



دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)

دانشکده مهندسی برق

پایان نامه کارشناسی ارشد

عنوان پایان نامه:

بهبود عملکرد روش HMM در دیکدر ATP گفتار پیوسته فارسی

نگارش : ساناز علیزاده

استاد راهنما : آقای دکتر صیادیان

بهمن ۱۳۸۵



تاریخ :
شماره :

فرم اطلاعات پایان نامه
کارشناسی ارشد و دکترا

دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)

معاونت پژوهشی
فرم پروژه تحصیلات تکمیلی ۷

مشخصات دانشجو

نام و نام خانوادگی : ساناز علیزاده ✓ دانشجوی آزاد بورسیه معادل
شماره دانشجویی: ۸۳۱۲۳۱۱۸ دانشکده : برق رشته تحصیلی: مخابرات سیستم

نام و نام خانوادگی استاد راهنما : دکتر ابوالقاسم صیادیان

عنوان به فارسی: بهبود عملکرد روش HMM در دیکدر ATP گفتار پیوسته فارسی

عنوان به انگلیسی: Improvement of HMM performance in Persian ATP for continuous speech

نوع پروژه: کارشناسی ارشد ✓ کاربردی بنیادی ✓ توسعه ای نظری
دکترا

تاریخ شروع: ۲۷ تیر ۱۳۸۳ تاریخ خاتمه : ۱۵ بهمن ۱۳۸۵ تعداد واحد: ۶

سازمان تامین کننده اعتبار : مرکز تحقیقات مخابرات ایران

واژه های کلیدی به فارسی : بازشناسی گفتار پیوسته - مبدل سیگنال شنیداری به معادل آن

واژه های کلیدی به انگلیسی : HMM – Continuous speech recognition - ATP

نظرها و پیشنهادهای به منظور بهبود فعالیت های پژوهشی دانشگاه:

استاد راهنما:

دانشجو:

امضاء استاد راهنما : تاریخ:

نسخه ۱: معاونت پژوهشی
نسخه ۲: کتابخانه و به انضمام دو جلد پایان نامه به منظور تسویه حساب با کتابخانه و مرکز اسناد و مدارک علمی

چکیده پایان نامه :

تغییرات مشخصه آکوستیکی واج‌ها تحت متن‌های مختلف موجب شده است که در پیاده‌سازی سیستم‌های بازشناسی گفتار، از واحدهای گفتاری وابسته به متن مانند هجا و نیم‌هجا که اثرات آواهای مجاور را نیز در نظر می‌گیرند، استفاده شود. با توجه به اینکه زبان فارسی از دسته زبان‌هایی می‌باشد که دارای ساختار هجایی ساده‌ای است، در این تحقیق واحد گفتاری نیم‌هجا برای مدلسازی طیفی مورد توجه قرار گرفته است و آزمون‌های متعددی برای تصدیق مطلب فوق در طی انجام پروژه صورت گرفته است. به علت فقدان پایگاه داده فارسی مبتنی بر نیم‌هجا، تلاش زیادی جهت طراحی متن و جملات مورد نیاز برای پایگاه داده در طی این تحقیق انجام پذیرفته است و داده‌های گفتاری مربوط به دو گوینده زن و دو گوینده مرد برای ارزیابی مدلها بیان و ضبط شده و به صورت با سرپرستی در سطح واکه و نیم‌هجا برچسب زده شده است.

با توجه به اینکه تشخیص نیم‌هجاها، ابتدا از آشکارسازی سکوت و واکه آغاز می‌گردد، در اولین قدم تمام توجه ما به آشکارسازی واکه‌ها معطوف شده است. در بازشناسی واکه‌ها، از ترکیب مدل آماری مارکوف پنهان با پارامترهای آکوستیکی مانند انرژی میانگذر استفاده شده است. در این پروژه محدوده واکه‌ها با بهره‌مندی از ویژگی‌های مدل آکوستیکی مانند سادگی، سرعت و ناوابسته بودن آن به گوینده‌ها، مشخص شده است. سپس با ترکیب نتایج حاصل از پارامترهای آکوستیکی و مدل آماری مارکوف پنهان به نتایج بسیار مناسبی در بازشناسی واکه‌ها دست یافتیم. در این پروژه در بهترین حالت‌ها، در گفتار پیوسته به خطای ۸/۹۸٪ و در گفتار گسسته به خطای ۲/۸۷٪ دست یافتیم.

فهرست مطالب

فصل اول - مقدمه

- ۱-۱ پیش گفتار..... ۱
- ۲-۱ مسائل مطرح در سیستم های بازشناسی گفتار..... ۲
- ۳-۱ مقایسه سیستم بازشناسی گفتار پیوسته و گسسته..... ۴
- ۴-۱ طراحی یک سیستم بازشناسی گفتار..... ۵
- ۵-۱ رویکرد این پروژه در جهت اعمال بهبود در بازشناسی گفتار..... ۸
- ۶-۱ ساختار پایان نامه..... ۱۰

فصل دوم-روش های رایج در بازشناسی گفتار

- ۱-۲ پیش گفتار..... ۱۱
- ۲-۲ پیش زمانی پویا..... ۱۱
- ۳-۲ شبکه های عصبی..... ۱۳
- ۴-۲ ماشین های بردار پشتیبان..... ۱۶
- ۵-۲ مدل قطعه بندی نرم..... ۱۹
- ۱-۵-۲ الگوریتم آموزش مدل..... ۲۰
- ۶-۲ مدل آماری مارکوف پنهان..... ۲۲
- ۱-۶-۲ سه مساله اساسی در HMM..... ۲۶
- ۱-۶-۲ ارزیابی (محاسبه احتمال)..... ۲۷
- ۲-۶-۲ دکدینگ..... ۳۰
- ۳-۶-۲ آموزش..... ۳۲
- ۲-۶-۲ مسائل پیاده سازی مدل های مارکوف پنهان..... ۳۵
- ۳-۶-۲ مزایای کاربرد مدل HMM در مدلسازی گفتار..... ۳۹
- ۴-۶-۲ معایب کاربرد مدل HMM در مدلسازی گفتار..... ۴۰
- ۵-۶-۲ نسخه های دیگری از مدل HMM..... ۴۱
- ۷-۲ خلاصه فصل..... ۴۶

