

بهبود عمل کرد سامانه بازشناسی گفتار پیوسته با ویژگی‌های استخراج شده از مانیفولدهای گفتاری در فضای بازسازی شده فاز

یاسر شکفته و فرشاد الماس گنج

آزمایشگاه پردازش گفتار، دانشکده مهندسی پزشکی، دانشگاه صنعتی امیرکبیر، تهران، ایران

چکیده

یکی از رویکردهای مؤثر در بهبود کارایی سامانه‌های بازشناسی گفتار، طراحی روش‌های متنوع استخراج ویژگی از سیگنال گفتار و ترکیب اطلاعات به دست آمده از آنهاست. تحقیقات اخیر نشان می‌دهد که سیگنال گفتار رفتار غیرخطی و آشوبی دارد؛ ولی از این مشخصه سیگنال گفتار در سامانه‌های بازشناسی پیوسته گفتار استفاده نمی‌شود. یکی از حوزه‌های مناسب برای نمایش مشخصه‌های پویا و غیرخطی سیگنال آشوبی، فضای بازسازی شده فاز (RPS) است، از این رو در این مقاله یک روش جدید استخراج ویژگی مبتنی بر RPS (LLRPS) پیشنهاد شده است. این ویژگی‌ها از امتیاز شباهت تراژکتوری سیگنال گفتار جاسازی شده در RPS با مجموعه‌ای از مانیفولدهای واجی از پیش تعیین شده محاسبه می‌شوند. سپس مقادیر احتمال پسین واجی به وسیله ساختار شبکه عصبی TMLP از روی ویژگی‌های LLRPS تخمین زده می‌شود. ساختار شبکه عصبی استفاده شده، به صورتی است که علاوه بر توانایی استخراج اطلاعات پویا، قابلیت پیاده‌سازی روش‌های متنوع ترکیب خروجی دارد. نتایج آزمایش‌ها بر روی مجموعه داده‌گان گفتاری فارسی نشان می‌دهد که ترکیب غیرخطی خروجی سامانه‌های بازشناسی، شامل ویژگی‌های متداول کپستروم MFCC و ویژگی‌های پیشنهادی LLRPS، به ترتیب منجر به بهبود ۳/۹۴ درصد در دقت بازشناسی قاب و ۴/۰۲ درصد در دقت بازشناسی واج نسبت به عمل کرد سامانه بازشناسی پایه شده است.

واژگان کلیدی: بازشناسی گفتار پیوسته، استخراج ویژگی، فضای بازسازی شده فاز، مانیفولدهای واجی، امتیاز درست‌نمایی، شبکه عصبی.

۱- مقدمه

سامانه تولید گفتار انسان شامل یک فرآیند پویای چندمتغیره است که خروجی آن به طور معمول به شکل یک مشاهده تک‌متغیره^۱ یا تک‌بعدی تحت عنوان سیگنال زمان-گسسته گفتار توسط میکروفون ضبط می‌شود. در فرآیند تولید سیگنال گفتار، به علت بروز برخی عوامل مانند ارتعاش غیرخطی تارهای صوتی و یا حرکات لایه‌ای هوا در مجرای صوتی، تظاهرات آشوب‌گونه^۲ در شکل سیگنال گفتار مشاهده شده است (Awrejcewicz, 1990; Berry, 1994; Herzfel, 1995; Jiang, 2002; Jiang, 2006). از این رو با توجه به خواص سیگنال‌های آشوب‌گونه، سیگنال گفتار را به صورت

یک سری زمانی با توصیف چندمتغیره، در فضای حالت^۳ یا فضای فاز^۴ می‌توان بیان نمود (Kokkinos, 2005; Haggmuller, 2006; Sun, 2007). این انتقال سیگنال به فضای حالت می‌تواند به صورتی انجام شود که رفتار تراژکتوری متناظر با آن، از لحاظ هندسی معادل با تراژکتوری واقعی آن باشد (Kantz, 1997).

در همین زمینه، تکنز روش‌های جاسازی فضایی^۵ در حوزه‌ای مشابه با فضای حالت (که فضای بازسازی شده فاز^۶ (RPS) نامیده می‌شود) را به عنوان یکی از تکنیک‌های مطالعه مشخصه‌های پویای غیرخطی و آشوبی معرفی کرد

³ State Space

⁴ Phase Space

⁵ Embedding Space

⁶ Reconstructed Phase Space

¹ Univariate

² Chaotic

(Takens, 1980). بر مبنای این روش‌ها، ثبت یک متغیر زمانی از رفتار سامانه پویا می‌تواند برای کسب اغلب خواص پویایی^۱ آن سامانه "کافی" باشد (Kantz, 1997). در نتیجه با انتقال سیگنال گفتار به حوزه RPS، می‌توان امید داشت که اطلاعات بیشتر و متمایزکننده‌تری از آن نسبت به دیگر روش‌های متداول استخراج کرد.

از همین اواخر تحقیقات متعددی در حوزه پردازش سیگنال گفتار بر مبنای استفاده از روش جاسازی سیگنال گفتار در RPS انجام شده است (Narayanan, 1995; Lindgren, 2004). اساس این تحقیقات، نمایش و استخراج اطلاعات مربوط به حوزه RPS سیگنال گفتار است که می‌تواند حاوی اطلاعات متمایزکننده‌تری نسبت به روش‌های متداول آنالیز سیگنال گفتار (مانند روش‌های استخراج ویژگی MFCC و PLP که مبتنی بر طیف هستند Davis, 1990; Hermansky, 1980) باشد. بخشی از این اطلاعات متمایز متأثر از مشخصه‌هایی است که به واسطه حذف اطلاعات فاز در روش‌های استخراج ویژگی مبتنی بر اندازه طیف سیگنال از بین رفته‌اند. با اینکه برخی از تحقیقات الهام گرفته شده از دستگاه شنوایی انسان ادعا می‌کنند که دستگاه شنوایی انسان به اطلاعات فاز سیگنال حساس نیستند؛ اما برخی تحقیقات منتشر یافته اخیر، نشان دهنده مؤثر بودن اطلاعات فاز در سامانه بازشناسی گفتار است (Paliwal, 2005; Hegde, 2007; Alsteris, 2007).

