





دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)

دانشکده مهندسی برق
گروه مخابرات - سیستم

پایان نامه کارشناسی ارشد

بهبود بازشناسی واجه‌های مصوت در گفتار پیوسته
به کمک روش ترکیب GMM و SVM

نگارش:

محمد نظری

استاد راهنما:

دکتر ابوالقاسم صیادیان

خرداد ماه ۱۳۸۷

سپاس و ستایش خدایی را که آفریننده قلم است و بدان سوگند یاد می کند
و هیچ کاری بدون لطف و اراده او غایت و سرانجامی نخواهد داشت.

نویسنده این اثر بر خود لازم می داند

تا از زحمات اساتیدی که بی شک بدون کمک و راهنمایی شان انجام این کار تحقیقی ممکن نبود،
خصوصاً استاد راهنمای محترم این پروژه جناب آقای دکتر صیادیان که با صبر و بردباری بسیار از
هیچ کمکی کوتاهی نکردند و همچنین جناب آقای دکتر احدی و جناب آقای دکتر همایون پور که در
ارزیابی و بهبود این اثر تاثیر فراوان داشته اند، تشکر نماید.

چکیده

با پیشرفت سریع در کاربردهای پردازش صوت اعم از سیستمهای بازشناسی گفتار زمان واقعی، سیستمهای بازشناسی گوینده، تطبیق گوینده و ... اهمیت وجود یک سیستم بازشناسی واجهای مصوت (واکه ها) هر چه بیشتر هویدا می شود. این اهمیت تا آنجاست که بخشی از تحقیقات امروز دانشمندان در زمینه پردازش صوت بر این موضوع متمرکز شده است. ما در این پایان نامه یک سیستم بازشناسی واکه ها را پیشنهاد کرده و افزایش قدرت بازشناسی آن را نسبت به سیستمهای موجود نشان داده ایم.

در این تحقیق، توجه عمده ما بر روی مدلسازی آکوستیکی و آماری واکه ها و در حقیقت ترکیب یک مدل مولد با یک مدل تمایزی برای کاهش خطای بازشناسی و تعیین مرزهای دقیقتر واکه ها معطوف شده است. در سیستم پیشنهادی جهت بازشناسی واکه ها، از تلفیق مدل مخلوط گوسی (GMM) با پنجره نرم و پارامترهای آکوستیکی مانند انرژی میان گذر با مدل تمایزی ماشین بردار پشتیبان احتمالی (PSVM) بهره گرفته ایم.

این سیستم دارای دو بخش اصلی شامل آشکارسازی و طبقه بندی واکه ها می باشد. مهمترین دستاورد این تحقیق در قسمت آشکارسازی حل مساله قدیمی مقدار آستانه برای مرزهای واکه می باشد. در سیستم پیشنهادی مرز واکه به صورت هوشمندانه و بسیار دقیق نسبت به داده های آموزشی اولیه انتخاب می شوند و سعی شده است سیستم تا حد ممکن نسبت به نویز محیطی مقاومت داشته باشد. در حقیقت در این گام از یک طبقه بندی کننده ترکیبی برای جداسازی واکه از غیر واکه استفاده می شود. در گام طبقه بندی نیز از روشی همانند قسمت آشکارسازی کمک می گیریم با این تفاوت که طبقه بندی کننده قبلی را به حالت چند کلاسی تعمیم می دهیم.

قابل ذکر است که برای افزایش سرعت در آشکارسازی و همچنین تخمین اولیه توان نویز محیطی از یک آشکارساز فعالیت گفتار (VAD) ساده استفاده شده است. در انتهای پایان نامه نتایج با الگوریتمهای دیگر که از تلفیق پارامترهای آکوستیکی و مدل های آماری بر روی پایگاه داده موجود پیاده سازی شده اند، مقایسه گردیده و نتایج مقایسه نشان دهنده کارایی بالاتر نسبت به سایر روشهای موجود می باشد.

کلمات کلیدی:

پردازش صوت، بازشناسی واج مصوت، گفتار پیوسته، مدل مخلوط گوسی (GMM)، ماشین بردار پشتیبان احتمالی (PSVM).

