیف واجی یک کلمه پیشوند کلمات دیگر می‌باشد. یال تهی به عنوان یک شاخه برای اتمام کلمه پیشوند عمل می‌کند.

³⁰ Lexical Tree

³¹ Prefix Tree

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 مؤسسه ملی اطلاع‌رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان‌سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

وقتی در واژگان کلماتی با توصیف واجی یکسان وجود دارد، مانند دو کلمه “two” و “to”. یالهای تهی برای اشاره به این کلمات استفاده می‌شود.



شکل ۸ نمونه‌ای کوچک از یک درخت واژگان

یک مزیت عمده استفاده از درخت واژگان این است که می‌تواند فضای جستجو را به طرز چشمگیری محدود کند.

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 ژورنال اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

۸ الگوریتم‌های جستجو برای بازشناسی گفتار

۱-۸ مقدمه

در سیستم‌های بازشناسی گفتار، دیکدر، اساساً یک فرایند جستجو برای آشکار کردن دنباله کلماتی $\hat{W} = w_1 w_2 \dots w_m$ است که بیشترین احتمال پسین $P(W|X)$ ³² را برای دنباله مشاهدات داده شده $X = x_1 x_2 \dots x_n$ داشته باشد.

$$\hat{W} = \arg \max_w P(W | X) = \arg \max_w \frac{P(W)P(X | W)}{P(X)} = \arg \max_w P(W)P(X | W) \quad (۹۳)$$

یک راه بدیهی آن است که همه دنباله کلمه‌های ممکن را جستجو کرده و سپس کلمه با بیشترین مقدار $P(W)P(X | W)$ را انتخاب کنیم.

وقتی مدل‌های کلمات موجود باشد، بازشناسی گفتار به یک مسأله جستجو تبدیل می‌شود که هدف آن یافتن دنباله‌ای از مدل‌های کلمات است که به بهترین صورت، شکل موج ورودی را در برابر مدل‌های لغات موجود بیان می‌کند. با توجه به اینکه تعداد کلمات و واج‌ها و همچنین مرزهای بین آنها در شکل موج ورودی مشخص نیست، استفاده از روش جستجویی مناسب که از عهده این الگوهای پویا با طول متغیر برآید، اهمیت زیادی دارد.

وقتی از HMM در سیستم بازشناسی گفتار استفاده می‌شود، حالت‌های HMM ها فضای جستجو را تشکیل می‌دهند ولی علاوه بر این مؤلفه‌های دیگری هم در میزان اندازه فضای جستجو نقش دارند. در ادامه ابتدا با مؤلفه‌های فضای جستجو در بازشناسی گفتار آشنا می‌شویم و پس از آن روش‌های جستجو در بازشناسی گفتار مطرح می‌شوند. اگرچه روش‌هایی که بیان خواهد شد بر مبنای HMM توضیح داده می‌شوند اما همه این روش‌ها را می‌توان در سیستم‌هایی که از تکنیک‌های مدل‌سازی دیگری استفاده کرده‌اند نیز به کار برد. درحقیقت بسیاری از این روش‌های جستجو قبل از اینکه از HMM برای بازشناسی گفتار استفاده شود، کشف شده بودند.

³² Maximum posterior probability

	عنوان پروژه:		 شورای ملی اطلاع رسانی	
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی			
	عنوان زیر پروژه:			
	تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی	کد زیر پروژه: پیک متن فارسی - ۲ - خ	ویرایش: ۱/۰	تاریخ: ۱۳۸۸/۰۳/۱۹

8-2 جستجوی ویتربی

وقتی از HMMها برای مدل آکوستیک استفاده می‌شود، از امتیاز³³ مدل آکوستیک طبق تعریف احتمال پیش‌رو³⁴ استفاده می‌شود. بنابراین همه دنباله‌ها باید ملاحظه شوند:

$$P(X | W) = \sum_{\text{all possible } S_0^T} P(X, S_0^T | W) \quad (94)$$

که در رابطه فوق مجموع همه حالت‌های ممکن S برای یک دنباله کلمه مورد نظر W، حساب می‌شود. الگوریتم پیش‌رو احتمال اینکه یک HMM، یک دنباله مشاهده را تولید کند، با جمع کردن احتمال‌های همه مسیرهای ممکن حساب می‌کند اما بهترین مسیر (یا دنباله حالت) را مشخص نمی‌کند، درحالی‌که در بسیاری از کاربردها یافتن چنین مسیری مطلوب است. واقعیت این است که یافتن مسیر بهینه (دنباله حالت بهینه) اساس جستجو در بازشناسی گفتار پیوسته می‌باشد. از آنجا که دنباله حالت در ساختار HMM پنهان است (مشاهده نمی‌شود) و با توجه به اینکه هدف دیکدینگ آشکار کردن بهترین دنباله کلمه است، می‌توان مجموع را با بیشینه‌ای برای یافتن بهترین دنباله حالت تقریب زد. معیاری که اغلب استفاده می‌شود، یافتن دنباله حالتی است که هنگام تولید دنباله مشاهده بیشترین احتمال را به دست آورد. طبق قانون بیز، معادله (94) به صورت زیر درمی‌آید:

$$\hat{W} = \arg \max_w P(W) P(X | W) \cong \arg \max_w \{ P(W) \max_{S_0^T} P(X, S_0^T | W) \} \quad (95)$$

معادله (95) اغلب با عنوان «تقریب ویتربی» شناخته می‌شود که می‌توان آن را چنین بیان کرد: «محتمل‌ترین دنباله کلمه» با «محتمل‌ترین دنباله حالت» تقریب زده می‌شود. اگرچه اصولاً نتایج جستجو به روش احتمال پیش‌رو و روش تقریب ویتربی می‌توانند متفاوت باشند، اما در عمل این وضعیت به ندرت اتفاق می‌افتد.

به الگوریتم ویتربی می‌توان به عنوان یک الگوریتم برنامه‌ریزی پویا که به HMM اعمال می‌شود و یا یک الگوریتم پیش‌روی اصلاح شده نگاه کرد. الگوریتم ویتربی به جای جمع کردن احتمال‌های مسیرهای مختلفی که به یک حالت نهایی می‌رسند، بهترین مسیر (بهترین دنباله حالت) را برمی‌گزیند و به خاطر می‌سپارد. این الگوریتم دنباله حالتی $S = (s_1, s_2, \dots, s_T)$ را جستجو می‌کند که مقدار $P(S, X | \Phi)$ را ماکسیمم سازد. این مسأله شباهت زیادی به مسأله یافتن مسیر بهینه در برنامه‌ریزی پویا دارد.

برای یافتن بهترین دنباله حالت $S = (s_1, s_2, \dots, s_T)$ برای دنباله مشاهده داده شده

³³ score
³⁴ forward

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 شورای عالی اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

$X = (x_1, x_2, \dots, x_t)$ لازم است کمیت $d_t(i)$ را تعریف کنیم:

$$d_t(i) = \max_{s_1, s_2, \dots, s_t} P[s_1, s_2, \dots, s_{t-1}, s_t = i, x_1, x_2, \dots, x_t | \Phi] \quad (96)$$

در طول یک مسیر بالاترین امتیازی است که برای اولین t مشاهده که در حالت i پایان می‌یابد، به دست می‌آید. با این مقدمه داریم:

$$d_{t+1}(i) = \max_i [d_t(i) a_{ij}] b_j(x_{t+1}) \quad (97)$$

برای بازیابی دنباله حالت، لازم است برای هر t و z آرگومانی که $d_t(i)$ را ماکسیمم می‌کند، نگه داریم. برای ذخیره آرگومان به آرایه‌ای $y_t(i)$ در الگوریتم نیاز است.

رویه استقراء برای الگوریتم ویتربی به صورت زیر می‌باشد:

مرحله اول: مقداردهی اولیه

$$\begin{aligned} d_t(i) &= p_i b_i(x_1) \quad 1 \leq i \leq N \text{ (states)} \\ y_1(i) &= 0 \end{aligned} \quad (98)$$

مرحله دوم: تکرار

$$d_t(j) = \max_{1 \leq i \leq N} [d_{t-1}(i) a_{ij}] b_j(x_t) \quad 2 \leq t \leq T, \quad 1 \leq j \leq N \quad (99)$$

$$y_t(j) = \text{Arg Max}_{1 \leq i \leq N} [d_{t-1}(i) a_{ij}] \quad (100)$$

مرحله سوم: اتمام

$$\text{بهترین امتیاز} = \max_{1 \leq i \leq N} [d_t(i)] \quad (101)$$

$$S_T^* = \text{Arg Max}_{1 \leq i \leq N} [B_T(i)] \quad (102)$$

مرحله چهارم: مرحله پسرو

$$S_t^* = y_{t+1}(S_{t+1}^*) \quad t = T-1, T-2, \dots, 1 \quad (103)$$

بهترین دنباله می‌باشد $S^* = (S_1^*, S_2^*, \dots, S_T^*)$

الگوریتم ویتربی در هر مرحله همه فرضیه‌ها را به طور موازی مورد پردازش قرار می‌دهد و به مرحله بعد

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 گروه کارپایه اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

بسط می‌دهد. بنابراین یک جستجوی همزمان ۳۵ محسوب می‌شود که قبل از رفتن به زمان $t+1$ همه محاسبات مربوط به زمان t را انجام می‌دهد. برای زمان t ، هر حالت با بهترین امتیاز از بین امتیازات همه حالات در زمان $t-1$ (به جای مجموع امتیازات همه مسیرهایی که وارد می‌شوند) به روز ۳۶ می‌شود. به همین دلیل است که معمولاً به این روش «جستجوی ویتربی همزمان» می‌گویند. وقتی به روز رسانی انجام می‌شود، اشاره گر backtracking نیز ذخیره می‌شود تا حالت ورودی با بیشترین احتمال را به خاطر داشته باشیم. در پایان جستجو، بهترین دنباله حالت را می‌توان با ردگیری این اشاره گر پیدا کرد. در الگوریتم ویتربی چون همه فرضیه‌ها متعلق به یک دنباله‌ی ورودی هستند و از لحاظ زمانی با هم برابرند، می‌توان آنها را به طور مستقیم با هم مقایسه کرد. این الگوریتم راه حلی بهینه برای مسایل پیچش زمانی غیرخطی بین مدل‌های HMM و مشاهدات آکوستیکی، تشخیص محدوده کلمه و تشخیص کلمه در بازشناسی گفتار پیوسته ارائه می‌کند.

3-8 جستجو در بازشناسی کلمات مجزا

در بازشناسی کلمات مجزا محدوده کلمات مشخص است. اگر مدل‌های HMM کلمات موجود باشد، احتمال مدل آکوستیک $P(X|W)$ را می‌توان با الگوریتم پیش‌رو³⁷ حساب کرد. جستجو به یک مسأله ساده تشخیص الگو تبدیل می‌شود و کلمه‌ای \hat{W} که بیشترین احتمال پیش‌رو را داشته باشد به عنوان کلمه بازشناسی شده انتخاب می‌شود. وقتی از مدل‌های زیرکلمه‌ای استفاده می‌شود، مدل‌های HMM کلمات با اتصال مدل‌های HMM زیرکلمه‌های متناظر بدست می‌آید.

4-8 جستجو در بازشناسی گفتار پیوسته

جستجو در بازشناسی گفتار پیوسته نسبتاً پیچیده است، حتی برای یک مجموعه لغات کوچک، چرا که الگوریتم جستجو بایستی احتمال شروع هر کلمه در هر فریم زمانی را در نظر بگیرد. در برخی از سیستم‌های بازشناسی گفتار پیوسته، برای بازشناسی گفتار پیوسته یک روش دو مرحله‌ای در نظر گرفته می‌شد، ابتدا محدوده‌های کلمات فرض می‌شد و سپس از تکنیک‌های انطباق الگو برای بازشناسی الگوهای قسمت شده استفاده می‌شد. اما این قسمت‌بندی و تعیین محدوده لغات چندان قابل اعتماد نبود.

در اینجا با مثال ساده‌ای (شکل ۹) نشان می‌دهیم که چگونه می‌توان روش‌های جستجوی کلمات مجزا

³⁵ Synchronous

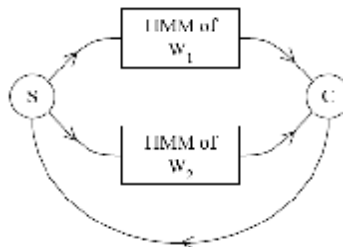
³⁶ update

³⁷ forward

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 گروه کا اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

را به بازشناسی گفتار پیوسته گسترش داد. این سیستم تنها شامل دو کلمه می‌باشد و فرض می‌کنیم که مدل زبانی استفاده شده در اینجا یک تک‌گرام یکسان³⁸ باشد $P(W_1)=P(W_2)=1/2$.

فراهم کردن ساختارهای زبان در همان چارچوب HMM اهمیت دارد. در شکل ۹ یک حالت شروع S و یک حالت جمع کننده³⁹ C اضافه کرده‌ایم. حالت شروع، یک انتقال پوچ⁴⁰ به حالت ابتدایی هریک از کلمات دارد که احتمال این انتقال‌ها براساس مدل زبانی مشخص می‌شود (در این مثال $1/2 =$). حالت خروجی هریک از کلمات نیز یک انتقال پوچ به حالت C دارد. همچنین حالت C یک انتقال پوچ به حالت شروع S دارد تا امکان تکرار وجود داشته باشد. مشابه با حالت جاسازی HMM های فونم‌ها در HMM کلمه برای بازشناسی کلمات مجزا، می‌توان HMM های کلمات W_1 و W_2 را در یک HMM جدید متناظر با ساختار شکل ۹ تعبیه نمود.



شکل ۹ مثال ساده‌ای از عملیات بازشناسی گفتار پیوسته با دو کلمه W_1 و W_2 . حالت S حالت شروع و حالت C یک جمع کننده حالت است تا پیوند کاملی بین لغات وجود داشته باشد.

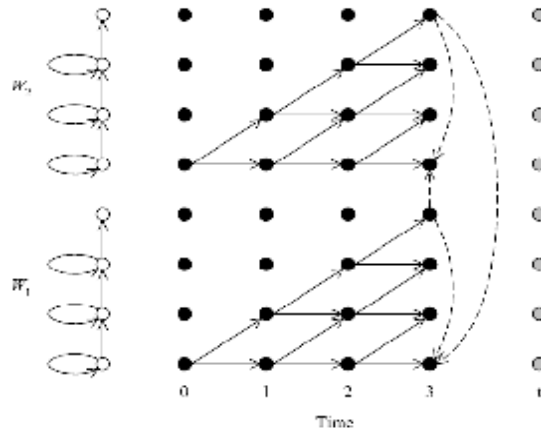
ساختار HMM ترکیبی که در شکل ۹ نشان داده شده است را می‌توان به عنوان یک شبکه آماری حالت نهایی⁴¹ با احتمالات انتقال و توزیعات خروجی در نظر گرفت. الگوریتم جستجو اساساً یک انطباق بین مشاهدات X و مسیر (در اینجا یعنی دنباله حالات و انتقالها) در شبکه حالت نهایی تولید می‌کند. برخلاف بازشناسی گفتار مجزا، در بازشناسی گفتار پیوسته لازم است دنباله کلمه بهینه \hat{W} مشخص شود. یک انتخاب بدیهی برای چنین کاری، الگوریتم ویتربی می‌باشد که دنباله حالت بهینه \hat{S} به دنباله کلمه بهینه \hat{W} متناظر می‌شود. شکل ۱۰ محاسبه شبکه ویتربی HMM را برای بازشناسی گفتار پیوسته دو کلمه شکل ۱۰ نشان می‌دهد.

³⁸ uniform unigram

³⁹ collector state

⁴⁰ null

⁴¹ stochastic finite state



شکل ۱۰ شبکه HMM برای مثال بازشناسی گفتار پیوسته در شکل ۸. وقتی حالت نهایی مدل یک کلمه ملاقات شد، یک کمان پوچ (که در شکل با خطوط منقطع مشخص شده است) از آن حالت به حالت ابتدایی همه کلمات پیوند می‌شود.