در حوزه پردازش سیگنال گفتار، برخلاف انجام تحقیقات متعدد با روش جاسازی، متأسفانه تاکنون روش‌های محدودی برای استفاده از قابلیت آن در سامانه‌های بازشناسی گفتار پیوسته ارائه شده است. عمده فعالیت‌های صورت گرفته در این حوزه، شامل استخراج ویژگی‌های آشوب گونه مانند بُعد همبستگی^۲، بُعد فرکتال^۳، بُعد فرکتال تعمیم یافته^۴ و نماهای لیاپانوف^۵ برای رده‌بندی واجی است (Yu, 2006; Pitsikalis, 2002&2009; Vaziri, 2010). در دو مشخصه جدید از سیگنال جاسازی شده در RPS معرفی شده است. این ویژگی‌ها مقادیر عددی هستند که مقدار جابه‌جایی و دایروی بودن مسیر تراژکتوری سیگنال در فضای فاز را بازنمایی می‌کنند.

برای اولین بار در سال ۲۰۰۲، از توزیع آماری سیگنال گفتار جاسازی شده در RPS با استفاده از روش هیستوگرام، برای رده‌بندی واج‌های مجزاً^۶ در مجموعه دادگان TIMIT استفاده شد (Ye, 2002). سپس در (Ye, 2003) استفاده از آنالیز مؤلفه‌های اساسی (PCA) در متعامد کردن و کاهش بُعد سیگنال جاسازی شده در فضای فاز بررسی شد. در ادامه با استفاده از تحلیل پویای تولید گفتار در RPS، روشی برای رده‌بندی واژه‌ها بر مبنای فاصله اندازه‌گیری شده بین جاذب‌های^۷ هر واج در فضای فاز مطرح شد (Liu, 2003). همچنین لیندگرن با استفاده از ویژگی‌های استخراج شده از سیگنال گفتار جاسازی شده در فضای RPS، یک سامانه بازشناسی گفتار واج مجزاً بر روی مجموعه دادگان TIMIT معرفی کرد (Lindgren, 2003). مشابه این سامانه، از یک مجموعه مدل مخلوط گوسی^۸ (GMM) برای مدل‌سازی جداگانه توزیع هر واج استفاده شده است (Povinelli, 2004&2006). در (Indrebo, 2006) یک روش جدید ترکیبی با استفاده از فضای بازسازی شده فاز اعمال شده بر روی زیرباندهای به دست آمده از تجزیه فیلتر بانکی سیگنال، برای طبقه‌بندی سیگنال گفتار ارائه شد. از طرف دیگر در (Narayanan 2012; Thasleema 2012) نیز طبقه‌بندی واحدهای CV گفتاری با توسعه روشی مبتنی بر هیستوگرام (SSPD) و طبقه‌بندی کننده‌های شبکه عصبی مصنوعی (ANN) و ماشین بردار پشتیبان (SVM) پیشنهاد شده است. در (جعفری، ۱۳۸۹) و (Jafari, 2010; Jafari, 2012) چند نوع بردار ویژگی ترکیبی پیشنهاد شده است که عناصر آن مخلوطی از ویژگی‌های متداول MFCC و ویژگی‌های خاص استخراج شده از RPS هستند. در این روش، ویژگی‌های مبتنی بر RPS از پارامترهای GMM یادگرفته شده از توزیع سیگنال جاسازی شده در RPS و یا قطع پوانکاره^۹ آن به دست آمده‌اند. سپس با انتخاب برخی پارامترهای مدل GMM تعلیم یافته، بردار ویژگی نهایی تولید شده است. متأسفانه هزینه محاسباتی بالا در پیاده‌سازی این روش‌ها یکی از نقاط ضعف اساسی آنها به شمار می‌رود.

در همین اواخر در (Shekofteh, 2013) روشی پیشنهاد شد تا ویژگی‌های مناسبی از RPS استخراج شود که علاوه بر هزینه محاسباتی و کارایی مناسب آن، قابلیت

1 Dynamic
2 Correlation Dimension
3 Fractal Dimension
4 Generalized Fractal Dimension
5 Lyapunov Exponents

6 Isolated Phoneme
7 Attractor
8 Gaussian Mixture Model
9 Poincare Section

بعلاوه در این مقاله عمل کرد روش پیشنهادی در حالت استفاده از اطلاعات تمامی مانیفولد‌های گفتاری واج بررسی شده است. در نهایت بوسیله سامانه بازشناس مبتنی بر شبکه عصبی و تخمین خروجی‌های احتمال پسین آن، روش‌های متنوعی از ترکیب خروجی سامانه‌ها و ساختار TANDEM در جهت بهبود کارایی ASR بررسی شده است.

ادامه مقاله به این صورت شکل یافته است: در بخش دوم تئوری جاسازی سیگنال در RPS معرفی خواهد شد. بخش سوم مدل سازی مانیفولد‌های گفتاری را شرح می‌دهد. در بخش چهارم روش استخراج ویژگی پیشنهادی آورده شده است. بخش پنجم و ششم به ترتیب روش بازشناسی و مجموعه دادگان مورد استفاده را معرفی می‌کنند. در بخش هفتم نتایج آزمایش‌های مربوط به سامانه بازشناس پایه ارائه شده است. بخش‌های هشتم و نهم به ترتیب شامل معرفی روش‌های ترکیب و ارائه نتایج آزمایش‌ها است. در انتها نیز نتیجه‌گیری مقاله آورده شده است.

