

به کارگیری اطلاعات زبانی در یک سیستمِ بازشناسی گفتار پیوسته فارسی

محمد بحرانی

دانشگاه صنعتی شریف (آزمایشگاه پردازش گفتار)

حسین صامتی

دانشگاه صنعتی شریف (آزمایشگاه پردازش گفتار)

چکیده

در این مقاله یک سامانه بازشناسی گفتار پیوسته برای زبان فارسی معرفی می‌شود و نقش مدل آوایی و مدل زبانی در آن مورد بررسی قرار می‌گیرد. مدل‌های آوایی با روش‌های مستقل از بافت و وابسته به بافت در این سامانه به کار رفته و نتایج به کارگیری آن‌ها ارائه شده است. مدل زبانی سه کلمه‌ای نیز با روش‌های مبتنی بر کلمه، مبتنی بر مقوله نحوی و مبتنی بر طبقه، با استفاده از پیکره متنی زبان فارسی استخراج و در سامانه بازشناسی به کار گرفته شده است. همچنین مدل زبانی دستوری مبتنی بر دستور ساخت-گروهی تعمیم یافته در این سامانه پیاده‌سازی شده و نیز در ترکیب با مدل زبانی آماری به کار رفته است. نتایج حاصل نشان می‌دهد که مدل آوایی وابسته به بافت، مطابق انتظار، بهترین عملکرد را دارد. همچنین مدل زبانی سه کلمه‌ای مبتنی بر کلمه، نسبت به سایر روش‌های استخراج مدل زبانی آماری برتری دارد. در ضمن ترکیب مدل زبانی دستوری با مدل زبانی آماری منجر به بهبود نتایج بازشناسی می‌شود. سامانه بازشناسی گفتار معرفی شده در این مقاله، اولین سامانه بازشناسی برای گفتار پیوسته فارسی بوده و با پشتوانه فعالیت‌های تحقیقاتی متعددی که برای پیاده‌سازی آن انجام شده است، قابلیت استفاده به صورت کاربردی را یافته است.

کلیدواژه‌ها: بازشناسی گفتار پیوسته، مدل سازی آوایی، مدل سازی زبانی، مدل زبانی آماری، مدل زبانی دستوری.

۱. مقدمه

مسئلهٔ بازشناسی گفتار و یا تبدیل گفتار به متن از مسائل مهم تحقیقاتی و کاربردی مرتبط با زبان‌شناسی است. یک موتور بازشناسی گفتار توسط پژوهشگرانی با تخصص‌های متفاوت، شامل مهندسی کامپیوتر، پردازش علائم دیجیتال، هوش مصنوعی، آواشناسی و زبان‌شناسی، طراحی و پیاده‌سازی می‌شود. تخصص‌های درگیر برای کار روی یک سامانهٔ بازشناسی گفتار نیاز به ارتباط نزدیک با یک‌دیگر برای مدل‌سازی و پیاده‌سازی اطلاعات گفتاری و زبانی دارند. گفتار مورد بازشناسی می‌تواند به صورت یک فایل صوتی باشد و یا از طریق میکروفون، از طریق خط تلفن و یا فرمان از راه دور باشد. در واقع هدف نهایی بازشناسی گفتار، ساخت ماشین‌هایی است که بتوانند مانند انسان بشنوند و عکس‌العمل مناسب نشان دهند.

اولین تلاش‌ها برای ساخت سیستم بازشناسی گفتار از اوایل دههٔ چهارم میلادی آغاز شد و تا اکنون این تلاش‌ها ادامه دارند. اگرچه رویکردهای مختلفی برای بازشناسی گفتار وجود دارند اما موفق‌ترین آن‌ها رویکرد تشخیص الگو^۱ است، که تقریباً همهٔ سیستم‌های موفق امروزی براساس آن عمل می‌کنند. در این رویکرد، گفتار به کمک تعدادی واحد آوایی^۲ (کلمه^۳، هجا^۴، سه واجی^۵ یا واج^۶) مدل می‌شود و برای بازشناسی نیز از تشخیص این واحدها و کنار هم قرار دادن آن‌ها، متن متناسب با گفتار، تشخیص داده می‌شود. این مدل به نوعی متناسب با نظریه‌های زبان‌شناسی چامسکی است که براساس آن، هر انسانی در برخورد با گفتاری از زبانی آشنا، آن‌را به رشته‌ای از کوچک‌ترین واحدهای زیرکلمه^۷ ممکن تبدیل می‌کند و سپس در مغز خود به پردازش این رشته می‌پردازد تا هجاها، کلمات، جملات و نهایتاً گفتار را بازشناسی کند. اگرچه امروزه سیستم‌های موفق مختلفی براساس این رویکرد ارائه شده‌اند و در عمل نیز از آن‌ها استفاده می‌شود ولی همگی آن‌ها به گونه‌ای برخی از محدودیت‌های ساده‌کننده را یک می‌کشند. حذف این محدودیت‌ها می‌تواند به طور قابل ملاحظه‌ای بر پیچیدگی آن‌ها بیفزاید. در واقع، یک سیستم بازشناسی کامل، که بتواند مشابه انسان عمل کند باید بتواند:

- قادر به بازشناسی گفتار پیوسته^۸ و محاوره‌ای^۹ باشد.
- گفتار افراد مختلف، حتی با لهجه‌های متفاوت را بازشناسی نماید.
- در محیط‌های شلوغ و نوفه‌ای هم جواب‌گو باشد.

1. pattern recognition
2. acoustic unit
3. word
4. syllable
5. triphone
6. phoneme
7. subword
8. continuous
9. spontaneous

- به صورت بی درنگ^۱ عمل کند.
 - قادر به فراگیری اطلاعات جدید، نظیر کلمات و قواعد زبانی باشد.
- سامانه‌های موجود بازشناسی گفتار معمولاً دارای محدودیت‌هایی هستند و مشخصات فوق را، به طور کامل برآورده نمی‌کنند. این سامانه‌ها ممکن است گفتار محاوره‌ای را به خوبی گفتار رسمی و کتابی تشخیص ندهند. همچنین عملکرد این سامانه‌ها در محیط‌های واقعی به مراتب ضعیف‌تر از محیط آکوستیک است. علاوه بر این، نوعاً محدود به واژگان تعریف شده و قوانین زبانی آموزش داده شده هستند و در مواجهه با واژگان جدید و قوانین زبانی جدید دچار خطا می‌شوند. میزان پیچیدگی یک سیستم بازشناسی گفتار به عوامل متعددی بستگی دارد. در این سیستم‌ها، برای کاربردی کردن، سعی می‌شود با ایجاد محدودیت‌هایی میزان پیچیدگی را کاهش دهند. این کار باعث تعریف دقیق سیستم برای یک کاربرد مشخص و کاهش دامنه کاربرد آن (با کارایی بهتر) می‌شود. مهم‌ترین این محدودیت‌ها که توانمندی یک سیستم بازشناسی گفتار را مشخص می‌کنند موارد زیر هستند.
- میزان وابستگی یا استقلال از گوینده:** سیستم‌هایی که تنها به یک و یا چند گوینده خاص پاسخ می‌گویند، وابسته به گوینده^۲، و آن‌هایی که به تمام گویندگان یک زبان پاسخ می‌گویند مستقل از گوینده^۳ نامیده می‌شوند.
- پیوسته یا گسسته بودن گفتار:** سیستم‌های بازشناسی گفتار ممکن است محدودیت‌هایی بر نحوه تلفظ کلمات توسط گوینده اعمال کنند. میزان پیوستگی یا گسستگی کلمات به سه دسته سیستم‌های بازشناسی کلمات مجزا^۴، سیستم‌های بازشناسی گفتار متصل^۵ و سیستم‌های بازشناسی گفتار پیوسته^۶ پیوسته^۶ طبقه‌بندی می‌شوند. در حالت مجزا، گفتار گوینده به صورت کلمه به کلمه و کاملاً مجزا از یکدیگر است. در سیستم‌های کلمات متصل، گوینده دنباله‌ای از کلمات محدود را تولید می‌کند. سیستم‌های بازشناسی پیوسته، گوینده را چندان به پیروی از قوانین خاصی در بیان گفتار مجبور نمی‌کنند بلکه گوینده، جملات را به طور پیوسته و طبیعی بیان می‌کند، هر چند در این حالت نیز فرض کتابی بودن گفتار وجود دارد. نوعی دیگر، که پیوسته اما کاملاً غیرماشینی است، گفتار محاوره‌ای^۷ است است که در آن گوینده می‌تواند گفتار را به صورت کاملاً طبیعی بیان کند. در چنین گفتاری، جملات ناقص، سرفه، تپق، مکث‌های طولانی و امثال آن وجود دارد.

1. real time
2. Speaker Dependent (SD)
3. Speaker Independent (SI)
4. Isolated Word Recognition (IWR)
5. connected word recognition
6. Continuous Speech Recognition (CSR)
7. spontaneous speech

حجم واژگان: حجم واژگان^۱ و تعداد کلمات مورد استفاده در یک سیستم‌بازشناسی، از عوامل مؤثر در دقت و سرعت سیستم است. بعضی از سیستم‌های بازشناسی فقط برای تشخیص تعداد محدودی کلمه طراحی شده‌اند درحالی‌که بعضی دیگر از سیستم‌ها قادرند مجموعه بزرگی از کلمات را تشخیص دهند.

محدودیت‌های زبانی: یکی از مهم‌ترین بخش‌های یک سیستم‌بازشناسی گفتار، مدل زبانی^۲ است، که درواقع بیانگر محدودیت‌های زبانی است. مدل زبانی یک زبان طبیعی مرکب از چهار جزء نمادها، دستور^۳، معنی‌شناسی^۴ و کاربردشناسی^۵ است. نمادهای زبان، واحدهایی هستند که پیام‌ها را تشکیل می‌دهند و درواقع کلمات یا واحدهایی کوچک‌تر از کلمه، نظیر هجاها یا واج‌ها هستند. دستور زبان مرکب از محدودیت‌های واژگانی^۶ و نحوی^۷ است که بیانگر نحوه شکل گرفتن کلمات از واحدهایی واحدهایی کوچک‌تر از کلمه و نیز شکل گرفتن جملات از کلمات است. جنبه معنایی مرتبط با نحوه ترکیب کلمات، برای شکل دادن پیغام‌های با معنا است. به‌عنوان مثال: جمله **صندلی غذا می‌خورد** از لحاظ نحوی درست ولی از لحاظ معنایی نادرست است. جنبه کاربردی یک زبان، در بالاترین سطح جای دارد و بیانگر وابستگی گفتار به گوینده‌ها و محیط است. محدودیت‌های معنایی و جنبه کاربردی به‌ندرت در سیستم‌های بازشناسی گفتار استفاده می‌شوند، زیرا که این محدودیت‌ها را به‌دشواری می‌توان به‌شکل فرمول بیان کرد. ولی محدودیت‌های دستوری تقریباً در تمامی سیستم‌های بازشناسی گفتار پیوسته، به‌صورت محدودیت‌های واژگانی و نحوی، مورد استفاده قرار می‌گیرند و تعداد جملات مجاز برای بازشناسی را کاهش می‌دهند و به‌عبارت‌دیگر، فضای مورد جستجو را کوچک‌تر می‌کنند. میزان محدودیتی که توسط مدل زبانی، درون یک سیستم‌بازشناسی ایجاد می‌شود، سرگشتگی^۸ حاصل از آن مدل زبانی نامیده می‌شود و سعی محققین، کاهش هرچه‌بیشتر سرگشتگی با استفاده از اطلاعات زبانی است.