8-5 نقش اشاره‌گر backtracking

در این جا لازم است نقش اشاره‌گر backtracking در جستجوی ویتربی برای بازشناسی گفتار پیوسته روشن شود. بطور کلی در بازشناسی گفتار ما به دنبال یافتن بهترین دنباله حالت نمی‌باشیم^{۴۲} در عوض ما فقط می‌خواهیم دنباله کلمه بهینه که با معادله ۹۳ مشخص می‌شود را پیدا کنیم. بنابراین ما فقط به‌خاطر یادآوری سابقه کلمه برای مسیر جاری از اشاره‌گر backtrack استفاده می‌کنیم و به این ترتیب دنباله کلمه بهینه در پایان جستجو قابل کشف می‌باشد. برای توضیح بیشتر، وقتی ما به حالت پایانی یک کلمه رسیدیم، یک نود سابقه^{۴۲} ایجاد می‌کنیم که شامل شناسه کلمه و شاخص زمان جاری^{۴۳} می‌باشد و این نود به اشاره‌گر فعلی backtracking ضمیمه می‌شود. سپس این اشاره‌گر به نود قبلی‌اش پاس می‌شود و به همین ترتیب پیش می‌رود تا دنباله مشخص شود. یک مزیت جانبی برای نگهداری این اشاره‌گر این است که دیگر ما نیازی به نگهداری کل شبکه در زمان جستجو نداریم بلکه فقط به فضایی برای نگهداری دو قسمت زمانی متوالی (ستونها) در محاسبه شبکه نیاز داریم (تکه زمانی قبلی و تکه زمانی جاری) چرا که همه اطلاعات backtracking اکنون در اشاره‌گر backtracking نگهداشته می‌شود. این ویژگی مزیت ممتازی برای پیاده‌سازی یک جستجوی همزمان ویتربی می‌باشد.

^{۴۲} البته با وجود اینکه ما به دنبال یافتن دنباله حالت‌های بهینه در ASR نیستیم اما آنها برای تعیین قطعه‌بندی فونتیک کاملاً مورد استفاده می‌باشند و می‌توانند اطلاعات مهمی برای توسعه سیستم‌های ASR فراهم کنند.

^{۴۲} history

^{۴۳} time index

	عنوان پروژه:			
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی			
	عنوان زیر پروژه:			
تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		کد زیر پروژه: پیکرمتن فارسی - ۲ - خ	ویرایش: ۱/۰	تاریخ: ۱۳۸۸/۰۳/۱۹

9 مقاوم سازی بازشناسی گفتار

1-9 مقدمه

دقت سیستمهای بازشناسی گفتار به هنگام عمل در محیطهای آکوستیکی مختلف که با یکدیگر ناسازگارند، به شدت کاهش می‌یابد. این کاهش را می‌توان ناشی از عدم تطبیق میان داده‌های گفتار آموزشی و داده‌های آزمایش در محیطهای آکوستیکی مختلف دانست. این عدم تطبیق می‌تواند ناشی از چند عامل باشد: آلودگی با نویز (جمع‌پذیر، انعکاسی⁴⁴ و کانال)، تغییر شیوه صحبت (سرعت بیان گفتار و اثر لومبارد⁴⁵) و یا تغییر گوینده (تغییر کیفیت صدا، فرکانس اصلی، جنسیت و لهجه گوینده). از این رو در سالهای اخیر روشهای متعددی برای مقاوم سازی بازشناسی گفتار و کاهش عدم تطبیق میان شرایط آموزش و آزمایش مطرح گردیده‌اند که از آن جمله می‌توان به بهبود سیگنال گفتار⁴⁶، استفاده از ویژگیهای مقاوم⁴⁷ و نرمال سازی ویژگیها، استفاده از آرایه‌های میکروفنی، استفاده از خصوصیات شنوایی و گوش انسان و تطبیق و دگرگونی مدل گفتار اشاره کرد.

افزایش و بهبود دقت بازشناسی بوسیله روشهای فوق با محدودیتهایی مواجه است که بخشی از این محدودیتها به علت نارسایی مدل ریاضی است که برای مشخص کردن تخریب ویژگیهای آکوستیک بکار می‌روند. به این منظور باید مدلی ساده یافت که چگونگی تاثیر محیط را بر پارامترهای مشخصه سیستمهای بازشناسی گفتار و ورودی آنها نشان دهد.

روشهای پیشنهادی برای توصیف اثر تخریب محیط بر گفتار را می‌توان به دو دسته کلی تقسیم کرد: شیوه‌های مبتنی بر داده⁴⁸ و شیوه‌های مبتنی بر مدل⁴⁹. روشهای مبتنی بر داده سعی می‌کنند که تاثیر محیط بر مشخصه‌های گفتار را مشخص نمایند و خود سیگنال گفتار یا پارامترها و ویژگیها را بهبود بخشند و اصلاح نمایند که برخی از این روشها به این منظور نیازمند به وجود هر دو نوع سیگنال گفتار تمیز و نویزی هستند. شیوه‌های مبتنی بر مدل تلاش می‌کنند که مدل ریاضی محیط را اصلاح نمایند و

۴۴ Reverberation
 ۴۵ Lombard
 ۴۶ Speech enhancement
 ۴۷ Robust features
 ۴۸ Data-based methods
 ۴۹ Model-based methods

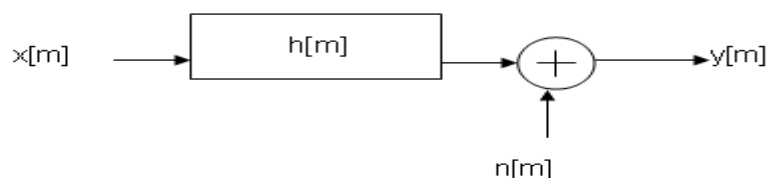
	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 گروه کا اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

آن را به نحوی تغییر دهند که با محیط تطبیق یابد. در این شیوه داده‌های مشاهده شده تغییر نمی‌کنند و فرض و تغییری در مورد سیگنال گفتار لازم نیست.

در این بخش پس از مروری کلی بر مهمترین روشهای مبتنی بر داده و مبتنی بر مدل که برای مقابله با تاثیر نویز بردقت سیستمهای بازشناسی گفتار، در سالهای اخیر بکار رفته‌اند، کارآیی برخی شیوه‌های مبتنی بر داده با یکدیگر مقایسه می‌شود و مدلی از تاثیر محیط بر گفتار ارائه می‌گردد. سپس شیوه‌های مبتنی بر داده و روشهای مبتنی بر مدل بررسی می‌شوند و در نهایت آزمایشها و نتایج کلی ارائه می‌شوند.

9-2 مدلی از تاثیر محیط بر سیگنال گفتار

در بسیاری از کاربردهای بازشناسی گفتار، دو منبع مهم نویز وجود دارد: نویز جمع‌پذیر و نویز فیلتر خطی. از جمله نویزهای جمع‌پذیر می‌توان به تاثیر جریان هوا و صداهای پس زمینه موجود در محیط اشاره کرد. نویزهای ناشی از کانال انتقال همانند خط تلفن یا اثرات میکروفن و یا نویزهای ناشی از انعکاس و بازتابهای سطحی صدا از اجسام نیز از جمله نویزهای فیلتر خطی محسوب می‌شوند. نویز موجود بر روی خط تلفن نیز از این نوع است. اگر چه انواع دیگری از نویز نیز وجود دارد، اما اثر محیط را معمولا با این دو نوع نویز مدل می‌کنند.



شکل ۱۱ مدل تاثیر نویز بر گفتار

در حالتی که نویز به این دو نوع محدود باشد، می‌توان مدلی همانند شکل ۱۱ برای اثر تخریب محیط ارائه کرد [۱۳]. در این مدل سیگنال تمیز $x[m]$ ابتدا از فیلتری با تابع تبدیل $h[m]$ عبور می‌کند و سپس خروجی آن با نویز جمع‌پذیر $n[m]$ که مستقل از $x[m]$ است، جمع می‌گردد و سیگنال نویزی $y[m]$ را شکل می‌دهد. رابطه زیر میان چگالی طیف توانهای سیگنال تمیز، سیگنال نویزی و نویز برقرار است:

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 ژورنال اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

$$P_y(w) = P_x(w) |H(w)|^2 + P_n(w) \quad (104)$$

این رابطه پس از اعمال لگاریتم در حوزه کپسترال به صورت زیر در می آید:

$$y = x + h + \log(1 - e^{n-x-h}) = x + h + s(x, n, h) \quad (105)$$

که در این رابطه y, h, n, x به ترتیب بیانگر $P_y(w), |H(w)|^2, P_n(w), P_x(w)$ می باشند. این رابطه را می توان به این صورت بازنویسی کرد:

$$y = x + f(x, n, h) \quad (106)$$

در این رابطه y را می توان بیانگر ویژگیها در محیط نویزی و x را نشانگر ویژگیها در محیط تمیز دانست و $f(x, n, h)$ که تابع محیط نامیده می شود، اثر نویزهای جمع پذیر و فیلتر خطی را بر بردارهای ویژگی در محیط تمیز نشان می دهد. یکی از روشها در جبران تاثیر تخریب محیط آن است که با عمل در فضای ویژگیها و با در دست داشتن ویژگیهای نویزی y ، تخمینی از ویژگیهای تمیز x بدست آید.

3-9 روشهای مبتنی بر داده

چنانکه گفته شد روشهای مبتنی بر داده تلاش می کنند که تاثیر محیط بر مشخصه های سیگنال گفتار را از طریق تاثیر آنها بر خود سیگنال و یا ویژگیها مشخص نمایند و این تاثیر را به نحوه ای بهبود بخشند و اصلاح نمایند که بتوان بازشناسی را تا حد امکان مستقل از تاثیر محیط انجام داد. این روشها را می توان به چند دسته کلی تقسیم کرد. اولین گروه روشهایی هستند که به صورت مستقیم بر روی سیگنال گفتار نویزی عمل می کنند و سیگنال گفتار تمیز را از سیگنال گفتار نویزی تخمین می زنند. از جمله شیوه های این گروه، می توان به شیوه تفاضل طیف و فیلتر کردن گفتار نویزی اشاره کرد [۲۸، ۱۳].

گروه دوم روشهای استخراج ویژگیهای مقاوم هستند. در این روشها در اولین مرحله پردازش سیگنال گفتار تبدیلی بکار برده می شود تا بتوان بردارهای ویژگی مشابهی برای نمونه های گفتار آموزشی و آزمایشی بدست آورد و تفاوت آنها را کمتر نمود. از جمله این روشها می توان از تفاضل میانگین ضرایب کپسترال⁵⁰ (CMS)، [۱۳] آنالیز پیشگویی خطی و درکی گفتار⁵¹ (PLP) [۲۶] و اعمال فیلتر

⁵⁰ Cepstral Mean Subtraction

⁵¹ Perceptual Linear Predictive analysis of speech

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 شورای ملی اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

RASTA بر آن [۲۲] و آنالیز ZCPA⁵² [۱۵] نام برد.

روشهای دیگری نیز وجود دارند که تلاش می‌کنند از خصوصیات شنوایی و گوش انسان در برخورد با نویز محیط استفاده کنند. از جمله این روشها استفاده از تکنیک پوشش نویز⁵³ است که برای آن الگوریتمهای متفاوتی ارائه گردیده است [۱۳].

قسمتهای بعدی این بخش به بررسی برخی از روشهای گروههای فوق که کاربرد بیشتری دارند و در کار حاضر مورد آزمایش قرار گرفته‌اند، می‌پردازد.

۹-۳-۱ تفاضل میانگین ضرایب کپسترال

تفاضل میانگین ضرایب کپسترال که در گروه ویژگیهای مقاوم جای می‌گیرد، یکی از ساده ترین و موثرترین الگوریتمهای موجود است که غالباً در بازشناسی گفتار با کتاب لغت وسیع بکار گرفته می‌شود. این الگوریتم متوسط بردارهای ویژگی کپسترال را برای یک نمونه از سیگنال گفتار محاسبه می‌کند و سپس این مقدار میانگین را از هر یک از بردارها کم می‌کند. این کار سبب می‌شود که تغییرات داده‌ها و ویژگیها کاهش یابد و نوعی نرمال سازی انجام گیرد. از این رو این شیوه را نرمال سازی ضرایب کپسترال با استفاده میانگین آنها⁵⁴ (CMN) نیز می‌نامند. اگر یک نویز از نوع ضرب شونده فرکانسی⁵⁵ نامتغیر با زمان موجود باشد، این مجموعه بردارهای ویژگی از تغییرات نویز تأثیر نمی‌پذیرند. در مورد نویز جمع‌پذیر نیز کم کردن این میانگین به مقاوم‌سازی سیستم کمک می‌کند. اما اگر هر دو نوع نویز مذکور به همراه هم وجود داشته باشند، به سختی می‌توان رفتار سیستم را پیش‌بینی کرد [۱۳، ۱۸].

۹-۳-۲ نرمال سازی ضرایب کپسترال با استفاده از میانگین و واریانس

در این روش که در گروه ویژگیهای مقاوم جای می‌گیرد، با توجه به تأثیر نویز بر میانگین و واریانس ضرایب کپسترال، نوعی نرمال‌سازی انجام می‌شود که در آن بردار ویژگی نه تنها از میانگین کل بردارهای ویژگی کم می‌شود، بلکه بر واریانس بردارهای ویژگی نیز تقسیم می‌شود. در این حالت کم کردن

⁵² Zero Crossing with Peak Amplitude

⁵³ Noise Masking

⁵⁴ Cepstral Mean Normalization

⁵⁵ Convolution

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 ژورنال مطالعات زبانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

میانگین همانند یک فیلتر کردن بالاگذر خطی و تقسیم بر واریانس همچون کنترل بهره خودکار⁵⁶ تعبیر می‌شود [۱۶،۱۷].

برای این کار ابتدا یک بافر با امکان ذخیره‌سازی N بردار کپسترال در نظر گرفته می‌شود و در هر لحظه ضرایب نرمال شده با توجه به بردارهای ویژگی که در داخل بافر قرار گرفته‌اند، محاسبه می‌شوند. هنگامی که نیمی از بافر با بردارهای کپسترال پر شود، ضرایب نرمال شده محاسبه می‌شوند و اولین بردار بافر نرمال سازی می‌گردد و بردار جدید در بافر جایگزین می‌شود و پس از آن بردار دوم بافر نرمال سازی می‌گردد و این روند تا پر شدن بافر با بردارهای نرمال شده ادامه می‌یابد و سپس بافر خالی شده و N بردار بعدی در بافر قرار می‌گیرند [۱۷].

آزمایشهای عملی نشان داده است که طول بافر (N) باید به نحوی انتخاب شود که در حدود ۴۰ میلی ثانیه از گفتار را پوشش دهد. به این ترتیب کارایی مناسب به دست می‌آید. علاوه بر اینکه اگر طول بافر از حدی بیشتر انتخاب شود، درصد بازشناسی از حد مشخصی فراتر نخواهد رفت و اصطلاحاً منحنی بازشناسی گفتار اشباع می‌شود [۱۶].

۳-۳-۹ آنالیز پیشگویی خطی و درکی گفتار و فیلتر RASTA

آنالیز پیشگویی خطی و درکی گفتار که به آنالیز PLP معروف است، از این ایده بهره می‌گیرد که قدرت تفکیک فرکانسی و حساسیت نسبت به تغییر فرکانس، در گوش انسان در فرکانسهای مختلف، یکسان نیست و حساسیت گوش نسبت به شدت و انرژی صوت نیز در فرکانسهای مختلف متفاوت است [۲۶]. در روش PLP ابتدا طیف توان زمان کوتاه سیگنال گفتار با استفاده از تبدیل فوریه محاسبه می‌گردد و سپس طیف میان باندهای بحرانی⁵⁷ تقسیم می‌شود. این تقسیم بر اساس معیار بارک⁵⁸ صورت می‌گیرد. پس از این مرحله، عمل کانولوشن طیف توان توزیع شده در باندهای بحرانی، با تابعی موسوم به تابع پوشش باندهای بحرانی⁵⁹ صورت می‌گیرد. این تابع مدلی از منحنی نامتقارن پوشش گوش در هر یک از باندهای بحرانی را ارائه می‌کند. حاصل این مرحله در تابع برابرسازی بلندی⁶⁰ ضرب می‌شود. این تابع تقریبی از نحوه حساسیت گوش نسبت به فرکانسهای مختلف را مدل می‌کند و حساسیت شنوایی را در

⁵⁶ Automatic Gain Control

⁵⁷ Critical bands

⁵⁸ Bark

⁵⁹ Critical-band masking curve

⁶⁰ Equal-loudness

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 مؤسسه ملی اطلاع‌رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان‌سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

سطح 40 dB نیز شبیه‌سازی می‌نماید. مرحله بعدی در عملیات PLP، اعمال قانون ریشه سوم شنوایی می‌باشد. این قانون بیان می‌دارد که در گوش انسان میزان احساس بلندی صوت با ریشه سوم انرژی آن متناسب است یعنی در گوش انسان برای تبدیل انرژی یا توان صوت به بلندی صوت، انرژی صوت به توان $0/33$ رسانده می‌شود. در مرحله آخر پس از اعمال تبدیل فوریه معکوس برحاصل، با اعمال یک مدل تمام قطب، ضرایب PLP محاسبه می‌شوند [۲۶].