۲- جاسازی سیگنال در فضای بازسازی شده فاز

جاسازی سیگنال تک‌متغیره زمانی به وسیله روش محورهای تأخیری^۶ بر مبنای تئوری‌های مطرح شده در (Takens, 1980) و (Sauer, 1991) است. در این روش، پس از تعیین پارامترهای مناسب تأخیر زمانی و بُعد جاسازی، نمونه‌های سیگنال زمانی اولیه را می‌توان به گونه‌ای به فضای چند بُعدی RPS انتقال داد که مشخصه پویای واقعی سامانه تولید کننده سیگنال، به صورت مناسب در آن فضا نمایش داده شود (Kantz, 1997). در ادامه، روش جاسازی سیگنال یک بُعدی در فضای با بُعد بالاتر بیان می‌شود:

فرض کنیم که سری زمانی $s = \{s[1], \dots, s[N]\}$ (نمونه‌های سیگنال زمانی یک قاب گفتاری) با تعداد N نمونه، سیگنال زمانی تک‌متغیره مورد بحث باشد. متناظر با این سری زمانی، یک مجموعه متوالی از نقاط جاسازی شده در فضای چندبُعدی RPS تعریف می‌شود. به عنوان مثال می‌توان $S[i]$ را به عنوان نمونه i ام جاسازی شده در فضای RPS با رابطه (۱) به صورت یک بردار d بُعدی نشان داد:

$$S[i] = [s[i], s[i + \tau], \dots, s[i + (d - 1)\tau]], \quad (1)$$

⁶ Combination

⁷ Delay-Coordinat

استفاده در سامانه‌های بازشناسی گفتار پیوسته را داشته باشند. این روش بر مبنای جاسازی سیگنال زمانی گفتار در فضای RPS (که علاوه بر اطلاعات دامنه و فاز سیگنال حاوی اطلاعات مانیفولد‌های گفتاری است) بود که در آن بایستی یک سری مدل‌های از پیش تعلیم یافته گفتاری مرجع آماده شوند. این مدل‌ها که می‌توان آنها را "مانیفولد" یا "جاذب" گفتاری تلقی کرد، وظیفه مدل سازی توزیع نمونه‌های هر یک از واحدهای گفتاری واجی را در RPS برعهده دارند. سپس برای هر قاب گفتاری، ویژگی‌های PPRPS، بر مبنای رابطه بیز و محاسبه مقدار شباهت (درست‌نمایی شرطی^۱) نمونه‌های قاب گفتاری جاسازی شده در RPS با هر یک از مدل‌های مرجع محاسبه شده است. بنابراین بُعد بردار ویژگی PPRPS برابر با تعداد کل مدل‌های مرجع گفتاری است. در (Shekofteh, 2013) نشان داده شد که استفاده از این ایده منجر به بهبود عمل کرد سامانه بازشناسی واج مجزا در مقایسه با روش ارائه شده در (Povinelli, 2006) می‌شود. همچنین کاربرد این روش برای استفاده در سامانه بازشناسی گفتار پیوسته با مدل مخفی مارکوف^۲ (HMM) و با مجموعه مجموعه محدودی از مانیفولد‌های انتخاب شده واجی بررسی شده است. برای انجام این کار نگاشت خطی آنالیز متمایزگر (LDA) بر روی تعدادی ویژگی اعمال شده است.

در (Furui, 1986) نشان داده شده است که با توجه به ماهیت پویای گفتار، استخراج اطلاعات دینامیک از توالی بردارهای ویژگی گفتاری می‌تواند بهبود چشم‌گیری در کارایی سامانه‌های بازشناسی گفتار داشته باشد. متأسفانه یکی از مشکلات روش پیشنهادی در (Shekofteh, 2013) عدم تولید ویژگی‌های مناسب پویا در مقایسه با ویژگی‌های پویای به دست آمده از روش‌های متداولی مانند MFCC (ضرایب دلتا و دلتا دلتا) است. از این رو در این مقاله برای استخراج مناسب تر اطلاعات پویا، استفاده از ساختار شبکه عصبی با ورودی قطعه‌ای^۳ پیشنهاد شده است که شامل توالی چندین بردار ویژگی از قاب‌های گفتاری است. همچنین ساختار شبکه عصبی مورد استفاده منجر به یادگیری الگوهای زمانی^۴ و ترکیب غیرخطی تمامی ویژگی‌های ورودی شده و تخمین مناسبی از احتمال پسین طبقه‌های واج^۵ گفتاری در خروجی خود تولید خواهد کرد.

¹ Conditional Likelihood

² Hidden Markov Model

³ Segmental

⁴ Temporal Pattern

⁵ Phoneme Classes

کلاس‌های واجی رفتاری بین موارد اشاره شده امکان‌پذیر است (Banbrook, 1996; Indrebo, 2006).

۳- مدل سازی مانیفولدهای گفتاری در فضای بازسازی شده فاز

در (Povinelli, 2006) یک سیستم بازشناسی واج مجزاً^۳ براساس مدل سازی آماری توزیع نقاط تراژکتوری سیگنال گفتار در RPS توسط پایونلی پیشنهاد شده است. توالی زمانی تراژکتوری یک سیگنال در حوزه RPS منجر به تشکیل "مانیفولد گفتاری" می‌شود.

سیستم بازشناسی واج مجزای پیشنهاد شده توسط پایونلی شامل دو مرحله تعلیم و آزمون بود. در مرحله تعلیم، برای هر یک از واج‌های گفتاری در RPS، یک توزیع چندمتغیره پارامتری با استفاده از GMM آموزش داده می‌شد. برای این منظور در ابتدا متناظر با هر واج، کلیه سیگنال‌های مربوط به آن واج (از مجموعه داده‌های تعلیم) در فضای RPS جاسازی شده و سپس برای توزیع نقاط تشکیل شده در RPS یک مدل کلی GMM تعلیم داده می‌شد. استفاده از GMM باعث می‌شد که ساختار هندسی مانیفولد شکل یافته در RPS به صورت نرم مدل شود به عنوان مثال تصویری از مدل سازی GMM مانیفولد واج واکدار /u/ در RPS دو بعدی در شکل (۲) نشان داده شده است.