فهرست مطالب

پیشگفتار	۱
فصل اول - الگوهای گفتاری و روشهای بازشناسی آنها	۴
۱-۱- مقدمه	۵
۱-۲- الگوها و واحدهای گفتاری	۶
۱-۳- واکه	۹
۱-۴- مسایل اصلی در مدل سازی آکوستیکی گفتار	۱۰
۱-۴-۱- تنوع بردارهای ویژگی	۱۱
۱-۴-۲- دنباله غیر ایستان بردارهای ویژگی	۱۱
۱-۴-۳- اثر کشش زمانی	۱۱
۱-۴-۴- توالی قطعات	۱۱
۱-۴-۵- توزیع آماری پیچیده	۱۲
۱-۴-۶- تبدیل تدریجی قطعات به هم	۱۲
۱-۴-۷- قدرت تعمیم پذیری	۱۲
۱-۴-۸- تعداد پارامترها	۱۲
۱-۴-۹- حجم محاسبات و زمان بازشناسی	۱۳
۱-۵- روشهای مورد استفاده در بازشناسی الگوهای گفتاری	۱۳
۱-۵-۱- پیچش زمانی پویا	۱۳
۱-۵-۲- روش کدینگ منبع	۱۵
۱-۵-۳- مدل مارکف پنهان (HMM)	۱۶
۱-۵-۴- مدل قطعه ای اتفاقی (SSM)	۲۰
۱-۵-۵- مدل قطعه ای نرم اتفاقی SSSM	۲۳
۱-۵-۶- مدل مسیر میانگین مشروط (CMT-SSM)	۲۳
۱-۵-۷- ماشین بردار پشتیبان (SVM)	۲۵
۱-۵-۸- شبکه های عصبی	۲۵
۱-۶- جمع بندی	۲۷