کارایی در حضور نوفه و در محیط‌های کاربردی مختلف: بعضی از سیستم‌های بازشناسی تنها در شرایط محیطی مساعد و با نوفه کم قادر به تشخیص گفتار با دقت مطلوب هستند. کارایی این سیستم‌ها با تغییر محیط به‌شدت کاهش می‌یابد. این‌که سیستمی بتواند با تغییر محیط و حضور نوفه،

-
1. vocabulary size
 2. language model
 3. grammar
 4. semantics
 5. pragmatics
 6. lexical
 7. syntactic
 8. perplexity

دقت بازشناسی خود را حفظ کند، یکی از معیارهای مهم، به‌ویژه در کاربردی بودن آن، است. در نظر گرفتن این مهم بر پیچیدگی مسئله می‌افزاید.

ابهام آکوستیکی^۱ و میزان اشتباه^۲ بین کلمات: کلماتی که شکل نوشتاری متمایزی دارند ولی از لحاظ گفتاری مانند هم تلفظ می‌شوند، ابهام آکوستیکی ایجاد می‌کنند. کلماتی نظیر «سمر» و «شمر». همچنین کلماتی که تلفظ آن‌ها به یک‌دیگر شباهت دارد، مانند «دو» و «ته» که ممکن است به‌جای یک‌دیگر بازشناسی شوند و دقت را پایین بیاورند. هرچه تعداد این‌گونه کلمات در واژگان بیشتر شود، دقت سیستم بازشناسی پایین‌تر می‌آید. لذا برای جبران این مسئله، سیستم بازشناسی باید از مدل زبانی در سطوح دستوری و معنایی کمک بگیرد.

تحقیقات بی‌شماری برای طراحی و پیاده‌سازی موتورهای و سامانه‌های بازشناسی گفتار انجام شده است. چون مسائل زبانی و زبان‌شناسی از بخش‌های اساسی بازشناسی گفتار هستند، زبان یک سامانه بازشناسی گفتار، آن‌را به‌صورت کلی، از سایر سامانه‌ها مجزا می‌کند. برای بازشناسی گفتار در زبان فارسی نیز فعالیت‌های زیادی انجام شده و مقالات و پایان‌نامه‌های زیادی منتشر شده است. از مهم‌ترین آن‌ها می‌توان به (احدی، ۱۹۹۹)، (غلامپور، ۱۳۷۹)، (الماس‌گنج و دیگران، ۱۳۸۰)، (الماس‌گنج و دیگران، ۱۳۸۳)، (همایون‌پور، ۱۳۸۳)، (ولی، ۱۳۸۵)، (صامتی و دیگران، ۱۳۸۳)، (باباعلی و صامتی، ۲۰۰۴)، (صامتی و دیگران، ۲۰۰۸)، و (صامتی و دیگران، ۲۰۰۹) اشاره کرد. در این مقاله موتور بازشناسی گفتار به‌صورت اجمالی معرفی می‌شود و بعضی پژوهش‌های انجام یافته برای به‌کارگیری مدل زبانی زبان فارسی شرح داده می‌شود. در بخش ۲، ساختار کلی سیستم بازشناسی گفتار مورد بحث در این مقاله، ارائه و واحدهای آن به‌طور مختصر تشریح می‌شود. بخش ۳ به تشریح نحوه مدل‌سازی آوایی، اعم از مستقل‌ازبافت^۳ (CI) و وابسته‌به‌بافت^۴ (CD) اختصاص دارد. در بخش ۴ نحوه ساخت مدل زبانی برای موتور بازشناسی در زبان فارسی، و مسائل مربوط به آن، مورد بحث قرار می‌گیرد. بخش ۵ به توضیح در مورد آزمایش‌ها و ارائه نتایج حاصل از به‌کارگیری مدل‌های مختلف زبانی اختصاص دارد. نهایتاً در بخش ۶، خلاصه و نتیجه‌گیری ارائه می‌شود.

۲. معرفی کلی سیستم بازشناسی گفتار

موتور بازشناسی گفتار پیوسته فارسی، مورد بحث در این مقاله، حاصل به‌کارگیری آخرین روش‌های شناخته‌شده برای پیاده‌سازی واحدهای مختلف یک سامانه بازشناسی گفتار و انطباق آن با اطلاعات واج‌شناسی و خواص دستوری و واژگانی زبان فارسی است. شکل ۱، نمودار جعبه‌ای این موتور

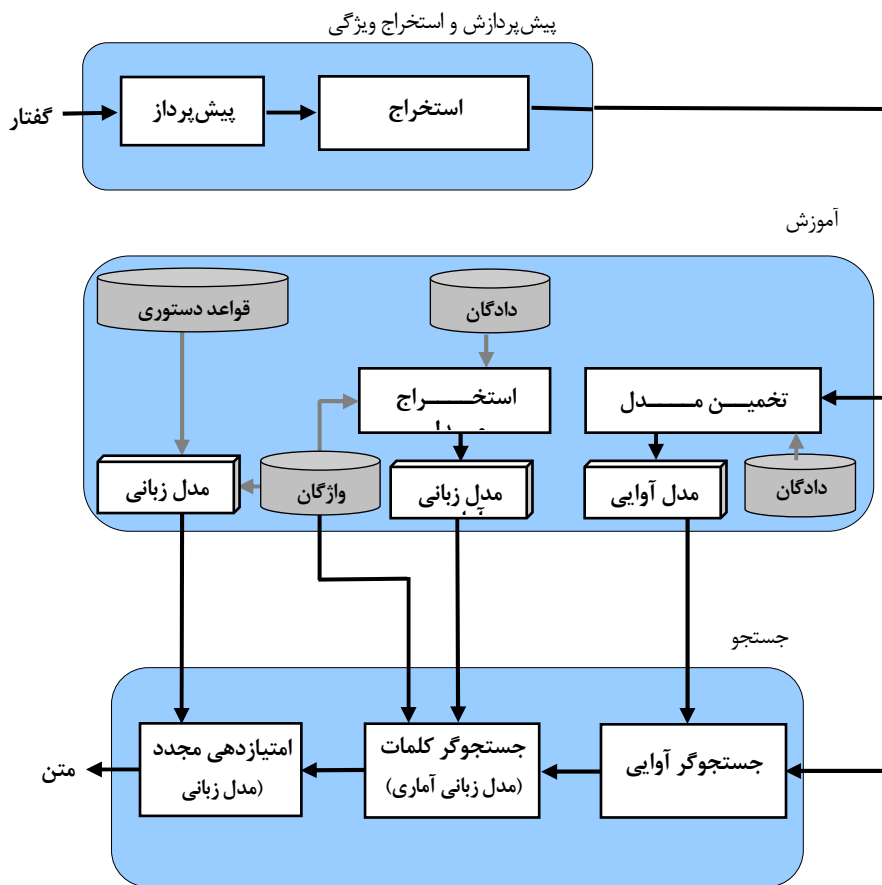
1. acoustic ambiguity
2. confusability
3. context independent
4. context dependent

را نشان می‌دهد. موتور به‌صورتی طراحی شده است که واحدهای مختلف آن به‌راحتی قابل تعویض و یا اصلاح باشند. بنابراین با ادامه کار گروه تحقیقاتی و به‌دست‌آمدن نتایج جدید، در مورد هر یک از واحدهای سامانه‌بازشناسی گفتار، به‌راحتی می‌توان اصلاحات موردنظر را در این سامانه وارد کرد. این سامانه شامل واحدهای مقاوم‌سازی نسبت به نوفه و شرایط محیطی است که وظیفه پیش‌پردازش، جداسازی سکوت از گفتار و حذف اولیه نوفه و استخراج مشخصه‌های اولیه برای بازشناسی را برعهده دارد. سپس این مشخصه‌ها با روش‌های متعدد و مؤثر، مقاوم‌سازی می‌شوند.

این سامانه هر دو نوع مدل وابسته به بافت و مستقل از بافت را به کار می‌گیرد. این مدل‌ها از نوع مدل مخفی مارکوف با مشاهدات پیوسته^۱ (CDHMM) هستند و شامل حالت‌هایی با تلفیق توزیع‌های گوسی در حوزه کپسترال هستند. در این مدل‌ها، پُرش به جلو، پُرش از روی حالت‌ها، و پُرش به خود حالت، مجاز است و ماتریس‌های کوواریانس^۲ از نوع قطری هستند (باباعلی و صامتی، ۲۰۰۴)، (صامتی و دیگران، ۲۰۰۸)، (صامتی و دیگران، ۲۰۰۹). پارامترهای مربوط به احتمال مشاهدات، با روش شباهت بیشینه^۳ (ML) آموزش داده شده، و مدل اولیه با قطعه‌بندی یک‌نواخت به‌دست می‌آید. هر دور از فرایند آموزش، شامل قطعه‌بندی با برنامه‌ریزی پویا (الگوریتم ویتربی)^۴ و تخمین مجدد پارامترهای مدل است که منجر به روش کا-مینز قطعه‌ای^۵ می‌شود (رایبیر^۶ و یوانگ^۷، ۱۹۹۳). در واحد جستجوگر آوایی^۸، یک جستجوی ویتربی با هرس از نوع شعاعی و هیستوگرام اجرا می‌شود. این بخش بخش واحدهای صوتی شناسایی شده را، به‌عنوان فرضیات اولیه، در اختیار واحد استخراج کلمات قرار می‌دهد. واحد استخراج، وظیفه جستجو در درخت واژگان را برعهده دارد. مدل‌های زبانی آماری و دستوری می‌توانند در هر دو واحد استخراج کلمات و امتیازدهی مجدد به کار روند. در این سامانه، از مدل زبانی آماری در واحد استخراج کلمات استفاده می‌شود و مدل‌های زبانی دستوری در بخش امتیازدهی مجدد به کار گرفته می‌شوند.