پس از تعلیم مدل‌های GMM که در بردارنده اطلاعات هندسی مانیفولد هر واج در RPS است؛ پایونلی از معیار بیشینه درست نمایی^۴ برای طبقه بندی در مرحله آزمون استفاده کرده است که در رابطه (۳) آورده شده است:

$$\hat{c} = \arg \max_{i=1, \dots, K} p(\mathbf{X} | c_i),$$

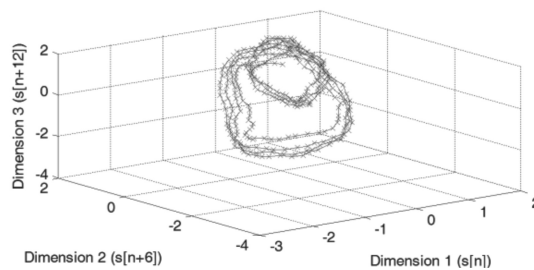
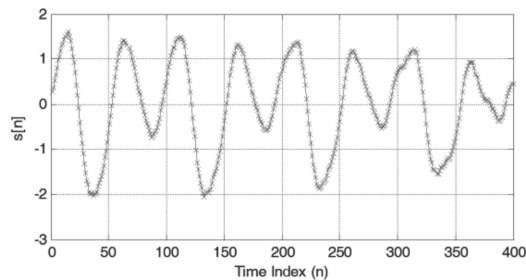
$$p(\mathbf{X} | c_i) = \frac{1}{L} \sum_{n=1}^L \log p(\mathbf{x}_n | c_i), \quad (3)$$

که در آن K تعداد کل مدل‌های واجی GMM (که هر کدام متناظر با یک کلاس واجی c_i هستند) و L تعداد نقاط جاسازی شده در RPS است. ماتریس X نیز شامل تعداد L نمونه جاسازی شده \mathbf{x}_n در فضای RPS است که معمولاً حاوی اطلاعات ماتریس تراژکتوری S و مشتقات زمانی آن، ΔS می‌باشد (Lindgren, 2003). همچنین مقدار درست

به طوری که پارامتر d، بُعد جاسازی و پارامتر τ معادل با زمان تأخیر است. واضح است که برای یک قاب با طول N، تعداد نقاط جاسازی شده در حوزه RPS به صورت زیر است:

$$L = N - (d - 1)\tau, \quad (2)$$

روش‌های مختلفی برای انتخاب مقدار بهینه بُعد جاسازی و تأخیر زمانی وجود دارند (Kennel 1992; Abarbanel, 1996; Povinelli, 2004; Johnson, 2005). یک انتخاب مناسب برای سیگنال گفتار میکروفونی با نرخ نمونه برداری شانزده کیلوهرتز، مقادیر $d=8$ و $\tau=6$ است (Jafari, 2010). در شکل (۱) نمونه‌ای از یک قاب سیگنال گفتار (شامل واج /u/) و همچنین تراژکتوری متناظر با آن را که در فضای سه بعدی RPS جاسازی شده‌اند، نشان می‌دهد.



شکل (۱) - واکه /u/ در نمایش معمولی تک بعدی (شکل بالا) و نمایش تراژکتوری سه بعدی آن (شکل پایین) با روش جاسازی در RPS.

رفتار تراژکتوری سیگنال گفتار در فضای RPS برای اکثر واج‌های واکدار مشابه با فرآیند قبض^۱ در سیگنال‌های آشوبی است. به عنوان مثال در شکل (۱) نمونه‌ای از این رفتار برای تراژکتوری واج واکدار /u/ نشان داده شده است. نتیجه این رفتار تولید مانیفولدهایی با الگوهای خاص و مشخص در RPS خواهد بود. از طرف دیگر برای واج‌های گفتاری انفجاری (مانند واج های /b/ و /t/) رفتار بسط^۲ سیگنال آشوبی در فضای فاز قابل مشاهده است. برای سایر

³ Isolated Phoneme Recognition

⁴ Max-Likelihood

¹ Folding Or Squeezing

² Stretching

۴- محاسبه تخمین احتمال پسین واجی در بازشناسی گفتار به وسیله ویژگی‌های مبتنی بر RPS

در این مقاله، مبنای روش استخراج ویژگی مبتنی بر RPS، براساس ایده‌ای در توسعه روش پاونیلی است. متأسفانه روش پاونیلی به‌انحصار قابلیت استفاده در کاربرد بازشناسی واج مجزا دارد. اما با روش توسعه‌ای مطرح شده در (Shekofteh, 2013)، توانایی استفاده از این روش در کاربرد بازشناسی گفتار پیوسته حاصل شده است.

همان‌طور که در رابطه (۳) در بخش (۳) دیدیم، در مرحله آزمون سامانه پاونیلی، به تعداد کلاس‌های واجی موجود، امتیاز آکوستیکی درست‌نمایی تولید خواهد شد. سپس مدل واجی که بالاترین امتیاز درست‌نمایی را به دست آورده باشد، به‌عنوان کلاس واج برنده انتخاب می‌شود. بنابراین در بخش تصمیم‌گیری روش پاونیلی، از اطلاعات فضای RPS هر مدل واجی به‌طور مستقل در تعیین طبقه واجی برنده استفاده می‌شود.

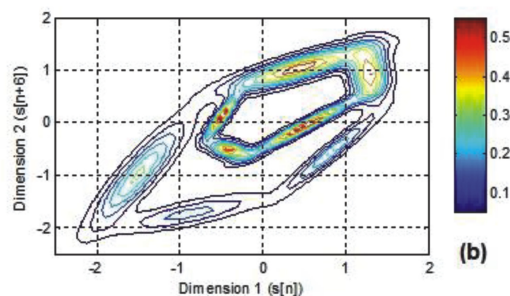
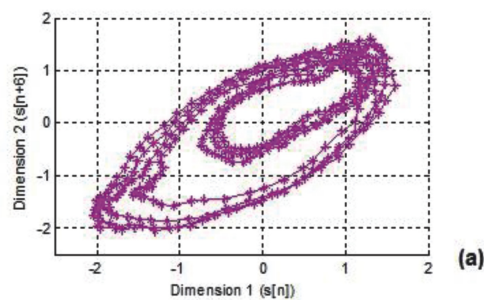
ایده مطرح‌شده در توسعه این روش، پیشنهاد استفاده هم‌زمان از زیرمجموعه‌ای از امتیازهای درست‌نمایی به دست آمده از مانیفولدهای واجی، به‌عنوان یک بردار ویژگی مناسب است. با این روش می‌توان به‌راحتی از اطلاعات مانیفولدهای گفتاری در RPS در پیاده‌سازی سامانه‌های بازشناسی گفتار پیوسته استفاده کرد.