فصل دوم - ماشینهای بردار پشتیبان احتمالی (PSVM) ۲۸

۲-۱-۱- مقدمه ۲۹

۲-۲- آموزش آماری ۲۹

۲-۲-۱- طبقه بندی خطی ۳۱

۲-۲-۲- طبقه بندی غیرخطی ۳۲

۲-۳- ماشینهای بردار پشتیبان ۳۲

۲-۳-۱- حداقل کردن ریسک عملی ۳۳

۲-۳-۲- حداقل کردن ریسک ساختاری ۳۳

۲-۳-۳- ماشینهای بردار پشتیبان خطی در حالت دو کلاسی با الگوی قابل جداسازی ۳۴

۲-۳-۴- ماشینهای بردار پشتیبان خطی در حالت دو کلاسی با الگوی جدایی ناپذیر ۳۵

۲-۳-۵- ماشینهای بردار پشتیبان غیر خطی و مفهوم هسته ۳۶

۲-۴- ماشینهای بردار پشتیبان احتمالی ۳۷

۲-۴-۱- کار های اخیر ۳۸

۲-۴-۲- ایده اصلی کار ۴۰

۲-۴-۵- برازش سیگموئید بعد از آموزش SVM ۴۱

۲-۵- جمع بندی ۴۲

فصل سوم - مدل مخلوط گوسی و روشهای ترکیب آن با SVM ۴۳

۳-۱- مقدمه ۴۴

۳-۲- طبقه بندی آماری و قانون بیز ۴۴

۳-۳- مدل مخلوط گوسی ۴۵

۳-۳-۱- آموزش مدل مخلوط گوسی ۴۶

۳-۳-۲- آزمون نسبت شباهت ۴۷

۳-۳-۳- هنجارسازی امتیاز مدل مخلوط گوسی ۴۸

۳-۳-۴- محدودیتهای مدل مخلوط گوسی ۴۸

۳-۴- ترکیب مدل مولد و مدل تمایزی ۴۸

۳-۴-۱- نحوه ترکیب مدل مخلوط گوسی و ماشین بردار پشتیبان ۴۹

۳-۴-۲- ابر بردار مدل مخلوط گوسی ۴۹

۳-۴-۳- روش ترکیب پیشنهادی مدل مخلوط گوسی و مدل PSVM ۵۱

۳-۵- جمع بندی ۵۲

فصل چهارم- بررسی الگوریتم پیشنهادی در آشکارسازی و طبقه‌بندی واکه‌ها .. ۵۳

۱-۴- مقدمه ۵۴

۲-۴- پایگاه داده ۵۴

۳-۴- سیستم پیشنهادی برای بازشناسی واکه ۵۵

۴-۴- آشکارسازی واکه و تخمین مرزهای آن ۵۶

۱-۴-۴- پیش پردازش و استخراج ویژگی ۵۷

۲-۴-۴- انرژی میان گذر در هر فریم ۵۸

۳-۴-۴- روش پیشنهادی برای ترکیب مدل GMM و PSVM ۶۱

۴-۴-۴- حوه آموزش مدل مخلوط گوسی با پنجره نرم در آشکارسازی ۶۳

۵-۴-۴- انتخاب بهترین تعداد مخلوط گوسی ۶۳

۶-۴-۴- نحوه انتخاب داده های آموزشی برای آموزش مدل PSVM ۶۷

۵-۴- نحوه آموزش مدل مخلوط گوسی برای طبقه بندی واکه ها ۶۷

۶-۴- نتایج پیاده سازی الگوریتمها ۶۸

۱-۶-۴- نحوه انتخاب مجموعه آزمون و آموزش ۶۸

۲-۶-۴- نتایج آشکارسازی واکه در گفتار پیوسته ۶۸

۳-۶-۴- نتایج طبقه بندی واکه در گفتار پیوسته ۶۹

۷-۴- مقایسه نتایج با مدل های دیگر ۷۰

۸-۴- نتایج ارزیابی مقاوم بودن دقت بازشناسی نسبت به نویز ۷۱

۱-۸-۴- نویزهای غیرپریودیک ۷۱

۲-۸-۴- نویزهای پریودیک ۷۳

۹-۴- جمع بندی ۷۴

فصل پنجم: نتیجه گیری و پیشنهادات ۷۵

۱-۵- نتیجه گیری ۷۶

۲-۵- پیشنهادات ۷۶

مراجع ۸۰

لیست واژه‌های انگلیسی

.....سطح اولیه	Basic-Level
.....بیز	Bayes
.....تابع تمایز	Discriminant Function
.....تجربی	Empirical
.....سطح ورودی	Entry-Level
.....تشخیص واکه	Vowel Detection
.....قابلیت عمومیت دادن	Generalization
.....ابر صفحه	Hyper-Plane
.....توابع پایه	Kernel
.....تصویر کردن خطی	Linear Projection
.....تابع تلف	Loss Function
.....بردارهای مشاهده	Observation Vectors
.....رد	Reject
.....ارائه	Representation
.....کمینه‌سازی ساختاری خطا	Structural Risk Minimization
.....تحت نظارت	Supervised
.....ماشین بردار پشتیبان	(SVM) Support Vector Machines

لیست واژه‌های اختصاری

Maximum A posteriori Probability.....	MAP
Gaussian mixture models	GMM
Support Vector Machines.....	SVM
Probability Support Vector Machines	PSVM
Principle Component Analysis.....	PCA
Linear Discriminant Analysis.....	LDA
Kullback-Leibler.....	KL
Maximum Likelihood.....	ML
Radial Basis Function.....	RBF
Least Squares.....	LS
Receiver Operator Characteristics.....	ROC

پیشگفتار

زبان از اجزای بسیار زیادی تشکیل شده است که کوچکترین جزء یا واحد آن "واج" است. همه زبانهای بشری به واج تجزیه می شوند؛ بنابراین در نخستین مرحله لازم است واجها با هم ترکیب شوند تا واحدهای بزرگتری مانند تکواژ و یا واژه ساخته شوند.

در مقوله پردازش گفتار، شناسایی واجهای مصوت (واکه‌ها)^۱ به دلیل خصوصیات ویژه آنها در گفتار پیوسته، نقش کلیدی بازی می کنند. برای مثال می توان به کاربردهای بازشناسایی واجهای مصوت برای استفاده در سیستمهای بازشناسایی خودکار گفتار (ASR)^۲ و شناسایی گوینده از روی گفتار اشاره کرد. این امر باعث گردیده تا طیف وسیعی از مطالعات امروز دانشمندان این رشته به روشهای جدید برای بهبود روشهای بازشناسی واجهای مصوت در گفتار پیوسته معطوف گردد [۲ و ۳].

گفتار واکنار نتیجه تحریک لوله صوتی توسط جریان پررودیک هوا در دهانه حنجره می باشد. واکه ها از مهمترین گروه اصوات می باشند که در همه زبانها موجود و دارای اهمیت زیادی در بازشناسی گفتار می باشند. این دسته اصوات پایدارترین گروه از نظر فرکانسی و زمانی می باشند و معمولاً نسبت به همخوانها دارای طول خوبی هستند؛ از این رو راحت تر و دقیق تر بازشناسی می شوند. در محیط های با نسبت سیگنال به نویز بالا، انرژی واکه ها، ۱۰dB تا ۲۰dB بالاتر از متوسط انرژی سیگنال گفتار است. همچنین بخشهای مصوت گفتار به دلیل دارا بودن ساختار هارمونیک، نسبت به شرایط نویزی مقاوم تر می باشند و تخمین محل وقوع آنها نسبت به بی واک ها آسان تر می باشد. به همین دلیل بسیاری از سیستمهای بازشناسی برای دستیابی به عملکرد بهتر، بر بازشناسی دقیق واکه ها تأکید دارند. مهمترین ویژگی آکوستیکی که برای تشخیص محل واکه استفاده می شود، انرژی می باشد؛ زیرا در گفتار، تمامی واکه ها صرف نظر از گوینده، دارای انرژی بالاتری نسبت به همخوانهای^۳ مجاور می باشند.

