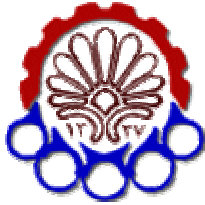


بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه صنعتی امیرکبیر

(پلی تکنیک تهران)

دانشکده مهندسی برق

گروه مخابرات سیستم

پایان نامه کارشناسی ارشد

بهبود عملکرد مدل ماشین‌های بردار پشتیبان
در دیکدر **ATP** گفتار پیوسته فارسی

دانشجو : ملیحه قیدی

استاد راهنما: آقای دکتر صیادیان

بهمن ۸۵

بسمه تعالی

تاریخ :
شماره :



فرم اطلاعات پایان نامه
کارشناسی ارشد و دکترا

دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)

معاونت پژوهشی
فرم پروژه تحصیلات تکمیلی 7

مشخصات دانشجو
نام و نام خانوادگی : ملیحه قیدی

✓ دانشجوی آزاد

بورسیه

معادل

شماره دانشجویی: 83123119

دانشکده : برق

رشته تحصیلی: مخابرات سیستم

نام و نام خانوادگی استاد راهنما : ابوالقاسم صیادیان

عنوان به فارسی: بهبود مدل ماشین های بردار پشتیبان در دیکدر ATP گفتار پیوسته فارسی

عنوان به انگلیسی: Improvement in Support Vector Machines Model in Persian Continuous Speech ATP Decoder

نوع پروژه: کارشناسی ارشد ✓
دکترا

کاربردی ✓

بنیادی

توسعه ای ✓

نظری

تاریخ شروع: 1384/4/27

تاریخ خاتمه: 1385/11/15

تعداد واحد: 6 واحد

سازمان تامین کننده اعتبار : مرکز تحقیقات مخابرات

واژه های کلیدی به فارسی : ماشین های بردار پشتیبان، نیم هجا، گفتار پیوسته، واکه

واژه های کلیدی به انگلیسی : Demi-syllable, coetaneous speech, support vector machine, vowel

نظرها و پیشنهادهای به منظور بهبود فعالیت های پژوهشی دانشگاه:

استاد راهنما:

دانشجو:

امضاء استاد راهنما :

تاریخ:

نسخه 1: معاونت پژوهشی
نسخه 2: کتابخانه و به انضمام دو جلد پایان نامه به منظور تسویه حساب با کتابخانه و مرکز اسناد و مدارک علمی

چکیده

در سیستم‌های بازشناسی گفتار انتخاب واحد گفتاری مناسب، از اهمیت ویژه ای برخوردار است. جهت انتخاب واحد آکوستیکی مناسب، در نظر گرفتن ساختار و ویژگی‌های زبان مورد استفاده، بسیار مهم است. با توجه به این که ساختار هجایی زبان فارسی تقریباً همانند زبان‌های هندی، چینی و ژاپنی ساده و محسوس است، در این تحقیق، واحد زیرکلمه نیم هجا به عنوان واحد آکوستیکی مناسب در زبان فارسی مورد توجه ما قرار گرفته است. با توجه به اینکه پایگاه داده گفتاری مبتنی بر نیم هجاها در زبان فارسی موجود نمی‌باشد، تلاش‌های زیادی جهت طراحی و پیاده سازی پایگاه داده گفتاری مناسب مبتنی بر نیم هجاها در طی انجام این تحقیق، صورت گرفته است. برای ارزیابی مدل‌ها، داده‌های گفتاری مربوط به دو گوینده زن و دو گوینده مرد ضبط شده و به صورت با سرپرستی در سطح واکه و نیم هجا برچسب زده شده است.

در این پایان نامه، به عنوان اولین قدم جهت تشخیص واحدهای نیم هجا در سیگنال گفتار، به آشکارسازی و بازشناسی واکه‌ها پرداخته شده است. سعی شده با ترکیب روش ماشین‌های بردار پشتیبان و روش استفاده از ویژگی‌های آکوستیکی، کارایی سیستم در این بخش تا حد ممکن بهبود داده شود. از پارامترهای آکوستیکی نظیر انرژی میان گذر به دلیل ویژگی‌های مناسبی چون سادگی محاسبات و ناوابسته بودن به گوینده، به منظور تشخیص اولیه محل واکه‌ها استفاده شده است. سپس از قدرت متمایز سازی خوب ماشین‌های بردار پشتیبان جهت طبقه بندی واکه‌ها و تعیین مرز دقیق تر آنها، بهره مند شدیم و به نتایج بسیار مناسبی دست یافتیم. در این تحقیق، برای دادگان گفتار گسسته در صد خطای کل ۱/۶۸٪ و برای گفتار پیوسته در صد خطای کل ۴/۸۶٪ حاصل شد. در مقایسه با نتایج حاصل از دو روش دیگر یعنی مدل مارکوف پنهان و مدل قطعه بندی نرم بر روی همین پایگاه داده، در می‌یابیم که ماشین‌های بردار پشتیبان در کاربردهای طبقه بندی بسیار کارآمد هستند. البته دست یابی به دقت بالا با استفاده از این روش، مستلزم صرف هزینه محاسباتی بیشتر و زمان آموزش طولانی تر خواهد بود.