از آنجا که هر یک از مؤلفه‌های بردار ویژگی به‌دست آمده از RPS دارای اطلاعات هم‌بسته با یکدیگر هستند، انتخاب یا ترکیب مناسب این ویژگی‌ها برای تولید اطلاعات غنی‌تر مفید به‌نظر می‌رسد. یک روش مناسب برای دستیابی به این هدف، تبدیل بردارهای ویژگی اولیه به تخمین مناسبی از احتمالات پسین واجی است که در همین‌اواخر به‌طور گسترده در کاربرد بازشناسی گفتار مطرح‌شده است (Park, 2011; Pinto, 2011; Ikbali, 2012). در اینجا استفاده از ساختار یک شبکه عصبی که قابلیت یادگیری الگوهای زمانی را دارد، برای تخمین احتمالات پسین پیشنهاد شده است. در شکل (۳) روندنمای روش پیشنهادی آورده شده است:

نمایی $p(\mathbf{x}_n | c_i)$ برای تعداد M گوسین مدل مخلوط گوسی (GMM) از رابطه (۴) محاسبه می‌شود:

$$p(\mathbf{x}_n | c_i) = \sum_{m=1}^M \omega_m^i N(\mathbf{x}_n; \mu_m^i, \Sigma_m^i), \quad (4)$$

که در آن $N(\mathbf{x}_n; \mu_m^i, \Sigma_m^i)$ نشان‌دهنده درست‌نمایی مؤلفه m ام توزیع گوسی با وزن ω_m^i از کلاس واج i ام (c_i) است.



(شکل ۲) - (a) توزیع تراژکتوری مانیفولد واج /u/ نشان‌داده شده در شکل (۱) در فضای RPS دوبعدی. (b) مدل‌سازی تابع چگالی احتمال آن به‌وسیله یک مدل GMM با تعداد هشت مؤلفه گوسی.

جانسن در (Johnson, 2005) نشان داد که قدرت تمایزگری بین واج‌ها که از فضای RPS به‌دست می‌آید، کمتر از ویژگی‌های متداول کپستروم MFCC است؛ اما با استفاده از روش ترکیب جمع وزن‌دار در سطح خروجی رده‌بندی‌کننده‌ها (ادغام امتیاز درست‌نمایی به‌دست آمده از ویژگی‌های حوزه RPS با امتیاز درست‌نمایی حاصل شده از ویژگی‌های MFCC) کارایی سامانه بازشناسی واج مجزا بر روی مجموعه داده TIMIT افزایش خواهد یافت. این افزایش کارایی مؤید هم‌افزا بودن اطلاعاتی است که توسط روش استخراج ویژگی MFCC و ویژگی‌های حوزه RPS به‌دست آمده‌اند.

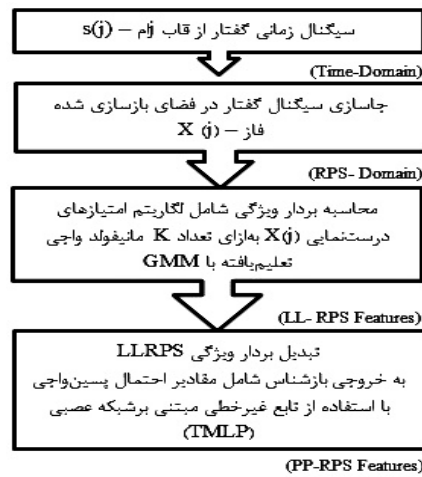
(حجم محاسبات) و دقت سامانه بازشناسی خواهد بود. پس از محاسبه امتیازات درست‌نمایی، بردار ویژگی اولیه $LLRPS(j)$ برای قاب گفتاری j به صورت زیر حاصل می‌شود:

$$LLRPS(j) = \begin{bmatrix} p(\mathbf{X}(j) | c_1) \\ p(\mathbf{X}(j) | c_2) \\ \vdots \\ p(\mathbf{X}(j) | c_K) \end{bmatrix} \quad (5)$$

از آنجا که ماهیت ویژگی‌های تولیدشده مبتنی بر لگاریتم درست‌نمایی^۲ (LL) است، به اختصار آن را LLRPS می‌نامیم. مرحله سوم: تخمین مقادیر احتمال پسین واجی با استفاده از ساختار شبکه عصبی؛

در این مرحله از یک ساختار شبکه عصبی جلوسو برای ترکیب غیرخطی تمامی ویژگی‌های LLRPS و تخمین مقادیر احتمال پسین^۳ واجی استفاده شده است. همچنین ساختار شبکه مورد استفاده، به گونه‌ای طراحی می‌شود که ورودی آن قطعه‌ای و شامل چند بردار ویژگی متوالی باشد تا بتواند به طور هم‌زمان اطلاعات پویای گفتاری و الگوهای زمانی مربوط به آن را از روی توالی بردارهای ویژگی ایستا^۴ متوالی استخراج کند.