¹ Vowels

² Automatic Speech Recognition

³ Consonant

در دنیای واقعی هنگامی که واجهای مصوت یکسان توسط افراد مختلف با جنسیت‌های مختلف، سنین مختلف و در شرایط مختلف ادا می‌شود و یا حتی هنگامی که یک واج مشخص توسط یک فرد خاص در زمانهای مختلف و یا در شرایط مختلف تلفظ می‌گردد، با یکدیگر تفاوت چشمگیر دارد. در این صورت باید تشخیص واجهای مصوت به درستی و با دقت زیاد انجام شود. بنابراین در اینجا مسئله برای بازشناسی واجهای مصوت در گفتار پیوسته مطرح می‌شود. هنگامی ماشین توانایی انجام درست این امر را دارد که توسط یک مدل آموزش پذیر مناسب آموزش دیده باشد. الگوریتمهای مختلف برای بازشناسی واجهای مصوت در گفتار پیوسته وجود دارد که بخش مهمی از الگوریتم‌های موجود در این زمینه مبتنی بر مدل‌های آماری می‌باشد.

از دیدگاه سیستمی، در یک سیستم بازشناسی الگو، پس از فرایندهای پیش پردازش نمونه‌های زمانی گفتار، ویژگی‌های مناسب از نمونه‌های زمانی استخراج می‌شود. سپس بردار ویژگی‌های استخراج شده به مدلی عرضه می‌شود تا تطبیق آن با مدل بررسی شود. میزان این تطبیق و مقایسه با تطابق با دیگر الگوها، به عنوان معیار بازشناسی شناخته می‌شود. همانند بسیاری از سیگنال‌های طبیعی دیگر، الگوهای مورد بازشناسی که دارای معانی یکسان می‌باشند، بسیار متنوع هستند و همچنین مرز دسته الگوهای متفاوتی که سیستم بازشناسی، وظیفه طبقه بندی آنها را بر عهده دارد، به خوبی روشن نیست. اگرچه استخراج ویژگی‌های مناسب می‌تواند پیچیدگی این مسأله را کاهش دهد، اما در حال حاضر اکثر سیستم‌های بازشناسی گفتار بر پایه مدل‌سازی آماری هر الگو بنا شده‌اند؛ به طوریکه پارامترهای مدل، در حین آموزش (مشاهده نمونه‌های مناسب از الگوی مورد نظر) تخمین زده می‌شوند.

در این پایان نامه، هدف ارایه روش جدیدی برای بهبود بازشناسی واجهای مصوت در گفتار پیوسته به کمک روش‌های ترکیبی GMM و SVM می‌باشد. در اینجا مسأله بازشناسی واجهای مصوت در گفتار پیوسته مورد بررسی قرار می‌دهیم و از یک مدل آماری برای حل این مسأله استفاده می‌کنیم و سپس به بهبود روش‌های موجود می‌پردازیم.

ابزار کلی در حل مسأله استفاده از دو گام اصلی می‌باشد. در ابتدا، گام اول که شامل یک گام آموزشی بر روی نمونه‌های گرفته شده است، انجام می‌پذیرد. این گام آموزشی بر روی یک مدل ترکیبی مدل مخلوط

گوسی (GMM) و ماشین بردار پشتیبان (SVM) اجرا می گردد. در گام دوم که گام بازشناسی واجه‌های مصوت در گفتار پیوسته می‌باشد، ورودی اعمال می‌شود تا واجه‌های مصوت مورد نظر شناسایی شود. در این پروژه ضمن پیاده سازی روشهای مختلف سعی می‌کنیم با استفاده از روشی ترکیبی که مورد طراحی قرار می‌گیرد عمل بازشناسی واجه‌های مصوت در گفتار پیوسته را بهبود دهیم.

این سیستم دارای دو بخش اصلی شامل آشکارسازی و طبقه‌بندی واکه‌ها می‌باشد. مهمترین دست‌آورد این تحقیق در قسمت آشکارسازی حل مساله قدیمی مقدار آستانه برای مرزهای واکه می‌باشد. در سیستم پیشنهادی مرز واکه به صورت هوشمندانه و بسیار دقیق نسبت به داده‌های آموزشی اولیه انتخاب می‌شوند و سعی شده است سیستم تا حد ممکن نسبت به نویز محیطی مقاومت داشته باشد. در حقیقت در این گام از یک طبقه‌بندی کننده ترکیبی برای جداسازی واکه از غیر واکه استفاده می‌شود. در گام طبقه‌بندی نیز از روشی همانند قسمت آشکارسازی کمک می‌گیریم با این تفاوت که طبقه‌بندی کننده قبلی را به حالت چند کلاسی تعمیم می‌دهیم.

به این منظور، در فصل اول به بیان ویژگی‌های یک سیستم بازشناسی الگو در کاربردهای پردازش گفتار پرداخته و در ادامه، مدل‌های معمول جهت بازشناسی، نقاط ضعف و قوت آنها را مورد بررسی قرار داده ایم. فصل دوم به بررسی مدل ماشین بردار پشتیبان (SVM) اختصاص دارد. در این فصل ایده اصلی، مدل‌سازی تحلیلی و الگوریتم‌های آموزش و بازشناسی این مدل بررسی می‌شوند و در خاتمه با تعمیم آن، مدل ماشین بردار پشتیبان احتمالی (PSVM) را معرفی می‌کنیم. فصل سوم به بررسی انواع مدل‌های ترکیب مدل مخلوط گوسی (GMM) و مدل ماشین بردار پشتیبان (SVM) اختصاص دارد. در فصل چهارم الگوریتم‌های پیشنهادی و نتیجه پیاده سازی آنها آورده شده است. در فصل پنجم، با نتیجه گیری نهایی و پیشنهاداتی برای ادامه کار، به مطالب پایان نامه خاتمه می‌دهیم.