1. Continuous Density Hidden Markov Models
2. covariance
3. Maximum Likelihood
4. Viterbi
5. segmental K-means
6. L. Rabiner
7. B. H. Juang
8. acoustic decoder

به‌کارگیری اطلاعات زبانی در یک سیستم بازشناسی ...



شکل ۱. ساختار سامانه بازشناسی گفتار

۳. مدل‌سازی آوایی

برای این‌که گفتار به‌طور دقیق و مناسب مدل‌سازی شود، لازم است گفتار به بخش‌های کوچک‌تری به نام «واحد آوایی» تقسیم و هربخش جداگانه مدل‌سازی شود. بسته به نوع و کاربرد سیستم بازشناسی (پیوسته یا منفصل بودن بازشناسی، حجم واژگان و زبان مورد استفاده)، واحدهای آوایی

مختلفی ممکن است به کار روند. واحدهای آوایی رایج عبارت‌اند از: واج^۱، هجا^۲، دیاد (نیم‌هجا)^۳، کلمه و واحدهای آوایی وابسته به بافت (مانند دوواج^۴، سه‌واج^۵ و واج‌گونه^۶).

متداول‌ترین ساختارهایی که برای مدل‌سازی گفتار در بازشناسی مورد استفاده قرار می‌گیرند، عبارت‌اند از: مدل مخفی مارکوف (HMM)، مدل انطباق زمانی پویا^۷ (DTW)، شبکه عصبی مصنوعی^۸ مصنوعی^۸ (ANN) و مدل‌های ترکیبی^۹ (ترکیب HMM با ANN). از بین مدل‌های مذکور، مدل مخفی مارکوف به دلیل توانایی بالای آن در مدل‌سازی تغییرات زمانی سیگنال گفتار، بیش از سایر مدل‌ها در بازشناسی گفتار پیوسته با مجموعه واژگان بزرگ به کار می‌رود (رابینر و یوانگ، ۱۹۹۳). در این مقاله نیز از مدل مخفی مارکوف برای مدل‌سازی آوایی استفاده شده است.

برای آموزش مدل‌های آوایی نیاز به دادگان گفتاری داریم. زیرا برای مدل‌سازی لازم است از هر واحد آوایی زبان، به تعداد کافی نمونه ضبط شده در دست باشد. در عمل تعداد این واحدها بسیار زیاد است و معمولاً به مهم‌ترین آن‌ها اکتفا می‌شود. با توجه به کاربردهای مختلف بازشناسی گفتار، دادگان‌های گفتاری مختلفی طراحی شده‌اند. از آن‌جاکه در بسیاری از سیستم‌های بازشناسی امروزی، واج‌ها به‌عنوان کوچک‌ترین واحدهای بازشناسی انتخاب می‌شوند، در طراحی دادگان‌های گفتاری به واج‌ها توجه خاصی می‌شود. در این مقاله از دو دادگان گفتاری فارس‌دات کوچک (بی‌جن‌خان و دیگران، ۱۹۹۴)، و فارس‌دات بزرگ (شیخ‌زادگان و بی‌جن‌خان، ۱۳۸۵) برای مدل‌سازی آوایی استفاده کرده‌ایم. برای مدل‌سازی آوایی، دو رویکرد عمده را به کار برده‌ایم: مدل‌سازی مستقل ازبافت و مدل‌سازی وابسته به بافت. در مدل‌سازی مستقل ازبافت، هر واحد آوایی به‌طور مستقل و بدون در نظر گرفتن واحدهای آوایی مجاور، مدل‌سازی می‌شود. در این مقاله از واج به‌عنوان واحد آوایی پایه استفاده شده است و به‌ازای هر کدام از ۲۹ واج زبان فارسی، یک مدل مخفی مارکوف آموزش داده شده است. یک مدل مخفی مارکوف نیز برای مدل‌سازی سکوت و نوفه در نظر گرفته شده است. بنابراین در مدل‌سازی مستقل ازبافت جمعاً ۳۰ مدل داریم.

در مدل‌سازی وابسته به بافت، هر واحد آوایی (مثلاً هر واج)، با توجه به واحدهای آوایی مجاور آن، مدل‌سازی می‌شود. واحدهای آوایی در هنگام مجاورت با هم، اثر متقابل بر یک‌دیگر می‌گذارند و تغییراتی در طرز تلفظ آن‌ها ایجاد می‌شود، که این تغییرات بسته به نوع واحدهای مجاور متفاوت است.

-
1. phoneme
 2. syllable
 3. dyad
 4. diphone
 5. triphone
 6. allophone
 7. Dynamic Time Warping
 8. Artificial Neural Network
 9. hybrid

مجاورت واج‌ها در گفتار پیوسته، باعث ایجاد «هم‌تولیدی»^۱ و «همگونی»^۲ می‌شود. «هم‌تولیدی» به این معناست که مجاورت واج‌ها با یکدیگر، در ناحیه مرزی بین واج‌ها تأثیر می‌گذارد. به عبارت دیگر، حالت شروع و پایان یک واج کاملاً وابسته به واج‌های اطراف آن در گفتار است. پدیده «همگونی» نیز در اثر وضعیت واحدهای تکلمی دخیل در تلفظ واج‌ها پیش می‌آید، زیرا وضعیت این واحدها هنگام گذار از یک واج به واج بعدی، به یک‌باره تغییر نمی‌کند و این باعث می‌شود که نحوه تلفظ یک واج، به واج یا واج‌های قبلی و یا واج‌های بعدی وابسته شود (به تلفظ‌های گوناگون یک واج در گفتار، «واج‌گونه» می‌گویند). به عنوان مثال واج [ŋ]، که یک واج بی‌واک است، در کلمه نهار [v/ŋ/p] در مجاورت واکه‌ها به صورت واک‌دار تلفظ می‌شود، یا در کلمه موم [μμμ]، واج [v] در مجاورت واج خیشومی [μ] به حالت خیشومی درمی‌آید. همچنین در کلمه حدس [ŋαδσ]، واج [δ] در مجاورت واج بی‌واک [σ] «واک‌رفته» شده و شبیه [τ] تلفظ می‌شود. اثر واج‌های مجاور گاهی نیز باعث می‌شود که یک واج کاملاً به واج دیگری تبدیل گردد، مانند واج [v] در کلمه شنبه [α.vβε] که در تلفظ به واج [μ] تبدیل می‌شود.

بنابر آنچه گفته شد، خصوصیات طیفی یک واج در گفتار، کاملاً به بافت اطراف آن وابسته است. برای مدل‌سازی دقیق و مناسب گفتار، به منظور بازشناسی، باید دگرگونی‌های حاصل از بافت را نیز در نظر گرفت و در واقع بافت آوایی را به نحوی در مدل‌سازی دخالت داد. تجربه نشان داده است که اطلاعات وابسته به بافت، نقش مهمی در بازشناسی دارند و با به کارگیری آن‌ها میزان خطای بازشناسی به مقدار قابل توجهی کاهش می‌یابد. برای این منظور ساختارهای وابسته به بافت متنوعی برای مدل‌سازی اثرات بافت آوایی به کار رفته‌اند. در این مقاله از دو روش مدل‌سازی وابسته به بافت استفاده شده است: مدل‌سازی مبتنی بر سه‌واج و مدل‌سازی مبتنی بر دسته‌بندی واج‌ها.

۳.۱. مدل‌سازی مبتنی بر سه‌واج

واحد آوایی سه‌واج یکی از زیرکلمه‌های وابسته به بافت بسیار پرکاربرد و مؤثر در بازشناسی گفتار است. در سه‌واج‌ها، هر واج با در نظر گرفتن واج قبلی و واج بعدی‌اش در گفتار به صورت یک واحد آوایی مستقل در نظر گرفته شده و مدل می‌شود. برای یک واج مانند P ، مجموعه سه‌واج‌های مرتبط به آن را به صورت P_L-P-P_R نمایش می‌دهیم.

مزیت استفاده از مدل‌سازی مبتنی بر سه‌واج در این است که برای هر واج، اثرات هر دو واج مجاور در مدل‌سازی دخالت داده می‌شود. در عوض، مشکل کمبود داده آموزشی در این‌جا بیشتر مشهود است، زیرا تعداد سه‌واج‌ها، به لحاظ نظری، برابر با تعداد واج‌ها به توان سه می‌شود، که تعداد بسیار زیادی

1. coarticulation
2. assimilation

است (برای مثال در زبان فارسی که حدود ۳۰ واج دارد، تعداد سه‌واج‌ها ۳۰^۳ می‌شود). اگرچه در عمل تعداد زیادی از سه‌واج‌ها در زبان گفتاری به کار نمی‌روند ولی بازهم تعداد باقی‌مانده بسیار زیادند. زیادبودن تعداد سه‌واج‌ها باعث می‌شود که در یک دادگانِ گفتاریِ محدود، بسیاری از آن‌ها خیلی به‌ندرت تلفظ شوند و به‌همین دلیل، نمی‌توان مدلی دقیق‌تری برای آن‌ها آموزش داد. همچنین ممکن است بعضی از سه‌واج‌های موجود در زبان، در دادگانِ آموزشی اصلاً وجود نداشته باشند. بنابراین نمی‌توان مدلی برای این سه‌واج‌ها تشکیل داد (به این سه‌واج‌ها، «سه‌واج مشاهده نشده»^۱ می‌گویند). در این حالت اگر در دادگانِ آزمون با این سه‌واج‌ها مواجه شویم، برای بازشناسی آن‌ها هیچ مدلی در اختیار نداریم، در نتیجه، سیستم بازشناسی در این حالت نمی‌تواند کاملاً مستقل از مجموعهٔ واژگان باشد.