فهرست مطالب

فصل اول – مقدمه

مقدمه	۱
۱-۱- دو مبحث مهم در طراحی سیستم‌های بازشناسی گفتار	۱
۲-۱- کارهای انجام شده در زمینه تقطیع و برچسب زنی اتوماتیک گفتار	۷
۱-۲-۱- استفاده از مدل مارکوف پنهان	۸
۲-۲-۱- استفاده از شبکه‌های عصبی	۹
۳-۲-۱- استفاده از ماشین‌های بردار پشتیبان	۱۰
۴-۲-۱- مدل قطعه بندی نرم	۱۰
۵-۲-۱- استفاده از ویژگی‌های آکوستیکی	۱۱
۶-۲-۱- روش‌های ترکیبی	۱۱
۳-۱- هدف و ساختار پایان نامه	۱۲

فصل دوم- ماشین‌های بردار پشتیبان

مقدمه	۱۵
۱-۲- ماشین‌های بردار پشتیبان	۱۵
۱-۱-۲- اصل مینیمم سازی ریسک ساختاری	۱۶
۲-۱-۲- ماشین‌های بردار پشتیبان خطی	۲۰
۱-۲-۱-۲- حالت جدایی پذیر خطی	۲۰
۲-۲-۱-۲- حالت جدایی ناپذیر	۲۴
۳-۱-۲- ماشین‌های بردار پشتیبان غیر خطی	۲۵
۱-۳-۱-۲- شرایط Mercer	۲۷
۴-۱-۲- ماشین‌های بردار پشتیبان در حالت چند کلاسه	۲۸
۱-۴-۱-۲- روش یکی در مقابل همه	۲۹
۲-۴-۱-۲- روش یکی در مقابل یکی	۲۹
۲-۲- مقایسه با شبکه‌های عصبی و سایر طبقه بندی کننده‌های غیر خطی	۳۰
۳-۲- مقایسه عملکرد SVM با طبقه بندی کننده RBF کلاسیک	۳۱
۴-۲- نکات عملی پیاده سازی SVM ها	۳۳
خلاصه فصل	۳۷

فصل سوم- انتخاب واحد گفتاری

مقدمه	۳۸
۱-۳- واحدهای گفتاری	۳۸
۱-۱-۳- واج	۳۸

۳۹ ۱-۱-۱-۳-۱-۱-۱-۳ واکه
۳۹ ۲-۱-۱-۳-۱-۱-۳ همخوان ها
۴۰ ۲-۱-۳-۱-۳ هجا
۴۲ ۳-۱-۳-۱-۳ نیم هجا
۴۳ ۲-۳-۱-۳-۱-۳ دادگان مناسب برای بازشناسی گفتار پیوسته فارسی
۴۳ ۱-۲-۳-۱-۳ ملاحظات طراحی متن دادگان
۴۳ ۱-۲-۳-۱-۳ کلمه‌های با ساختار CVC
۴۴ ۲-۲-۳-۱-۳ کلمه‌های با ساختار CVCC
۴۴ ۳-۲-۳-۱-۳ کلمه‌های با ساختار CV
۴۵ ۴-۲-۳-۱-۳ کلمه‌های با ساختار CV تحت متن
۴۵ ۵-۲-۳-۱-۳ تعداد کل کلمه‌های طراحی شده
۴۶ خلاصه فصل

فصل چهارم- نتایج پیاده سازی روش پیشنهادی

۴۷ مقدمه
۴۷ ۱-۴-۱-۴ انرژی
۴۸ ۱-۴-۱-۴ ملایم سازی منحنی انرژی
۴۸ ۲-۴-۱-۴ انرژی میان گذر
۵۱ ۲-۴-۱-۴ نتایج پیاده سازی روش پیشنهادی بر روی دادگان گفتار گسسته
۵۱ ۱-۴-۲-۴ مشخصات دادگان مورد استفاده
۵۱ ۲-۴-۲-۴ استخراج ویژگی
۵۱ ۳-۴-۲-۴ الگوریتم استفاده از انرژی میان گذر
۵۵ ۴-۲-۴ ترکیب مدل SVM و روش استفاده از ویژگیهای آکوستیکی
۵۶ ۱-۴-۲-۴ مرحله آموزش
۵۶ ۱-۴-۲-۴ نرمالیزه کردن داده ها
۵۶ ۲-۴-۲-۴ انتخاب نوع طبقه بندی کننده چند کلاسه
۵۷ ۳-۴-۲-۴ انتخاب نوع تابع هسته
۵۷ ۴-۴-۲-۴ انتخاب پارامترهای تابع هسته و پارامتر C
۵۸ ۵-۴-۲-۴ نحوه انتخاب داده های آموزش
۵۹ ۲-۴-۲-۴ مرحله تست
۵۹ ۱-۴-۲-۴ نحوه ارزیابی خروجی SVM
۶۰ ۲-۴-۲-۴ نحوه انتخاب مجموعه دادگان تست
۶۰ ۲-۴-۲-۴ بهبود نتایج حاصل از انرژی میان گذر با استفاده از SVM ها
۶۳ ۵-۴-۲-۴ مقایسه نتایج روش SVM با روش های HMM و قطعه بندی نرم در گفتار گسسته
۶۳ ۱-۵-۲-۴ تحلیل نتایج مدل ترکیبی