فصل اول

الگوهای گفتاری و روشهای بازشناسی آنها

۱-۱ - مقدمه

شنوایی، به عنوان یکی از حسهای اصلی انسان، نقش بسیار مهمی در به دست آوردن اطلاعات محیطی بر عهده دارد. گوش انسان با شنیدن گفتار اطلاعات فراوانی از قبیل محتوای گفتار، جنس، سن، وضعیت احساسی و حتی هویت فرد گوینده به دست می‌آورد. این توانایی، دانشمندان را بر آن داشته است تا با تحقیقات فراوان سعی در مدل کردن این توانایی بشری با الگوهای موجود نمایند تا هر چه بیشتر بتوانند در مسیر نزدیکی انسان و ماشین قدمی بردارند. این ایده باعث شده است تا دانشمندان علوم اعصاب شناختی، علوم ریاضی و علوم کامپیوتر دست در دست هم به سمت این مدل‌سازی‌ها حرکت کنند.

اولین و مهمترین زمینه در شناسایی الگوهای گفتاری، بازشناسی خودکار گفتار (ASR¹) می‌باشد. بازشناسی خودکار گفتار به عنوان یکی از آرزوهای پژوهشگران سیستمهای پردازشی، هم اکنون به یکی از شاخه‌های اصلی پردازش سیگنال مبدل شده است. سیستمهای مترجم خودکار، سیستمهای محاوره ماشینی، منشی‌های ماشینی و رابط‌های گفتاری، کاربردهای مختلف در این زمینه می‌باشند. اگرچه بسیاری از سیستمهای بازشناسی گفتار به موفقیت‌هایی دست یافته‌اند ولی همواره فاصله انسان و بهترین روشهای بازشناسی گفتار حفظ شده است. به طوریکه در عمل، نه تنها سیستمهای کاربردی نتوانسته‌اند در این زمینه از تواناییهای انسان پیشی بگیرند، بلکه هیچگاه جانشین مطلوبی نیز محسوب نشده‌اند. از این رو گستره وسیعی از پژوهش برای بهبود کارایی سیستمهای ماشینی بازشناسی گفتار گشوده شده است.

در سیستمهای گفتاری اعم از بازشناسی و یا سنتز، گفتار به صورت رشته‌ای از اجزای تشکیل دهنده اش مورد بررسی قرار می‌گیرد. این اجزای تشکیل دهنده می‌توانند جملات، عبارات، کلمات، هجاها و یا واجها باشند. برای مثال در بازشناسی لازم است میزان شباهت بین سیگنال ورودی و گفتار (یا مدل‌های گفتار) موجود در لغت نامه سنجیده شود. این مقایسه می‌تواند در سطح هر کدام از اجزای تشکیل دهنده گفتار انجام شود و یا در سنتز گفتار، باید با استفاده از این اجزای تشکیل دهنده بیان مورد نظر ساخته شود.

¹ *Automatic Speech Recognition*

در این فصل ابتدا به صورت مختصر به الگوهای گفتاری و مخصوصاً واحدهای زبانی می‌پردازیم سپس خواص واکه‌ها و مزایای بازشناسی دقیق آنها را مورد بررسی قرار می‌دهیم. آنگاه ویژگیهای یک سیستم بازشناسی الگوی گفتار را بررسی می‌کنیم. سپس در مورد مدلهای مختلف و خصوصاً آماری به کار رفته برای بازشناسی الگوی گفتار توضیحاتی داده شده و نقاط قوت و ضعف این مدلها مورد بررسی قرار خواهد گرفت.

۱-۲- الگوها و واحدهای گفتاری

یک سیستم بازشناسی الگوی گفتاری متشکل از اجزای زیادی است که به صورت متوالی به کار رفته‌اند و دقت خروجی سیستم به کارایی کلیه اجزای آن بستگی دارد. اما حساس‌ترین بخش این سیستم، استفاده از مدلی است که تنوع بردارهای ویژگی استخراج شده را تحمل می‌کند، به واحدهای زبانی متفاوت، امتیاز تطابق متناسب می‌کند و آنها را از لحاظ تطابق با ورودی رتبه بندی می‌کند.