به‌منظور فائق آمدن بر مشکل کمبود دادهٔ آموزشی برای مدل‌های وابسته به متن، معمولاً سه‌واج‌های مشابه (یا حالت‌های مربوط به مدل آن‌ها) به‌هم گره زده می‌شوند. به‌این معنا که حالت‌های مشابه در سه‌واج‌های مختلف در یک دسته قرار می‌گیرند و با هم یکی در نظر گرفته می‌شوند. دو روش مشهور برای گره‌زدن حالت‌ها در مدل‌های مبتنی بر سه‌واج، دسته‌بندی مشتق‌شده از داده^۲ (یانگ^۳ و وودلند^۴، ۱۹۹۳) و دسته‌بندی مبتنی بر درختِ تصمیم‌گیری^۵ (یانگ و دیگران، ۱۹۹۴) است. در این مقاله، از روش دوم، که رایج‌تر است، استفاده کرده‌ایم. در هر دو روش این روش‌ها رویهٔ کلی آموزش به‌این‌صورت است که همهٔ سه‌واج‌های موجود در دادگانِ آموزش، ابتدا با استفاده از همان تعداد دادهٔ آموزشی موجود، آموزش داده می‌شوند، و سپس حالت‌های مربوط به سه‌واج‌های با بافت مشابه، دسته‌بندی می‌شوند و حالت‌های شبیه‌تر در یک دسته قرار می‌گیرند. در نهایت، حالت‌های موجود در یک دسته به‌عنوان یک حالتِ واحد در نظر گرفته، و به اصطلاح به هم گره زده می‌شوند. در مرحلهٔ بعد، حالت‌های گره‌زده‌شده، مجدداً تحت آموزش قرار می‌گیرند.

در روش مبتنی بر درختِ تصمیم‌گیری، که در این مقاله استفاده شده است، رویهٔ دسته‌بندی حالت‌ها با استفاده از درختِ تصمیم‌گیری صورت می‌گیرد. در این روش، تمام حالت‌های مربوط به بافت‌های مشابه (واجِ میانیِ یک‌سان در سه‌واجی‌ها) و همچنین مکان مشابه در مدل، در ریشهٔ درخت قرار می‌گیرند. سپس تعدادی پرسش در مورد بافتِ چپ و راستِ حالت‌ها مطرح می‌شود. با توجه به پرسش مطرح شده، حالت‌ها به دو دسته تقسیم می‌شوند، و دو گرهٔ جدید در درخت ایجاد می‌کنند. از بین پرسش‌های مطرح‌شده، پرسشی در نهایت انتخاب می‌شود که دو دستهٔ حاصل از آن بیشترین

-
1. unseen triphone
 2. tie
 3. data driven clustering
 4. S. J. Young
 5. P. C. Woodland
 6. tree-based state tying decision

افزایش در میزان شباهت^۱ را نسبت به دسته اولیه، داشته باشند. این روند ادامه پیدا می‌کند و درخت به تدریج رشد می‌کند تا در نهایت به جایی برسد که افزایش در میزان شباهت از حد آستانه خاصی کمتر شود و یا داده آموزش، برای حالت‌های موجود در گره‌های درخت، به حد کافی نباشد. پس از تشکیل درخت، هر برگ از درخت نشان‌دهنده یک طبقه از حالت‌هاست. حالت‌های موجود در هر کلاس به هم گره زده می‌شوند و یک حالت را تشکیل می‌دهند.

در این مقاله برای آموزش مدل‌های مبتنی بر سه‌واج از دادگان‌های فارس‌دات کوچک و فارس‌دات بزرگ استفاده کرده‌ایم. فرهنگ لغت به‌کاربرده شده برای دادگان فارس‌دات کوچک شامل ۱۱۴۷ کلمه و برای فارس‌دات بزرگ شامل ۴۷۸۰۲ کلمه است. با استفاده از فرهنگ لغات می‌توان تمام انواع سه‌واج‌های دیده‌شده در دادگان آموزش را مشخص کرد. با توجه به این فرهنگ لغات، در مجموع، ۴۳۶۲ نوع سه‌واج درون کلمه‌ای^۲ در دادگان فارس‌دات کوچک و ۱۶۳۶۴ نوع سه‌واج درون کلمه‌ای در دادگان فارس‌دات بزرگ وجود دارد. در روند آموزش، تعداد حالت‌ها برای هر مدل، ۵ مورد و تعداد گوسی‌ها در هر حالت، ۱ مورد در نظر گرفته شد. در نتیجه تعداد کل حالت‌ها در مجموعه دادگان فارس‌دات کوچک، $4362 \times 5 = 21810$ و در مجموعه دادگان فارس‌دات بزرگ، $16364 \times 5 = 81820$ است.

در مجموع، ۱۳۰ پرسش خودکار با استفاده از روش مذکور، در مقاله سینگ^۳ (سینگ و دیگران، ۱۹۹۹)، برای ساخته شدن درخت تصمیم‌گیری به کار برده شده است، که ۶۵ عدد از این پرسش‌ها در خصوص بافت چپ و ۶۵ عدد بقیه در خصوص بافت راست مطرح می‌شوند. با توجه به تعداد واج‌ها (۳۰ واج)، و تعداد حالت‌ها در هر مدل (۵ حالت)، تعداد $30 * 5 = 150$ درخت تصمیم‌گیری تشکیل می‌شود. بعد از انجام رویه دسته‌بندی و پیداشدن حالت‌های مشابه، داده‌های آموزش مربوط به همه این حالت‌ها برای آموزش یک حالت عمومی (گره زده شده) به کار برده می‌شود. به این حالت عمومی، زنون^۴ گفته می‌شود. تعداد زنون‌ها را می‌توان با تغییر مقادیر آستانه در الگوریتم دسته‌بندی، کنترل کرد. در مرحله آموزش زنون‌ها می‌توان تعداد گوسی‌ها را در هر حالت افزایش داد. در این مقاله دسته‌بندی و آموزش سه‌واج‌ها را با تعداد زنون‌های مختلف و تعداد گوسی‌های مختلف در هر زنون انجام دادیم. آزمایش‌ها نشان می‌دهند که تعداد ۵۰۰ زنون در فارس‌دات کوچک و ۴۰۰۰ زنون در فارس‌دات بزرگ و همچنین ۸ گوسی در هر زنون، بهترین دقت بازشناسی را نتیجه می‌دهند. نتایج بازشناسی در بخش ۵ آمده است.

1. likelihood
2. within-word triphone
3. R. Singh
4. senone

۲.۳. مدل‌سازی وابسته‌به‌بافت با استفاده از دسته‌بندی واج‌ها

در این روش، از ایده واحدهای آوایی چندگانه^۱ (مدل‌سازی واج‌گونه‌ها^۲)، برای مدل‌سازی وابسته‌به‌بافت استفاده کرده‌ایم. مزیت استفاده از این روش این است که می‌توان با توجه به حجم داده آموزش در دسترس، تعداد واحدهای آوایی وابسته‌به‌بافت را انتخاب کرد.

اندیشه کلی واحدهای آوایی چندگانه این است که یک واج در کلمه‌ها و متن‌های مختلف، به حالت‌های گوناگونی تلفظ می‌شود. بنابراین می‌توان گفت که خصوصیات طیفی یک واج (واحد آوایی مستقل از بافت) در بافت‌های مختلف دچار تغییرات و دگرگونی‌هایی می‌شود. این تغییرات را می‌توان در چند دسته کلی جای داد. به عبارت دیگر، هر واحد آوایی مانند P ، بسته‌به‌این که در چه بافتی قرار دارد، چند حالت گوناگون خواهد داشت (مانند P_1, P_2, P_3). پس می‌توان به جای یک واحد آوایی مستقل از بافت، انواع گوناگون وابستگی به بافت آن را به عنوان واحدهای بازشناسی به کار برد.

روش‌های مختلفی برای پیدا کردن خودکار حالات گوناگون یک واج وجود دارد. یک روش مناسب برای دسته‌بندی حالات گوناگون یک واج، استفاده از تعدادی قواعد آوایی است. با استفاده از این قواعد می‌توان نوع تأثیرپذیری یک واج را از واج‌های اطرافش تعیین کرده و براساس آن، واج را در یک دسته خاص جای داد (غلامپور، ۱۳۷۹). برای به کارگیری این روش نیاز به قواعد آواشناسی داریم. در این مقاله به دلیل عدم دسترسی به قواعد آوایی زبان فارسی، دسته‌بندی حالات گوناگون واج‌ها را به صورت بی‌نظارت و با استفاده از الگوریتم دسته‌بندی کا-مینز^۳ انجام داده‌ایم (بحرانی و صامتی، ۱۳۸۴). دلیل انتخاب این الگوریتم سادگی و سهولت پیاده‌سازی آن است.

برای دسته‌بندی واج‌ها باید هر نمونه واج از گفتار را با یک دنباله ویژگی^۴ نمایش داد. برای این منظور لازم است دادگان گفتاری، تقطیع واجی شده باشد، به این معنا که واج‌های تلفظ‌شده در گفتار و مرز آن‌ها در سیگنال گفتار دقیقاً مشخص باشد.

ابتدا سیگنال گفتار، پیش‌پردازش می‌شود. به این صورت که سیگنال گفتار به تعدادی قاب^۵ با طول مساوی (و با مقداری همپوشانی بین قاب‌ها) تقسیم شده و از هر قاب، پس از طی مراحل، تعدادی ضریب بازنمایی (در این سیستم، ضرایب مل-کپستروم) به عنوان بردار ویژگی استخراج می‌گردد. چون دادگان گفتاری، تقطیع شده هستند، قاب‌های ابتدایی و انتهایی هر واج را در سیگنال گفتار می‌توان مشخص کرد. بنابراین هر واج از تعدادی قاب تشکیل شده که هر قاب نیز با یک بردار ویژگی^۶ n بُعدی مشخص می‌گردد. پس می‌توان برای هر واج یک دنباله ویژگی به صورت دنباله‌ای از l بردار ویژگی قاب

1. multiple phone units
2. allophones
3. clustering
4. K-means
5. feature sequence
6. frame

تعریف کرد که l مشخص کننده تعداد قاب‌های تشکیل دهنده واج است. در این مقاله، تعداد قاب تشکیل دهنده دنباله ویژگی یک واج را به عنوان «طول دنباله ویژگی» می‌شناسیم. تعداد قاب تشکیل دهنده هر واج بسته به میزان کشیدگی زمانی آن واج در گفتار متفاوت است، بنابراین دنباله ویژگی واج‌ها دارای طول‌های متفاوتی خواهند بود.