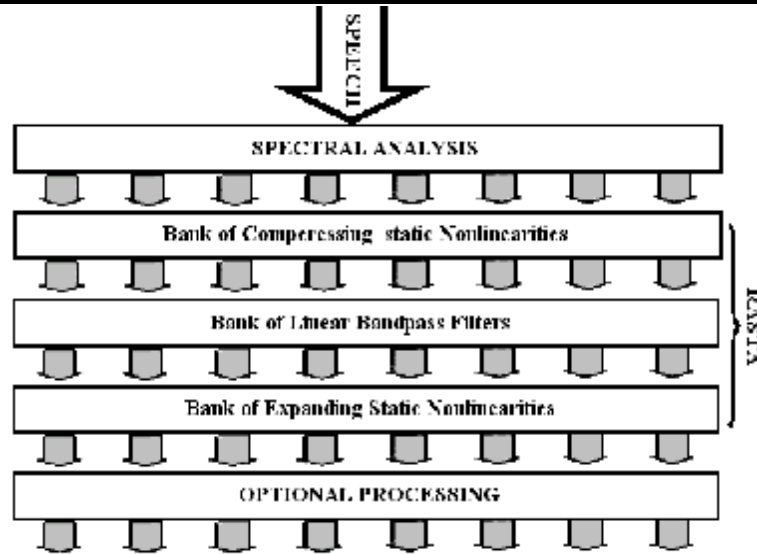
روش PLP بر اساس محاسبه طیف زمان کوتاه گفتار صورت می‌گیرد و بنابراین. مانند دیگر تکنیک‌های مبتنی بر طیف زمان کوتاه، به دلیل تاثیر کانال‌های مخابراتی بر مقادیر طیف زمان کوتاه، آسیب‌پذیر است. بنابراین ضرایب PLP ویژگی‌های مقاومی محسوب نمی‌شوند. بنابراین باید ضرایب PLP را در مقابل اثرات تخریبی کانال ارتباطی مقاوم نمود. بکارگیری فیلتر RASTA از جمله شیوه‌هایی است که این مقاوم‌سازی را انجام می‌دهد. نام این فیلتر از واژه طیف نسبی⁶¹ گرفته شده است [۲۲].

آنالیز طیف نسبی با این هدف انجام می‌شود که تغییرات آهسته یا تندی که در خارج از محدوده زبان شناسی سیگنال گفتار قرار می‌گیرند، از این سیگنال حذف گردند. در واقع فیلتر RASTA براساس عملکرد گوش انسان عمل می‌کند. گوش انسان نسبت به مؤلفه‌های فرکانسی که شدت آنها در طول زمان با فرکانس حدود 4 kHz تغییر می‌کنند، در مقایسه با مؤلفه‌هایی که شدت آنها در طول زمان سریعتر یا آهسته‌تر از 4 kHz تغییر می‌کنند، حساسیت بیشتری دارد.

مراحل کلی آنالیز طیف نسبی مطابق شکل ۱۲ است. در عملیات RASTA ابتدا یک تبدیل غیرخطی فشرده‌کننده بر طیف توان گفتار اعمال می‌شود. سپس یک فیلتر میان‌گذر بر طیف توان تبدیل یافته اثر می‌کند. در نهایت یک عمل تبدیل منبسط‌کننده که معکوس عمل فشرده‌سازی است، بر نتیجه اعمال می‌شود. مجموع عملیات فشرده‌سازی طیف، اعمال فیلتر و منبسط‌سازی طیف، آنالیز طیف نسبی یا RASTA نامیده می‌شود [۵].

^{۱۵} Relative Spectra

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 گروه تحقیقات اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

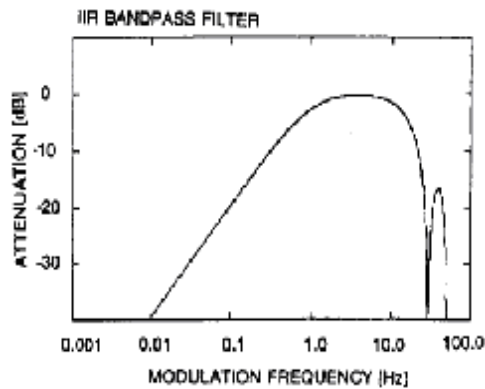


شکل ۱۲ مراحل کلی آنالیز طیف نسبی

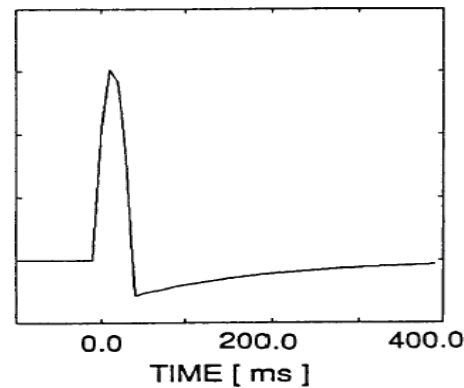
پس از آنالیز طیف نسبی، هر آنالیز دیگری مانند PLP می تواند انجام شود. اگر ضرایب PLP پس از آنالیز RASTA بدست آید، ضرایب حاصل ویژگیهای مقاومی هستند که ضرایب RASTA-PLP نامیده می شوند.

فیلتر بکار گرفته شده در آنالیز طیف نسبی می تواند شکل‌های مختلفی داشته باشد. شکل ۱۳ پاسخ فرکانسی و پاسخ ضربه یک فیلتر RASTA را نشان می دهد [۲۲]. تبدیل فشرده ساز در آنالیز RASTA می تواند به صورت $y = \ln(n)$ باشد. در این حالت فیلتر در حوزه لگاریتم طیف عمل می کند و آنالیز RASTA را log-RASTA می نامند.

هنگامی که عملیات RASTA در حوزه لگاریتم صورت گیرد، نویز فیلتر خطی یا تاثیر کانال که در حوزه لگاریتم طیف به صورت مولفه‌های طیفی جمع پذیر ظاهر می شود، به طور مؤثر حذف می شود. اما در مورد نویز جمع پذیر غیر وابسته به سیگنال، پس از اعمال تبدیل غیر خطی لگاریتم، نویز وابسته به سیگنال می شود و به طور مؤثر با آنالیز RASTA حذف نخواهد شد. جهت حل این مسأله، استفاده از تبدیل فشرده ساز $y = \ln(1+Jx)$ به جای $y = \ln(x)$ پیشنهاد شده است که در آن J یک مقدار ثابت مثبت وابسته به سیگنال است. در صورتی که J بسیار کوچکتر از یک اختیار شود، تبدیل حاصله یک تبدیل شبه خطی است.



(ب) پاسخ فرکانسی



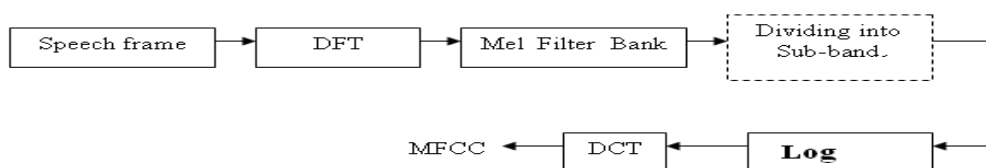
(الف) پاسخ ضربه

شکل ۱۳ فیلتر RASTA.


در صورتی که مقدار J بسیار بزرگتر از یک باشد، تبدیل شبه لگاریتمی خواهد بود. عملیات منبسط کننده متناظر با فشرده ساز به صورت $x = (e^y - 1)/J$ می باشد. به آنالیز RASTA با چنین فشرده ساز و منبسط کننده ای lin-log-RASTA گفته می شود و گاه آن را J-RASTA نیز می نامند [۲۲]. فیلتر J-RASTA کارآیی یک سیستم بازشناس را هم در حضور نویز جمع شونده و هم در حضور نویز فیلتر خطی، افزایش می دهد [۱۳].

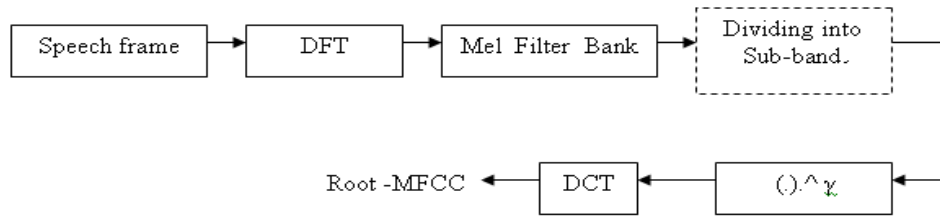
۹-۳-۴ ضرایب کپسترال ریشه ای

استفاده از دیگر انواع توابع فشرده سازی همانند تابع ریشه بجای تابع لگاریتم در استخراج ضرایب مل کپستروم، روش دیگری برای مقابله با حضور نویز به ویژه در حضور نویز جمع پذیر است [۲۳]. به این ترتیب میزان فشرده سازی انرژی زیرباندهای مل در مقایسه با یکدیگر و در حضور نویز، تغییر می کند و موثرتر از تابع لگاریتم می گردد. تفاوت ویژگیهای مل کپسترال ریشه ای و ویژگیهای مل کپسترال عادی در شکل ۱۴ قابل رویت است.



(الف) روش قراردادی استخراج ضرایب مل کپسترال (MFCC)

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 گروه کارشناسی اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	



(ب) ضرایب کپسترال ریشه ای (RMFCC) جاییکه $1 < \gamma < -1$

شکل ۱۴ روشهای استخراج MFCC و RMFCC

۹-۳-۵ ضرایب خود همبستگی فاز^{۶۲} (PAC)

تبدیل فوریه دنباله خودهمبستگی یک سیگنال در حوزه زمان، معادل با طیف توان آن سیگنال در حوزه فرکانس است. دنباله خود همبستگی را می توان به صورت یک ضرب داخلی مابین بردارهای گفتاری نیز در نظر گرفت. اخیراً معیار دیگری برای خود همبستگی بنام خود همبستگی فاز (PAC) مطرح شده است که از زاویه میان بردارهای گفتار به عنوان معیار همبستگی استفاده می کند [۹]. ایده استفاده از زاویه بر این اساس شکل گرفته است که در یک ضرب داخلی بین بردارها، زاویه کمتر از اندازه از نویز تاثیر می پذیرد [۲۷]. ضرایب خود همبستگی فاز را می توان به صورت زیر محاسبه کرد:

$$P[k] = q_k = \text{Arc cos} \left(\frac{R[k]}{|x|^2} \right) \quad (107)$$

در این رابطه، $R[k]$ و $|x|^2$ به ترتیب بیانگر ضرایب خود همبستگی عادی و انرژی یک قاب از سیگنال هستند. طیف فرکانسی محاسبه شده از ضرایب PAC بنام طیف خود همبستگی فاز یا طیف PAC نامیده می شود. مشابه با ویژگیهایی که از طیف عادی استخراج می شوند، ویژگیها می توانند از طیف خود همبستگی فاز نیز استخراج شوند. ویژگیهای MFCC که از طیف خود همبستگی فاز استخراج می شوند، بنام PAC-MFCC شناخته می شوند. نتایج آزمایشها در [۶] و [۹] نشان داده اند که ویژگیهای PAC-MFCC در حضور نویز مقاوم هستند، اما کارایی بازشناسی گفتار تمیز را کاهش می دهند.

^{۶۰} Phase AutoCorrelation

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 گروه کا اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

۹-۳-۶ ویژگیهای مبتنی بر تاخیر گروه^{۶۳}

تابع تاخیر گروه به صورت منفی مشتق فاز سیگنال صحبت تعریف می شود. این تابع را می توان اینگونه محاسبه کرد [۸،۱۰،۲۴]:

$$t_p(w) = - \frac{X_R(w)Y_R(w) + X_I(w)Y_I(w)}{|X(w)|^2} \quad (10.8)$$

در این رابطه اندیسهای I و R بیانگر بخشهای موهومی و حقیقی هستند و $X(\omega)$ و $Y(\omega)$ به ترتیب نشانگر تبدیل فوریه دنباله های $x(n)$ و $n.x(n)$ هستند. در [۸]، تابع اصلاح شده ای بر اساس رابطه (۱۰۶) برای محاسبه تاخیر گروه پیشنهاد شده است که مخرج این رابطه را حذف کرده است. زیرا مخرج این رابطه سبب ایجاد پرشهای ناگهانی^{۶۴} در مقدار تابع تاخیر گروه می شود که مطلوب نیست. رابطه جدید موسوم به طیف ضربی^{۶۵} است که به صورت زیر تعریف می شود:

$$Q(w) = |X(w)|^2 t_p(w) = X_R(w)Y_R(w) + X_I(w)Y_I(w) \quad (10.9)$$

در کار حاضر، ما ویژگیهای MFCC را از این طیف ضربی استخراج می کنیم. آزمایشها نشان داده است که این ویژگیها نسبت به حضور نویز جمع پذیر برای بازشناسی کلمات و واجها مقاوم است [۸،۵]. این ویژگیها را در کار حاضر با نام GDF مشخص می نماییم.

4-9 روشهای مبتنی بر مدل

چنانکه گفته شد شیوههای مبتنی بر مدل تلاش می کنند که بجای اصلاح سیگنال یا پارامترها، مدل آکوستیک محیط را در مرحله بازشناسی اصلاح نمایند. مزیت این روشها آن است که در آنها دادههای مشاهده شده تغییر نمی کنند و هیچ نوع فرض یا تصمیم گیری قبلی درباره سیگنال گفتار ضروری نیست و از این رو نیازمند پایگاه داده استریو نیستند. برخی از روشهای مبتنی بر مدل عبارتند از: تجزیه مدل مخفی مارکوف^{۶۶} [۲۵]، ترکیب موازی مدلهای^{۶۷} (PMC) [۱۸] و بازگشت خطی با بیشترین شباهت^{۶۸}

^{۵۶} Group Delay

^{۵۷} Spikes

^{۵۸} Product Spectrum

^{۶۶} HMM decomposition

^{۶۷} Parallel Model Combination

^{۶۸} Maximum Likelihood regression

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 ژورنال اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

(MLLR) [۱۴].

در تجزیه مدل مخفی مارکوف، گفتار تمیز و نویز به صورت جداگانه و به عنوان دنباله‌ای از حالات مدل مخفی مارکوف مدل می‌شوند. مدل‌های گفتار و نویز در یک مدل ادغام می‌شوند و به این ترتیب امکان انتقال از حالات گفتار به حالات نویز و بر عکس به وجود می‌آید. به این ترتیب، به هنگام بازشناسی الگوریتم ویتربی محتمل‌ترین مسیر را در یک فضای گسترده جستجو می‌کند. این شیوه نه تنها برای سیگنال و نویز، بلکه برای بازشناسی دو سیگنال همزمان با یکدیگر نیز کاربرد دارد [۲۵]. ایراد این شیوه، هزینه بالای محاسباتی جستجو در فضای گسترده است.

شیوه ترکیب موازی مدلها تلاش می‌کند تا تاثیر محیط بر توزیع کپستروم گفتار تمیز را در مدل محیط منعکس نماید و با فرض وجود دانش کامل از نویز، بردار میانگین و ماتریس کوواریانس توزیعهای گاوسی مدل مخفی مارکوف را به نحوی تغییر دهد که به توزیعهای ایده‌آل کپستروم گفتار نویزی نزدیک شوند. چندین شیوه مختلف برای تغییر بردارهای میانگین و ماتریسهای کوواریانس در الگوریتم ترکیب موازی مدلها موجود است [۱۸].

به هر حال، در تمامی فرمهای الگوریتم PMC، به دانش قبلی از نویز و بردارهای کانال نیاز است و نمونه‌های مجزا از نویز باید وجود داشته باشند تا بتوان پارامترهای PMC را تخمین زد [۱۹].

روش MLLR در اصل برای تطبیق گوینده طراحی شده است، اما برای جبران اثر محیط و نویز نیز می‌تواند بکار رود. این الگوریتم بردارهای میانگین و ماتریسهای کوواریانس گفتار تمیز مدل شده بوسیله مدل مخفی مارکوف را تغییر می‌دهد و برای این تغییر از داده‌های آموزشی کمی استفاده می‌کند. این روش در نهایت یک ماتریس تبدیل می‌یابد که احتمال مشاهده بردار کپستروم نویزی در مدل را بیشینه کند [۱۴].

روش MLLR از یک مدل صریح از محیط استفاده نمی‌کند، بلکه تنها فرض می‌کند که بردارهای میانگین توزیع کپستروم گفتار تمیز، تحت تاثیر محیط دوران یافته و انتقال می‌یابند [۱۹].