- ۶۴.....SVM از مدل SVM استفاده از مدل SVM ۶-۲-۴ تعیین مرز واکه ها تنها با استفاده از مدل SVM
- ۶۵.....نتایج پیاده سازی روش پیشنهادی در دادگان گفتار پیوسته فارسی ۳-۴-۳
- ۶۵.....۱-۳-۴ مشخصات دادگان مورد استفاده ۳-۴-۱
- ۶۶.....۲-۳-۴ استخراج ویژگی ۳-۴-۲
- ۶۶.....۳-۳-۴ الگوریتم مورد استفاده ۳-۴-۳
- ۶۹.....۵-۳-۴ ترکیب مدل SVM و روش استفاده از اطلاعات آکوستیکی ۳-۴-۵
- ۷۰.....۱-۵-۳-۴ بهبود نتایج حاصل از انرژی میان گذر با استفاده از SVM ها ۳-۴-۵-۱
- ۷۱.....۲-۵-۳-۴ مقایسه نتایج روش SVM با روش های HMM و قطعه بندی نرم در گفتار پیوسته ۳-۴-۵-۲
- ۷۱.....۳-۵-۳-۴ تحلیل نتایج مدل ترکیبی ۳-۴-۵-۳

نتیجه گیری و پیشنهادات ۷۳

مراجع ۷۶

فهرست شکل‌ها

- شکل ۱-۱-۱- بلوک دیاگرام کلی سیستم بازشناسی گفتار بر اساس واحدهای آکوستیکی..... ۳
- شکل ۲-۱- مثالی از مسئله تفکیک دو کلاس توسط تخمین ML که سطح آستانه تصمیم‌گیری ناشی از آن ایده آل نیست و در ناحیه خاکستری رنگ احتمال خطا وجود دارد. ۶
- شکل ۳-۱- ساختار سلسله‌مراتبی HMM ها برای بازنمایی واج‌ها، کلمات و جملات..... ۹
- شکل ۴-۱- بلوک دیاگرام کلی سیستم ترکیبی تشخیص واژه..... ۱۲
- شکل ۱-۲- تفاوت طبقه‌بندی کننده حاصل از ریسک تجربی و ریسک ساختاری، هر سه ابرصفحه C_0, C_1, C_2 ریسک تجربی را حداقل می‌کنند، اما تنها ابرصفحه C_0 بهینه است و ریسک ساختاری را حداقل می‌کند. ۱۸
- شکل ۲-۲- اصل مینیمم‌سازی ریسک ساختاری در واقع تلاش برای پیدا کردن نقطه بهینه در منحنی باند ریسک واقعی میباشد..... ۱۸
- شکل ۳-۲- حداکثر سه نقطه در فضای دو بعدی توسط خطوط جهت‌دار مستقیم، شکسته می‌شود..... ۱۹
- شکل ۴-۲- ابرصفحه جداکننده دو کلاس در حالت داده‌های جدایی‌پذیر خطی، بردارهای پشتیبان با حلقه‌ای دور نقاط مشخص شده‌اند..... ۲۰
- شکل ۵-۲- ابرصفحه جداکننده در حالت داده‌های جدایی‌ناپذیر، بردارهای پشتیبان با حلقه‌ای دور نقاط مشخص شده‌اند..... ۲۵
- شکل ۶-۲- ماشین‌های بردار پشتیبان در حالت غیر خطی، قدرت SVM ها در این است که می‌توانند داده‌ها را به فضای با دیمانسیون بالا منتقل کنند و در آن فضا یک طبقه‌بندی کننده خطی بسازند..... ۲۷
- شکل ۷-۲- بردارهای پشتیبان (مراکز خوشه‌ها) که به طور خوارک توسط SVM ها بدست آمده، با دایره اضافی روی نمونه‌ها مشخص شده است. مراکز خوشه‌ها که با استفاده از روش RBF کلاسیک حاصل شده، با ضربدر مشخص شده است..... ۳۲
- شکل ۸-۲- اثر استفاده از مجموعه ارزیابی بر عملکرد SVM..... ۳۵
- شکل ۹-۲- هیستوگرام خروجی SVM برای نمونه‌های مثبت و منفی با استفاده از روش cross-validation..... ۳۶
- شکل ۱۰-۲- یک نمونه از تخمین با استفاده از تابع توزیع سیگموئیدی برای یک طبقه‌بندی کننده خاص..... ۳۶
- شکل ۱-۴- منحنی انرژی برای جمله "دستی از غیب او را یاری داد."، این منحنی با پنجره همپینگ به طول ۱۳ هموار شده است..... ۵۰
- شکل ۲-۴- منحنی انرژی میان‌گذر برای جمله "دستی از غیب او را یاری داد."، این منحنی با پنجره همپینگ به طول ۱۳ هموار شده است..... ۵۰
- شکل ۳-۴- پنجره اعمال شده برای محاسبه انرژی میان‌گذر..... ۵۲
- شکل ۴-۴- منحنی انرژی میان‌گذر برای کلمه "جت"، به همراه مرز واژه حاصل از تقطیع دستی و مرز واژه حاصل از انرژی میان‌گذر..... ۵۵
- شکل ۵-۴- منحنی انرژی میان‌گذر برای کلمه "نیم"، به همراه مرز واژه حاصل از تقطیع دستی و مرز واژه حاصل از انرژی میان‌گذر، در این حالت استفاده از انرژی میان‌گذر منجر به خطای درج شده است..... ۶۲
- شکل ۶-۴- بهبود حاصل از اعمال مدل SVM برای کلمه "نیم"، خطای درج حاصل از روش انرژی میان‌گذر از بین رفته است..... ۶۲

شکل ۴-۷- منحنی انرژی میان گذر برای جمله " دستی از غیب او را یاری داد "، به همراه مرز واکه حاصل از تقطیع دستی و مرز واکه حاصل از انرژی میان گذر..... ۶۹