نحوه تعلیم شبکه عصبی هم به گونه‌ای است که در خروجی آن تخمینی از مقدار احتمال پسین (PP) طبقه واجی متناظر با بردار ویژگی ورودی (LLRPS) از قاب گفتاری میانی قطعه ورودی به شبکه عصبی را تولید کند. از این رو بردار تولیدشده در خروجی شبکه را PPRPS می‌نامیم. در این مقاله برای به دست آوردن تخمین احتمال پسین، از مدل شبکه عصبی^۵ TMLP استفاده شده است که ساختار آن الهام گرفته شده از خاصیت تونوتوپیک دستگاه شنوایی انسان است (Chen, 2005). ساختار این شبکه باعث می‌شود که توالی زمانی هر یک از ویژگی‌های ورودی در لایه‌های پایینی شبکه به طور مستقل پردازش شده و سپس اطلاعات به دست آمده از آنها در لایه‌های بالاتر آن ترکیب شوند. این ساختار شبکه عصبی به یادگیری بهتر الگوهای زمانی از هر یک از ویژگی‌های ورودی به شبکه کمک می‌کند. در (شکفته، ۱۳۸۶) نشان داده شد که قدرت یادگیری این مدل تمایزی^۶ از شبکه‌های عصبی متداول MLP و یا TDNN



(شکل ۳) - روندنمای روش پیشنهادی

روندنمای شکل (۳) شامل مراحل زیر است:

مرحله اول: جاسازی سیگنال زمانی گفتار در RPS؛

در اینجا برای فرآیند جاسازی سیگنال در فضای RPS، از مقادیر $d=8$ و $\tau=6$ استفاده شده است. به طور متداول نمونه‌های سیگنال قبل از جاسازی در RPS، بهنجار می‌شوند (Povinelli, 2004). اعمال روش‌های هنجارسازی باعث می‌شود که اثر شدت یا همان اندازه دامنه سیگنال گفتار ضابط شده کم‌رنگ شده و تنها شکل هندسه سیگنال مورد بررسی قرار گیرد. از طرفی دیگر، اعمال روش‌های هنجارسازی منجر به قرارگیری مانیفولدهای گفتاری در محدوده نسبتاً مشخصی از RPS خواهند شد که باعث مقاوم‌سازی^۱ ویژگی‌های استخراج شده از آن می‌شود. در این مقاله نمونه‌های هر قاب گفتاری به گونه‌ای بهنجار می‌شوند که مقدار میانگین نمونه‌ها صفر و انحراف معیار آن‌ها برابر با مقدار واحد شود.

مرحله دوم: محاسبه بردار ویژگی اولیه شامل امتیاز درست‌نمایی متناظر با هر مدل مانیفولد واجی مبتنی بر GMM؛

در این مرحله، از رابطه (۳) برای محاسبه امتیازات درست‌نمایی از هر مدل مانیفولد واجی برای سیگنال قاب j گفتاری جاسازی شده در RPS، $\mathbf{X}(j)$ استفاده خواهد شد. بنابراین بایستی مشابه با روش پاوینلی، در یک مرحله جداگانه تعلیم، به تعداد واج‌های مجموعه دادگان ($K=30$) مدل GMM از روی داده‌های تعلیم آموزش داده شود. در (Shekofteh, 2013) نشان دادیم که استفاده از GMM با ۴ مؤلفه گوسی ($M=4$) مصالحه مناسبی بین زمان پردازش

² Log Likelihood

³ Posterior Probability

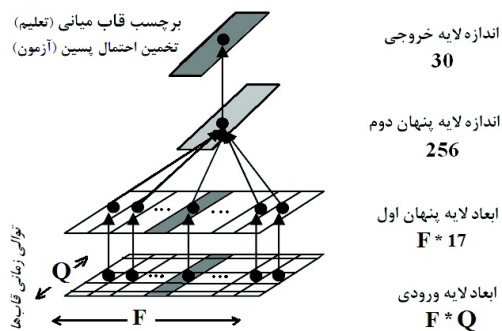
⁴ Static

⁵ Tonotopic Multi Layer Perceptron

⁶ Discriminative

¹ Robustness

(رحیمی نژاد، ۱۳۸۲) در کاربرد بازشناسی گفتار بیشتر است. بنابراین در این مقاله از ساختار شبکه عصبی TMLP استفاده شده است که در شکل (۴) نمایش داده شده است.



(شکل ۴) - ساختار شبکه عصبی TMLP با توالی $Q=21$ قاب گفتاری و بردار ویژگی F بعدی ورودی برای هر قاب گفتاری.

در ساختار شبکه عصبی TMLP مورد استفاده، یک قطعه از بردارهای ویژگی متوالی (در اینجا تعداد Q قاب متوالی گفتاری) به ورودی شبکه اعمال خواهند شد تا طبقه واجی متناظر با قاب میانی ورودی آن به وسیله یک برچسب باینری در خروجی شبکه عصبی پیش‌بینی شود. همچنین هر یک از بردارهای ویژگی ورودی دارای بُعد F (در اینجا $F=K=30$ و معادل با تعداد واج‌های دادگان می‌باشد) است. ساختار شبکه عصبی TMLP، برای توالی $Q=21$ قاب از بردارهای ویژگی ورودی LLRPS (با بُعد ۳۰) به صورت زیر در نظر گرفته شده است:

$$30 \cdot (21 \cdot 17) - 30 \cdot (17 \cdot 256) - 256 \cdot 30$$

که در برگزیده اطلاعات پویا از مجموعه ۲۱ قاب گفتاری متوالی (بیش از ۲۰۰ میلی‌ثانیه) است. انتخاب اندازه لایه‌های پنهان شبکه براساس بیشینه‌سازی دقت بازشناسی بر روی داده‌های اعتبارسنجی^۱ و رعایت نسبت تعداد پارامترهای شبکه به تعداد پارامترهای داده تعلیم (برای جلوگیری از برازش بیش از حد^۲ شبکه) انجام شده است.

در اینجا برای تعلیم شبکه TMLP، از برچسب‌دهی گذشته باینری نوع سخت (One-Hot) استفاده شده است (رحیمی نژاد، ۱۳۸۲؛ ولی، ۱۳۸۵). در این نوع برچسب‌دهی، یک خروجی ۳۰ نرونی (به تعداد طبقه‌های واجی) به عنوان خروجی مطلوب شبکه تعریف می‌شود که یک نرون آن حاوی مقدار یک (متناظر با برچسب طبقه واجی مربوط به قاب میانی مجموعه بردار ورودی) و بقیه نرون‌های آن مقدار

صفر دارند (Tohidypour, 2012). اگر الگوریتم تعلیم شبکه براساس کمینه‌سازی میانگین مجذور خطا^۳ (MSE) باشد، تخمینی از مقدار احتمالاتی پسین طبقه‌های واجی در خروجی شبکه به شرط قطعه قاب‌های گفتاری ورودی اعمال شده به آن تولید خواهد شد (White, 1989; Zavaliagos, 1994).