۹-۴-۱ معیار تصویردهی وزن دار

تئوری معیار تصویردهی وزن دار (WPM) بر اساس این مشاهده *Juang* و *Mansour* استوار است که اندازه بردارهای ویژگی کپسترال در حضور نویز جمع پذیر سفید کاهش می‌یابد [۲۷]. طبق این خصوصیت در [۲۷] یک معیار محاسباتی بر پایه عمل تصویر کردن معرفی شد که کارایی بازشناسی گفتار بوسیله روش DTW را در حضور نویز به طرز قابل توجهی بهبود بخشید. *Celemnnt* و *Carlson* این معیار تصویردهی را گسترش دادند و آن را در یک سیستم بازشناسی مبتنی بر مدل مخفی مارکوف با چگالی پیوسته

	عنوان پروژه:		 گروه کار اطلاع رسانی
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		
	عنوان زیر پروژه:		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	کد زیر پروژه: پیک متن فارسی - ۲ - خ

(CDHMM) بکار گرفتند. آنها یک عامل مقیاس را در توزیع حالات CDHMM و به عبارت دیگر در توزیع گاوسی احتمالات شباهت دخالت دادند تا به این وسیله کاهش اندازه بردار کیسترال جبران گردد. ویژگیهای مورد استفاده آنان ضرایب مل کیستروم (MFCC) بود. نتایج در [۲۱،۵،۷،۲۹] نشان داده است که بکار بردن این معیار تصویردهی وزن دار، نرخ بازشناسی کلمات مجزا را به طرز قابل توجهی در حضور انواع مختلف نویز از جمله نویز سفید و نویز رنگی بهبود می بخشد. رابطه جبران تعریف شده آنها برای توزیع گاوسی را می توان به صورت زیر بیان کرد:

$$b_{j,i}(c_t) = N(c_t, I_{j,i,t}, m_{j,i}, C_{j,i}) = \frac{\exp\left(-\frac{1}{2}(c_t - I_{j,i,t} m_{j,i})^T C_{j,i}^{-1} (c_t - I_{j,i,t} m_{j,i})\right)}{(2p)^{\frac{n}{2}} |C_{j,i}|^{\frac{1}{2}}} \quad (110)$$

که پارامترهای موجود در آن این گونه تعریف می شوند:

c_t : بردار مشاهده کیسترال برای قاب t

n : بعد بردار مشاهده

$m_{j,i}$: بردار میانگین مخلوط گاوسی زام در حالت i

$C_{j,i}$: ماتریس کوواریانس مخلوط گاوسی زام در حالت i

$I_{j,i,t}$: عامل مقیاس برای قاب t در مخلوط گاوسی زام در حالت i

$b_{j,i}(c_t)$: احتمال تولید شده برای بردار مشاهده c_t توسط مخلوط گاوسی زام در حالت i

با گرفتن لگاریتم از رابطه (۱۰۵) و محاسبه لگاریتم احتمال شباهت، می توان به رابطه ای برای الگوریتم ویتربی دست یافت که بیانگر معیار تطبیق میان بردار مشاهده و بردارهای میانگین مخلوطهای گاوسی است. این رابطه را می توان اینگونه نوشت:

$$\log b_{j,i}(c_t) = (c_t - I_{j,i,t} m_{j,i})^T C_{j,i}^{-1} (c_t - I_{j,i,t} m_{j,i}) + \log |C_{j,i}| + N \log(2p) \quad (111)$$

با گرفتن مشتق از رابطه (۱۱۱) نسبت به $I_{j,i,t}$ می توان مقدار بهینه ای برای $I_{j,i,t}$ بدست آورد که به واسطه آن احتمال مشاهده $b_{j,i}(c_t)$ بیشینه شود. این مقدار بهینه به صورت زیر بدست می آید:

$$I_{j,i,t} = \frac{c_t^T C_{j,i}^{-1} m_{j,i}}{m_{j,i}^T C_{j,i}^{-1} m_{j,i}} \quad (112)$$

با جایگذاری این مقدار برای $I_{j,i,t}$ در رابطه (۱۱۱)، مقداری برای لگاریتم احتمال شباهت بدست می آید

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 مؤسسه ملی اطلاع‌رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان‌سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

که *WPM* نامیده می‌شود. رابطه (۱۱۱) و (۱۱۲) برای تمامی مخلوطهای گاوسی موجود در یک حالت مدل مخفی مارکف بکار می‌رود و احتمال شباهت نهایی با استفاده از مخلوطهای گاوسی تطبیق یافته با معیار تصویردهی وزن دار محاسبه می‌شود.

5-9 شیوه‌هایی برای تطبیق با گوینده جدید

یکی از دلایل عدم تطبیق میان داده‌های گفتار آموزشی و داده‌های آزمایش، تغییر گوینده یا به عبارتی شناسایی کلام یک گوینده جدید است. برای آنکه سیستم بتواند با گوینده جدید تطبیق یابد، از شیوه‌های تطبیق گوینده استفاده می‌شود تا بتوان با در دست داشتن کمترین نمونه‌های گفتاری گوینده، گفتار او را بازشناسی کرد. در این راستا، شیوه‌های متعددی مطرح شده‌اند، همانند: نرمال سازی طول جهاز صوتی⁶⁹ (VTN)، خوشه‌بندی گوینده⁷⁰، روش MLLR و انواع اصلاح شده آن، شیوه تخمین بیشینه پارامترهای پسین⁷¹ (MAP) و انواع تغییر یافته آن [۱۲،۱۴].

شیوه نرمال سازی جهاز صوتی بر اساس این ایده استوار است که طول جهاز صوتی در گویندگان مختلف متفاوت است و این امر بر فرکانسهای فرمانت تاثیر می‌گذارد و سبب تغییراتی می‌شود و بنابراین باید داده‌ها بر اساس یک مقیاس فرکانسی خطی نرمال سازی شوند تا تاثیر تغییر طول جهاز صوتی در فرکانسهای فرمانت به حساب آورده شود و به این ترتیب به نوعی سیستم با تغییر گوینده تطبیق یابد [۱۲،۱۴].

در خوشه‌بندی گوینده، تلاش می‌شود که گروهی از گویندگان یافت شوند که از لحاظ آکوستیکی به گوینده جدید نزدیک‌ترند و از ویژگیهای این گویندگان برای تخمین و اصلاح مجدد پارامترهای مدل گفتار استفاده می‌شود. ساده‌ترین شیوه این خوشه بندی بر اساس جنسیت گوینده است و توجه به پارامترهایی همانند فرکانس اصلی که در دو جنس متفاوت است [۱۴].

شیوه MLLR نیز چنانچه گفته شد، عمل می‌کند اما این بار ماتریس تبدیل را چنان می‌یابد که احتمال مشاهده بردارهای کپستروم مربوط به گوینده جدید را بیشینه کند.

شیوه MAP نیز یک تخمین پسین برای پارامترهای مدل در نظر می‌گیرد و با استفاده از این تخمین در مرحله آموزش مدل مخفی مارکوف تغییری ایجاد می‌کند تا بتواند با داده‌های آموزشی کم مدل را با

⁶⁹ Vocal Tract Normalization

⁷⁰ Speaker Clustering

⁷¹ Maximum A Posteriori

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 مؤسسه ملی اطلاع‌رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان‌سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

گوینده جدید تطبیق دهد [۱۲،۱۴].

علاوه بر ابزارهای پردازش و بازشناسی گفتار که برای سیستم مورد نظر پیاده‌سازی شده‌اند، ابزارهای آماده دیگری نیز برای اهداف کاربردی و تحقیقاتی در مراکز تحقیقاتی جهان تهیه شده‌اند که استفاده از آنها برای همگان آزاد است و اطلاعات و منابع مربوط به آنها به صورت آزاد بر روی شبکه جهانی موجود است و برای کاربران آنها لیستهای پستی موجود است. از جمله این ابزارها که کار بر روی آن بیش از یک دهه است ادامه دارد، HTK⁷² است. HTK ابزاری برای ساخت مدل مخفی مارکوف و پردازش سیگنال گفتار بوسیله آن است که در دانشگاه کمبریج انگلستان تولید شده و توسعه یافته است. علاوه بر این امکانات دیگر پردازش گفتار نظیر استخراج ویژگیها، گرامر زبان و تطبیق گوینده نیز در آن وجود دارد. از شیوه‌های تطبیق گوینده، روشهای MAP و MLLR در HTK پیاده‌سازی گردیده‌اند [۱۲].

در این گزارش ابتدا مقدمه‌ای در مورد سیستمهای بازشناسی گفتار و روشهای متداول آن ارائه شد. سپس بحث پیش‌پردازش و بررسی ویژگیهای بکار رفته در بازشناسی بررسی شد. در ادامه مدل مخفی مارکوف بطور کامل تشریح شد و آنرا به عنوان روشی برای بازشناسی تحلیل کردیم. در فصل بعدی مسائل مهم در پیاده‌سازی عملی سیستم بازشناسی گفتار دیده شدند. در ادامه نیز روشهای جستجوی معمولی، مؤلفه‌های فضای جستجو و روشهای جستجوی ویتربی بکار رفته در سیستم بازشناسی گفتار ارائه شدند. در فصل آخر متدهای مختلف مقاوم سازی بازشناسی گفتار و نتایج بدست آمده از این متدها آورده شد.

سیستم بازشناسی گفتار به عنوان یکی از سرویسهای بکار رفته در پروژه حاضر از اهمیت خاصی برخوردار است که بدلیل بحثهای پردازش سیگنال به صورت گزارش مستقل آورده شد.

⁷²Hidden Markov Model Toolkit

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 ژورنال اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

10 تطبیق گوینده

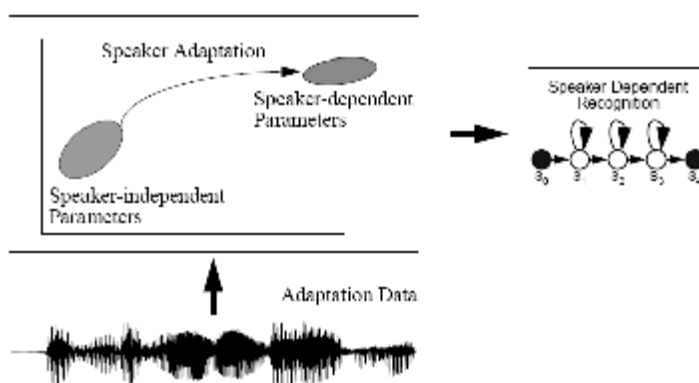
10-1 مقدمه

مبحث تطبیق گوینده هم برای سیستمهای بازشناسی گفتار و هم برای سیستمهای سنتز مطرح است. تطبیق گوینده در سیستمهای بازشناسی گفتار به منظور افزایش راندمان بازشناسی و در سیستمهای سنتز برای تغییر و تطبیق صدای سنتز شده با صدای فرد مورد نظر صورت می‌گیرد.

10-2 تطبیق گوینده برای سیستمهای بازشناسی گفتار


سیستمهای بازشناسی گفتار می‌توانند از نوع وابسته به گوینده (SD^{73}) یا مستقل از گوینده (SI^{74}) باشند. سیستمهای وابسته به گوینده معمولا دارای راندمان بالاتری برای گوینده‌ای که بازنه داده‌های گفتاری او آموزش دیده‌اند می‌باشند. چنانچه سیستم بازشناسی موجود از نوع مستقل از گوینده باشد با استفاده از روشهای تطبیق گوینده و به کمک داده‌های گفتاری یک گوینده خاص سیستم را به سیستم وابسته به گوینده تبدیل و راندمان آنرا افزایش داد.

ایده اصلی تطبیق گوینده استفاده از مقدار کمی اطلاعات تطبیقی به منظور تغییر سیستم بازشناسی است به طوری که تا حد ممکن بتواند اطلاعات مربوط به گوینده جدید را مدل کند. شکل ۱۵ عمل تطبیق گوینده جدید با مدل‌های موجود در سیستم تشخیص با استفاده از داده‌های تطبیقی کم را نشان می‌دهد.



⁷³ Speaker Dependent

⁷⁴ Speaker Independent

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 مؤسسه ملی اطلاع‌رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

شکل ۱۵ نمایش سطح بالای عمل تطبیق گوینده

با توجه به اینکه سیستم بازشناسی استفاده شده مبتنی بر HMM است روشهای تطبیق گوینده مبتنی بر HMM مورد بررسی قرار می‌گیرند. تکنیکهای تطبیق گوینده برای سیستمهای بازشناسی مبتنی بر HMM^{۷۵} به دو دسته اصلی تقسیم شوند. دسته اول شامل تکنیکهایی هستند که گفتار ورودی گوینده جدید را به فضای برداری مشترک با گفتارهای آموزشی تبدیل می‌کنند. این تکنیکها با نام روشهای نگاشت طیفی^{۷۶} شناخته می‌شود. دسته دوم شامل روشهایی هستند که با تغییر پارامترهای مدل باعث تطبیق بهتر آن مدلها با ویژگیهای داده های تطبیقی می‌شوند. این تکنیکها را تکنیکهای نگاشت مدل^{۷۷} می‌نامند.

10-3 انواع تطبیق ها

دو موضوعی که هنگام بحث در مورد روشهای نگاشت مدل باید مشخص شود، مد آموزش^{۷۸} (با نظارت^{۷۹} در مقابل بی نظارت^{۸۰}) و مد تطبیق (افزایشی^{۸۱} در مقابل دسته ای^{۸۲}) است. در مد آموزشی با نظارت، سیستم تشخیص، آوانویسی درست کلمات را در اختیار داشته و تنها گفتار کاربر را با آوانویسی های موجود مقایسه و هم تراز^{۸۳} می‌کند. در تطبیق بدون نظارت تشخیص دهنده آوانویسی درست کلمات را در اختیار نداشته و در نتیجه ممکن است گفتار کاربر را با یک آوای غلط مقایسه کند و به عبارتی تشخیص سیستم ممکن است با خطا همراه باشد. در نتیجه، مد با نظارت اگر قابل دسترسی باشد به مد بی نظارت ترجیح داده می‌شود.

10-4 روش نگاشت طیفی

ایده اصلی روش نگاشت طیفی افزایش بازده سیستم بازشناسی گفتار بوسیله نگاشت بردار ویژگی بدست آمده از گوینده جدید است. در این روش ویژگیهای گوینده را با یک تبدیل مناسب به ویژگی گوینده مدل آموزشی نگاشت می‌کنند.

⁷⁵ HMM-based recognition systems

⁷⁶ spectral mapping techniques

⁷⁷ model mapping approaches

⁷⁸ training mode


⁷⁹ supervise

⁸⁰ unsupervise

⁸¹ incremental

⁸² batch

⁸³ align

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 ژورنال اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

$$y_n = P(x_n) = A^T x_n \quad (113)$$

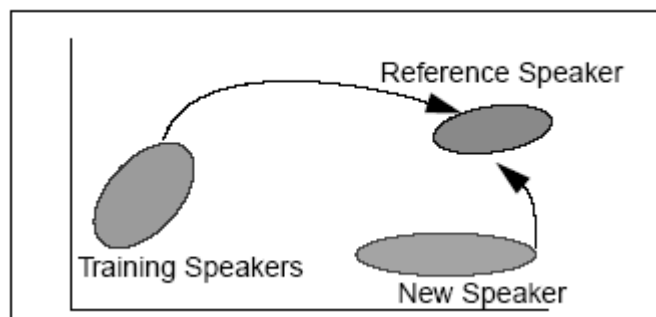
که $P(\cdot)$ تبدیلی است که بر روی ویژگی گوینده x_n ، ورودی اعمال می‌شود. A^T ماتریس تبدیل و y_n نتیجه تبدیل یافته بردار ویژگی ورودی است. ماتریس تبدیل بصورتی تخمین زده می‌شود تا مقدار MSE^4 بین ویژگی بدست آمده y_n و بردار ویژگی گوینده مرجع x_r مینیمم گردد.

$$D = E[(x_r - y_n)^T \cdot (x_r - y_n)] = E[(x_r - A^T x_n)^T \cdot (x_r - A^T x_n)] \quad (114)$$

نتیجه حاصل بصورت زیر است:

$$\min_A(D) \Rightarrow A = (E[x_n x_n^T])^{-1} \cdot E[x_n x_r^T] \quad (115)$$

عمل نگاشت طوری طراحی می‌شود که اختلاف میان مجموعه بردار مرجع و مجموعه بردار نگاشت شده حداقل باشد. این اختلافات به واسطه اختلافات طیفی میان سیستمهای تولید گفتار گوینده ها است. این روش در شکل ۱۶ نمایش داده شده است و معمولاً با نام روش نرمال کردن گوینده شناخته می‌شود.