شکل ۴-۸- بلوک دیاگرام کلی روش پیشنهادی..... ۷۲

فهرست جداول

- جدول ۳-۱- واژه های زبان فارسی ۳۹
- جدول ۴-۱- نتایج حاصل از پیاده سازی روش استفاده از انرژی میان گذر بر روی دادگان گفتار گسسته ۵۳
- جدول ۴-۲- نتایج پیاده سازی مدل ترکیبی در گفتار گسسته بر مجموعه آموزش ۶۱
- جدول ۴-۳- نتایج پیاده سازی مدل ترکیبی در گفتار گسسته بر مجموعه آزمون ۶۱
- جدول ۴-۴- مقایسه نتایج روش SVM با روش های HMM و قطعه بندی نرم در گفتار گسسته بر مجموعه آزمون ۶۳
- جدول ۴-۵- نتایج پیاده سازی مدل SVM در گفتار گسسته بر مجموعه آموزش ۶۵
- جدول ۴-۶- نتایج پیاده سازی مدل SVM در گفتار گسسته بر مجموعه آزمون ۶۵
- جدول ۴-۷- نتایج حاصل از پیاده سازی روش استفاده از ویژگیهای آکوستیکی بر روی دادگان گفتار پیوسته ۶۸
- جدول ۴-۸- نتایج پیاده سازی مدل ترکیبی در گفتار پیوسته بر مجموعه آموزش ۷۰
- جدول ۴-۹- نتایج پیاده سازی مدل ترکیبی در گفتار پیوسته بر مجموعه آزمون ۷۰
- جدول ۴-۱۰- مقایسه نتایج روش SVM با روش های HMM و قطعه بندی نرم در گفتار پیوسته بر مجموعه آزمون ۷۱

فصل اول

مقدمه

مقدمه

سیستم‌های بازشناسی گفتار پیوسته، یکی از پرکاربردترین سیستم‌های کاربردی جهت ارتباط طبیعی انسان و ماشین می‌باشند. هدف اصلی این سیستم‌ها، تبدیل سیگنال صحبت به متن می‌باشد. اغلب سیستم‌های بازشناسی گفتار، شامل دو زیر مرحله کلی می‌باشند: (۱) تبدیل سیگنال صحبت به سمبل یا در واقع تبدیل آکوستیک به فونتیک^۱ (۲) تبدیل سمبل به متن. در مرحله اول سیگنال صحبت به قطعات کوچکتر (واحدهای آکوستیکی) تقسیم می‌شود و به هر یک از این واحدهای آکوستیکی، یک سمبل یا فونتیک اختصاص داده می‌شود. در این پایان نامه هدف، ایجاد بهبود در این مرحله از بازشناسی گفتار می‌باشد. در مرحله دوم، این رشته سمبل‌ها با استفاده از فرهنگ لغات و در نظر گرفتن قواعد دستوری و معنایی زبان، به متن تبدیل می‌شوند.

۱-۱- دو مبحث مهم در طراحی سیستم‌های بازشناسی گفتار

هر سیستم بازشناسی گفتار پیوسته مبتنی بر واحدهای آکوستیکی، از قسمت‌های اساسی نشان داده شده در شکل ۱-۱ تشکیل شده است. در بلوک اول، ویژگی‌های مناسب برای توصیف خواص سیگنال گفتار استخراج می‌شود. در مرحله بعد واحد گفتاری مناسب برای سیستم انتخاب می‌شود، سپس مدل آکوستیکی-آماری مناسب برای بازشناسی واحدهای آکوستیکی، آموزش داده می‌شود. در بلوک سوم با استفاده از این مدل‌ها، واحدهای آکوستیکی بازشناسی می‌شوند. پس از بازشناسی واحدهای آکوستیکی، با استفاده از فرهنگ لغات بر مبنای تلفظ-های واقعی، تطبیق در سطح کلمه انجام می‌شود. در نهایت با در نظر گرفتن مدل‌های دستوری و معنایی زبان، جمله متناسب با سیگنال گفتار ورودی، بازشناسی می‌شود.

در این پروژه صرفاً "بهبود عملکرد بلوک‌های (۲) و (۳) مورد توجه و تحقیق می‌باشد. برای انجام وظایف بلوک شماره (۱) از الگوریتم‌های آماده استفاده می‌شود.

¹ Acoustic to Phonetic (ATP)