۵- روش بازشناسی گفتار

با توجه به استفاده روش پیشنهادی از مدل شبکه عصبی به عنوان یک تخمین‌زننده مقدار احتمال پسین واجی، در این حالت سامانه بازشناسی گفتار می‌تواند در چارچوب‌های متنوعی پیاده‌سازی شود که در ادامه به شرح هر یک از آنها خواهیم پرداخت:

۵-۱- بازشناسی گفتار به وسیله مقدار احتمالات پسین واجی تخمین زده شده در خروجی مدل شبکه عصبی

در این حالت ارزیابی عمل کرد سامانه بازشناسی گفتار با دو معیار دقت بازشناسی در سطوح آوایی قاب (frame) و واج (phoneme) انجام خواهد شد. دقت بازشناسی قاب، به طور مستقیم از خروجی احتمال پسین به دست آمده از هر قاب گفتاری قابل محاسبه است؛ به این صورت که طبقه واجی متناظر با نرونی از خروجی شبکه که بیشترین مقدار احتمال پسین را داشته باشد، به عنوان برچسب واج بازشناسی شده برای قاب متناظر با آن در ورودی شبکه در نظر گرفته می‌شود. از این رو معیار درصد دقت بازشناسی قاب به صورت نسبت تعداد قاب‌های گفتاری درست تشخیص داده شده به تعداد کل قاب‌های گفتاری مورد استفاده تعریف می‌شود (Tohidypour, 2012).

برای محاسبه دقت بازشناسی در سطح آوایی واج، از اعمال الگوریتم ویتربی بر روی بردارهای متوالی خروجی شبکه عصبی که حاوی تخمینی از مقدار احتمال پسین هستند مشابه با (Pinto, 2011) استفاده خواهد شد. در این روش فرض می‌کنیم که هر واج با یک مدل چندحالتی HMM^۴ مدل شده است که مقدار درست‌نمایی انتشاری^۴ تمامی حالات^۵ آن با یکدیگر یکسان و برابر با مقدار خروجی مقیاس^۶ شده، احتمال پسین به دست آمده از نرون خروجی

³ Mean Square Error

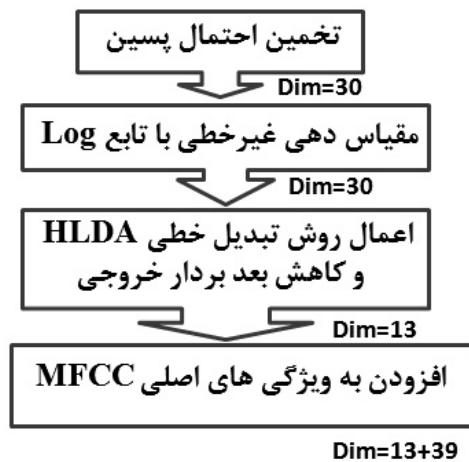
⁴ Emission Likelihood

⁵ States

⁶ Sclaed

¹ Validation

² Over-Training



(شکل ۵) - بلوک دیاگرام ساختار TANDEM در تولید بردار ویژگی گفتاری مورد استفاده در مدل سازی صوتی GMM/HMM.

در اینجا پیاده سازی مدل سازی صوتی GMM/HMM توسط جعبه ابزار HTK انجام شده است (HTK, V.3.4). همچنین از واحدهای آوایی سه واجی^۴ گره زده شده برای مدل سازی صوتی استفاده شده است (Sameti, 2011). در این حالت ارزیابی عمل کرد سامانه در سطح آوایی واج و با معیار NIST (رابطه ۶) انجام خواهد شد.

۶- مجموعه دادگان

دادگان گفتاری مورد استفاده در این مقاله از مجموعه "دادگان فارس دات کوچک میکروفونی" با حجم زمانی حدود شش ساعت است که نرخ نمونه برداری آن به 16kHz تقلیل یافته است (FARSDAT). این دادگان شامل جملات پیوسته بیان شده از ۳۰۴ نفر گوینده زن و مرد است که هر نفر تعداد بیست جمله را به زبان فارسی بیان کرده اند. جملات دادگان شامل برچسب دهی آوایی^۵ با ۴۳ برچسب مختلف (بدون سکوت) است. در کار حاضر با تجمیع تعدادی از برچسبها با یکدیگر، از تعداد سی برچسب واجی (۲۹) واج متداول فارسی به همراه یک برچسب سکوت (K=30) استفاده شده است (Bijankhan, 2003). از داده های مربوط به ۲۲۰ گوینده اول در بخش تعلیم مدل بازناس، داده های ۵۴ گوینده آخر در بخش آزمون و برای بخش توسعه (Development) و اعتبارسنجی مدل شبکه عصبی و الگوریتم ویتربی نیز از داده های سی گوینده باقیمانده فارس دات استفاده شده است.

شبکه عصبی متناظر با آن واج است. از این رو می توان این روش را مشابه با پیاده سازی روش های پدید HMM/ANN تلقی کرد (Bourlard, 1994).