شکل ۱۶ نگاشت بردارهای ویژگی گوینده های آموزشی و گوینده جدید به یک فضای مشترک

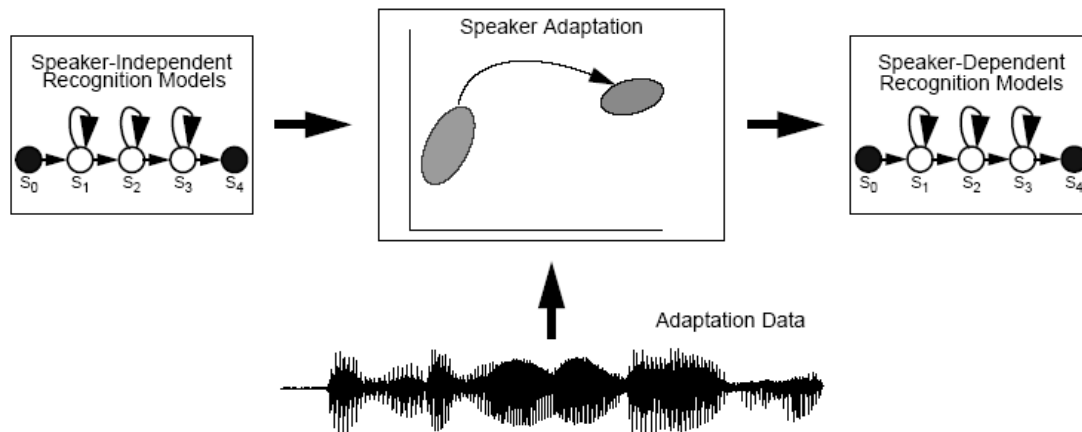
5-10 روش نگاشت مدل

هدف تکنیکهای نگاشت طیفی بهبود تطبیق میان گوینده مرجع و گوینده جدید است. اما این روشها به صورت واضح سعی نمی کنند تا دقت مدلها را برای گوینده جدید افزایش دهند. به عبارت دیگر این روشها به طور کامل از مزایای داده‌های تطبیقی استفاده نمی کنند. به همین دلیل از دسته دیگری از روشها با نام روشهای نگاشت مدل استفاده می شود که به جای نگاشت تمام گوینده ها به یک فضای

⁸⁴ Mean-Squared Error

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 شورای ملی اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

پارامترهای مدل را با بهترین ارائه گوینده جدید تنظیم می‌کند. بلوک دیاگرام کلی این روشها در شکل ۱۷ دیده می‌شود.



شکل ۱۷ روش نگاشت مدل برای تطبیق مدل‌های مخفی مارکوف

اکثر روشها، پارامترهای تابع چگالی HMM را بوسیله ماکزیمم کردن احتمال یک مدل (I) برای یک مجموعه مشاهدات (O) تخمین می‌زند.

$$P(I | O) = \frac{P(O | I)P_p(I)}{P(O)} \quad (116)$$

و با فرض اینکه I از مقدار توزیع قبلی $P_p(I)$ بدست آمده است.

از روشهای معروف و شناخته شده که برای تطبیق مدل مخفی مارکوف با گوینده موردنظر استفاده می‌شوند می‌توان به موارد زیر اشاره نمود:

^{۸۵}MAP (۱)

^{۸۶}MLLR (۲)

^{۸۷}CAT (۳)

در ادامه روش MLLR شرح داده می‌شود.

⁸⁵ Maximum A Posteriori

⁸⁶ Maximum Likelihood Linear Regression

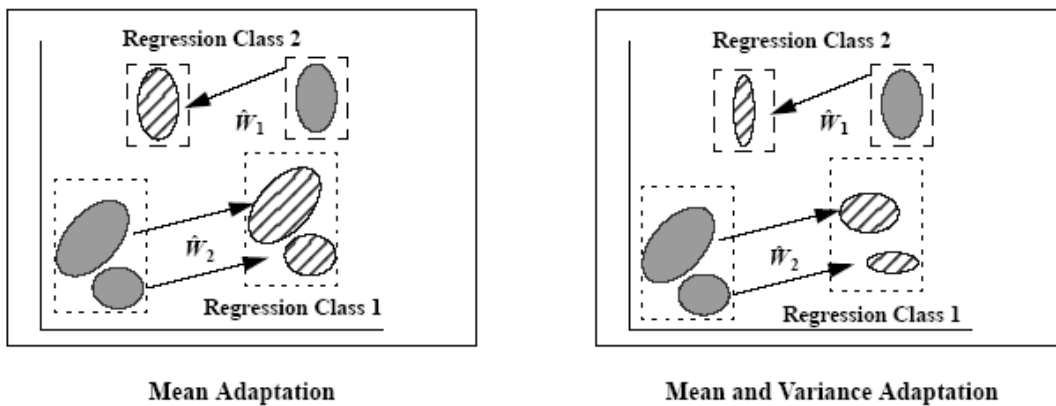
⁸⁷ Cluster Adaptive Training

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 گروه محاسبات و پردازش زبان
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

6-10-6-10 تطبیق مدل با روش MLLR

در روش MLLR، مجموعه‌ای از تبدیلات محاسبه می‌شوند تا عدم شباهت بین مدل اولیه و داده تطبیقی را کاهش دهند. MLLR یک روش تطبیق مدل است که در آن پارامترهای میانگین و واریانس برای مدل مخفی مارکوف دارای چگالی گوسی بوسیله تبدیلهای خطی تخمین زده می‌شوند. نتیجه این تبدیلات این است که میانگین‌های اولیه شیفت و واریانس اولیه به کلی تغییر می‌کند.

MLLR قادر است تبدیلات تطبیقی قوی⁸⁸ بسازد که حتی برای مدل‌هایی که داده تطبیقی برای آنها وجود ندارد نیز با استفاده از تبدیل مشترک⁸⁹ عمل تطبیق را انجام دهد. این مساله به ما کمک می‌کند تا مشکل داده‌های تطبیقی محدود را برطرف کنیم. با وجود مقدار کمی داده، تنها یک تبدیل سراسری⁹⁰، می‌تواند برای تمام مدلها به کار رود. هر چه داده بیشتری در دسترس باشد، تبدیلات دقیقتر و بیشتری می‌تواند به کار گرفته شود. در شکل ۴ اساس روش MLLR را آمده است که یک مدل مستقل از گوینده (بیضی توپر) را گرفته و از یک تبدیل استفاده می‌کند تا فضای مدل را به سمت یک به مدل وابسته به گوینده (بیضی راه راه) انتقال دهد. معمولاً تنها تطبیق، روی میانگین انجام می‌شود، به خاطر اینکه فرض می‌شود اختلاف اصلی میان گوینده‌ها به خاطر موقعیت میانگین آواها در فضای صوتی است. تطبیق روی کواریانس معمولاً کمتر استفاده می‌شود، چون تاثیر کمتری نسبت به تطبیق میانگین دارد. در شکل ۱۸ می‌بینیم که تطبیق میانگین، موقعیت میانگین مدل‌ها را در فضا جابجا می‌کند، در حالیکه تطبیق کواریانس شکل توزیع را تغییر می‌دهد.



شکل ۱۸ تطبیق مدل مستقل از گوینده به مدل وابسته به گوینده با استفاده از روش MLLR

⁸⁸ robust adaptation transforms

⁸⁹ transform sharing

⁹⁰ global transform

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 ژورنال اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

یک ماتریس تبدیل برای یافتن مقدار جدید میانگین استفاده می‌شود تا میانگین جدید بوسیله این تبدیل تخمین زده شود:

$$\hat{m} = Wx \quad (117)$$

که W یک ماتریس $n \times (n+1)$ و n بعد داده است. x نیز بردار میانگین گسترش یافته است که یک مقدار w به آن اضافه شده است.

$$x = [w, m_1, m_2, \dots, m_n]^T \quad (118)$$

می‌توان W را بصورت یک ماتریس A و یک بردار بایاس b بصورت $W = [b \ A]$ تجزیه نمود. سپس ماتریس تبدیل W با حل مسأله ماکزیمم کردن بوسیله روش ماکزیمم کردن امید ریاضی^{۹۱} (EM) بدست می‌آید. اگر $g_q(t)$ را احتمال تولید o_t در حالت q و در زمان t ، با دنباله مشاهدات O و مدل I باشد، ماتریس W_q با حل معادله زیر بدست می‌آید:

$$\sum_{i=1}^T \sum_{r=1}^R g_{qr}(t) U_{qr}^{-1} W_q x_{qr} x'_{qr} \quad (119)$$

که در آن x' ترانهاده ماتریس میانگین و U واریانس چگالی گوسی می‌باشند.

در ادامه تخمین ماتریس تبدیل برای واریانس فقط برای سیستمهای با کوواریانس قطری انجام می‌پذیرد. کوواریانس گوسی بصورت زیر تبدیل می‌یابد.

$$\hat{\Sigma}_m = B_m^T H_m B_m \quad (120)$$

که در آن H_m تبدیل خطی ای است که باید تخمین زده شود و B_m معکوس ضریب چولسکی^{۹۲} Σ_m^{-1} است. بنابراین:

$$\Sigma_m^{-1} = C_m C_m^T \quad (121)$$

$$B_m = C_m^{-1} \quad (122)$$

بعد از بازنویسی تابع کمکی، ماتریس تبدیل H_m بصورت زیر تخمین زده می‌شود:

⁹¹ Expectation Maximisation

⁹² Choleski

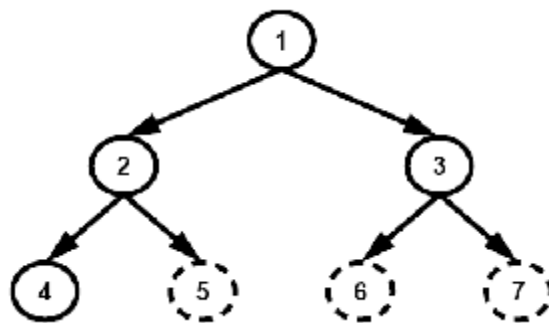
	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 ژورنال اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

$$H_m = \frac{\sum_{r=1}^{R_c} C_{m_r}^T [L_{m_r}(t)(o(t) - m_{m_r})(o(t) - m_{m_r})^T] C_{m_r}}{L_{m_r}(t)} \quad (123)$$

7-10 تبدیل اشتراک

به طور ایده ال، ما می خواهیم برای هر مدل گوسی در یک سیستم بر پایه HMM، از یک تبدیل تطبیقی استفاده کنیم تا تمام اختلافات میان مدل مستقل از گوینده و مدل وابسته به گوینده را به طور دقیق پیدا کنیم. هر چند، این کار در عمل نیاز به داده های تطبیقی زیادی دارد تا به طور دقیق، مدل های مناسب را تخمین بزنند. به همین خاطر، معمولاً یک فرم از تبدیل اشتراک به کار گرفته می شود که مجموعه مدل های گوسی که با هم کار می کنند را به وسیله یک تبدیل یکسان، تطبیق می دهد. در این طرح، مدل هایی که داده های تطبیقی کمی دارند و اصلاً داده تطبیقی ندارند، توسط داده های تطبیقی مدل های مشابه تطبیق داده می شوند.

یک روش معمول برای تبدیل ادغام^{۹۳} استفاده از یک درخت بازگشت باینری^{۹۴} است که در شکل ۱۹ نیز نمایش داده شده است. همانطور که در شکل دیده می شود برگ های ۴، ۵، ۶ و ۷ کلاس های پایه هستند. دایره های خط چین گره هایی را نشان می دهند که داده آموزشی کافی در اختیار ندارند و بنابراین این گره ها با گره های سطح بالاتر شریک می شوند. دایره های یک پارچه گره هایی هستند که داده آموزشی کافی برای تخمین انتقال دارند.



شکل ۱۹ یک درخت بازگشت باینری برای ادغام اجزای تطبیق شده

⁹³ transform pooling

⁹⁴ binary regression tree

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 شورای عالی اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

یک الگوریتم ادغام نیاز به تعیین اجزایی دارد که در هر گره با هم ادغام شده اند. در اغلب موارد از یک الگوریتم جدا کننده مرکز ثقل^{۹۵}، برای جدا کردن اجزای یک گره به دو مجموعه از اجزا استفاده می شود. این عمل طوری انجام می شود که به ما اطمینان می دهد، اجزایی که به یکدیگر نزدیکترند، در فضای احتمال در یک گره با هم ادغام شوند.

⁹⁵ Centroid splitting algorithm

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 مؤسسه ملی اطلاع‌رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان‌سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

11 بازشناسی هویت از طریق گفتار

1-11 مقدمه:

بازشناسی گفتار یکی از مهمترین مباحث تحقیقات گفتاری است که در طول پنج دهه گذشته، یعنی از زمانی که چنین تحقیقاتی در اوایل دهه ۵۰ میلادی آغاز شد، پیشرفتهای زیادی داشته است.

فرایند بازشناسی گوینده^{۹۶} که در اینجا به شاخه عمومی مسائلی از قبیل تعیین هویت گوینده^{۹۷}، تصدیق گوینده^{۹۸} و طبقه بندی گوینده اطلاق می شود، اولین بار توسط آتال^{۹۹} در این زمینه مطرح شد. این اولین قدم در معرفی این زمینه تحقیقاتی بود که البته در ابتدا از موفقیت کمی برخوردار بود.

تعیین هویت گوینده به این صورت تعریف می شود که از میان N مدل گوینده مرجع، آن مدل گوینده ای که نزدیکترین و بیشترین شباهت را به گوینده نامشخصی ورودی دارد پیدا می کند. از آنجایی که الگوهای گفتار با تک تک مدل‌های مرجع مقایسه می شوند و همچنین از آنجایی که برای هر تصمیم گیری نادرست احتمال معینی برای هر مقایسه وجود دارد، بنابراین واضح است که احتمال تصمیم گیری کلی تابعی از N بوده و در نتیجه هر چه تعداد مدل‌ها بیشتر باشد احتمال خطا یا تصمیم گیری نادرست در تعیین هویت گوینده بیشتر می شود.

مسئله تصدیق گوینده به این صورت مطرح می شود که گفتار یک گوینده نامشخص و مدل گوینده ای که وی ادعا می نماید داده شده است و بایستی مشخص شود که آیا گفتار این گوینده به اندازه کافی به مدل گوینده ادا شده شباهت دارد یا نه. بنابراین در این حالت تعداد مقایسه یکی است و معمولاً مستقل از تعداد جمعیت مدل مرجع است.


در یک تقسیم بندی، بسته به اینکه مدلی که برای شناسایی مورد آزمون قرار می گیرد جزو مدل‌های مجموعه باشد یا نباشد، دو تعریف زیر ارائه می گردد:

^{۹۶} Speaker Recognition

^{۹۷} Speaker Identification

^{۹۸} Speaker Verification

^{۹۹} Atal

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 ژورنال اطلاع رسانی
	عنوان زیرپروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	کد زیرپروژه: پیک-متن-فارس - ۲ - خ	ویرایش: ۱/۰	

۱. تعیین هویت در مجموعه بسته: که در این حالت مدل مورد آزمون جزو N مدل موجود در شبکه می باشد.

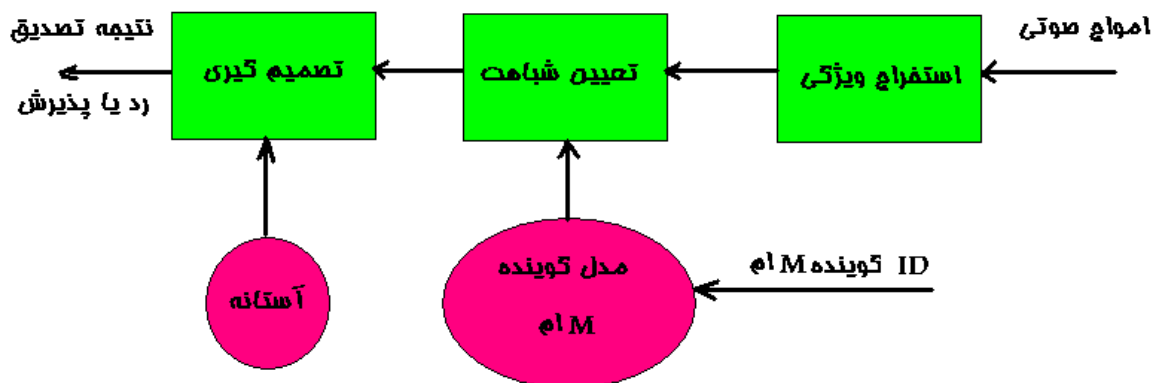
۲. تعیین هویت در مجموعه باز: که در این حالت گفتار ورودی نامشخص متعلق به هیچ کدام از اعضای مدل‌های گویندگان مرجع موجود در سیستم نمی باشد.