در سیستم‌های بازشناسی گفتار، انتخاب واحد آکوستیکی در بلوک (۲)، از اهمیت ویژه‌ای برخوردار است. برای سیستم‌های بازشناسی گفتار گسسته با تعداد لغات کم، معمولاً از واحد کلمه استفاده می‌شود. با افزایش دایره لغات، استفاده از کلمه به عنوان واحد آکوستیکی، سبب ایجاد محدودیت‌هایی در آموزش، ذخیره و جستجوی واحد کلمه می‌شود. برای حل این مشکل به واحدهای زیر کلمه مانند واج^۱، هجا^۲ و نیم هجا^۳ روی آورده شده است [1],[2],[3],[4]. در فصل سوم راجع به این واحدهای آکوستیکی بیشتر توضیح داده خواهد شد. واحد آکوستیکی مرسوم در بلوک (۲) معمولاً "واج می‌باشد. واج‌ها تحت متن‌های مختلف از نظر مشخصه آکوستیکی به شدت تغییر می‌کنند و در واقع باید اثرات واج قبل و بعد از آنها نیز لحاظ شود. نشان داده شده است که مدلسازی وابسته به متن^۴، بازدهی بازشناسی بسیار بهتری دارد [5],[6],[7],[8]. به این ترتیب روش‌های دو آوایی^۵ و سه آوایی^۶ مطرح شدند و تحقیقات بسیاری در این زمینه انجام شد [9],[10],[11]. اما یکی از مشکلات این روش‌ها این است که تعداد حالت‌های آنها زیاد است. بنابراین حجم جستجو و پردازش بالا می‌رود. هدف ما مدلسازی مناسب واج‌ها تحت متن می‌باشد. یکی دیگر از مدل‌های وابسته به متن، استفاده از واحد آکوستیکی هجا می‌باشد. با توجه به اینکه ساختار هجایی زبان فارسی همانند زبان‌های هندی، چینی و ژاپنی نسبتاً ساده و محسوس است، استفاده از هجا به عنوان واحد زیر کلمه برای این نوع زبان‌ها مناسب می‌باشد. اما یکی از مشکلات استفاده از این واحد آکوستیکی نیز تعداد حالت‌های زیاد آن می‌باشد. یک روش برای حل این مشکل، می‌تواند استفاده از واحدهای نیم هجا باشد. بنابراین با توجه به پوشش مناسب و تعداد حالت‌های کم نیم هجا، استفاده از این واحد آکوستیکی کارآیی سیستم‌های بازشناسی گفتار پیوسته را بالا می‌برد. لذا طی دو دهه اخیر تحقیقاتی در زمینه بازشناسی گفتار مبتنی بر نیم هجاها در زبانهای دیگر انجام گرفته است [6],[7],[12],[13],[14]. در زبان فارسی در این قسمت، تحقیق و توجه کمتری شده است. این امر انگیزه‌ای برای انجام تحقیق در زمینه بازشناسی گفتار مبتنی بر نیم هجاها در زبان فارسی شد.

¹ Phoneme

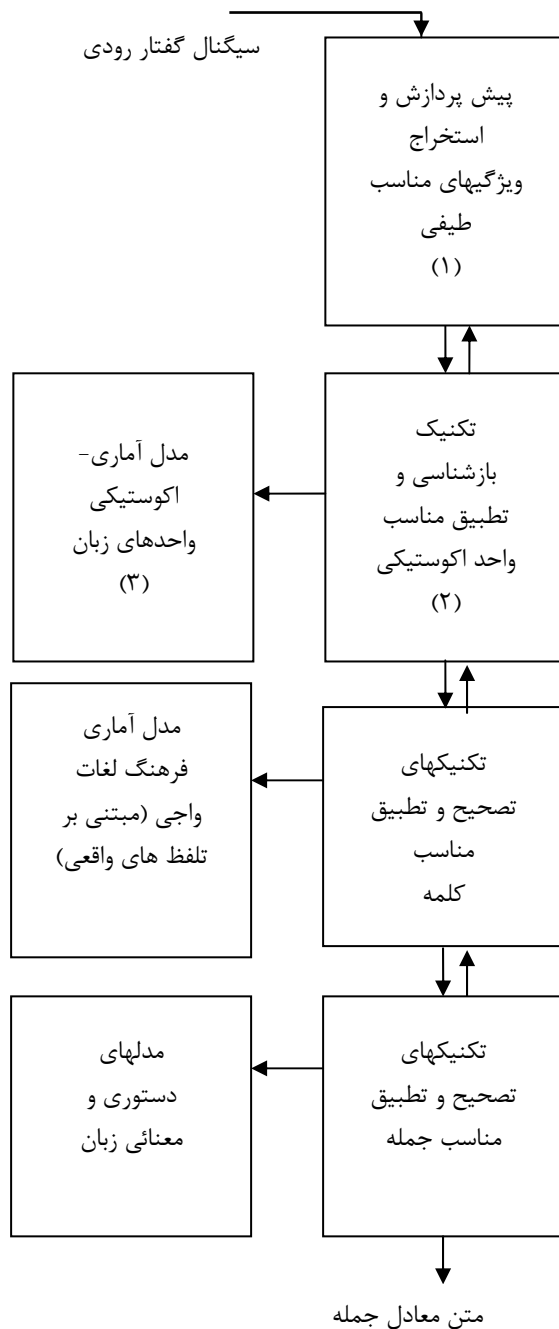
² Syllable

³ Demi-syllable

⁴ Context-dependent

⁵ Diphone

⁶ Triphone



شکل ۱-۱-۱- بلوک دیاگرام کلی سیستم بازشناسی گفتار بر اساس واحدهای آکوستیکی

با توجه به این که واکه‌ها^۱ از نظر مشخصه فرکانسی و زمانی نسبت به سایر واج‌ها پایدارتر هستند، در بازشناسی گفتار بسیار پراهمیت هستند. ثابت شده است که در یک سیگنال گفتار واکه‌ها برای بازشناسی خودکار

^۱Vowels

گفتار کاملاً مناسب می‌باشند [1], [15], [16], زیرا معمولاً به اندازه کافی طولانی بوده و انرژی کافی برای تشخیص در سیگنال صوتی را دارند. یکی دیگر از ویژگی‌ها واکه‌ها پریودیک بودن آنها می‌باشد. با استفاده از این ویژگی‌ها می‌توان با محاسبات نسبتاً ساده به محل تقریبی واکه‌ها در سیگنال گفتار دست یافت. بررسی‌ها نشان می‌دهد که این ویژگی واکه‌ها تقریباً مستقل از گوینده و نیز مقاوم در برابر شرایط محیطی (از جمله نویز) می‌باشد. لذا با استفاده از آنها می‌توان نقاط قابل اعتمادی در سیگنال گفتار پیوسته پیدا کرد.