فصل سوم - پایگاه داده

۴۷	۱-۳ پیش گفتار.....
۴۸	۲-۳ واحدهای گفتاری برای بازشناسی گفتار.....
۵۱	۳-۳ ویژگی های دادگان طراحی شده.....
۵۱	۱-۳-۳ دادگان مناسب برای آموزش گوینده ها.....
۵۲	۲-۳-۳ دادگان مناسب برای بازشناسی گفتار پیوسته فارسی سیستم نا وابسته به گوینده.....
۵۲	۴-۳ ملاحظات طراحی متن دادگان.....
۵۴	۱-۴-۳ تعداد کل کلمه های طراحی شده.....
۵۴	۲-۴-۳ تعداد جمله های طراحی شده.....
۵۵	۵-۳ خلاصه فصل.....

فصل چهارم - پیاده سازی مدل پیشنهادی

۵۶	۱-۴ پیش گفتار.....
۵۷	۲-۴ نتایج استفاده از پارامتر انرژی میانگذر.....
۵۷	۱-۲-۴ نتایج استفاده از انرژی میانگذر در پایگاه داده گسسته.....
۶۲	۲-۲-۴ نتایج استفاده از انرژی میانگذر در پایگاه داده پیوسته.....
۶۵	۳-۴ قطعه بندی و طبقه بندی واکه به کمک مدل آماری مارکوف پنهان.....
۶۵	۱-۳-۴ نتایج پیاده سازی مدل آماری در پایگاه داده گسسته.....
۶۷	۲-۳-۴ نتایج پیاده سازی مدل آماری در پایگاه داده پیوسته.....
۶۸	۴-۴ مقایسه نتایج پیاده سازی مدل های آماری.....
۶۹	۵-۴ پیشنهادات جهت ادامه کار.....
۷۲	مراجع.....

۱-۱ پیش گفتار

سیستمهای بازشناسی گفتار^۱، از پرکاربردترین سیستمهای کاربردی جهت ارتباط طبیعی انسان و ماشین می‌باشند. اگرچه بسیاری از سیستم‌های بازشناسی گفتار در کاربردهای محدود به موفقیت‌هایی دست یافته‌اند، اما با توجه به اهمیت و کارایی این سیستم‌ها، تلاش برای بهبود و توسعه محدوده عملکرد آنها ادامه دارد. بازشناسی گفتار، فرآیند تبدیل صحبت به متن و بطورکلی تر تبدیل صحبت به رشته‌ای از برچسب‌ها می‌باشد. با توجه به تعریف بازشناسی گفتار، در نگاه اول بازشناسی گفتار فرآیندی ساده به نظر می‌رسد و دو مرحله برای اجرای آن در ذهن تداعی می‌شود:

- تصمیم‌گیری در مورد نوع آوای تلفظ شده
- جستجوی کلمه حاصل از توالی آواهای شناخته شده، در فرهنگ لغاتی که از کلمات بیان شده تهیه شده است.

سیستم‌های بازشناسی گفتار که جهت بازشناسی ارقام و واژه‌ها در دهه ۵۰ طراحی شده بودند، بازده چشمگیری داشتند [۱]. بازشناسی در این سیستم‌ها بر اساس دو مرحله عنوان شده در بالا انجام می‌شد. به نظر می‌آید بتوان با تعمیم تکنیکهای به کاررفته در این سیستم‌ها، این نتایج را برای سیستم‌های پیچیده‌تری تعمیم

^۱ Speech recognition

داد. اما متأسفانه کاربرد تکنیک های استفاده شده در سیستم های با عملکرد و کاربردهای محدود نتوانست نتایج و کارایی مناسبی در سیستمهای پیچیده تر ارائه دهد.

سیستمهایی بازنشاسی گفتار با قابلیت های بهتر دارای پارامترهای بیشتری نسبت به سیستم های ساده می- باشند. دسته بندی این سیستم ها بر اساس پارامترهای گفتار ورودی انجام می شود و تغییر هر یک از آنها در گفتار ورودی موجب پیچیده شدن فرآیند بازنشاسی گفتار می شود. تنوع در این فاکتورها می تواند ناشی از تعدد و تنوع گوینده ها، تنوع در حالت های ادای گفتار، تنوع در لهجه ها و تنوع در محیط پیرامون و یا مربوط به گستردگی دامنه لغات باشد. در حال حاضر هیچ روشی وجود ندارد که بتواند تنوع در تمامی منابع این پارامترها را پوشش دهد. سیستم های طراحی شده با هدف پوشش برخی از این عوامل و برای کاربردهای خاص طراحی شده اند. در بخش بعدی، ۵ عامل اصلی تنوع در گفتار ورودی که موجب پیچیدگی بازنشاسی گفتار می شوند را بررسی خواهیم نمود.