در سیستم های تعیین هویت در مجموعه باز نتایج خروجی $N+1$ حالت می باشد که N تعداد مدل‌های مرجع است. در واقع در این نوع سیستم، ترکیبی از دو نوع سیستم تعیین هویت مجموعه بسته و تصدیق هویت می باشد. در این سیستم ابتدا نزدیکترین مدل به گفتار ورودی پیدا شده و سپس صحت تعلق گفتار ورودی به آن مدل مورد بررسی قرار می گیرد. بنابراین کارایی این سیستم نسبت به تعیین هویت مجموعه بسته کمتر است.

تصدیق هویت گوینده راحتی و اطمینان بیشتری را برای بسیاری از فعالیتهای روزمره آدمی که نیاز به امنیت دارد ایجاد می نماید. این تکنولوژی بدلیل اینکه باز خصوصیات ذاتی موجود در صدای انسان استفاده می کند، در مقابله با مسئله تقلید و کلاه برداری بسیار مقاوم بوده و از لحاظ ویژگی بازشناسی هویت از راه دور نسبت به برخی روشها از جمله اثر انگشت کارایی بسیار بیشتری دارد.


11-2 روشهای پیاده سازی سیستم های تصدیق گوینده:

شمای کلی یک سیستم تصدیق گوینده در شکل زیر آمده است:



شکل ۲۰ شمای کلی یک سیستم تصدیق گوینده

سیگنال گفتار گوینده به عنوان ورودی به سیستم انتقال داده می شود. سپس پیش پردازش و نهایتاً استخراج ویژگی ها بر روی آن اعمال می گردد و به این ترتیب به یک فضای جدید منتقل می گردد. ویژگی های استخراج شده از گفتار گوینده با مدل مرجع او مقایسه می شود. عمل مقایسه به این صورت

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 مؤسسه ملی اطلاع‌رسانی
	عنوان زیرپروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

انجام می‌گیرد که میزان شباهت یا اختلاف این ویژگی‌ها با مدل مرجع به دست می‌آید. میزان شباهت به دست آمده با میزان آستانه مقایسه می‌شود و با توجه به این مقایسه خروجی سیستم مبنی بر تعلق و یا عدم تعلق گفتار ورودی به گوینده ادعا شده مشخص می‌شود.

قبل از معرفی روشهای پیاده سازی سیستم، به بررسی نکاتی در طراحی سیستم های بازشناسی (اعم از تعیین هویت و یا تصدیق هویت) می پردازیم:

۱- انتخاب ویژگی ها: برای استخراج ویژگی های گوینده بایستی چند عامل مد نظر قرار گیرد؛ از جمله: این ویژگی ها بایستی به نوبت محیط انتقال مانند کانال های مخابراتی حساس باشند، همچنین این ویژگی ها بایستی به حالات روانی و فشار های محیط که گوینده در آن در حال صحبت کردن است وابسته باشد. در کل روش بهینه ای برای استخراج ویژگی ها وجود ندارد و معمولاً از طریق تجربه این عمل صورت می گیرد.

لازم به ذکر است که در بازشناسی گوینده از ویژگی های دینامیک که بستگی زیادی به حالات گوینده دارد استفاده می شود، در صورتی که در روشهای دیگر خصوصیات فیزیکی استاتیک و پایدار مورد بررسی قرار می گیرند. بنابراین یک سری محدودیتهای ذاتی در مورد استفاده از این سیگنالها وجود دارد. برای درک این محدودیتهای بایستی اطلاعات تمایزدهنده گوینده ها و نحوه قرار گرفتن آنها در سیگنال گفتار مورد بررسی قرار گیرد. سیگنال گفتار از طریق حرکت اندامهای تولید گفتار بوجود می آید و توسط حنجره و سیستم عصبی کنترل می گردد. بنابراین دو منع اطلاعات گوینده در سیگنال صحبت وجود دارد یکی مربوط به خصوصیات فیزیکی و ساختاری مجرای گفتار و دیگری اطلاعات کنترلی از مغز و ماهیچه ها اندام گویایی است. این اطلاعات همراه با اطلاعاتی مربوط به هنگام حرکت دادن مفصل های اندامهای تولید گفتار، وارد سیگنال صحبت می شود. در کل اطلاعات سیگنال گفتار را به دو دسته سطح بالا (مانند: لحن، محتوای گفتار و استیل گفتار یعنی طریقه استفاده گرامری و نحوی از کلمات) و سطح پایین (مانند: خصوصیات سیگنال گفتار از قبیل دامنه طیف، فرکانی گام واکدار، فرکانس فرمانت، پهنای باند و خصوصیات تناوبی گفتار واکدار) تقسیم بندی می کنند. اطلاعات سطح بالا عملاً در بازشناسی گوینده توسط انسان کاربرد داشته و در عوض سیستم های بازشناسی گوینده اتوماتیک از ویژگی های سطح پایین سیگنال استفاده می شود.

۲- مدل های گویندگان: در سیستم های بازشناسی گوینده، مدل هر گوینده شامل خصوصیات آماری او است. معمولاً در مدل های گویندگان از یکی از دو مدل پارامتری (مانند مدل گوسی) و غیر پارامتری (مانند مدل مراکز خوشه ها یا حالات چند گانه) استفاده می شود.

۳- طول گفتار آموزشی: که با افزایش این زمان نتایج مطلوب تری حاصل می گردد.

۴- انتخاب گفتار مناسب: پیشنهاد می گردد سکوت و بخشهای غیر گفتاری حذف گردد. همچنین

	عنوان پروژه:		 ژورنال اطلاع رسانی	
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی			
	عنوان زیرپروژه:			
تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		کد زیرپروژه: پیکرمتن فارسی - ۲ - خ	ویرایش: ۱/۰	تاریخ: ۱۳۸۸/۰۳/۱۹

گفتارهای کم انرژی و بی واک در تمایز بین گوینده ها کم اثر بوده و پیشنهاد می گردد که از اطلاعات مدل حذف گردند.

۵- محیط نویزی: نویز پشت زمینه یک مشکل عمومی است و بهتر است برای تخمین آن از نویز پشت زمینه در زمان سکوت استفاده و با اعمال آن به مدل، تخمینی از مدل بدون نویز داشته باشیم.

۶- تنوع گوینده: اغلب سیگنال گفتار یک شخص در حالت های خوشحالی، ناراحتی و خستگی متفاوت و به این ترتیب کارایی سیستم را کاهش می دهند. بنابراین خصوصیات آماری یک فرد ثابت نیست.

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 مؤسسه ملی اطلاع‌رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان‌سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

12 سنتز گفتار (Speech Synthesis)

1-12 مقدمه

سنتز گفتار یک فناوری است که بوسیله آن متن به گفتار تبدیل می‌شود (TTS). یکی از کاربردهای سنتز گفتار تبدیل email و Fax به صدا و اخیراً تبدیل محتویات صفحات وب به صوت می‌باشد. نرم افزارهایی که در این زمینه طراحی شده‌اند، بسیار زیاد هستند. با وجود آمدن گفتار مصنوعی نرم افزارهایی برای خواندن کتابها طراحی شده‌اند.

عروض (prosody)

عروض یکی از فاکتورهای اصلی برای بدست آوردن یک گفتار مصنوعی با کیفیت زیاد می‌باشد. مفهوم عروض، زیر و بم کردن صدا و ریتم گفتار که باعث تلفظ و برداشت مفهوم‌های مختلفی از گفتار می‌گردد، می‌باشد. همچنین اطلاعاتی راجع به وضعیت روانی سخنگو به ما می‌دهد. صدای هر شخص دارای یک فرکانس اصلی می‌باشد. که بر اساس آن تارهای صوتی به لرزه در می‌آیند. می‌توان گفت فرکانس اصلی یکی از فاکتورهای جدانشونده صدا می‌باشد و آن قابل تغییر است. بطور ساده، فرکانس اصلی بوسیله اصلاح شکل موج قابل دستکاری کردن می‌باشد. با استفاده از تکنیک پردازش سیگنال، امواج صوتی می‌توانند از این تکنیک‌ها استفاده کنند. صدای مصنوعی کیفیت صدای طبیعی را ندارد چندین ایده را هنگام انجام کار می‌توان مد نظر قرار داد. یکی از مشکلات اساسی و لاینفک در این موضوع مسأله کنترل وقفه‌های زمانی بوسیله سیستم کنترل امواج می‌باشد که به منظور تولید یک صدای خوب و قابل فهم توسط انسان بر این مشکل بایستی غلبه گردد. هنگامیکه سنتز گفتار در یک سیستم بزرگ محاوره‌ای بکار گرفته شود، استفاده الگوهای عروضی گفتار اهمیت بحرانی داشته و استفاده نادرست از آن منجر به شکست خواهد شد. غلبه بر مشکلات عروضی گفتار بصورت دستکاری کردن شکل موج می‌باشد. بعضی از محققان با جلوگیری کردن از طولانی شدن متن‌ها تلاش در کم کردن این مشکلات را دارند.

12-2 انواع متدهای تبدیل متن به گفتار

تکنیک‌های سنتز گفتار در سه مقوله زیر قابل بررسی می‌باشند:

۱- سنتز شمرده به شمره لغات (Articulatory Synthesis)

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 ژورنال اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

۲- سنتز فرمانت (Formant Synthesis)

۳- سنتز اتصالی (Concatinative Synthesis)

۱۲-۲-۱- سنتز شمرده به شمرده لغات (Articulatory Synthesis)

در این روش از یک مدل کامپیوتری - مکانیکی برای تولید گفتار استفاده می شود. مثل حنجره که بوسیله ارتعاش تارهای صوتی باعث تولید صدا و تحریکات حلقی می شود. زبان بوسیله فشردن هوا در محوطه دهان باعث جدا کردن صداها از همدیگر می شود. مسیر عبور هوا بوسیله زبان با فشار بر روی سقف دهان یا دندانها تغییر می کند. لبها برای ایجاد هماهنگی و کمک به فرایند Articulation به زبان کمک می کنند. در حالت ایده آل سنتز کننده های شمرده به شمرده می توانند بوسیله ماهیچه های مصنوعی و کنترل آنها به شکل زبان ، لب و دهان شبیه سازی شوند. این کار مستلزم حل معادلات دیفرانسیل درجه سه وابسته به زمان می باشد. در هر صورت محاسبه ایجاد گفتار در خروجی بسیار مشکل می باشد و نتیجه آن در اکثر موارد تولید یک صدای غیر طبیعی است.

۱۲-۲-۲- سنتز فرمانت (Formant Synthesis)

سنتز فرمانت از یک سری از قوانین برای کنترل یک شکل موج که گفتار انسان را شبیه سازی می کند ، استفاده می کند. این روش باعث تولید یک گفتار ماشینی صحیح ، شبیه صدای روبات می باشد. اکثر نرم افزارهای TTS از این شکل سنتز برای تولید گفتار استفاده می کنند. روش سنتز فرمانت باعث تولید گفتاری قابل فهم و واضح می گردد ، اما صدای آن همانند صدای طبیعی نمی باشد. از مزایای این روش استفاده معقول از حافظه و محاسبات کامپیوتری می باشد. سنتز فرمانت بطور گسترده برای پاسخگویی به کاربر مورد استفاده قرار می گیرد.

۱۲-۲-۳- سنتز اتصالی (Concatinative Synthesis)

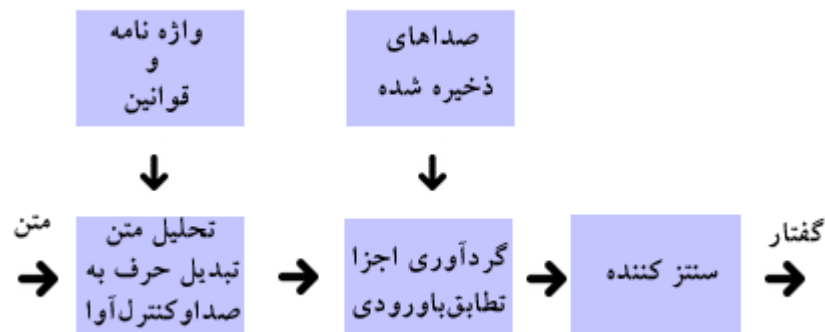
در این رویکرد گفتار ذخیره شده طبیعی بصورت تکه تکه در کنار هم قرار می گیرند تا تولید یک گفتار خروجی کنند. به عنوان مثال ، برای اینکه به کاربر مقدار حساب بانکی اش را بیان کنیم باید لغات زیر را با اتصال به یکدیگر بیان کنیم.

"Your + total + account + balance + is + five + thousand + three + hundred + forty + nine + dollars + and + twenty + six + cents"

به علت اینکه فاکتورهای زیادی می توانند بر روی ویژگی های قطعات گفتاری در یک زبان طبیعی تاثیرگذار باشند در سیستمهای امروزی قطعات گفتاری زیادی در نوعها و قالبهای مختلف ذخیره شده است. همچنین برای زبانهای مختلف ، قطعات جداگانه ای ذخیره شده است. قطعه های گفتار می توانند به شکل موج خام و یا مجموعه ای از پارامترهای مشتق شده از انواع گوناگون موج ذخیره شوند. اگرچه

	عنوان پروژه:		 مؤسسه ملی اطلاع‌رسانی
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		
	عنوان زیر پروژه:		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	کد زیر پروژه: پیک-متن-فارس - ۲ - خ

ذخیره قطعه‌ها بصورت مجموعه‌ای از پارامترها فضای کمتری را اشغال می‌کند، اما روند حرکت این تکنولوژی بر اساس ذخیره موج خام بوده است. هدف از این کار تولید صدایی شبیه به صدای انسان می‌باشد. هنگام سنتز یک جمله، قطعه‌ها از بانک اطلاعات استخراج می‌شود و به هم متصل می‌گردند و به منظور داشتن فاصله زمانی مناسب و تلفظ مناسب اصلاح می‌شوند. یکی از کارهای مشکل و پیچیده این روش، شناخت و انتخاب قطعه مناسب بر اساس ساختار جمله گفتاری مورد نظر می‌باشد. اجزای گفتار مثل واجها شامل نیمه اول یک صدا و نیمه دوم صدای دیگر می‌باشند. بعضی از سنتز کننده‌های الحاقی که به آنها demissyllables گفته می‌شود از diphon برای مقیاس زمانی سیلابها استفاده می‌کنند. معنای diphon زمان گذرای بین واجها و داده‌ها و واجهای ضبط شده می‌باشد و معنای دمی سیلاب زمان گذرای نیم سیلاب می‌باشد (این زمان می‌تواند از زمان شروع تا وسط سیلاب و یا از وسط سیلاب تا پایان سیلاب در نظر گرفته شود). در شکل ۲۱ بلوک دیاگرام TTS الصاقی نشان داده شده است.



شکل ۲۱ فرایند سنتز گفتار

بلوک اول ماژول تحلیل متن را بصورت کد اسکی دریافت کرده و آن را بصورت یک سری از سمبلهای تلفظ و عروض، شامل فرکانس اصلی و فاصله‌های زمانی و دامنه موج در می‌آورد. ماژول تحلیل متن از یک سری از ماژولهای جداگانه تشکیل شده است. متن ورودی در ابتدا تحلیل می‌شود و سمبلهای غیر از حروف الفبا و حروف اختصاری به لغات کامل تبدیل می‌شوند. برای مثال در جمله:

Dr. Kay lives at 5201 Peak Dr.