در این پایان‌نامه، نیم‌هجاها و واکه‌ها به عنوان بهترین واحد زیر کلمه برای گفتار پیوسته فارسی انتخاب شده است. به این ترتیب که برای هر واکه دو بخش ایستان و گذرا در نظر گرفته می‌شود. بخش ایستان واکه تقریباً نشان‌دهنده مشخصات خود واکه می‌باشد و اثر آواهای قبل و بعد از آن نسبتاً ناچیز است. بنابراین ابتدا قسمت ایستان واکه‌ها که تشخیص آنها ساده‌تر و مطمئن‌تر است، در سیگنال گفتار شناسایی می‌شوند، سپس همخوان^۱ (یا همخوان‌ها) و قسمت گذرای واکه به عنوان نیم‌هجا مدل سازی می‌شود. به این ترتیب اثرات متقابل واکه‌ها و همخوان‌ها در نظر گرفته می‌شود. البته در این پروژه فقط قسمت اول یعنی آشکارسازی محل واکه‌ها و طبقه‌بندی آنها انجام شده است. سعی شده با استفاده از ترکیب روش استفاده از ویژگی‌های آکوستیکی و روش آماری، کارایی سیستم در این بخش، تا حد ممکن بهبود داده شود. انجام قسمت دوم جزء کارهای آتی می‌باشد.

یکی دیگر از مباحث مهم در طراحی سیستم‌های بازشناسی گفتار، نحوه مدل‌سازی آکوستیکی-آماری سیگنال گفتار در بلوک (۳) می‌باشد. مدل‌های مارکوف پنهان^۲ با چگالی مشاهدات مخلوط گوسی^۳ روش غالب در سیستم‌های بازشناسی گفتار می‌باشد. این سیستم‌ها معمولاً از قدرت بازنمایی^۴ خوب برای مدل سازی استفاده می‌کنند [17]. یکی از علل مهم موفقیت HMM‌ها در بازشناسی گفتار، توانایی آنها در مدل کردن حالت‌های زمانی سیگنال گفتار، توسط فرایند مارکوف می‌باشد. تابع توزیع احتمالی که برای هر حالت در مدل HMM در نظر گرفته می‌شود، تغییرات سیگنال گفتار ناشی از گوینده‌های مختلف و یا ناشی از متون مختلف آوایی را مدل می‌کند. این تابع توزیع معمولاً مخلوطی از توابع گوسی در نظر گرفته می‌شود.

¹Consonant

²Hidden Markov Model (HMM)

³Gaussian Mixture Model (GMM)

⁴Representation

برای تخمین پارامترهای HMM، تخمین بر اساس ماکزیمم کردن درستنمایی^۱، یکی از مؤثرترین روش-هاست. به هر حال مشکلاتی در رابطه با تخمین ML برای کاربردهایی همچون بازشناسی گفتار وجود دارد. تخمین ML سعی در بازنمایی بهتر یک کلاس دارد ولی از نظر ایجاد تمایز میان دو کلاس ضعیف می‌باشد، در حالی که قدرت تمایز بالا برای سیستم‌های بازشناسی گفتار، یک نیاز اساسی می‌باشد [17], [18], [19], [20]. در تخمین ML پارامترهای مدل تنها براساس داده‌های مربوط به کلاس مورد نظر^۲ و بدون در نظر گرفتن داده‌های خارج از کلاس مورد نظر^۳، تخمین زده می‌شوند [18]. شکل ۱-۲ یک مثال ساده از این مشکلات را نشان می‌دهد. دو کلاس با تابع توزیع یکنواخت به صورت کاملاً جدا از هم نشان داده شده است. از تخمین ML برای بازنمایی این دو کلاس به صورت گوسی و از قاعده بییز^۴ برای طبقه بندی داده‌ها استفاده شده است. همان‌طور که مشاهده می‌شود سطح آستانه تصمیم‌گیری که از تخمین ML به دست آمده است، با سطح تصمیم‌گیری بهینه فاصله دارد و در این محدوده احتمال خطای قابل ملاحظه‌ای وجود دارد. این مساله برای تمامی داده‌هایی همانند سیگنال گفتار که در فضای ویژگی‌ها همپوشانی دارند و مرز کلاس‌ها در همسایگی یکدیگر قرار دارد، وجود دارد [18].

در این مثال دیدیم که آموزش ML یک مدل گوسی، هیچگاه به طبقه بندی کامل دست نمی‌یابد. اما نکته مهم این است که مدل‌های گوسی لزوماً انتخاب نامناسبی نیستند، بلکه این مطلب نشان دهنده نیاز به روش‌های با قدرت تمایز^۵ بالا برای ایجاد مدل‌های مقاوم‌تر و دقیق‌تر می‌باشد [18].