پیچیدگی این مدل‌سازی به توانایی بخش‌های دیگر سیستم وابسته است. به عنوان مثال، در صورتی که ویژگی‌های استخراج شده از گفتار آنقدر تفکیک‌کننده باشد که هر دنباله بردارهای ویژگی مستقیماً به یک واحد زبانی اشاره کند، نیازی به این مدل‌سازی نیست. همچنین در صورتی که مجموعه هدف خروجی آنقدر کوچک باشد که اشتباه در بازشناسی اکوستیکی با مدل زبانی جبران شود، مدل‌سازی اکوستیکی نقش حساس خود را از دست می‌دهد. اما در عمل، نه ویژگی‌های موجود تفکیک‌پذیری کافی را ارائه می‌دهند و نه کاربردهایی با مجموعه‌ای چنین محدود مورد علاقه هستند. لذا پیچیدگی زیادی به بخش مدل‌سازی اکوستیکی واحدهای زبانی تحمیل می‌گردد. در واقع، مدل‌سازی اکوستیکی، باید آنچه را که توسط بخش‌های دیگر بازنمایی نشده است، به شیوه‌ای آماری مدل کند. این مسایل شامل تنوع در گوینده، محیط گفتار، وابستگی به متن گفتار، تلفظ، لهجه و احساسات بیانی گوینده و حتی تنوع حاصل از دو بیان یک واحد گفتاری توسط یک گوینده می‌شود [۱].

واحدهای زبانی معمول در بازشناسی گفتار شامل آوا، واج، هجا و اجزای آن، (ترکیب‌های نیم هجایی) و کلمه است که از لحاظ اکوستیکی به ترتیب ساده به پیچیده ارائه شده‌اند. آوا کوتاه‌ترین واحد اکوستیکی در

بازشناسی گفتار پیوسته است که شکل لوله صوتی در طول آن تقریباً ثابت می‌ماند و واج کوتاه‌ترین واحد مورد استفاده زبان‌شناسان است.

در بازشناسی گفتار پیوسته، ترکیب‌های مجاز گفتاری بسیار متنوع است و معمولاً از واحدهای زبانی کوتاه‌تر و کم‌تنوع‌تر استفاده می‌شود تا با الفبایی محدود، تمام ترکیب‌های مجاز قابل استخراج و مدل‌سازی باشد. برای استفاده از مزیت مدل‌های مورد استفاده زبان‌شناسان، معمولاً از آواهای متناظر با واحدهای پایه به عنوان واحد آکوستیکی استفاده می‌شود، اگرچه در بعضی پژوهش‌ها برای مدل‌سازی تاثیر متقابل واحدها به هم از ترکیب‌های دو آوایی و سه آوایی نیز بهره گرفته شده است. رویکرد اخیر مستلزم مدل‌سازی تعداد زیادی واحد آکوستیکی است که خود باعث پیچیدگی بازشناسی آکوستیکی و نیاز به بانک اطلاعاتی بسیار وسیع برای تخمین پارامترهای سیستم می‌شود. واضح است که بیان چنین حجم پیچیدگی، تنها به صورتی آماري قابل مدل‌سازی است.

در بازشناسی گفتار، انتخاب واحد گفتاری برای سیستم، از اهمیت ویژه‌ای برخوردار است. برای سیستم‌های با تعداد لغات کم، معمولاً از واحد کلمه استفاده می‌شود. استفاده از کلمه به عنوان واحد گفتار، برای سیستم‌های بازشناسی کلمات گسسته و بازشناسی کلمات بدنبال هم به کار می‌رود. از مزایای انتخاب واحد کلمه به عنوان واحد گفتار این است که نمایش آکوستیکی کلمات به خوبی معین شده و تغییرات آوایی عمدتاً در ناحیه شروع و انتهای کلمه قرار دارد. مزیت دیگر اینکه، استفاده از واحد کلمه، نیاز به فرهنگ لغت را مرتفع می‌سازد و بنابراین ساختار بازشناسی ذاتاً ساده‌تر می‌شود. اما معایب استفاده از مدل‌های گفتاری کلمه برای بازشناسی گفتار پیوسته هنگامی نمایان می‌گردد که در فاز آموزش نیاز به انواع حالات و ترکیبات کلمه در جمله به وجود می‌آید، که این امر آموزش در سیستم‌های عملی را غیر ممکن می‌سازد.

تعدادی از واحدهای آکوستیکی عبارتند از آوا، دوآوایی، سه آوایی، هجا و نیم هجا، که در ادامه در مورد آنها توضیحاتی داده می‌شود.

- آوا: کوچکترین واحد آکوستیکی است که به دو گروه واکه^۱ و همخوان^۲ تقسیم می‌شود. بر طبق [۱] زبان فارسی ۸ واکه و ۲۳ همخوان دارد. این ۸ واکه شامل ۶ واکه اصلی و ۲ واکه مرکب است.