پس از پیش‌پردازش می‌توان به‌ازای هر یک از واج‌های زبان، دنباله ویژگی تمام نمونه‌های آن را در دادگان گفتاری جمع‌آوری کرد. اگر زبان دارای N واج باشد، N مجموعه دنباله ویژگی داریم که هر مجموعه مربوط به یکی از واج‌های زبان خواهد بود. حال می‌توان هر کدام از این مجموعه‌ها را به چند دسته تقسیم کرد به طوری که هر دسته بیان‌گر یکی از گوناگونی‌های آن واج باشد. به دلیل غیرهم‌طول بودن دنباله‌های ویژگی واج‌ها، برای دسته‌بندی آن‌ها با استفاده از الگوریتم k -مینر با دو مشکل مواجه هستیم. مشکل اول این که نمی‌توان برای محاسبه مرکز دسته‌ها از میانگین‌گیری معمولی دنباله‌ها استفاده کرد و مشکل دوم این که برای محاسبه فاصله بین دنباله‌ها، توابع فاصله معمولی (مانند فاصله اقلیدسی) را نمی‌توان به کار برد. برای حل این مشکلات از روش انطباق زمانی پویا^۱ (DTW) استفاده کرده‌ایم. پس از تعیین تعداد دسته مناسب برای هر واج، که معمولاً بین ۱ تا ۵ دسته است، الگوریتم k -مینر را برای دسته‌بندی دنباله‌های ویژگی مربوط به نمونه‌های مختلف آن واج به کار می‌بریم. به جای تابع فاصله از روش انطباق زمانی پویا استفاده می‌کنیم و برای محاسبه مرکز دسته‌ها، روش «انطباق و میانگین‌گیری» (بحرانی و صامتی، ۱۳۸۴) را به کار می‌بریم.

پس از اعمال الگوریتم دسته‌بندی، واج مورد نظر به K دسته تقسیم می‌شود که هر دسته را می‌توان به عنوان یکی از حالات گوناگون آن واج در نظر گرفت. در واقع هر دسته بیان‌گر یک واحد آوایی وابسته به بافت است و دنباله‌های ویژگی‌ای که در یک دسته قرار دارند داده‌های آموزش را برای مدل‌سازی آن واحد آوایی فراهم می‌سازند. به عنوان مثال اگر واج $[a]$ به سه دسته تقسیم شود، سه واحد آوایی $[a1]$ ، $[a2]$ و $[a3]$ تولید می‌شود که هر یک بیان‌گر یکی از گوناگونی‌های واج $[a]$ در متن‌های مختلف خواهد بود. حال به جای آن که یک مدل برای واج $[a]$ ساخته شود، برای هر کدام از سه نوع گوناگون واج $[a]$ ، مدل جداگانه‌ای ساخته می‌شود. مدل‌سازی می‌تواند به طور مستقیم با استفاده از دنباله‌های ویژگی موجود در هر دسته صورت گیرد. تعداد کل مدل‌هایی که باید آموزش داده شوند برابر با تعداد کل واحدهای آوایی وابسته به بافت و به عبارت دیگر برابر با تعداد کل دسته‌های مربوط به همه واج‌ها است. آزمایش‌ها نشان می‌دهند که استفاده از مدل وابسته به بافت در سامانه بازشناسی گفتار پیوسته فارسی، نرخ خطای کلمات را به میزان ۲٪ کاهش می‌دهد.

1. Dynamic Time Warping

۴. مدل‌سازی زبانی

یکی از مؤثرترین راه‌های افزایش دقت سیستم‌های بازشناسی گفتار پیوسته، به‌کارگیری اطلاعات زبانی است (به‌صورت آماری، نحوی و معنایی). انسان‌ها به‌هنگام شنیدن گفتار، علاوه بر اطلاعات آوایی، از اطلاعات زبانی هم استفاده می‌کنند. آن‌ها از احتمالات رخداد یک کلمه، که به‌نوعی در ذهن آن‌ها نقش بسته است، استفاده می‌کنند. اگر بخواهیم اطلاعات زبانی را، که انسان‌ها به‌طور ناخودآگاه به‌کار می‌برند، سطح‌بندی کنیم، در اولین سطح، اطلاعات واژگانی و در سطوح بعد اطلاعات نحوی و معنایی قرار دارند (آلن^۱، ۱۹۹۵). در کاربردهای بازشناسی گفتار، سعی می‌شود که این سطوح از اطلاعات، مدل‌سازی شوند. هدف از مدل‌سازی زبان این است که سیستم بازشناسی، به‌سمت جمله‌هایی متناسب با کاربرد موردنظر سوق داده شود و از جمله‌های بی‌معنا اجتناب گردد. به‌این‌منظور مدل‌های زبانی متنوعی پیشنهاد شده‌اند که هر یک سعی می‌کنند تمام یا بخشی از سطوح اطلاعاتی مذکور را پوشش دهند.

برای استخراج مدل‌های زبانی در هر زبان خاص، نیاز به حجم عظیمی از داده‌های متنی آن زبان داریم. هرچه حجم داده‌های مورد استفاده بیشتر باشد، تخمین بهتری از این مدل‌ها به‌دست می‌آید. داده‌های متنی مورد استفاده برای استخراج مدل‌های زبانی، به‌صورت دادگان‌هایی به‌نام پیکره متنی^۲ در دسترس هستند. پیکره‌های متنی برای هر زبان، شامل متن‌های مختلف آن زبان در زمینه‌های مختلف و در حجم بالاست (در حدود چند میلیون کلمه)، که معمولاً با برچسب‌های مختلف زبانی و گرامری مشخص شده‌اند. در این مقاله از نسخه اولیه^۳ "پیکره متنی زبان فارسی" (بی‌جن‌خان و دیگران، ۲۰۰۹)، برای استخراج مدل‌های زبانی آماری استفاده کرده‌ایم.

۴.۱. مدل زبانی چندکلمه‌ای

مدل زبانی چندکلمه‌ای^۳ از متداول‌ترین مدل‌های زبانی مورد استفاده در سیستم‌های بازشناسی گفتار پیوسته است. مدل‌های چندکلمه‌ای احتمال شرطی رخداد یک کلمه در زبان، پس از هر رشته $n-1$ کلمه‌ای از آن زبان را مشخص می‌سازند. در این روش احتمال دنباله کلمات $W = w_1 w_2 \dots w_T$ از رابطه زیر محاسبه می‌شود:

$$p(W) = p(w_1 w_2 \dots w_T) \cong \prod_{i=1}^T p(w_i | w_{i-n+1} \dots w_{i-1}) \quad (1)$$

1. J. Allen
2. text corpus
3. n-gram

به‌کارگیری اطلاعات زبانی در یک سیستم بازشناسی ...

چون همیشه میزان داده‌های متنی در دسترس برای آموزش مدل‌های چندکلمه‌ای محدود است، معمولاً n برابر با ۱، ۲، یا ۳ انتخاب می‌شود، و مدل‌های آماری استخراج‌شده به ترتیب یک‌کلمه‌ای^۱، دوکلمه‌ای^۲ و سه‌کلمه‌ای^۳ نامیده می‌شوند. همان‌طور که گفتیم در این مقاله از "پیکره متنی زبان فارسی" برای آموزش مدل‌های زبانی استفاده شده است. متون "پیکره متنی زبان فارسی" از منابع مختلف شامل روزنامه‌ها، مجلات و کتاب‌های مختلف فراهم شده است. این پیکره شامل حدود ۱۰۰ میلیون کلمه است ولی نسخه اولیه آن، که در این مقاله استفاده شده، شامل حدود ۱۰ میلیون کلمه است که به‌زای هر کلمه آن یک برچسب مشخص شده است. این برچسب‌ها شامل مقوله نحوی یا اجزای کلام (POS)^۴، و در صورت لزوم شامل زیرمقوله‌های نحوی - معنایی هر کلمه هستند. در کل ۸۸۲ نوع برچسب در پیکره متنی موجود است. با استفاده از این پیکره متنی، مدل‌های چندکلمه‌ای مختلفی شامل مدل‌های چندکلمه‌ای مبتنی بر کلمه، مبتنی بر مقوله نحوی و مبتنی بر طبقه را استخراج کرده‌ایم. مدل‌های چندکلمه‌ای مبتنی بر مقوله نحوی و مبتنی بر طبقه، از آن جهت استخراج شدند که در وهله اول به‌نظر می‌رسید میزان داده‌های متنی برای آموزش مدل‌های مبتنی بر کلمه کافی نباشد. برای استخراج مدل‌های چندکلمه‌ای مختلف از پیکره متنی، با دو مشکل مواجه بودیم. مشکل اول یک‌دست‌نبودن پیکره از لحاظ املائی، و مشکل دوم زیادبودن تعداد برچسب‌ها در پیکره متنی بود. مشکل اول از سیستم نگارش فارسی و استفاده از الفبای عربی در نگارش فارسی ناشی می‌شود. در نگارش فارسی بعضی از واژگها^۵ (مانند پسوندها و پیشوندها) اجازه دارند که به‌صورت پیوسته یا جدا از کلمه اصلی نوشته شوند. در حالت جدانویسی نیز دو گزینه داریم: استفاده از فاصله به‌عنوان جداکننده و استفاده از نویسه ZWNJ^۶ یا نیم‌فاصله. به‌عنوان مثال سه شکل ممکن برای اتصال پسوند جمع "ها" و پیشوند استمراری "می" به کلمات، در جدول ۱ نشان داده شده است. همه این اشکال املائی در پیکره متنی استفاده شده‌اند.

جدول ۱. صورت‌های نگارشی مختلف برای اتصال پیشوند "می" و پسوند "ها" به کلمه اصلی

متصل	جدا بدون فاصله	جدا با فاصله
کتابها	کتاب‌ها	کتاب ها
می‌روند	می‌روند	می روند

1. monogram
2. bigram
3. trigram
4. parts of speech
5. morphemes
6. Zero-Width Non-Joiner

مشکل دیگر، عدم وجود استاندارد خاص برای نگارش کلمات فارسی است که باعث می‌شود یک کلمه خاص چند شکل املائی مختلف داشته باشد. به‌عنوان مثال کلمه **مسئولیت** می‌تواند به شکل‌های **مسوولیت** و **مسؤولیت** نیز نوشته شود. برای حل این مشکلات با بررسی‌های فراوانی که بر روی متون پیکره انجام گرفت، همه شکل‌های مختلف نگارشی یک کلمه به‌شکل واحدی تبدیل شدند. مثلاً همه پسوندها و پیشوندها با استفاده از نویسه نیم‌فاصله به کلمه اصلی متصل شدند و کلمات دارای چند شکل املائی، همگی به‌شکل استاندارد فرهنگستان زبان و ادب فارسی (صادقی و زندی‌مقدم، ۱۳۸۵) تبدیل شدند.