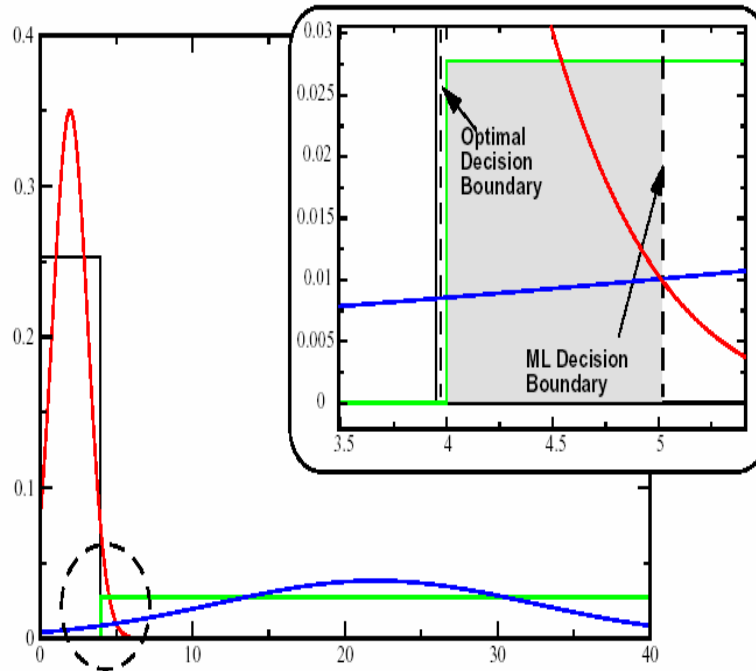
¹ Maximum Likelihood (ML)

² In-class data

³ Out-of-class data

⁴ Bayes' rules

⁵ Discriminative



شکل ۱-۲ - مثالی از مسئله تفکیک دو کلاس توسط تخمین ML که سطح آستانه تصمیم گیری ناشی از آن ایده آل نیست و در ناحیه خاکستری رنگ احتمال خطا وجود دارد [18].

یکی دیگر از موارد مطرح در بحث بازشناسی، قدرت تعمیم دهی^۱ سیستم است. عملکرد سیستم‌های HMM در یک سیستم حلقه بسته^۲ بسیار خوب است اما در یک سیستم حلقه باز^۳ کارایی آن به طور قابل ملاحظه ای کاهش می‌یابد. منظور از یک سیستم حلقه بسته این است که شرایط آموزش^۴ و آزمون^۵ یکسان باشد و همین طور، منظور از یک سیستم حلقه باز این است که شرایط آموزش و آزمون متفاوت باشد. همچنین سیستم‌های HMM در یک سیستم وابسته به گوینده به مراتب بهتر از سیستم ناوابسته به گوینده عمل می‌کنند [17].

¹Generalization

²Closed-loop

³Open-loop

⁴Training

⁵Test

در این پایان نامه یک روش نسبتاً جدید بر اساس اصل مینیم سازی ریسک ساختاری^۱ بررسی می-شود که از یک قالب با قدرت تمایز بالا بر اساس ماشین‌های بردار پشتیبان^۲ برای بازشناسی گفتار استفاده می‌کند. ماشین‌های بردار پشتیبان توانایی همزمان قدرت تعمیم دهی و قدرت تمایز مدل‌های آکوستیکی را دارند [17], [18], [19], [20].

ماشین‌های بردار پشتیبان در چند سال گذشته در بسیاری از کاربردهای طبقه بندی موفقیت‌های چشمگیری داشته اند [18]. مجموعه ای که باعث شد، ماشین‌های بردار پشتیبان در اوایل دهه ۹۰ بسیار برجسته شوند، داده های رقمی سرویس پستی آمریکا بود که SVM ها بهترین ارقام را گزارش دادند [21]. ابتدا تلاش‌های زیادی برای به کارگیری SVM ها برای بازشناسی گوینده شد [22], [23]. اما این تلاش‌ها به دلیل کمبود روش‌های پیاده سازی موثر برای تخمین SVM ها در آن زمان، موفقیت‌های محدودی داشتند. بعدها با گسترش روش های بهینه سازی، ماشین‌های بردار پشتیبان برای طبقه بندی های وسیع تر مثل طبقه بندی متون، تشخیص چهره و غیره استفاده شدند. موفقیت های چشمگیر ماشین های بردار پشتیبان در بسیاری از کاربردهای طبقه بندی موجب شد که از آن ها در سیستم های بازشناسی گفتار استفاده شود. اما همچنان که ملاحظه می‌شود، تمامی کاربرد هایی که ذکر شد، یک ویژگی مشترک دارند، این که همه این طبقه بندی ها استاتیک هستند. SVM ها برای ساختار زمانی داده ها طراحی نشده اند و قادر نیستند حالت دینامیک گفتار را مدل کنند، برای رفع این مشکل معمولاً از SVM ها به صورت ترکیب با HMM در بازشناسی گفتار استفاده می‌شود [18], [19], [24].

۱-۲- کارهای انجام شده در زمینه تقطیع و برچسب زنی اتوماتیک گفتار

به طور کلی هر مبدل آکوستیک به فونتیک می‌تواند شامل دو قسمت زیر باشد:

¹ Structural-Risk Minimization

² Support Vector Machines

(۱) تقطیع^۱: در این مرحله سیگنال گفتار به نواحی گسسته در زمان تقسیم می‌شود که در این نواحی، خواص آوایی سیگنال، نمایش دهنده یک واج (یا چندین واج) می‌باشد. خطا در تعیین محل واج ها نه فقط به خاطر محدودیت الگوریتم ها بلکه به دلیل ابهام ذاتی در آنالیز سیگنال گفتار نیز ایجاد می‌شود و به طور کلی حدود دقیق یک قطعه را نمی‌توان تعیین کرد.