۱-۲ مسائل مطرح در سیستم های بازنشاسی گفتار

مهمترین مساله ای که در طراحی یک سیستم بازنشاسی گفتار باید در نظر گرفته شود، پیوسته و یا گسسته بودن گفتار ورودی می باشد. گفتار گسسته^۱ به گفتاری گفته می شود که در آن گوینده کلمات را بصورت مجزا و با فاصله زمانی مشخصی بیان می کند، در واقع گوینده پس از بیان هر کلمه مکث می کند. در حالیکه در گفتار پیوسته^۲ فاصله زمانی معین و از پیش تعیین شده ای بین کلمات متوالی وجود ندارد و گفتار بصورت رشته ای از لغات متوالی بیان می شود.

یکی دیگر از پارامترهای تعیین کننده در سیستم بازنشاسی گفتار، دامنه لغات قابل بازنشاسی توسط این سیستم می باشد [۲]. پیچیدگی یک سیستم بازنشاسی گفتار پیوسته و یا گسسته با پایگاه داده ای با تعداد لغات محدود، بسیار کمتر از پیچیدگی یک سیستم بازنشاسی گفتار پیوسته و یا گسسته با پایگاه داده ای با دایره لغات گسترده می باشد. هنگامیکه دادگان پایگاه داده محدود است، می توان کلمه را به عنوان واحد گفتاری در نظر

^۱ discrete

^۲ continuous

گرفت و برای هر کلمه یک مدل آکوستیکی طراحی نمود. اما با افزایش دایره لغات، طراحی مدل برای هر یک از آنها و نیز جستجو در چنین پایگاه داده‌ای، بسیار وقت گیر می‌باشد. اکثر سیستم‌های تجاری موجود دارای دامنه لغات محدود می‌باشند. اما نیاز به سیستم‌های با دایره لغات وسیع تر موجب شده است که مدل‌های گفتاری کوچکتری طراحی شوند. در این نوع سیستم‌ها، مدلها از واحدهای زیر لغوی^۱ به جای کلمه برای مدلسازی گفتاری بهره می‌برند. سیستم‌های بازشناسی مبتنی بر واج^۲، هجا^۳ و نیم هجا^۴ نمونه‌هایی از این سیستم‌ها هستند.

همچنین در سیستم‌های بازشناسی گفتار پیوسته از مدل‌های زبانی^۵ یا دستور زبان ساختگی برای محدود ساختن تعداد جملات قابل طراحی با کلمات پایگاه داده استفاده می‌شود. ساده‌ترین نوع مدل زبانی مدلی است که ترتیب خاصی از کلمات را مجاز می‌سازد و نوع رایج تر مدلی است که تقریباً مشابه دستور زبان طبیعی می‌باشد و بصورت قواعد وابسته به متن می‌باشد. ترکیب تاثیر دامنه لغات و دستور زبان می‌تواند با سرگشتگی^۶ اندازه‌گیری می‌شود [۱]. سرگشتگی بیانگر تعداد متوسط لغاتی است که می‌توانند در هر نقطه تصمیم‌گیری رخ دهد که با افزایش پیچیدگی سیستم بازشناسی گفتار افزایش می‌یابد.

تفاوت‌های گویندگان از نظر جنسیت، سن و آهنگ و لحن گفتار نیز موجب تنوع در سیگنال آکوستیکی می‌شود. سیستم‌های بازشناسی وابسته به گوینده^۷ تنها برای یک یا تعداد محدودی گوینده خاص طراحی و آموزش داده می‌شوند، درحالیکه سیستم‌های ناوابسته به گوینده^۸ فرآیند بازشناسی گفتار را مستقل از گوینده و برای تمامی گویندگان انجام می‌دهند. سیستم‌های نوع اول تنها قابلیت ایجاد تنوع در لحن و نرخ گفتار گوینده‌ای که با صدایش آموزش داده شده‌اند را دارند. در حالیکه سیستم‌های نوع دوم می‌توانند تنوع سنی، جنسی و حالت‌های گفتار گوینده‌های متفاوت را پوشش دهند.

¹ Sub-word

² phonem

³ syllable

⁴ Demi-syllable

⁵ Language modeling

⁶ perplexity

⁷ Speaker dependent

⁸ Speaker independent

عامل مهم دیگری که سبب ایجاد تغییر و تنوع در سیگنال گفتار ورودی و در نتیجه تضعیف عملکرد سیستم بازشناسی گفتار می‌شود، نویز می‌باشد. نویز می‌تواند ناشی از عوامل محیطی همچون ترافیک، باران و یا سر و صدای محیط پیرامون و یا ناشی از سرفه و یا بازدم گوینده باشد. نویزهای دیگری نیز که در بخش‌های مختلف سیستم بازشناسی گفتار به وجود می‌آیند موجب کاهش کیفیت سیستم می‌شوند.

۱-۳ مقایسه سیستم بازشناسی گفتار پیوسته و گسسته

هر چند جدا سازی کلمات در سیستم‌های بازشناسی گفتار گسسته به دلیل وجود وقفه بین کلمات به سادگی انجام می‌شود، اما مکث گوینده پس از بیان هر کلمه موجب می‌شود گفتار از حالت طبیعی خارج شود. بازشناسی گفتار پیوسته به سه دلیل بسیار مشکل‌تر از بازشناسی گفتار گسسته می‌باشد:

- نبودن فاصله زمانی معین بین کلمات
- تاثیر آواهای مجاور بر یکدیگر
- تغییر لحن و تلفظ کلمات به علت قرار گرفتن در ساختار جمله

از آنجاییکه در گفتار پیوسته فاصله زمانی معینی بین کلمات وجود ندارد، اولین مشکل، تعیین مرز ابتدایی و انتهایی کلمات به کار رفته در یک جمله می‌باشد تا با توجه به مرز تعیین شده بتوان نوع آنها را تعیین نمود.