مشکل دیگر زیاد بودن تعداد برچسب‌ها در پیکره متنی است. همان‌طور که گفته شد، ۸۸۲ نوع برچسب در پیکره متنی به‌کار رفته است. علاوه‌براین که این تعداد برچسب بسیار زیاد و بعضی از آن‌ها بیش‌ازحد جزئی هستند، تعداد قابل‌ملاحظه‌ای از آن‌ها نیز بسیار کم به‌کار رفته‌اند. بنابراین با بررسی برچسب‌ها و استخراج آمار هر برچسب، آن‌ها را براساس تشابهشان (از لحاظ نحوی)، به ۱۶۴ دسته کلی تقسیم کردیم و به هر دسته، یک برچسب کلی اختصاص دادیم. با این کار هم تعداد کل برچسب‌ها کاهش می‌یابد و هم برچسب‌های کم‌کاربرد در دسته‌های کلی‌تر جای می‌گیرند. به تعدادی از برچسب‌ها نیز که بسیار کم به‌کار رفته بودند، و درضمن در هیچ دسته‌ای جای نمی‌گرفتند، برچسب IGNORE نسبت داده شد. یک برچسب NULL نیز به‌عنوان نشانه پایان جمله و شروع جمله جدید در نظر گرفته شد. بنابراین در کل، ۱۶۶ برچسب به‌دست آمد که در استخراج مدل‌های مبتنی بر مقوله نحوی مورد استفاده قرار گرفتند (بحرینی و دیگران، ۲۰۰۶). پس از یک‌دست‌سازی املائی و کاهش تعداد برچسب‌ها، آمارهای زیر از پیکره متنی استخراج شدند (بحرانی و دیگران، ۱۳۸۵):

آمار یک‌کلمه‌ای‌ها یا تعداد رخداد هر کلمه در پیکره متنی (با استفاده از این آمار ۲۰.۰۰۰ کلمه پرکاربرد از پیکره متنی استخراج و به‌عنوان مجموعه واژگان در نظر گرفته شد)، آمار دوکلمه‌ای‌ها، آمار سه‌کلمه‌ای‌ها، آمار هر برچسب، آمار دو برچسب متوالی، آمار سه برچسب متوالی و تعداد رخداد هر برچسب به‌زای هر کلمه (آمار برچسب-کلمه).

لازم‌به‌ذکر است که آمارهای دوکلمه‌ای، سه‌کلمه‌ای و برچسب-کلمه برای ۲۰.۰۰۰ کلمه موجود در مجموعه واژگان، استخراج شدند. سایر کلمات، همگی به‌عنوان یک کلمه خارج از واژگان در نظر گرفته شدند.

با استفاده از آمار استخراج‌شده، مدل‌های چندکلمه‌ای مبتنی بر کلمه و مبتنی بر مقوله نحوی ساخته شدند. از آن‌جا که در مدل‌های دوکلمه‌ای و سه‌کلمه‌ای، تعداد زیادی پارامتر با مقدار صفر وجود

دارند، روش هموارسازی ویتن^۱ - بل^۲ (ویتن و بل، ۱۹۹۱) برای از بین بردن مقادیر صفر بر روی مدل‌ها اعمال گردید.

علاوه بر مدل‌های چندکلمه‌ای مبتنی بر کلمه و مبتنی بر مقوله نحوی، مدل‌های چندکلمه‌ای مبتنی بر طبقه نیز از پیکره متنی استخراج شد. در این روش، کلمات دسته‌بندی شدند و هر کلمه در یک یا چند طبقه جای گرفت. سپس مدل‌های چندکلمه‌ای، بر اساس دسته‌های کلمات (به‌جای خود کلمات) استخراج گردیدند. مدل‌های مبتنی بر طبقه، مانند مدل‌های مبتنی بر مقوله نحوی، مشکل کمبود داده‌های آموزشی را تا حد زیادی برطرف می‌سازند. برای دسته‌بندی کلمات، روش‌های خودکار زیادی پیشنهاد شده‌اند. دو روش مشهور در این زمینه، الگوریتم مارتین^۳ (مارتین و دیگران، ۱۹۹۸)، و الگوریتم براون^۴ (براون و دیگران، ۱۹۹۲) هستند. در این مقاله از الگوریتم مارتین برای دسته‌بندی کلمات استفاده شده است. در این الگوریتم ابتدا یک دسته‌بندی اولیه از کلمات به وجود می‌آید و سپس کلمات، بین طبقات، آن قدر جابه‌جا می‌شوند تا معیار سرگشتگی^۵ به حداقل برسد. پس از دسته‌بندی کلمات، آمار دوکلمه‌ای و سه‌کلمه‌ای، بین دسته‌های کلمات استخراج می‌گردد، و مدل‌های دوکلمه‌ای و سه‌کلمه‌ای مبتنی بر طبقه با استفاده از این آمار ساخته می‌شوند.

مدل‌های چندکلمه‌ای مذکور، در بخش جستجوی کلمات^۶، در سیستم بازشناسی به کار می‌روند (شکل ۱). در این روش امتیاز مدل زبانی و امتیاز مدل آکوستیک "در حین جستجو" با هم ترکیب می‌شوند (هارپر^۷ و دیگران، ۱۹۹۴). به‌کارگیری مدل چندکلمه‌ای، در حین جستجو به این صورت است که هنگامی که الگوریتم جستجو، فرضیه‌های مختلف را برای بازشناسی کلمات به پیش می‌برد، با شناسایی یک کلمه جدید، احتمال چندکلمه‌ای آن را نیز، همراه با امتیاز آوایی آن، در امتیاز فرضیه ضرب می‌کند. به این معنا که اگر امتیاز کنونی یک فرضیه پس از شناسایی کلمه w_n ، S_n باشد، و این فرضیه پس از بسط داده شدن، کلمه w_{n+1} را به عنوان کلمه بعدی شناسایی کند، امتیاز جدید فرضیه برابر است با:

$$S_{n+1} = S_n \cdot S_{AM}(w_{n+1}) \cdot S_{LM}(w_{n+1})^{LMW} \quad (2)$$

-
1. I. Witten
 2. T. Bell
 3. S. Martin
 4. P. Brown
 5. perplexity
 6. word decoding
 7. M.P. Harper

$S_{AM}(w_{n+1})$ امتیاز مدل آوایی کلمه w_{n+1} و $S_{LM}(w_{n+1})$ امتیاز مدل زبانی آن است. معمولاً به دلیل تفاوت در مقیاس‌های $S_{AM}(w_{n+1})$ و $S_{LM}(w_{n+1})$ ، یک پارامتر وزن (LMW)^۱ برای امتیاز مدل زبانی در نظر گرفته می‌شود. معمولاً برای پرهیز از به‌کار بردن اعداد خیلی کوچک، به جای خود امتیازها از لگاریتم آن‌ها استفاده می‌گردد:

$$\log S_{n+1} = \log S_n + \log S_{AM}(w_{n+1}) + LMW \cdot \log S_{LM}(w_{n+1}) \quad (۳)$$

بنابراین لگاریتم امتیاز زبانی کلمه جدید به صورت وزن دار، با لگاریتم امتیاز آوایی کلمه و همچنین با لگاریتم امتیاز فرضیه جمع زده می‌شود و امتیاز جدید فرضیه را می‌سازد. امتیاز فرضیه به همین روش، با استفاده از ترکیب امتیازهای آوایی و زبانی به دست می‌آید، و در نهایت پس از کامل شدن فرضیه‌ها، فرضیه با بالاترین امتیاز (در بین N دنباله خروجی الگوریتم جستجو)، خروجی بخش بازشناسی سیستم خواهد بود. در این روش، مدل چندکلمه‌ای، در واقع، رویه جستجو را برای یافتن دنباله کلمات محتمل‌تر هدایت می‌کند.

برای مدل‌های یک کلمه‌ای، دو کلمه‌ای و سه کلمه‌ای مبتنی بر کلمه، به ترتیب با روابط زیر محاسبه می‌شود:

$$S_{monogram}(w_{n+1}) = P(w_{n+1}) = \frac{N_{monogram}(w_{n+1})}{N_{total}} \quad (۴)$$

$$S_{bigram}(w_{n+1}) = P(w_{n+1} | w_n) = \frac{N_{bigram}(w_n w_{n+1})}{N_{monogram}(w_n)} \quad (۵)$$

$$S_{trigram}(w_{n+1}) = P(w_{n+1} | w_{n-1} w_n) = \frac{N_{trigram}(w_{n-1} w_n w_{n+1})}{N_{bigram}(w_{n-1} w_n)} \quad (۶)$$

در روابط بالا، $N_{monogram}(w_{n+1})$ ، $N_{bigram}(w_n w_{n+1})$ و $N_{trigram}(w_{n-1} w_n w_{n+1})$ به ترتیب آمارهای یک کلمه‌ای، دو کلمه‌ای و سه کلمه‌ای و N_{total} تعداد کل کلمات، در پیکره متنی است.

1. Language Model Weight

۲.۴. مدل زبانی دستوری

خروجی سیستم‌های بازشناسی گفتار پیوسته، معمولاً به صورت فهرستی از بهترین دنباله‌های کلمات (N-best) معادل با گفتار ورودی هستند. در این سیستم‌ها، بهترین دنباله‌های کلمات، معمولاً با در نظر گرفتن مدل‌های آوایی و مدل‌های زبانی آماری به دست می‌آیند، و امتیاز آن‌ها بر اساس این مدل‌ها محاسبه می‌شود و از میان آن‌ها، دنباله با بالاترین امتیاز به عنوان خروجی نهایی در نظر گرفته می‌شود. چون مدل آوایی و مدل زبانی آماری هیچ محدودیت دستوری بر روی خروجی‌ها اعمال نمی‌کنند، پس دنباله‌های کلمات حاصل ممکن است از لحاظ دستوری، درست یا نادرست باشند. بنابراین هیچ تضمینی برای درست بودن بهترین فرضیه از لحاظ دستوری نیست. در حالی که ما انتظار داریم که اگر جمله گفتاری ورودی از نظر دستوری درست باشد، سیستم بازشناسی، یک جمله درست دستوری به دست دهد. بنابراین لازم است خروجی‌های سیستم بازشناسی از لحاظ درست بودن دستوری نیز بررسی گردند. برای این کار احتیاج به یک مدل دستوری داریم. مدل دستوری، شامل الگوهای نحوی جملات و همچنین مشخصه‌های دستوری کلمات مورد استفاده است. با استفاده از این مدل، یک تحلیل گرنحوی، خروجی سیستم بازشناسی را پردازش می‌کند و دنباله‌های نادرست (از نظر دستوری) را رد کرده و در نهایت دنباله‌ای را به عنوان خروجی نهایی تولید می‌کند که هم از لحاظ دستوری درست باشد، و هم دارای بیشترین امتیاز (از لحاظ آوایی و آماری) در بین N دنباله خروجی باشد.