کلمه Dr. اول به کلمه doctor تبدیل می‌شود و کلمه دوم به کلمه Drive و پس از آن ۵۲۰۱ به صورت پنج - صفر - صفر - یک تبدیل می‌شود. پس از آن پارسر سنتز کننده که وظیفه آن تشخیص قطعه‌های گفتار برای هر لغت در جمله می‌باشد نشانه گذاری روی متن را انجام می‌دهد. ماژول بعدی دارای دو جزء اصلی است:

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 شورای ملی اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

۱. جزء نرمال ساز متن

۲. جزء تجزیه کننده متن

قوانین نرمال کننده متن ، متن ورودی را می خواند و بصورت رشته ای از قطعات گفتاری نمایش داده و لغات مخفف را بصورت حرف حرف در می آورد. برای مثال جمله St. John St. به جمله Saint john Street تبدیل می شود. یا 1 oz به one ounce تبدیل می شود. یا \$500 به پانصد دلار تبدیل می شود. کاربر می تواند بعضی از رفتارهای سیستم را با دستکاری بعضی از لغات در واژه نامه تغییر دهد. قوانین نرمالسازی ، زمان پایان جمله را تعیین می کنند. به مثال زیر که برای جمله Dr. Smith is 43 آمده است توجه کنید:

Sentence : | sent |

Inp : | D|r|.| | S|m|i|t|h| | |i|s| | |4| |3|.| |

Text : |d|o|c|t|o|r| | |s|m|i|t|h| | |i|s| | |f|o|r|t|y| | |t|h|r|e|e|.| |

خطوط عمودی برای هماهنگ سازی رشته متن در نقاط مناسب مورد استفاده قرار می گیرد. نتیجه حاصل از نرمال سازی جمله برای قوانین تجزیه کننده متن فرستاده می شود و در آنجا پردازش کلامی صورت می گیرد. قوانین تجزیه متن عمل خود را روی متن نرمال شده قبل انجام می دهند و آن را به عبارتهای طبیعی - لغات - سیلابها و diphone ها می باشد. در شکل قبل دومین بلوکی که یونیتهای مختلف را برای ایجاد یک صدای طبیعی گردآوری می کند ، نشان داده شده است. این اطلاعات به یک سنتز کننده ، داده شده که برای کاربر تولید یک شکل موج گفتار می نماید.

یکی از کارکردهای کلیدی سنتز گفتار ، گردآوری رشته صحیحی از صداها به همراه یک واژه نامه تلفظ کمکی می باشد. بنابراین در جمله Dr. Kay lives at 5201 Peak DR فعل live از جمله شناسایی می شود و همچنین در جمله (to Mr Wright to give him the right direction Write a letter) سنتز کلمات wright و right و write به راحتی انجام می شود. اگر واژه نامه لغات دچار ایراد شود قوانین کلی تبدیل حروف الفبا به صدا اعمال می شوند. در پایان بوسیله جملات نقطه گذاری شده و اطلاعات صوتی و نحوی موجود ، یک ماژول عروضی عملیات تولید یک عبارت را برای مقصد مورد نظر انجام می دهد.

برای مثال برای فرکانس اصلی ، فاصله زمانی بین واجها و دامنه توسط ماژول تعیین می شود. در سالهای اخیر TTS پیشرفتهای خوبی در زمینه تولید صدای طبیعی داشته است که نتایج آن در پورتالهای صوتی و سرویس بانکها و آژانسهای مسافرتی شنیده می شود ، بطوریکه برای مشتری های این شرکتها تشخیص صدای کاربر واقعی و مصنوعی قابل تشخیص نمی باشد. سنتز یک بخش انتخاب شده بوسیله عوامل زیر قابل تاثیر است. یکی از مهمترین دیدگاهها به این موضوع افزایش قدرت پردازشی و ذخیره سازی سرور می باشد که بطور مستقیم بر روی ساینز دستگاههای ذخیره سازی صوت تاثیر دارد.

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 مؤسسه ملی اطلاع‌رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان‌سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

امروزه سنتز اتصالی در بسیاری از موارد استفاده می‌شود. نقطه نظر دیگر نیز این است که تکنیک‌های کارای جستجو به ما اجازه می‌دهند که میلیون‌ها بخش صوتی متفاوت را بصورت بلادرنگ برای بدست آمدن یک گفتار بهینه مورد جستجو قرار دهیم. در پایان اینکه، می‌توانیم از برچسب زنده‌های اتوماتیک که باعث سرعت بخشی دسترسی بانک اطلاعات عروض و تلفظ می‌شوند، استفاده کنیم.

12-3 ارزیابی سنتز گفتار

سه فاکتور کلیدی برای داشتن یک گفتار مصنوعی موفق وجود دارد.

۱- مفهوم بودن (Intelligibility)


درجه ایست که شنونده گفتار را متوجه می‌شود. این میزان بوسیله اندازه‌گیری تعداد واجهای تولید شده بوسیله یک نرم افزار در مقابل تعداد واجهایی که توسط کاربر قابل تشخیص است، قیاس می‌شود.

۲- طبیعی بودن (Naturalness)

درجه ایست که شنونده گفتار سنتز شده مصنوعی را همانند گفتار طبیعی انسان می‌پندارد. اندازه‌گیری آن از لحاظ ذهنی بیشتر از قابلیت وضوح صدا می‌باشد ولی از نظر اهمیت در یک سطح قرار دارند.

۳- مطبوعیت (Pleasantness)

درجه ای است که بر اساس آن شنونده از صدا لذت می‌برد یا آن را قابل تحمل می‌داند. تعیین درجه‌های فوق بسته به نوع کاربران و نوع سرویس متفاوت است.

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 ژورنال اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

13 دادگان های گفتاری

1-13 مقدمه

استفاده از گفتار جهت برقراری ارتباط با ماشین‌ها در سال‌های اخیر مورد توجه پژوهشگران قرار گرفته است. امروزه با توجه به افزایش قدرت پردازشی پردازنده‌ها و حجم حافظه‌ها تمایل افراد به سمت ذخیره-ی اطلاعات گفتاری در رایانه‌ها بیشتر شده و این سبب شده است تا جستجو، تشخیص و پردازش‌های لازم مستقیماً روی فایل‌های گفتاری ذخیره شده در رایانه اعمال شود.


در تهیه‌ی دادگان‌های صوتی قواعد خاصی مورد نظر قرار می‌گیرد از جمله پایگاه داده باید شامل کلمات و عبارات مناسبی باشد که معرف مجموعه گفتاری زبان مورد نظر باشد و از طرفی گویندگان باید از جنسیت، سن و سطح سوادهای مختلف جهت پوشش دادن مشخصه‌های گفتاری مختلف، انتخاب شوند. در فصل حاضر مطالعه‌ی به نمونه‌هایی از پایگاه داده‌های گفتاری، خصوصیات آنها و کاربردشان اشاره می‌شود. برای تهیه دادگان گفتاری جدید خصوصیات دادگان‌های موجود باید مطالعه و بررسی شده و بسته به کاربرد مورد نظر دادگان جدید طراحی گردند.

13-2 دادگان گفتاری فارس‌دات

دادگان فارس‌دات یک بانک اطلاعاتی از گفتار فارسی است که توسط پژوهشکده هوشمند علائم ارائه گردیده است. این دادگان حاوی ۳۸۶ جمله (ساخته شده با ۱۰۰۰ کلمه شامل تمامی دو واجها) است که توسط ۳۰۰ گوینده و با ۱۰ لهجه مختلف بیان شده‌اند. هر گوینده ۲۰ جمله را ادا کرده است. در مجموع ۶۰۰۰ جمله متمایز در فارس‌دات موجود است. همچنین ۱۱۰۰۰ کلمه مجزا نیز توسط ۱۰۰ گوینده بیان شده است. کلیه واجهای موجود در جملات این بانک اطلاعاتی دارای برچسب زمانی هستند. این دادگان می‌تواند در سیستم‌های بازشناسی گفتار، بازشناسی گوینده و لهجه بکار گرفته شود.

13-3 دادگان گفتاری فارس‌دات بزرگ

دادگان فارس‌دات بزرگ یک بانک اطلاعاتی از گفتار فارسی است که توسط پژوهشکده هوشمند علائم

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 مؤسسه ملی اطلاع‌رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

ارائه گردیده است. این دادگان حاوی ۱۴۰ ساعت گفتار است که توسط ۳۰۰ گوینده و با ۱۰ لهجه مختلف بیان شده‌اند. هر گوینده ۴۰۰۰ کلمه را ادا کرده است. کلیه سیگنالها موجود در این بانک اطلاعاتی دارای برچسب در سطح کلمه هستند.

مناسب جهت کاربردهایی از قبیل بازشناسی گفتار پیوسته فارسی مستقل از گوینده با واژگان بزرگ، بازشناسی گوینده مستقل از متن و تحقیقات در زمینه واج‌شناسی آزمایشگاهی می باشد.

13-4 دادگان گفتاری فارسی‌دات تلفنی (مونولوگ)

دادگان فارسی‌دات تلفنی یک بانک اطلاعاتی از گفتار فارسی است که توسط پژوهشکده هوشمند علائم ارائه گردیده است. این دادگان توسط ۶۴ گوینده و با لهجه‌های مختلف بیان شده‌اند. این دادگان شامل ۷ ساعت سیگنال گفتاری است. کلیه سیگنالها موجود در این بانک اطلاعاتی دارای برچسب در سطح واج هستند. برای کاربردهایی از قبیل سیستم‌های رایانه‌ای تلفنی در بازشناسی گفتار و بازشناسی گوینده و بازشناسی زبان فارسی مناسب می باشد.

13-5 دادگان گفتاری فارسی‌دات تلفنی بزرگ (محاوره ای)

دادگان فارسی‌دات تلفنی بزرگ یک بانک اطلاعاتی از گفتار فارسی است که توسط پژوهشکده هوشمند علائم ارائه گردیده است. این دادگان حاوی ۲۵۲ مکالمه تلفنی دو طرفه پیرامون ۱۳ موضوع مختلف (شامل موضوعات فرهنگی، سیاسی، اجتماعی، اقتصادی، علمی، آموزشی، هنری، خانواده، ورزشی، صنعتی و فناوری، کشاورزی و دامداری، عمران و خدمات) است که توسط ۲۰۰ گوینده و با ۱۰ لهجه مختلف بیان شده‌اند. بخشی از سیگنالها موجود در این بانک اطلاعاتی دارای سه نوع برچسب زبانی در سطح کلمه، شامل برچسب‌واجی، آوایی و نوشتاری هستند. همچنین عناصر غیرزبانی از قبیل نوفه خط تلفن، نوفه محیطی، صدای تنفس، صدای برخورد لب‌ها صدای شک و تردید و کلماتی که ناقص تولید می‌شوند، برچسب‌دهی شده اند.

مناسب جهت کاربردهایی از قبیل بازشناسی گفتار پیوسته تلفنی مستقل از گوینده با واژگان بزرگ، بازشناسی گوینده تلفنی مستقل از متن، بازشناسی گفتار مکالمه‌ای روی خط تلفن به منظور خودکارسازی مراکز تلفنی، تحقیقات واج‌شناسی آزمایشگاهی و تجزیه و تحلیل کلام (گفتمان) می باشد.

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 گروه اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

6-13 دادگان اعداد و ارقام منفصل و پیوسته فارسی

دادگان اعداد متصل و پیوسته فارسی یک بانک اطلاعاتی از اعداد فارسی است که توسط پژوهشکده هوشمند علائم ارائه گردیده است. تعداد گویندگان اعداد پیوسته ۱۰۰ نفر و تعداد گویندگان اعداد متصل ۱۱۰ نفر با لهجه تهرانی می باشد. این بانک اطلاعاتی دارای برچسب در سطح عدد هستند.

7-13 دادگان گفتاری TIMIT

دادگان TIMIT یک بانک اطلاعاتی از گفتار پیوسته انگلیسی است که توسط شرکت TI¹⁰⁰ و دانشگاه MIT¹⁰¹ تهیه شده است و اداره استاندارد آمریکا (NIST¹⁰²) آن را تأیید کرده 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. در مجموع ۲۴۳۲ جمله متمایز در TIMIT موجود است که شامل ۲ جمله مشترک میان تمامی گویندگان، ۴۵۰ جمله مشترک میان گروههای ۷ نفری گویندگان و ۱۸۹۰ جمله تک گوینده است. کلیه کلمات و واجهای موجود در جملات این بانک اطلاعاتی دارای برچسب زمانی هستند. بانک اطلاعاتی TIMIT به دو بخش آموزش و آزمون تقسیم شده است که بخش آموزش شامل ۴۶۲ گوینده و بخش آزمون شامل ۱۶۸ گوینده است و گویندگان و جملات ادا شده در هر یک از این دو بخش با یکدیگر متفاوتند.

دادگان TIMIT عاری از نویز است و معمولاً برای ارزیابی نرخ بازشناسی واجها در بازشناسی گفتار پیوسته مورد استفاده قرار می‌گیرد. هر چند که با وجود برچسبهای زمانی برای کلمات و واجها، برای ارزیابی نرخ بازشناسی کلمات مجزا نیز قابل استفاده است. برای کاربرد این دادگان در ارزیابی روشهای بازشناسی گفتار در حضور نویز، باید نویز را بطور مصنوعی به این دادگان اضافه نمود.

8-13 دادگان گفتاری TIDIGITS

دادگان TIDIGITS یک بانک اطلاعاتی گفتار انگلیسی است که در شرکت TI تهیه شده است. هدف از تهیه این دادگان طراحی و ارزیابی الگوریتمهای بازشناسی مستقل از گوینده برای دنباله متصل اعداد¹⁰³ بوده است. این دادگان حاوی گفتار ۳۲۶ گوینده متشکل از ۱۱۱ مرد، ۱۱۴ زن، ۵۰ پسر بچه و ۵۱ دختر

¹⁰⁰ Texas Instruments

¹⁰¹ Massachusetts Institute of Technology

¹⁰² National Institute for Standards and Technologies

¹⁰³ Connected Digit Sequences

	عنوان پروژه:		 انستیتو مطالعات زبان و ادبیات	
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی			
	عنوان زیر پروژه:			
تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		کد زیر پروژه: پیک متن فارسی - ۲ - خ	ویرایش: ۱/۰	تاریخ: ۱۳۸۸/۰۳/۱۹

بچه است که هر یک از گویندگان ۷۷ دنباله متصل عددی را بیان کرده اند. دنباله‌های عددی بیان شده با استفاده از ۱۱ رقم زیر ساخته شده‌اند:

.zero, oh, one, two, three, four, five, six, seven, eight, nine

این دنباله‌های عددی بیان شده برای هر گوینده را می‌توان این گونه تقسیم کرد: ۲۲ بار بیان اعداد مجزا (۲ بار بیان برای هر عدد) و ۱۱ مرتبه بیان برای هر یک از رشته‌های متشکل از ۲، ۳، ۴، ۵ و ۷ عدد که در مجموع ۷۷ دنباله مذکور را برای هر گوینده شکل می‌دهند. تمامی این دنباله‌ها به دو زیر مجموعه دادگان آموزش و دادگان آزمایش تقسیم شده‌اند.

دادگان TIDIGITS نیز عاری از نویز است و برای ارزیابی روشهای بازشناسی گفتار در حضور نویز باید به آن بطور مصنوعی نویز اضافه نمود.

13-9 دادگان گفتاری AURORA2

دادگان گفتاری AURORA2 برای ارزیابی دقت بازشناسی سیستم‌های بازشناسی گفتار در محیط‌های نویزی و به ویژه برای ارزیابی روش‌های جبران ویژگی در شرایط نویزی متفاوت مورد استفاده قرار می‌گیرد. این دادگان با استفاده از گفتار بزرگسالان دادگان گفتاری TIDIGIT ساخته شده است، به این ترتیب که ابتدا نرخ نمونه‌برداری به ۸ کیلو هرتز کاهش یافته است و سپس فیلتری برای شبیه‌سازی تأثیر کانال انتقال تلفنی بر روی دادگان اعمال گردیده است. برای شبیه‌سازی اثر خط تلفن، دو نوع فیلتر استاندارد مورد استفاده قرار گرفته‌اند که با نامهای G.712 و MIRS شناخته می‌شوند.

گفتار حاصل از اعمال فیلترهای فوق، دادگان تمیز گفتاری را شکل می‌دهند. سپس به این دادگان تمیز نویزهای جمع‌پذیری با نسبت‌های سیگنال به نویز ۲۰، ۱۵، ۱۰، ۵، ۰ و ۵- دسی‌بل اضافه می‌شوند. نویزهای جمع‌پذیر بکار رفته عبارتند از: نویزهای مترو¹⁰⁴، همهمه¹⁰⁵، ماشین¹⁰⁶، نمایشگاه¹⁰⁷، رستوران¹⁰⁸، خیابان¹⁰⁹، فرودگاه¹¹⁰ و ایستگاه قطار¹¹¹.

در دادگان AURORA2 دو مجموعه مجزای آموزش و آزمایش مشخص شده‌اند. داده‌های آموزش خود

¹⁰⁴ Subway
¹⁰⁵ Babble
¹⁰⁶ Car
¹⁰⁷ Exhibition
¹⁰⁸ Restaurant
¹⁰⁹ Street
¹¹⁰ Airport
¹¹¹ Train Station

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 مؤسسه ملی اطلاع‌رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	کد زیر پروژه: پیک‌متن‌فارس - ۲ - خ	ویرایش: ۱/۰	

به دو گروه تقسیم می‌شوند: داده‌های آموزش تمیز و داده‌های آموزش مختلط (هم تمیز و هم نویزی) است. این دو مجموعه آموزش بوسیله فیلتر G.712 فیلتر شده‌اند و شامل ۸۴۴۰ جمله می‌باشند. در دادگان آموزش مختلط ۸۴۴۰ جمله مذکور به ۲۰ زیرمجموعه ۴۲۲ جمله‌ای دسته‌بندی شده‌اند. این زیرمجموعه‌ها در برگیرنده ۴ نوع نویز مترو، همهمه، ماشین و نمایشگاه با چهار نسبت سیگنال به نویز ۵، ۱۰، ۱۵ و ۲۰ دسی‌بل و همچنین چهار زیرمجموعه بدون نویز می‌باشند.