مشکل دیگر ناشی از تاثیر آواهای مجاور در کلمات متوالی بر هم می‌باشد که به عنوان اثرات هم تولیدی^۱ شناخته می‌شود و دارای ساختاری پیچیده است. در گفتار پیوسته تاثیر آواهای مجاور بر یکدیگر در درون کلمه و در بین کلمات متوالی بصورت تاثیر آوای انتهایی یک کلمه بر آوای ابتدایی کلمه مجاورش مشاهده می‌شود. اثرات هم تولیدی در گفتار موجب می‌شود تلفظ یک کلمه با توجه به موقعیت آن کلمه در جمله نسبت به کلمات مشابه دیگر تغییر کند. چنانچه در یک سیستم بازشناسی گفتار، تاثیر آواهای متوالی بر یکدیگر مدلسازی شود به این

^۱ co-articulation

سیستم یک سیستم بازشناسی گفتار وابسته به متن^۱ گفته می‌شود و در غیر این صورت یک سیستم بازشناسی گفتار غیر وابسته به متن^۲ می‌باشد.

وقتی کلمات در یک جمله به دنبال یکدیگر بیان می‌شوند، لحن و نوع تلفظ هر یک از کلمات نسبت به حالتیکه بصورت مجزا بیان می‌شوند تغییر می‌کند. همچنین در نوعی از گفتار پیوسته که به آن گفتار بی تکلف^۳ گفته می‌شود و در واقع همان گفتار رایج در میان مردم می‌باشد، ممکن است نوعی اختصار و یا حذف کلمات قرینه به اقتضای توالی آنها در جملات رخ دهد. اثرات دیگر ناشی از سرعت گفتار نیز در این بین خودنمایی می‌کنند که سبب ادغام آواها و تولید آوایی بینابین می‌شود. لازم به ذکر است که این نوع گفتار از قوانین دستوری پیروی نمی‌کند و به همین دلیل بازشناسی این نوع گفتار پیوسته نسبت به گفتار پیوسته نوشتاری که دارای قواعد دستوری مشخص می‌باشد، بسیار پیچیده‌تر است.

۱-۴ طراحی یک سیستم بازشناسی گفتار

سیستم های بازشناسی گفتار اولیه، از روش انطباق با مدل استفاده می‌کردند. همانطور که از نام این روش برمی‌آید، در این روش ابتدا یک قالب^۴ برای هر یک از لغات پایگاه داده طراحی می‌شود. سپس ورودی با کلیه قالب ها مقایسه می‌شود و با توجه به میزان شباهت ورودی با قالب ها، بهترین تطابق، معادل ورودی را مشخص می‌نماید. لازم به ذکر است که یک قالب می‌تواند از میانگین گیری بر روی نمونه های مختلف ضبط شده از تلفظ های متعدد یک کلمه تهیه شود. این روش برای بازشناسی گفتار گسسته با دایره لغات محدود دارای کارایی مناسبی می‌باشد.

اما متأسفانه تعمیم این روش به گفتار پیوسته و گفتار گسسته با پایگاه داده بزرگ به دلیل مشکلات ناشی از تهیه قالب برای این تعداد کلمات و نیز مشکل تعیین مرز بین کلمات امکان پذیر نیست. از آنجاییکه ارتباط بین

¹ Context dependent

² Context independent

³ spontaneously

⁴ prototype

کلمات موجود در پایگاه داده و کلمات به کار رفته در یک جمله از یک رابطه قطعی و یقینی پیروی نمی‌کند، باید از روشهای ریاضی به عنوان جایگزین قالب‌ها استفاده نمود.

در یک سیستم بازشناسی گفتار پیوسته، گفتار بیان شده از طریق میکروفن پس از نمونه برداری با فرکانس ۸ یا ۱۶ کیلو هرتز بصورت یک سیگنال دیجیتال به واحد پیش پردازش وارد می‌شود. سپس بردارهای ویژگی این نمونه‌های فریم بندی شده استخراج می‌گردند. بنابراین سیگنال آکوستیکی حاصل از بیان جمله W ، بصورت رشته‌ای از مشاهدات آکوستیکی درمی‌آید. نوع بردار مشاهده آکوستیکی یا همان بردار ویژگی^۱، بسته به تکنیک به کار رفته در استخراج بردارهای ویژگی می‌تواند از نوع LPC, MFCC و غیره باشد. در نهایت مدل آکوستیکی سیگنال ورودی بصورت زیر خواهد بود:

$$A = a_1, a_2, \dots, a_m \quad a_i \in A \quad (1-1)$$

یک جمله شامل یک رشته از n کلمه و یا بطور کلی تر شامل n واحد آکوستیکی مانند کلمه و یا زیر کلمه می‌باشد. در یک سیستم بازشناسی، پایگاه داده \mathcal{V} کلیه واحدهای آکوستیکی را که هر یک از دنباله‌ای از بردارهای ویژگی تشکیل شده‌اند، پوشش می‌دهد. w_i بیانگر سمبل ورودی در لحظه i می‌باشد.

$$\mathbf{W} = w_1, w_2, \dots, w_n \quad w_i \in \mathcal{V} \quad (2-1)$$

وظیفه یک سیستم بازشناسی گفتار، پیدا کردن محتمل‌ترین دنباله لغات و یا سمبل‌ها از روی دنباله مشاهدات آکوستیکی می‌باشد با کمک رابطه بیز می‌توان $P(W|A)$ را بصورت زیر نوشت:

$$\hat{W} = \arg \max_w P(W|A) \quad (3-1)$$

$$P(W|A) = \frac{P(W)P(A|W)}{P(A)}$$

احتمال مشاهده دنباله آکوستیکی $P(A)$ ، برای یک دنباله مشاهده مشخص به ازای تمام W های ممکن، مقداری ثابت می‌باشد. بنابراین $P(A)$ را در محاسبات مربوط به بیشینه سازی $P(W|A)$ وارد نمی‌کنیم. لذا خواهیم داشت:

¹ Feature vector

$$\hat{W} = \arg \max_w P(W)P(A|W) \quad (۴-۱)$$

حساس‌ترین بخش در بازشناسی گفتار، تعیین احتمال $P(A|W)$ می‌باشد. برای محاسبه این احتمال از روشهای آماری متعددی استفاده می‌گردد. برای محاسبه $P(W)$ نیاز به یک مدل زبانی داریم که با توجه به ماهیت گفتار، مدلی آماری می‌باشد. در نهایت برای پیدا کردن محتمل‌ترین دنباله سمبل‌ها از رابطه (۴-۱)، به یک الگوریتم جستجو نیاز داریم. لذا از روشهای ریاضی برای محاسبه $P(A|W)$ استفاده می‌شود. تمامی روشهای رایج برای محاسبه این احتمال به جز روش پیچش زمانی پویا^۱ روشهای آماری می‌باشند. در کلیه روشهای آماری، مستقل از مدل آکوستیکی انتخاب شده برای بازشناسی گفتار، هر واحد آکوستیکی با یک مدل آماری مدل‌سازی می‌شود.

با توجه به اینکه پارامترهای آماری یک واحد آکوستیکی در طول این واحد تغییر می‌کند و ثابت نمی‌ماند، برای مدل‌سازی آماری هر یک از این واحدها، آن را به قطعه^۲‌هایی تقسیم می‌کنند، به گونه‌ای که بتوان توزیع آماری عناصر آن قطعه را ثابت در نظر گرفت. انتخاب مدل آماری مناسب از چهار جهت حائز اهمیت می‌باشد که عبارتند از:

- مدل‌سازی توالی قطعات
- نمایش همبستگی بردارهای یک قطعه
- نمایش همبستگی بین قطعات و تاثیر قطعات متوالی بر یکدیگر
- در بر گیرندگی تنوع بردارهای ویژگی و تغییرات مشخصات آماری آنها در طول زمان

همچنین انتخاب مناسب طول قطعه به گونه‌ای که مدل آماری منتخب بتواند تنوع در مشخصات آماری و پارامترهای یک سیستم بازشناسی گفتار مانند تنوع در گوینده، لهجه و حالات مختلف گویش را پوشش دهد نیز

^۱ DTW

^۲ Segment

حائز اهمیت می‌باشد. رایج ترین مدل آماری مورد استفاده در بازشناسی گفتار مدل مارکوف پنهان^۱ (HMM) می‌باشد. در فصل دوم مدل مارکوف پنهان و سایر مدل‌های رایج را شرح می‌دهیم.

۱-۵ رویکرد این پروژه در جهت اعمال بهبود در بازشناسی گفتار

کیفیت و کمیت دادگان گفتاری تاثیر بسزایی در بهبود سیستم و الگوریتم های کاربردی بازشناسی گفتار پیوسته دارد. تا کنون تحقیقات وسیعی به منظور تهیه دادگان در زبان های مختلف انجام گرفته است. اکثر دادگان های رایج در این زبانها مبتنی بر واج هستند. فارس دات، مطرح ترین دادگان برای گفتار پیوسته در زبان فارسی نیز بر مبنای آموزش واج های فارسی طراحی شده است [۳].

همانگونه که در بخش ۱-۲ شرح داده شد واج ها به عنوان واحد گفتاری در سیستم های ناوابسته به متن به کار می‌روند. اما در سیستم های وابسته به متن که اثر واج های متوالی بر روی یکدیگر را در نظر می‌گیرند، از دو آوایی ها که اثر آوای سمت راست و یا اثر آوای سمت چپ را نیز لحاظ می‌کنند و یا از سه آوایی ها که اثرات آوای سمت چپ و سمت راست را نیز در مدلسازی دخیل می‌کنند استفاده می‌شود. هجا و نیم هجا واحدهای گفتاری وابسته به متن می‌باشند. زبان فارسی دارای ساختار هجایی ساده ای می‌باشد و طبق تعریف برخی زبان شناسان، دارای سه نوع هجا می‌باشد [۳]:

- ساختار همخوان - واکه (CV)
- ساختار همخوان - واکه - همخوان (CVC)
- ساختار همخوان - واکه - همخوان - همخوان (CVCC)

هجاها دارای اطلاعات خوبی از لحن، آهنگ و نرخ گفتار می‌باشند. در زبان فارسی حدود ۶۰۰۰ هجا وجود دارد. مشکل اصلی استفاده از دو آوایی ها، سه آوایی ها و هجاها، تعداد زیاد این واحدها در زبان فارسی می‌باشد که موجب پیچیده و زمان بر شدن فرآیند جستجو می‌شود. مشکل دیگر استفاده از این واحدها، کمبود داده های آموزشی به دلیل تنوع آنها و نیاز به داده های آموزشی زیاد می‌باشد.

¹ Hidden Markov Modeling

هر هجا به دو نیم‌هجای ابتدایی و انتهای تقسیم می‌شود. نیم‌هجای ابتدایی شامل همخوان ابتدایی و بخشی از هسته هجا و نیم‌هجای انتهایی شامل باقیمانده هسته تا انتهای هجا می‌باشد. تعداد نیم‌هجاهای موجود در زبان فارسی که به صورت CV، VC و VCC می‌باشند حدود ۷۰۰ نیم‌هجا می‌باشد. با توجه به پوشش مناسب و تعداد حالت‌های کم نیم‌هجاها، استفاده از این واحدهای آکوستیکی کارآیی سیستم‌های بازشناسی گفتار پیوسته را بالا می‌برد.