¹ Vowel

² Consonant

همخوان‌ها شامل دو زیر گروه واکدار^۱ و بی واک^۲ می‌باشند. در میان زیر کلمه‌ها تعداد آواهای زبان فارسی از سایر واحدها کمتر می‌باشد. با توجه به اثر آواها بر روی یکدیگر، نمی‌توان آنها را مستقل از هم در نظر گرفت و در واقع باید اثرات درون هجایی و برون هجایی آواها بر یکدیگر لحاظ شود. یعنی بسته به کلمه و محل قرار گرفتن آواها، باید مدل آوای مورد نظر را تغییر داد که در این صورت، کارایی سیستم بازشناسی به میزان زیاد افزایش پیدا می‌کند [۱ و ۲].

- دو آوایی^۳: هر آوا با توجه به آوای سمت راست یا سمت چپ به طور جداگانه مدل می‌شود.
- سه آوایی^۴: هر آوا با توجه به آوای سمت راست و سمت چپ، به طور جداگانه مدل می‌شود.
- هجا^۵: واحدهای آوایی - مانند واکه و همخوان - مواد خامی هستند که بر طبق قواعد و الگوهای معین گرد هم آمده و واحدهای پیچیده تر از خود، یعنی هجاها را تشکیل می‌دهند. هجا در فارسی عبارت است از یک رشته آوایی پیوسته که از یک واکه و یک تا سه همخوان تشکیل یافته است. منظور از رشته آوایی پیوسته این است که اجزاء سازنده هجا طی یک فرآیند تولیدی و بدون مکث تولید می‌گردند. واکه به منزله مرکز یا هسته هجا می‌باشد. بدین معنی که موجودیت هجا بستگی به وجود واکه دارد و اگر واکه را حذف کنیم دیگر هجایی باقی نمی‌ماند. در صورتیکه می‌توان از یک هجا، یک یا دو همخوان را حذف کرد بدون اینکه به موجودیت آن لطمه ای وارد شود.
- نیم هجا^۶: هر هجا به دو نیم هجای ابتدایی و انتهایی تقسیم بندی می‌شود. نیم هجای ابتدایی شامل همخوان ابتدایی و بخشی از هسته هجا و نیم هجای انتهایی شامل باقیمانده هسته تا انتهای هجا می‌باشد. نیم هجاهای موجود در زبان فارسی به صورت CV ، VC و VCC می‌باشند که تعداد آنها حدود ۷۰۰ است.

ساختار هجاها به این نحو می‌باشد که هر هجا از سه بخش شروع، هسته و انتها تشکیل می‌شود. هر چند که بسیاری از هجاها دارای هر سه قسمت می‌باشند، اما بخشهای با اهمیت هجا معمولاً شامل یک یا دو قسمت

1 Voiced
 2 Unvoiced
 3 Biphone
 4 Triphone
 5 Syllable
 6 Demi-Syllable

از بخشهای ذکر شده می باشد. اگر هجایی دارای یک جزء تکی باشد، الزاماً آن جزء هسته می باشد. بعضی از زبانهای دنیا [۵]، از جمله ژاپنی [۶] و چینی [۷]، از ساختار هجایی نسبتاً روشن و شفاف بر خوردارند. زبان فارسی از دسته زبانهایی می باشد که دارای ساختار هجایی ساده ای است. زبان فارسی دارای سه نوع مختلف هجا می باشد:

۱. ساختار همخوان - واکه (CV)

۲. ساختار همخوان - واکه - همخوان (CVC)

۳. ساختار همخوان - واکه - همخوان - همخوان (CVCC)

بررسی ها نشان می دهد که نوع اول بیشترین تعداد وقوع در گفتار فارسی را دارد و پس از آن نوع دوم و سوم به ترتیب بیشترین تعداد وقوع را دارند. نشان داده شده است که تغییرات تلفظی در سطح هجا بسیار منظم تر از سطح قطعات آوایی می باشد [۱].

درسیستمهای بازشناسی مبتنی بر نیم هجا، در صورتیکه ابتدا واکه، که هسته هجا می باشد، آشکار سازی شود، بازشناسی نیم هجای قبل و بعد از واکه با دقت بالاتر و مقاوم تر نسبت به شرایط محیطی انجام می - گردد. در زبانهای مبتنی بر هجا، برای افزایش دقت بازشناسی و تقطیع گفتار به هجاها از آشکارسازی واکه استفاده می شود [۸-۳]، در ادامه در مورد خواص واکه ها بیشتر توضیح می دهیم.

۱-۳ - واکه

گفتار واکدار نتیجه تحریک لوله صوتی توسط جریان پررودیک هوا در دهانه حنجره می باشد. واکه ها از مهمترین گروه اصوات می باشند که در همه زبانها موجود و دارای اهمیت زیادی در بازشناسی گفتار می - باشند. این دسته اصوات پایدارترین گروه از نظر فرکانسی و زمانی می باشند و معمولاً نسبت به همخوانها دارای طول خوبی هستند؛ از این رو راحت تر و دقیق تر بازشناسی می شوند. در محیط های با سیگنال به نویز پایین، انرژی واکه ها، ۱۰dB تا ۲۰dB بالاتر از متوسط انرژی سیگنال گفتار است. همچنین بخشهای مصوت گفتار به دلیل دارا بودن ساختار هارمونیک، نسبت به شرایط نویزی مقاوم تر می باشند و تخمین

محل وقوع آنها نسبت به بی واک ها آسان تر می باشد. به همین دلیل بسیاری از سیستمهای بازشناسی برای دستیابی به عملکرد بهتر، بر بازشناسی دقیق واکه ها تأکید دارند.