(۲) برچسب زنی^۲: در این مرحله به هریک از قطعه های به دست آمده از قسمت قبل، یک برچسب نسبت داده می‌شود.

از آنجاییکه هدف این پایان نامه ایجاد بهبود در مبدل آکوستیک به فونتیک گفتار پیوسته فارسی می‌باشد، در این بخش، به بررسی کارهای انجام شده در زمینه تقطیع و برچسب زنی اتوماتیک گفتار در مقالات می‌پردازیم.

۱-۲-۱- استفاده از مدل مارکوف پنهان

در این دسته از روش ها سیگنال گفتار به صورت آماری مدلسازی می‌شود. مدل های مارکوف پنهان، حالت های زمانی سیگنال گفتار را با استفاده از زنجیره مارکوف^۳، مدل می‌کنند. این مدل ها، احتمال مشاهده یک رشته زمانی^۴ را بدون دانستن رشته حالت هایی^۵ که موجب تولید این رشته مشاهدات^۶ شده اند، بدست می‌آورند. تابع توزیع احتمالی که برای هر حالت در مدل HMM در نظر گرفته می‌شود، تغییرات سیگنال گفتار ناشی از گوینده های مختلف و یا ناشی از متون مختلف آوایی را مدل می‌کند. این تابع توزیع معمولاً مخلوطی از توابع گوسی در نظر گرفته می‌شود. شکل ۱-۳ نشان می‌دهد که چگونه در نظر گرفتن حالت ها و انتقال حالات^۷ در یک مدل HMM می‌تواند به صورت سلسله مراتبی^۸، واجها، کلمات و جملات را نمایش دهد. از مدل های HMM به طور وسیعی برای مدلسازی واحد های گفتاری در سطوح مختلف استفاده شده است [12], [25], [26], [27]. همان-

¹ Segmentation

² Labeling

³ Markov chain

⁴ Temporal sequence

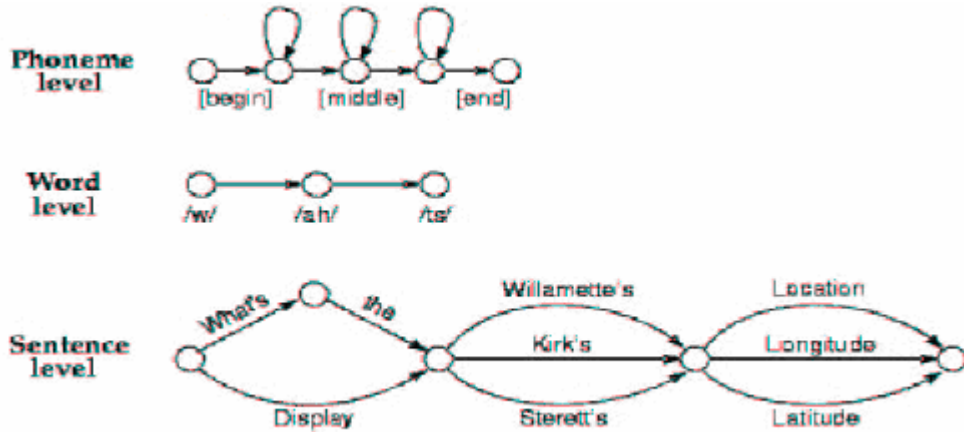
⁵ State sequence

⁶ Observation sequence

⁷ State transition

⁸ Hierarchically

طور که در بخش قبل توضیح داده شد، مدل HMM متداول‌ترین روش در سیستم‌های بازشناسی گفتار می‌باشد. این مدل دارای قدرت بازنمایی بالایی است. اما برای ایجاد مدل‌های مقاوم‌تر و دقیق‌تر نیاز به روش‌های با قدرت تمایز بالا برای سیستم‌های بازشناسی گفتار وجود دارد.



شکل ۱-۳- ساختار سلسله مراتبی HMM ها برای بازنمایی واج‌ها، کلمات و جملات [6]

۱-۲-۲- استفاده از شبکه‌های عصبی

در روش‌های تقطیع بر اساس شبکه‌های عصبی، معمولاً شبکه‌های عصبی برای تخمین احتمال وقوع واحد آکوستیکی مورد نظر، آموزش می‌بینند. در بخش برچسب زنی نیز، از خاصیت طبقه بندی و متمایز سازی خوب شبکه‌های عصبی استفاده می‌شود. برای مثال در مرجع [28] روشی برای تخمین محل شروع هجاها ارائه شده است. ابتدا ویژگی‌های مناسبی از سیگنال گفتار استخراج می‌شود و سپس طبقه بندی کننده از این ویژگی‌ها برای اندازه گیری احتمال وقوع هجا استفاده می‌کند. در مرجع [6] نیز برای بازشناسی واحدهای زیر کلمه با ساختار CV در زبان هندی، از شبکه‌های عصبی استفاده شده است. در این مرجع اساس کار پیدا کردن نقطه شروع واژه‌ها^۱ در سیگنال گفتار می‌باشد. پس از پیدا کردن نقطه شروع واژه، یک فاصله زمانی ثابت ۲۰۰ میلی ثانیه در اطراف این نقطه، به منظور بازشناسی واحد زیر کلمه CV در سیگنال گفتار مورد پردازش قرار می‌گیرد.

¹ Vowel onset point (VOP)