دادگان آزمایش نیز به سه مجموعه تقسیم می‌شوند که با نامهای A، B و C شناخته می‌شوند. دادگان آزمایش A شامل ۴۰۰۴ جمله فیلتر شده با فیلتر G.712 می‌باشند که به چهار زیرمجموعه ۱۰۰۱ جمله‌ای تقسیم شده‌اند. علاوه بر زیرمجموعه‌های جداگانه تمیز مذکور، به هر یک از آنها نیز نویزهای مترو، همهمه، ماشین و نمایشگاه با نسبت‌های سیگنال به نویز ۵-، ۰، ۵، ۱۰، ۱۵ و ۲۰ دسی‌بل افزوده گردیده‌اند. دادگان آموزش B نیز مشابه با دادگان آموزش A ساختاردهی شده‌اند با این تفاوت که در آنها از نویزهای جمع‌پذیر رستوران، خیابان، فرودگاه و ایستگاه قطار استفاده شده است. برای مجموعه B عدم تطبیق میان داده‌های آموزش و آزمایش حتی برای دادگان آموزش مختلط نیز وجود دارد، زیرا دادگان آموزشی مختلط تنها شامل نویزهایی است که به مجموعه A اضافه شده‌اند. دادگان آزمایش C شامل دو زیرمجموعه از چهار مجموعه ۱۰۰۱ جمله‌ای است. در این مجموعه داده‌های گفتار و نویز بوسیله فیلتر MIRS فیلتر شده‌اند و سپس نویزهای مترو و خیابان با نسبت‌های سیگنال به نویز ۵-، ۰، ۵، ۱۰، ۱۵ و ۲۰ دسی‌بل به گفتار افزوده شده‌اند. در این مجموعه نیز داده‌های تمیز بطور جداگانه وجود دارند. به این ترتیب مجموعه C برای آزمایش اثر کانال‌های متفاوت مناسب است.

	عنوان پروژه:		 مؤسسه ملی اطلاع‌رسانی
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		
	عنوان زیر پروژه:		
تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		کد زیر پروژه: پیک‌متن: فارس - ۲ - خ	ویرایش: ۱/۰
تاریخ: ۱۳۸۸/۰۳/۱۹			

14 نتیجه گیری

در این گزارش به بررسی اجمالی پردازش گفتار و معرفی پایگاه های گفتار فارسی موجود پرداختیم. چنانکه دیده شد، هنوز خلا داشتن پایگاه داده گفتاری در اکثر کاربردهای پردازش گفتار دیده می شود. در ادامه ویژگیهای مهمی که در طراحی پایگاه داده گفتار باید مدنظر قرار بگیرد آورده شده است.

۱- تعداد زبان

تعداد زبانهایی که در پایگاه داده ی گفتاری وجود دارد، مانند فارسی، انگلیسی.

۲- تعداد گویندگان

تعداد گویندگان نیز از اهمیت خاص خود برخوردار است.

۳- ویژگی های گویندگان

الف- یک زبانه بودن یا چند زبانه بودن

زبان مادری گوینده فارسی است و یا گوینده ای است که زبان فارسی بعنوان زبان دوم آن محسوب می شود. چراکه زبان مادری برای گوینده حل شده است و هیچ تردیدی در مورد آن ندارد، این امر باعث می شود که ادای جملات با فونوتیک درست صورت پذیرد.

ب- پراکندگی سنی

وجود داده گفتاری از هر گروه سنی می تواند به عملکرد بهتر سیستمهای گفتاری کمک کند.

ج- لهجه (تعداد گویش)

بسته به کاربرد می توان تنوع لهجه مانند خراسانی، یزدی، اصفهانی را در ساخت پایگاه داده ی گفتاری لحاظ نمود و یا می توان فقط از گویش کتابی استفاده نمود.

د- تجربه گوینده

داشتن تجربه قبلی ضبط صدا می تواند باعث کاهش لرزش صدا ناشی از استرس و

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 مؤسسه ملی اطلاع‌رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

اضطراب گردد. همچنین از تبق زدن نیز جلوگیری می شود. در این زمینه می توان از روزنامه نگاران، گویندگان اخبار، بازیگران زن و مرد، راویان و یا خوانندگان بهره جست.

- حالت‌های رفتاری

سیگنال گفتار گوینده در مدهای مختلف رفتاری دارای ویژگیهای متفاوت خواهد بود. عصبانیت، ناراحتی، خوشحالی و ... نمونه هایی از حالات گوینده می باشد که برای ضبط گوینده نبایستی در یکی از این حالات بسر ببرد.

۴- مدهای مکالمه

الف- قالب مکالمه

مکلمه می تواند براساس متن از قبل تنظیم شده باشد و یا مکالمه فی البداهه در شرایط طبیعی باشد.

ب- کیفیت صدا

صدا باید خوشایند، سازگار و دارای کیفیت یکسانی برای تمامی گویندگان در تمامی جلسات ضبط صدا باشد.

ج- رسائی صدا

باید برای ضبط از گویندگانی استفاده شود که دارای رسائی صدا هستند. صداهای گنگ، مبهم و ضعیف باعث کاهش کیفیت سیگنالهای گفتاری خواهد شد.

۵- ویژگیهای سیگنالینگ

الف- فرکانس نمونه برداری

فرکانس نمونه برداری می تواند از 8kHz تا 96kHz متغیر باشد.

ب- اندازه هر نمونه

تعداد بیت مورد نیاز برای نمایش هر نمونه که غالباً ۸ یا ۱۶ بیت می باشد.

۶- ابزار و شرایط ضبط صدا

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 شورای عالی اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

الف- نوع میکروفون

میکروفون می تواند بصورت رومیزی و یا بصورت هدست باشد. همچنین میکروفون هایی وجود دارد که مقداری از نویز را حذف می کنند که این امر ممکن است باعث ایجاد تخریب شیگنال گفتار شود. لذا در استفاده از این نوع میکروفونها بایستی دقت لازم را داشته باشیم.

ب- محیط نویزی یا بدون نویز (نرخ سیگنال به نویز)

دادگان گفتار را می توان در یک اتاق آکوستیک (استودیو ضبط صدا) انجام داد تا دادگان ضبط شده حاوی نویز نباشند. در این حالت در صورت نیاز به دادگان نویزی بطور مصنوعی به دادگان تمیز ضبط شده نویز اضافه می کنیم. همچنین دادگان را می توان در شرایط مختلف نویزی ضبط نمود که در این حالت امتیاز داشتن شرایط طبیعی نویزی را خواهیم داشت.

۷- ویژگی های پایگاه از لحاظ فونتیک

الف- گفتار پیوسته یا گسسته

ب- تعداد جملات / کلمات / فونمها و پراکندگی مناسب آنها


ج- سطح برجسب زنی و دقت آن

د- تنوع موضوع (مانند سیاسی، ورزشی و هنری)

در بالا به ویژگیهای یک پایگاه داده گفتار استاندارد اشاره شد. در ادامه به ویژگیهایی که پایگاه داده برای هر کاربرد بایستی داشته باشد اشاره خواهیم کرد.

۱- بازشناسی گفتار گسسته و پیوسته:

- داشتن دادگان آموزشی و تست بطور مجزا
- داشتن دادگان زیاد جهت آموزش و ارزیابی درست سیستم بازشناسی
- داشتن تنوع گویش، جنسیت، رده های سنی
- داشتن دو بخش دادگان تمیز و نویزی جهت ارزیابی سیستم بازشناسی در شرایط نویزی. این امر می تواند روشهای آموزش multi-condition و clean را فراهم کند.
- برجسب زنی دقیق در سطح واج، سه واج، آوا و کلمه یا برجسب دنباله واجها یا کلمات

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 ژورنال اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

- انتخاب مناسب متون بمنظور داشتن دادگان مناسب از هر نوع واج، سه واج یا کلمه

- داشتن فرکانس نمونه برداری و bit rate مناسب

۲- بازشناسی گفتار مستقل از گوینده

- تمامی موارد گفته شده برای بازشناسی گفتار گسسته و پیوسته

- داشتن گویندگان زیاد بمنظور داشتن تنوع زیاد در هنگام آموزش مدلها

۲- تطبیق گوینده

- تمامی موارد گفته شده برای بازشناسی گفتار گسسته و پیوسته

- داشتن دادگان به اندازه کافی از هر گوینده هم در بخش دادگان آموزش و هم در بخش دادگان تست.

۳- بازشناسی گوینده

- تمامی موارد گفته شده برای بازشناسی گفتار گسسته و پیوسته

- داشتن گویندگان زیاد هم در بخش آموزش و هم بخش تست. بهتر است که کسر بیشتری از گویندگان دو بخش آموزش و تست همپوشانی داشته باشند.

- برای دو روش کلی بازشناسی گوینده وابسته به متن و مستقل از متن تمهیدات لازم دیده شود.

۴- سنتز گفتار

- داشتن دادگان زیاد جهت بکارگیری از آن در سیستم بازشناسی

- برچسب زنی دقیق در سطح واج، سه واج، آوا و کلمه یا برچسب دنباله واجها یا کلمات

- انتخاب مناسب متون بمنظور داشتن دادگان مناسب از هر نوع واج، سه واج یا کلمه

- تنوع لهجه، گوینده و رده سنی انتخابی است (ضروری نمی باشد).

۵- تشخیص کلمات کلیدی

- برای این کاربرد می توان از دادگان تهیه شده برای بازشناسی گفتار بهره جست، بشرطی که بخوبی در سطح واج و کلمه برچسب زنی دقیق داشته باشند.

در این مجموعه تلاش شد تا در ابتدا به معرفی اجمالی شاخه های پردازش گفتار پرداخته شود، سپس نگاهی به ویژگیهای دادگان گفتاری داشته، در انتها نیز بسته با کاربرد پردازش گفتار ویژگی دادگان

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 شورای عالی اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

بررسی گردید.

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 مؤسسه ملی تحقیقات رایانه و فناوری اطلاعات
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

15 منابع و مآخذ

- [1] L. R. Rabiner. "A tutorial on hidden Markov models and selected applications in speech recognition". Proceeding of the IEEE.
- [2] X. D. Huang, Y. Ariki, and M. Jack. "Hidden Markov Models for Speech Recognition". Edinburgh University Press, 1990
- [3] Lalit R. Bahl, Fred Jelinek, and R. L. Mercer. "A maximum likelihood approach to continuous speech recognition". IEEE Transactions on Pattern Analysis and Machine Intelligence, 5(2):179-190.
- [4] Kao, Yu-Hung, "N-best search for continuous speech recognition using viterbi pruning for non-output differentiation states", Texas Instruments Incorporated; April 16, 2002.
- [5] Babak Nasersharif, Ahmad Akbari, Mohammad Mehdi Homayounpour, "Application of HMM adaptation and robust features to sub-band speech recognition in noise", *Proceeding of 11'th international CSI computer conference*, Vol. 2, pp. 215-219, Tehran, Iran, January 2006.
- [6] Babak Nasersharif, Ahmad Akbari, "Sub-band weighted projection measure for robust sub-band speech recognition", *Proceeding of EUROSPEECH*, pp. 945-948, Lisbon, 2005.
- [7] Babak Nasersharif, Ahmad Akbari, "Improved HMM entropy for robust sub-band speech recognition", *Proceeding of 12'Th European Signal Processing Conferences (EUSIPCO)*, Turkey, 2005.
- [8] Zhu, D., Paliwal, K., "Product of power spectrum and group delay function for speech recognition", *IEEE International Conference on Acoustic, Speech, and Signal processing*, vol. 1, pp. 125-128, 2004.
- [9] Ikbali, S., Misra, H., Bourlard, H., "Phase autocorrelation derived robust speech features", *IEEE International Conference on Acoustic, Speech, and Signal processing*, vol. 2, pp. 133-136, 2003.
- [10] Murthy, H.A., Gadde, V., "The modified group delay function and its application to phoneme recognition", *IEEE International Conference on Acoustic, Speech, and Signal processing*, vol. 1, pp. 68-71, 2003.
- [11] Tyagi, V. , McCown, I., Misra, H., Bourlard, H., "Mel-cepstrum modulation spectrum (MCMS) features for robust ASR", *IEEE workshop on Automatic Speech Recognition & Understanding*, pp. 399-404, 2003.
- [12] S. Young, G. Evermann, D. Kershaw, G. Moore, J. J. Odell, D. Ollason, V. Valtchev, P. C. Woodland, "The HTK Book, Version 3.1 ", Cambridge University Engineering Department, 2002.
- [13] X. Huang, A. Acero, H. Hon, " Spoken language processing: a guide to theory,

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 ژورنال اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

algorithm and system development", Prentice Hall, 2001.

[14] P.C.Woodland, "Speaker Adaptation for Continuous Density HMMs: A Review", Cambridge University, 2001.

[15] D.Kim, S.Lee, R.M.Kil, "Auditory Processing of Speech Signals for Robust Speech Recognition in Real-World Noisy Environment", IEEE Trans. Speech & Audio Processing, January 1999.

[16] J.Hakkien, J.Sountasta, R.Hariharan, M.Vasilache, K.Launilla, "Improved Feature Vector Normalization for noise Robust Connected Speech Recognition", EuroSpeech'99, vol.6, pp.1833-2836, 1999

[17] O.Vikki, K.Launilla, "Cepstral Domain Segmental Feature Vector Normalization", Speech Communication, No. 25, pp.133-147, 1998.

[18] M.J.F.Gales, "Model-based Techniques for Noise Robust Speech Recognition", Phd thesis, Cambridge University, 1996.

[19] P.J.Moreno, "Speech Recognition in Noisy Environment", Phd Thesis, Carnegie Mellon University, 1996.

[20] P.C.Woodland, M.J.F.Gales, D.Pye, "Improving Enviromental Robustness in Large Vocabulary Speech Recognition", Proc. ICASSP'96, Vol. 1, pp. 65-68, Atlanta, 1996.

[21] B.A. Carlson, M.A. Clements, "A projection-based likelihood measure for speech recognition in noise", IEEE Trans. on Speech and Audio Processing, vol. 2, no. 1, pp. 97-102, January 1994.

[22] H.Hermansky, N.Morgan, "RASTA Processing of Speech", IEEE Trans. Speech and Audio Processing, Vol.2, No. 4, pp. 587-589, October 1994.

[23] Alexandre, P., Lockwood, P., "Root cepstral analysis: A unified view. Application to speech processing in car noise environments", Speech Communication, Vol. 12, Iss. 3, pp. 277-288, July 1993.

[24] Murthy, H.A, Yegnanarayana, B. "Formant Extraction from group delay function", Speech Communication, vol. 10, pp.209-221, 1991.

[25] A.P.Varga, R.K.Moore, "Hidden Markov Model Decomposition of Speech and Noise", ICASSP'90, pp. 845-848, 1990.

[26] H.Hermansky, "Perceptual Linear Predictive Analysis of Speech", Journal of Acoustical Society of America, pp. 1738-1752, April 1990.

[27] D. Mansour, B. Juang, "A family of distortion measure based upon projection operation for robust speech recognition", IEEE Trans. on Acoustic, Speech and signal processing, vol. 37, pp.1659-1671, 1989.

[28] S.F.Boll, "Suppression of Acoustic Noise in Speech using Spectral Subtraction", IEEE Trans. Acoustics, Speech & Signal Processing, April 1979.

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		 شورای عالی اطلاع رسانی
	عنوان زیر پروژه: تحلیل ابعاد دادگان گفتاری و امکان سنجی تهیه موتورهای بازشناسی گفتار زبان فارسی		
	تاریخ: ۱۳۸۸/۰۳/۱۹	ویرایش: ۱/۰	

[۲۹] بابک ناصرشریف، احمد اکبری، محمد مهدی همایونیپور، "ترکیب روشهای مبتنی بر مدل و پردازش چنبدانندی گفتار برای مقاوم سازی بازشناسی گفتار نسبت به نویز"، سیزدهمین کنفرانس مهندسی برق ایران، دانشگاه زنجان، صفحات ۲۰۹-۲۱۴، اردیبهشت ۱۳۸۴.