جدول ۳-۱: واکه های زبان فارسی

مثال	علامت	آوا
n@m, نم	@	ا
mat, مات	a	آ
cek, چک	e	اِ
miz, میز	i	ای
kot, کت	o	اُ
nur, نور	u	او
lowh, لوح	ow	اوُ
\$eix, شیخ	ei	ایِ

در [۹] برای افزایش دقت بازشناسی هجاها در گفتار پیوسته، بر یافتن محل دقیق واکه ها تأکید شده و از ترکیب HMM، شبکه عصبی و انرژی برای این آشکارسازی استفاده شده است. مهمترین ویژگی آکوستیکی که برای تشخیص محل واکه استفاده می شود، انرژی می باشد؛ زیرا در گفتار، تمامی واکه ها صرف نظر از گوینده، دارای انرژی بالاتری نسبت به همخوانهای مجاور می باشند. هشت واکه موجود در زبان فارسی - شش واکه ساده و دو واکه مرکب - در جدول (۳-۱) آورده شده اند.

۱-۴- مسایل اصلی در مدل سازی آکوستیکی گفتار

گفتار، پدیده پیچیده ای است که علی رغم تنوع بسیار در ابعاد گوناگون می تواند معنای یکسانی را منتقل کند. اگرچه استخراج ویژگی مناسب می تواند این تنوع را کاهش دهد، اما طبیعت گفتار پیچیده تر از آن است که با دنباله ای از بردارهای ویژگی بازنمایی گردد. بنابراین یک مدل آکوستیکی موفق باید با شناخت طبیعت گفتار آنچنان طراحی گردد که تا حد ممکن به این طبیعت نزدیک شود. در این فصل به مسایلی پرداخته می شود که یک مدل سازی آکوستیکی موفق باید به شکلی آنها را تحمل کند. مدل سازی مسایل زیر که همگی ابعاد گوناگون سیگنال طبیعی گفتار هستند، مستقیماً بر روی نرخ بازشناسی یک سیستم بازشناسی گفتار موثر است.

۱-۴-۱- تنوع بردارهای ویژگی

مهمترین مشخصه بردارهای ویژگی یک گفتار، عدم یکسانی آنها به ازای بیان آنها در شرایط مختلف یا بردارهای لحظات مختلف یک قطعه گفتاری است. تا زمانی که نظم این تنوع به صورت معین قابل مدل سازی نباشد، بازنمایی این تنوع به تخمین توزیع چگالی احتمال بردارها می انجامد.

۱-۴-۲- دنباله غیر ایستان بردارهای ویژگی

بردارهای ویژگی یک گفتار نه تنها متنوعند، بلکه مشخصات آماری آنها در طول زمان تغییر می کند و مدل آکوستیکی باید بتواند این تغییر مشخصات را بازنمایی کند. این بازنمایی معمولاً با قطعه بندی دنباله بردارهای ویژگی یک واحد زبانی انجام می پذیرد.

۱-۴-۳- اثر کشش زمانی

یک ویژگی بارز سیگنال گفتار، تنوع در طول گفتارهایی است که همگی بیانگر یک واحد زبانی هستند. این کشش زبانی که معمولاً نشانه لهجه، تاکید، یا احساسات گفتاری گوینده است، نباید موجب اشتباه در نتیجه بازشناسی واحدهای زبانی گردد.

اثر کشش زمانی از دیدگاه کمی نیز قابل بررسی است. برای تحمل صحیح کشش زمانی هر قطعه در یک واحد زبانی، باید رفتار آماری کشش زمانی هر قطعه مورد بررسی قرار گیرد. مدل سازی مناسب طول هر قطعه آکوستیکی می تواند راه کار مناسبی برای این منظور باشد.

۱-۴-۴- توالی قطعات

قطعات مدل سازی شده در هر واحد گفتاری به طور طبیعی به همان توالی مدل سازی شده ظاهر می شوند، اما در بیان یک واحد زبانی ممکن است بعضی قطعات اصلاً بیان نشوند که مدل مناسب باید بتواند این پدیده را بازنمایی کند. چنانچه مشاهده خواهد شد، عدم مدل سازی این پدیده تاثیر جدی در نرخ بازشناسی ندارد و علاوه بر مدل های آماری، از مدل های معین نیز برای مدل سازی توالی قطعات استفاده می شود.