یکی از اهداف این پروژه در راستای بهبود کارآیی سیستم بازشناسی گفتار پیوسته فارسی، تهیه متن یک پایگاه داده مناسب بر پایه نیم‌هجاهای رایج در زبان فارسی می‌باشد. شرح دقیق ملاحظات به کاررفته در تهیه متن این پایگاه داده در فصل سوم ارائه خواهد شد. در این پروژه واحد نیم‌هجا به عنوان واحد گفتاری در نظر گرفته شده است. یکی دیگر از مزایای این واحد گفتاری نسبت به سایر واحدهای آکوستیکی به ویژگی‌های واکه‌ها برمی‌گردد. بررسی منحنی انرژی سیگنال گفتار نشان می‌دهد، واکه‌ها در اکثر موارد از نظر فیزیکی به خوبی تلفظ می‌شوند و اثر آکوستیکی خوبی از خود باقی می‌گذارند. همچنین از منحنی انرژی سیگنال صحبت مشاهده می‌شود که انرژی متوسط واکه‌ها، ۱۰dB تا ۱۵dB بیشتر از انرژی متوسط سیگنال در سایر بخش‌های گفتار می‌باشد. ویژگی دیگر واکه‌ها پریودیک بودن آنها برای تمامی گویندگان می‌باشد. این ویژگی‌ها سبب می‌شوند که تشخیص واکه‌ها و تخمین مرز آنها نسبت به همخوان‌ها آسان‌تر و دقیق‌تر صورت گیرد. لذا این واحد آکوستیکی همزمان از مزایای پوشش مناسب، تعداد حالت‌های کم نیم‌هجاها در زبان فارسی و پایداری نسبی واکه نسبت به شرایط محیطی و سادگی آشکارسازی آن بهره‌مند می‌شود.

بنابراین برای ساده‌سازی بازشناسی گفتار با استفاده از مدل آماری مارکوف پنهان، ابتدا با استفاده از سکوت بین کلمات در یک جمله، مرزهای کلمات را مشخص می‌نماییم. آنگاه با بهره‌مندی از ویژگی‌های واکه‌ها، با تعیین مرز تقریبی آنها در هر کلمه و سپس طبقه‌بندی آنها توسط مدل آماری، جمله را به زیر بخش‌های واحد آکوستیکی تقسیم می‌کنیم. در آخر هر یک از زیر بخش‌ها را به کمک مدل آماری مارکوف پنهان بازشناسی می‌کنیم. در فصل سوم به بررسی ویژگی‌های واکه‌ها که موجب سهولت و دقت در تشخیص و قطعه‌بندی واکه‌ها و نهایتاً منجر به انتخاب واحد آکوستیکی نیم‌هجا شده‌اند، می‌پردازیم.

۱-۶ ساختار پایان نامه

در فصل ۲ مدل‌های آماری رایج برای بازشناسی گفتار مورد بررسی قرار می‌گیرند. از آنجاییکه مدل آماری استفاده شده در این پروژه مدل مارکوف پنهان می‌باشد، پارامترهای این مدل و الگوریتم‌های رایج جهت پیاده‌سازی آن با جزئیات بیشتری نسبت به سایر مدلها شرح داده می‌شود. در ادامه نکات قوت و ضعف این مدل در بازشناسی گفتار و مدل‌هایی که جهت رفع این کاستی ارائه شده اند معرفی می‌شوند.

در فصل ۳ ابتدا به بررسی واحد‌های زبانی زبان فارسی و مزایا و مشکلات استفاده از هر یک از آنها برای کاربرد بازشناسی گفتار می‌پردازیم. سپس دلایل انتخاب واحد گفتاری نیم‌هجا برای بازشناسی گفتار پیوسته فارسی را مورد بررسی قرار می‌دهیم. در پایان نگاهی اجمالی خواهیم داشت بر ویژگی‌های متن پایگاه دادگان تهیه شده در آزمایشگاه تحقیقاتی پردازش اطلاعات که با هدف پوشش کلیه نیم‌هجا‌های زبان فارسی تهیه شده است.

فصل ۴ به ارائه نتایج استفاده از انرژی میانگذر و مدل ترکیبی آکوستیکی آماری اختصاص دارد. در این فصل ابتدا به بررسی ویژگی‌های واکه‌ها و مزایای استفاده از این ویژگی‌ها برای بازشناسی واکه‌ها در گفتار پیوسته و گسسته می‌پردازیم. سپس نتایج استفاده از انرژی میانگذر بر روی گفتار گسسته و پیوسته را ملاحظه می‌کنید. در ادامه نتایج پیاده‌سازی دو روش ترکیبی آکوستیکی و آماری مارکوف پنهان بر روی دادگان گفتاری گسسته و پیوسته ارائه شده است. در انتهای این فصل جمع‌بندی و پیشنهادات برای ادامه این پروژه نیز مطرح شده است.

۲-۱ پیش گفتار

همانطور که پیشتر گفته شده با افزایش پیچیدگی سیستم های بازشناسی گفتار، تهیه قالب برای هر یک از دادگان بانک اطلاعات کاری غیر عملی می باشد. لذا از روشهای ریاضی برای تهیه مدل ها استفاده می نماییم. در این فصل در ابتدا روش غیر آماری پیچش زمانی پویا^۱ و سپس روشهای آماری رایج در بازشناسی گفتار مطرح می شوند. با توجه به موفقیت مدل آماری مارکوف پنهان در بازشناسی گفتار، این مدل رایج ترین مدل آماری مورد استفاده برای بازشناسی گفتار و روش به کار رفته در این پروژه می باشد. بنابراین این مدل را با جزئیات بیشتری نسبت به سایر مدل ها شرح می دهیم.