در این تحقیق از نظریه ساخت-گروهی تعمیم یافته، یا GPSG^۲ (گزار^۳ و دیگران، ۱۹۸۵)، برای ارائه مدل دستوری مورد نظر استفاده شده است. از آن جاکه GPSG یک نوع دستور ساخت-گروهی است، تلاش می‌کند جملات زبان را به صورت ترکیبی از تعدادی گروه دستوری^۴ در نظر بگیرد و سپس این گروه‌های دستوری را نیز ترکیبی از گروه‌های دستوری کوچک‌تر، و همین‌طور الی آخر در نظر گرفته تا نهایتاً به مرز کلمات برسد. در انتخاب گروه‌های زبان فارسی، با اندکی انحراف، از نظریه X-تیره^۵ (ردفورد^۱، ۱۹۸۸) پیروی شده و اسم (N)، فعل (V)، صفت (ADJ)، قید (ADV) و حرف اضافه (P) به عنوان مقوله‌های نحوی پایه‌ای در نظر گرفته شده‌اند، که می‌توانند به عنوان هسته گروه‌های اسمی (\bar{N}, \bar{N})، فعلی (\bar{V}, \bar{V})، صفتی (\bar{ADJ}, \bar{ADJ})، قیدی (\bar{ADV}, \bar{ADV}) و حرف اضافه‌ای (\bar{P}, \bar{P}) قرار گیرند. سپس تلاش شده است تا ساخت نحوی هر یک از این گروه‌ها بر اساس مقوله‌های نحوی کوچک‌تر صورت‌بندی شود. به عنوان مثال، \bar{N} ترکیب اسم با همه وابسته‌های پسین آن و \bar{N} ترکیب

1. feature
2. Generalized Phrase Structured Grammar
3. G. Gazdar
4. phrase
5. x-bar Theory
6. A. Radford

\bar{N} با همه وابسته‌های پیشین اسم در نظر گرفته شده‌اند. البته در برخی موارد، در نظر گرفتن دو لایه برای تعریف قواعدی که بتوانند تمامی ساختارهای نحوی ممکن یک گروه را دربر بگیرند، کافی نبود. در این‌گونه موارد بین لایه‌های ۱ و ۲، لایه‌های میانی تعریف شدند. به‌عنوان مثال، در مورد گروه‌های اسمی، به دلیل وجود چهار نوع وابسته پسین برای اسم (وابسته صفتی، وابسته حرف اضافه‌ای، وابسته گروه‌اسمی و وابسته جمله)، بالاخص با توجه به امکان تکرار وابسته صفتی، یک لایه \bar{N}^+ بین لایه‌های \bar{N} و \bar{N} تعریف شد، و \bar{N} ترکیب اسم با همه وابسته‌های صفتی آن، و \bar{N}^+ نیز ترکیب \bar{N} با سایر وابسته‌های پسین اسم در نظر گرفته شدند.

یکی از قواعد ساخت گروهی \bar{N}^+ در زیر نشان داده شده است:

$$\bar{N}^+ \rightarrow * \bar{N} [\text{GEN} +, \text{PRO} -] \bar{N} (\bar{P}) (\text{S} [\text{COMP} +, \text{GAP}])$$

این قاعده، ساختار گروه اسمی را با وابسته‌های پسین آن، هنگامی که اسم دارای وابسته گروه اسمی (مضاف الیه) است، نشان می‌دهد. همان‌گونه که گفته شد، \bar{N} در سمت راست قاعده فوق، ترکیب اسم با وابسته‌های صفتی احتمالی آن را نشان می‌دهد. مشخصه دوارزشی GEN+، بیان می‌کند که ترکیب \bar{N} باید حتماً دارای کسره اضافه باشد و مشخصه دوارزشی - PRO، تأکید می‌کند که \bar{N} نمی‌تواند دارای هسته ضمیر باشد.

در حدود ۱۷۰ قاعده، مشابه با قاعده فوق استخراج شدند، که بسیاری از ساختارهای رایج گروه‌های اسمی، فعلی، حرف اضافه‌ای، صفتی، قیدی و بالاخره جمله را پوشش می‌دهند. جزئیات بیشتر در مورد دامنه پوشش قواعد دستوری استخراج شده، در مقاله حافظی و دیگران (۱۳۸۵) آمده است.

۵. آزمایش‌ها و نتایج

در این بخش، نتایج به‌کارگیری مدل‌های زبانی در سامانه بازشناسی گفتار پیوسته زبان فارسی ارائه می‌شود. سیستم برای ۲۹ واج زبان فارسی، و سکوت به‌عنوان واج سی‌ام آموزش داده شد. مدل‌های مخفی مارکوف، که به‌کار گرفته شدند، از نوع چپ‌به‌راست و شامل ۶ حالت و ۱۶ تلفیق گوسی در هر حالت هستند. حجم واژگان، ۱۰۰۰ کلمه است و از زیرمجموعه‌های آموزش و آزمایش دادگان فارس‌دات کوچک، به‌ترتیب برای آموزش و بازشناسی استفاده شد. برای ارزیابی مدل‌های زبانی به‌کاررفته با پارامترهای مختلف، از دو معیار سرگشتگی و نرخ خطای کلمه^۱ (WER) استفاده شد. جدول ۲ نتایج بازشناسی را روی زیرمجموعه آزمایش فارس‌دات کوچک نشان می‌دهد، که شامل ۱۴۰

1. Word Error Rate

به کارگیری اطلاعات زبانی در یک سیستم بازشناسی ...

جمله از ۷ گوینده است. برای همه مدل های زبانی، روش هموارسازی^۱ ویتن- بل استفاده شد. تعداد دسته های مورد استفاده برای آزمایش های مربوط به مدل زبانی مبتنی بر طبقه، ۲۰۰ دسته است. اولین نتیجه ای که از نتایج این جدول حاصل می شود آن است که سیستم پایه^۲، بدون مدل زبانی، دارای نرخ خطای کلمات بالایی است و اطلاعات موجود در مدل زبانی عملکرد سیستم را بهبود می بخشد. از آن جاکه مدل زبانی آماری مبتنی بر کلمه، بهبود بیشتری را در عملکرد سیستم، به نسبت سایر مدل های زبانی آماری ایجاد می کند، از این مدل زبانی برای آزمایش های بعد استفاده شد. همچنین اثر استفاده از اطلاعات دستوری در کاهش نرخ خطای کلمات در جدول ۲ مشهود است.

جدول ۲. عمل کرد موتور بازشناسی در سطح کلمه در شرایط بدون نوفه.

مدل زبانی	نرخ خطای کلمات (WER%)
بدون مدل زبانی	۳۸.۱۴
سه کلمه ای مبتنی بر مقوله نحوی	۲۴.۶۸
سه کلمه ای مبتنی بر طبقه	۲۳.۴۰
سه کلمه ای مبتنی بر کلمه	۲۱.۷۶
سه کلمه ای مبتنی بر مقوله نحوی به همراه قواعد دستوری	۱۸.۲۰

جدول ۳ میزان سرگشتگی محاسبه شده روی ۷۵۰ جمله (حدود ۱۰.۰۰۰ کلمه) از دادگان فارس دات بزرگ را براساس مدل چندکلمه ای مبتنی بر کلمه (براساس واژگان ۲۰.۰۰۰ کلمه ای) نشان می دهد. به منظور کاهش اندازه حافظه مورد نیاز برای مدل زبانی، چندکلمه ای های با فراوانی کم، از مدل خارج شدند و تعداد آستانه مورد نظر برای تعیین حذف یا عدم حذف چندکلمه ای ها، به عنوان پارامتر میان بر^۳ (کلارکسون^۴ و روزنفلد^۵، ۱۹۹۷) در نتایج آمد. نتایج مندرج در جدول ۳ نشان می دهد که چگونه میان برهای مربوط به دو کلمه ای و سه کلمه ای، روی اندازه حافظه مورد نیاز و سرگشتگی مدل زبانی سه کلمه ای تأثیر می گذارند. نتیجه خلاصه ای که از این جدول می توان گرفت این است که مقادیر میان بر، اندازه حافظه مورد نیاز برای مدل زبانی را کاهش می دهد، ولی افزایش معناداری در سرگشتگی ایجاد نمی کند. براساس نتایج این جدول برای تعداد دو کلمه ای، مقدار میان بر برابر با صفر و

1. smoothing
2. baseline
3. cutoffs
4. P. Clarkson
5. R. Rosenfeld

مجله زبان و زبان‌شناسی

برای تعداد سه کلمه‌ای، مقدار آن برابر با ۱ در نظر گرفته شد. همچنین در جدول ۴، نرخ خطای کلمات موتورِ بازشناسی روی فارسی‌داتِ کوچک و بزرگ، با استفاده از مدلِ واجیِ مستقل‌زبافت و وابسته‌به‌بافت آمده است. حجمِ واژگان در این آزمایش ۲۰.۰۰۰ کلمه است. نتایجِ مندرج در این جدول نشان می‌دهد که استفاده از مدلِ وابسته‌به‌بافت در همهٔ حالات منجر به کاهشِ نرخِ خطای کلمات می‌شود. همچنین بهترین نتایج، در حالتِ استفاده از دادگانِ فارسی‌داتِ بزرگ، برای آموزش و آزمایش و مدل‌سازیِ وابسته به بافت، حاصل می‌شود.

جدول ۳. اثر پارامتر میان‌بُر بر اندازه و سرگشتگی یک مدلِ زبانی سه کلمه‌ای.

میان‌بُر (سه کلمه‌ای)	میان‌بُر (دو کلمه‌ای)	سرگشتگی	اندازه (MB)
۰	۰	۱۳۴.۵۴	۳۶
۰	۱	۱۳۴.۷۶	۲۰
۰	۲	۱۳۵.۸۲	۱۷
۱	۱	۱۴۳.۱۸	۱۰
۱	۲	۱۴۳.۲۶	۷.۸

جدول ۴. نرخ خطای کلمات موتورِ بازشناسی روی فارسی‌داتِ کوچک و بزرگ با استفاده از مدلِ واجیِ مستقل‌زبافت (واج)، و وابسته‌به‌بافت (سه‌واج).