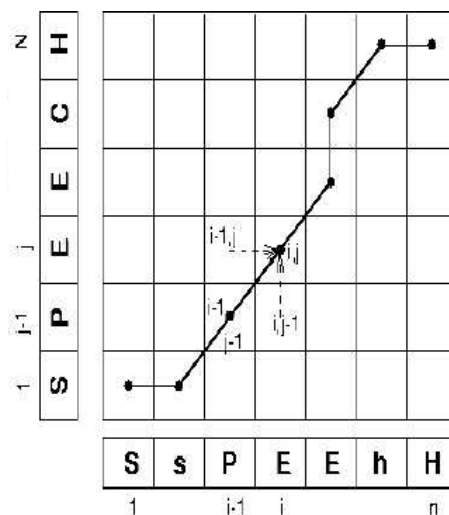
۲-۲ پیچش زمانی پویا

این روش یکی از اولین روش های مورد استفاده در بازشناسی گفتار و تنها مدل غیر آماری در این زمینه به شمار می رود [۱] و [۴]. پیچش زمانی پویا یک روش برنامه ریزی پویا می باشد که از مفهوم تطبیق الگو استفاده می کند. طول زمانی بیان های مختلف یک کلمه توسط افراد مختلف و حتی توسط یک فرد در بیان های گوناگون یکسان نیست. چنانچه گوینده ای کلمه ای را چند بار تکرار نماید، طول هر یک از واجهای ادا شده در کلمه مورد نظر با طول همان واج در تکرار دیگر از این کلمه متفاوت خواهند بود و ممکن است در برخی موارد گوینده یک واج را نسبت به حالت قبل کشیده تر و یا فشرده تر ادا کند. به همین دلیل اگر تعداد فریم های گویش اول M و

^۱ Dynamic Time Warping (DTW)

تعداد فریم های گویش دوم N باشد، نمی توان فریم i ام از گویش اولی را با یک تابع خطی با فریم z ام از گویش دوم متناظر کرد. بنابراین برای استفاده از این روش در بازشناسی گفتار، ابتدا باید گفتار آزمایشی از نظر زمانی نسبت به گفتار مرجع تنظیم شود.

لذا باید از یک تابع غیر خطی برای تطبیق الگو استفاده نمود. الگوریتم معروف دایکسترا^۱ برای یافتن کوتاه ترین مسیر بین دو گره یک گراف و تطبیق الگوها به کار می رود. شکل (۱-۲) نمونه ای از عملکرد این روش را برای منطبق کردن قسمتهای مختلف از دو گویش یک کلمه را نشان می دهد.



شکل (۱-۲). تطبیق دو بیان از یک کلمه با استفاده از DTW

این روش تعیین می کند که بردار ویژگی i_x از گویش اول کلمه به طور بهینه به کدام بردار ویژگی از گویش دوم این کلمه شبیه است و اختلاف کمتری با آن نسبت به سایر بردارهای ویژگی گویش دوم دارد. به عبارتی یک تابع پیچش زمانی با استفاده از تکنیک برنامه ریزی پویا، در فضای دو بعدی اندیس i_x از گویش اول را با اندیس i_y از گویش دوم متناظر می کند. در واقع روش DTW برای تطبیق زمانی دو کلمه و به دست آوردن میزان اختلاف آنها به کار می رود. از این روش می توان برای تطبیق بردارهای ویژگی مجموعه آزمون با بردارهای ویژگیهای مجموعه آموزش که به عنوان مرجع در نظر گرفته می شوند، استفاده نمود. یکی از نقاط ضعف DTW نیاز به داشتن مدل های مرجع برای مقایسه می باشد. با افزایش دامنه لغات از مدلهای زیر لغوی مانند آوا به جای

^۱ Dijkstra

کلمات استفاده می‌شود. اما تهیه مدل های مرجع مناسب از آواها با توجه تغییرات آنها وابسته به متن کاری مشکل می‌باشد. عیب اصلی این روش، کند بودن آن است. اگرچه روش‌های تکاملی برای افزایش سرعت DTW مورد استفاده واقع شده است، اما در هنگامی که تنوع گفتارهای ورودی یا تعداد واحدهای زبانی مورد بازشناسی زیاد باشد، استفاده از DTW در بازشناسی گفتار مناسب نیست.

۲-۳ شبکه های عصبی^۱

شبکه های عصبی با توجه به توان بالا در پردازش موازی و محلی، قابلیت یادگیری، تعمیم، طبقه بندی، به خاطر سپردن و به خاطر آوردن الگوها به صورت محتوایی خیزش وسیعی را در زمینه هایی از قبیل سیستم های هوشمند مصنوعی، شناسایی و کنترل تطبیقی سیستم ها، سیستم های شناسایی گفتار، تصویر و ... ایجاد کرده است [۴-۸].

شبکه های عصبی به جای انجام مراحل الگوریتمی و قدم به قدم، قوانین انجام کار را برای خود تولید می‌نمایند. به بیان دیگر شبکه با روش سعی و خطا، به تدریج نحوه انجام کار را یاد می‌گیرد. عملکرد شبکه های عصبی بر اساس شبیه سازی عملکرد مغز انسان می‌باشد. شبکه عصبی مجموعه‌ای است از عناصر ساده محاسباتی که دارای اتصالات زیادی با یکدیگر می‌باشند. هر عنصر مجموعه که دارای چندین ورودی و یک خروجی است را "نرون" و اتصالات ما بین آنها را "سیناپس" می‌نامند. فرض کنید که بردار ویژگیهای مورد نظر یا ورودی ها بصورت $x = [x_1, x_2, \dots, x_N]^T$ باشد، آنگاه در هر نرون که هسته اصلی یک شبکه عصبی می‌باشد، ورودی x و خروجی y با رابطه (۱-۲) به هم مربوط می‌باشند:

$$y = f\left(\sum_{i=1}^N w_i x_i - Q\right) \quad (1-2)$$

یک نرون در واقع یک ابر صفحه در فضای ویژگی است که دو ناحیه را با یک صفحه از هم جدا می‌کند. مقدار Q مقدار سطح آستانه نرون است و w_i ها بردار وزن هستند و بردار y صفحه را تشکیل می‌دهد. تابع f تابع