آموزش		تست	
دادگان	وابستگی به بافت	فارسی‌داتِ بزرگ	فارسی‌داتِ کوچک
فارسی‌داتِ کوچک	مستقل	۲۹.۶۰	۲۵.۷۷
	وابسته	۲۰.۵۱	۱۶.۷۹
فارسی‌داتِ بزرگ	مستقل	۶.۱۰	۳۷.۳۹
	وابسته	۵.۲۱	۲۶.۸۵

۶. خلاصه و نتیجه‌گیری

در این مقاله، خلاصه‌گزارش تحقیقات انجام‌شده برای طراحی و ساخت یک موتور بازشناسی گفتار پیوسته زبان فارسی ارائه شد. این سامانه حاصل سال‌ها فعالیت تحقیقاتی بوده و اولین موتور بازشناسی زبان فارسی است که در عمل به‌کارگیری شده است. پس از تشریح ساختار کلی سامانه، روش‌های به‌کاررفته برای مدل‌سازی واجی و طراحی مدل زبانی موتور بازشناسی، مطرح شد. در مدل‌سازی آوایی، برای به‌دست‌آوردن نتایج مطلوب، از مدل‌های وابسته به بافت مبتنی بر سه‌واج، استفاده شد. در ضمن برای کاهش تعداد مدل‌های لازم، روش دسته‌بندی واج‌ها به کار رفت، و پس از طراحی اولیه، بر روی دادگان مناسب پیاده‌سازی و اجرا شد. تأکید اصلی این مقاله بر روی مدل‌های زبانی به‌کاررفته بود. مدل‌های زبانی آماری و دستوری در موتور بازشناسی استفاده شدند. مدل‌های زبانی آماری، از سه‌نوع مبتنی بر کلمه، مبتنی بر مقوله نحوی و مبتنی بر طبقه، استفاده و نتایج عمل‌کرد آن‌ها با یک‌دیگر مقایسه شد. در این بین، از مدل زبانی آماری مبتنی بر کلمه، بهترین نتایج حاصل شد. در ضمن اثر پارامتر میان‌بر، که نقش آستانه تعداد چندکلمه‌ای مشاهده‌شده برای حذف یا لحاظ‌شدن در مدل زبانی آماری را ایفا می‌کند، بررسی شد و مقدار بهینه آن به‌دست آمد. همچنین مدل زبانی دستوری به‌کاررفته، شرح داده شد و در عمل مورد آزمایش قرار گرفت.

منابع

- الماس گنج، فرشاد و سیدعلی سیدصالحی و محمود بی‌جن‌خان و حسین صامتی و جواد شیخ‌زادگان (۱۳۸۰). "شنوا ۱ - سیستم بازشناسی گفتار پیوسته فارسی"، مجموعه مقالات کنفرانس بین‌المللی مهندسی برق ایران، ۶۳-۵۶.
- _____ و سیدعلی سیدصالحی و محمود بی‌جن‌خان و حسین رازی‌زاده و محمدرضا اصغری (۱۳۸۳). "نرم‌افزار بازشناسی گفتار پیوسته فارسی شنوا ۲". مجموعه مقالات اولین کارگاه پژوهشی زبان فارسی و رایانه، ۸۲-۷۷.
- بحرانی، محمد و حسین صامتی (۱۳۸۴). "استخراج و مدل‌سازی واحدهای آوایی وابسته به بافت برای بهبود دقت بازشناسی گفتار پیوسته با روش دسته‌بندی واج‌ها". نشریه مهندسی برق و مهندسی کامپیوتر ایران. سال ۳، شماره ۱، ۵۱-۴۵.
- _____ و حسین صامتی و نازیلا حافظی و سعیده ممتازی و حامد موثق (۱۳۸۵). "به‌کارگیری پیکره متنی زبان فارسی در ساخت مدل‌های زبانی آماری برای سیستم‌های بازشناسی گفتار پیوسته فارسی". مجموعه مقالات دومین کارگاه پژوهشی زبان فارسی و رایانه، ۱۰۹-۹۲.

حافظی، محمدمهدی و حسین صامتی و نیلوفر منصوری و نیلوفر منتظری و محمد بحرانی و حامد موثق (۱۳۸۵). "ارائه یک مدل دستوری برای بهبود دقت سیستم‌های بازشناسی گفتار پیوسته فارسی". مجموعه مقالات دومین کارگاه پژوهشی زبان فارسی و رایانه، ۸۰-۹۱.

شیخ‌زادگان، جواد و محمود بی‌جن‌خان (۱۳۸۵). "دادگان‌های گفتاری زبان فارسی". مجموعه مقالات دومین کارگاه پژوهشی زبان فارسی و رایانه، ۲۶۱-۲۴۷.

صادقی، علی‌اشرف و زهرا زندی‌مقدم (۱۳۸۵). فرهنگ املائی خط فارسی. انتشارات فرهنگستان زبان و ادب فارسی: تهران.

صامتی، حسین و حامد موثق و باقر باباعلی و محمد بحرانی و خسرو حسین‌زاده و امین فاضل‌دهکردی و حمیدرضا ابوطالبی و هادی ویسی و یاسمین مگری و نیلوفر منتظری و محمد نظامی رنجبر (۱۳۸۳). "سیستم بازشناسی گفتار پیوسته فارسی با واژگان بزرگ". مجموعه مقالات اولین کارگاه پژوهشی زبان فارسی و رایانه، ۶۹-۷۵.

غلامپور، ایمان (۱۳۷۹). "بازشناسی مستقل از گوینده واج‌های فارسی در صحبت پیوسته". پایان‌نامه دکترا. دانشکده مهندسی برق، دانشگاه صنعتی شریف.

ولی، منصور (۱۳۸۵). "بازشناسی مقاوم گفتار به‌منظور جبران‌سازی تنوعات گفتار میکروفونی - تلفنی توسط شبکه‌های عصبی". پایان‌نامه دکترا. دانشکده مهندسی پزشکی، دانشگاه امیرکبیر.

همایون‌پور، محمدمهدی (۱۳۸۳). "سیستم بازشناسی گفتار پیوسته به کمک هیبرید شبکه عصبی و مدل مخفی مارکوف و با استفاده از مدل‌های زبانی و روش‌های جستجو". مجموعه مقالات اولین کارگاه پژوهشی زبان فارسی و رایانه، ۱۸۳.

- Ahadi, S. M. (1999). "Recognition of Continuous Persian Speech Using a Medium-sized Vocabulary Speech Corpus". In *Proc. Eurospeech99*. 863-866.
- Allen, J. (1995). *Natural Language Understanding*. the Benjamin/Cummings Publishing Company, Inc: Redwood City, CA.
- Babaali, B. & H. Sameti (2004). "The Sharif Speaker-Independent Large Vocabulary Speech Recognition System". *The 2nd Workshop on Information Technology & Its Disciplines*. Kish Island: Iran.
- Bahrani, M. & H. Sameti & N. Hafezi & H. Movasagh (2006). "Building and Incorporating Language Models for Persian Continuous Speech Recognition Systems". In *Proc. the 5th international conference on Language Resources and Evaluation*. 2590-2593. Genoa: Italy.
- Bijankhan, M. & J. Seikhzadeghan & M. Roohani & Y. Samareh & K. Lucas & M. Tebyani (1994). "FARSDAT-The Speech Database of Farsi Spoken Language". In *Proc. the 5th Australian International Conference on Speech Science and Technology*. 826-831.
- _____ & J. Seikhzadeghan & M. Bahrani & M. Ghayoomi (2011). "Lessons from Creation of a Persian Written Corpus: Peykare". *Language Resources and Evaluation Journal*. Vol. 45, No. 2. 143-164.
- Brown, P. & V. Della Pietra & P. deSouza & J. Lai & R. L. Mercer (1992). "Class-based

- n-gram models of natural language". *Computational Linguistics*, 18(4). 467-479.
- Clarkson, P. & R. Rosenfeld (1997). "Statistical Language Modeling Using the CMU-Cambridge Toolkit". In *Proc. Eurospeech97*. vol. 5. 2707-2710.
- Gazdar, G. & E. H. Klein & G. K. Pullum & I. A. Sag (1985). *Generalized Phrase Structure Grammar*. Harvard University Press: Great Britain.
- Harper, M. P. & L. H. Jamieson & C. D. Mitchell & G. Ying & S. Potisuk & P. N. Srinivasan & R. Chen & C. B. Zoltowski & L. L. McPheters & B. Pellom & R. A. Helzerman (1994). "Integrating Language Models with Speech Recognition". *AAAI-94 Workshop on the Integration of Natural Language and Speech Processing*. 139-146.
- Martin, S. & J. Liermann & H. Ney (1998). "Algorithms for bigram and trigram word clustering". *Speech Communication*. vol. 24. 19-37.
- Rabiner, L. & B. H. Juang (1993). *Fundamentals of Speech Recognition*. New Jersey: Prentice Hall.
- Radford, A. (1988). *Transformational Grammar*. Cambridge University Press.
- Sameti, H. & H. Veisi & M. Bahrani & B. Babaali & K. Hosseinzadeh (2008). "Nevisa, a Persian Continuous Speech Recognition System". In *Communications in Computer and Information Science, Advances in Computer Science and Engineering*. Vol. 6. 485-492. Springer Berlin: Heidelberg.
- _____. & H. Veisi & M. Bahrani & B. Babaali & K. Hosseinzadeh (2009). "A Large Vocabulary Continuous Speech Recognition System for Persian Language". submitted to *IEICE - Transactions on Information and Systems*.
- Singh, R. & B. Raj & R. M. Stern (1999). "Automatic Clustering and Generation of Contextual Questions for Tied States in Hidden Markov Models". In *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP'99)*. 117-120.
- Witten, I. & T. Bell (1991). "The Zero-Frequency Problem: Estimating the Probabilities of Novel Events in Adaptive Text Compression". In *IEEE Transactions on Information Theory* 37(4). 1085-1094.
- Young, S.J. & P. C. Woodland (1993). "The Use of State Tying in Continuous Speech Recognition". In *European Conference on Speech Communication and Technology (EUROSPEECH'93)*. Berlin: Germany ISCA. 2203-2206.
- _____. & J. J. Odell & P. C. Woodland (1994). "Tree-based State Tying for High Accuracy Acoustic Modeling". In *Human Language Technology Conference*. Association for Computational Linguistics Morristown.