¹ Neural Networks

فعالیت واحد عصبی است که به عنوان مثال می‌تواند یک تابع پله، تانژانت هیپربولیک و یا یک تابع سیگموئید باشد. چند نمونه تابع فعالیت که دارای مقادیر بین صفر و یک هستند، در زیر نشان داده شده اند:

$$f(x) = \begin{cases} +1 & x \geq 0 \\ -1 & x < 0 \end{cases}, \quad f(x) = \operatorname{tgh}(\beta x) \quad \beta > 0, \quad f(x) = \frac{1}{1 + e^{-x}}$$

در فضای ویژگی از ابرصفحه‌ها برای نمایش مرزهای یک ناحیه با اطرافش استفاده می‌کنند و در واقع با تعدادی از واحدهای عصبی ناحیه‌ای از فضای بردار ویژگی را از سایر نواحی جدا می‌نمایند. چنانچه هر کلاس از کلاسهای ویژگی را با تعدادی نرون مدلسازی نماییم، عمل دسته بندی یا جدا کردن یک دسته از دسته‌های دیگر با تعداد زیادی از واحدهای عصبی صورت می‌گیرد.

یکی از ویژگی‌های عمده شبکه عصبی این است که با استفاده از آن می‌توان تابع پیچیده غیر خطی را پیاده سازی نمود و فضا را به نواحی مختلف تقسیم نمود. در این شبکه‌ها نیازی به شناخت الگوریتم تولید خروجی از روی داده ورودی وجود ندارد و چنانچه داده ورودی اختلاف اندکی با ورودی مناسب برای تولید یک خروجی مشخص داشته باشد، سیستم اختلاف را در نظر نگرفته و خروجی مطلوب را ایجاد می‌کند. در حقیقت سیستم یک نگاشت بین داده‌های ورودی و خروجی ایجاد می‌کند.

شبکه‌های عصبی را می‌توان به دو دسته پس‌خورد^۱ و پیش‌خورد^۲ تقسیم نمود. اگر خروجی نرونهای یک لایه به ورودی نرونهای لایه قبل متصل باشد شبکه عصبی از نوع پیش‌خورد خواهد بود. برای استفاده از شبکه‌های عصبی ابتدا باید آنها را آموزش داد. در فرآیند آموزش، با استفاده از یک سری داده آموزشی، وزنه‌های اتصالات نرونها تنظیم می‌شوند. در فرآیند آموزش اگر داده‌های آموزشی شامل زوجهای ورودی خروجی باشد شبکه از نوع یادگیری با سرپرستی می‌باشد. یکی از شبکه‌های عصبی معروف در پردازش گفتار شبکه عصبی پرسپترون^۳ نام دارد.

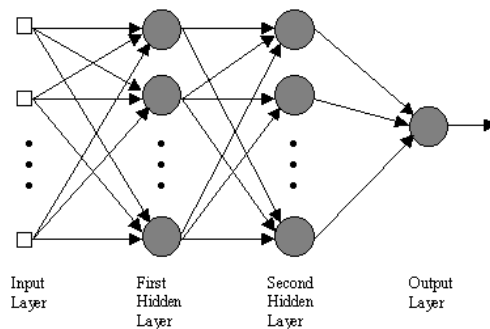
^۱ feedback

^۲ feed forward

^۳ perceptron

- شبکه پرسپترون چند لایه^۱ (MLP)

شبکه MLP یک شبکه متشکل از لایه های ورودی ، مخفی و خروجی است که به ترتیب قرار گرفتن واحدهای عصبی در آن در شکل (۲-۲) نشان داده شده است. گره های ورودی این شبکه به هر یک از المانهای بردار ویژگی متصل خواهند شد. گره های لایه های مخفی هر کدام بیانگر یک ابر صفحه می باشند و هر یک از گره های خروجی متناظر با یکی از کلاسهای فضای ویژگی هستند که در صورت تعلق بردار ویژگی به آن دسته، گره نظیر آن فعال می گردد. البته ممکن است شبکه MLP دارای چند لایه مخفی نیز باشد.



شکل (۲-۲). شبکه پرسپترون چند لایه

الگوریتم یادگیری و تصحیح وزنها در شبکه عصبی MLP، الگوریتم انتشار خطا به عقب^۲ می باشد که همواره در صدد مینیمم کردن خطاهای دسته بندی بردارهای ویژگی است. شبکه MLP استاندارد در حالت عادی یک شبکه استاتیک است و دینامیک بودن و تغییرات زمانی گفتار را در نظر نمی گیرد و این شاید تنها ایرادی باشد که بتوان بر آن گرفت. ولی یکی از مزایای شبکه MLP قدرت تمایز بین کلاسها است.

- شبکه عصبی با تاخیر زمانی (TDNN)

همانطور که قبلا اشاره شد سیگنال گفتار طبیعتی غیر ایستان دارد و مشکل اکثر روش ها آن است که طبیعت دینامیک سیگنال را در نظر نمی گیرند. شبکه عصبی با تاخیر زمانی راهی برای در نظر گرفتن

^۱ Multi Layer Perceptron

^۲ Back propagation