در خروجی الگوریتم ویتربی، دنباله ای از چند واج متوالی تولید خواهد شد که بالاترین احتمال وقوع را دارند. در اینجا برای گذر بین واج های مختلف، مقدار احتمال یکسان در نظر گرفته شده است که معادل با عدم استفاده از مقدار احتمالات مدل زبانی بایگرم^۱ واجی در سامانه بازناسی گفتار پیوسته است. همچنین مقدار بهینه جریمه درج^۲ واج در الگوریتم ویتربی، با داده های توسعه^۳ تعیین شده است. وظیفه جریمه درج، جلوگیری از پرش سریع خروجی الگوریتم بین واج های مختلف است که منجر به کاهش تعداد درج واج و در نتیجه افزایش مقدار دقت بازناسی خواهد شد. این فرآیند با اضافه کردن یک جریمه (در اینجا یک مقدار منفی) به مقدار لگاریتم درست نمایی انباشته شده واج هایی غیر از واج شناسایی شده فعلی انجام می شود. با استفاده از معیار NIST، درصد دقت بازناسی واج نهایی از رابطه زیر محاسبه می شود:

$$\%Phoneme Accuracy = \frac{A-D-I-S}{A} * 100 \quad (6)$$

که در آن A تعداد واج های موجود در برچسب داده مرجع، D تعداد واج های حذف شده، I تعداد واج های درج شده و S تعداد واج های جانشین شده است.

۵-۲- بازناسی گفتار با مدل مخفی مارکوف (HMM) در ساختار TANDEM

در روش TANDEM، از مقادیر تخمین احتمال پسین واجی که در خروجی شبکه عصبی تولید شده است، به عنوان بردار ویژگی خام و اولیه برای تولید ویژگی گفتاری مناسب در مدل سازی صوتی مبتنی بر GMM/HMM استفاده می شود (Hermansky, 2000; Zhu, 2004; Ikbal, 2012). در شکل (۵) ساختار TANDEM مورد استفاده در این مقاله برای تولید بردار ویژگی گفتاری پس پردازش شده نشان داده شده است.

¹ Bigram
² Insertion Penalty
³ Development

⁴ Triphone
⁵ Phonetic Labeling

۷- آزمایش‌های سامانه‌های پایه

در این بخش به بررسی نتایج بازشناسی به دست آمده از اعمال بردار ویژگی LLRPS به مدل بازشناسی شبکه عصبی خواهیم پرداخت. در جدول (۱) نتایج درصد دقت بازشناسی قاب و واج آن آورده شده است.

(جدول ۱) - مقایسه نتایج آزمون درصد دقت بازشناسی قاب و واج برای بردار ویژگی LLRPS، ویژگی متداول MFCC و روش TANDEM به وسیله مدل بازشناسی شبکه عصبی TMLP و یا مدل متداول HMM.

	ساختار مدل بازشناسی	بردار ویژگی ورودی	بُعد ویژگی	% دقت قاب	% دقت واج
۱	TMLP	LLRPS	۳۰	۷۴/۲۷	۶۲/۶۰
۲	TMLP	MFCC13	۱۳	۷۹/۲۳	۶۹/۶۷
۳	TMLP	MFCC39	۳۹	۸۰/۹۸	۷۲/۲۵
۴	HMM	MFCC39	۳۹	-	۷۵/۲۲
۵	HMM	TANDEM (PPRPS13)	۱۳	-	۶۰/۰۵
۶	HMM	TANDEM (PPRPS13 +MFCC39)	۵۲	-	۷۵/۸۴

همچنین در جدول (۱)، نتایج سامانه بازشناسی پایه، شامل ویژگی‌های MFCC (به‌عنوان یک روش استخراج ویژگی پرکاربرد گفتاری) آورده شده است. برای محاسبه این نتایج، بردار ویژگی ایستا MFCC، شامل سیزده ضریب کپستروم (MFCC13) و یا بردار ویژگی ایستا-پویا MFCC39، شامل سیزده ضریب کپستروم به همراه مشتقات زمانی اول و دوم آن (ضرایب دلتا و دلتادلتا) به ورودی شبکه عصبی TMLP اعمال شده است. سپس مشابه با روش الگوریتم ویتربی مطرح شده در بخش ۴، نتایج بازشناسی آن به وسیله مقدار احتمال پسین به دست آمده از خروجی شبکه (که در اینجا آنها را PPMFCC می‌نامیم) به دست آمده است.

در ردیف چهارم جدول (۱)، نتایج دقت بازشناسی واج به دست آمده از مدل‌سازی صوتی HMM با بردار ویژگی‌های MFCC39 (در حالت عدم استفاده از مدل زبانی بایگرم واجی در کدگشای آن) آورده شده است که به‌عنوان یک سامانه بازشناسی پایه، متداول و پرکاربرد مطرح است. همچنین در دو سطر آخر جدول ۱، نتایج دقت بازشناسی واج مدل‌سازی صوتی HMM با ویژگی‌های به دست آمده از روش TANDEM، که به ترتیب در حالات اضافه نشده یا

افزوده شده به ویژگی متداول MFCC39 هستند، آورده شده است.

همان‌طور که از نتایج جدول (۱) مشاهده می‌شود، کارایی مدل بازشناسی HMM (با ویژگی‌های متداول MFCC) از مدل بازشناسی مبتنی بر شبکه عصبی TMLP (با ویژگی‌های متداول MFCC) حدود سه درصد در دقت بازشناسی واج بهتر است. همچنین در بین نتایج به دست آمده از مدل بازشناسی شبکه عصبی، بردار ویژگی MFCC39 که حاوی اطلاعات ایستا و پویای گفتاری است، بالاترین درصد بازشناسی قاب و واج را به دست آورده است؛ این در حالی است که در این آزمایش بردار ویژگی‌های LLRPS و MFCC13 به‌انحصار شامل اطلاعات ایستای هر قاب گفتاری هستند. البته لازم به ذکر است که ساختار شبکه عصبی TMLP که در آن از مجموعه چندین قاب متوالی در ورودی خود استفاده می‌کند، قابلیت استخراج اطلاعات پویا را دارد، با این وجود نتایج به دست آمده از جدول (۱) نشان می‌دهد که استفاده مجزا از اطلاعات دینامیک در بردار ویژگی MFCC39، به ترتیب منجر به بهبود ۱/۷۵٪ و ۲/۵۸٪ در نتایج دقت بازشناسی قاب و واج نسبت به روش MFCC13 شده است. همچنین مدل بازشناسی TMLP با بردار ویژگی ورودی LLRPS دقت بازشناسی واج ۶۲/۶۰٪ را کسب کرده است که بیان‌گر کارایی پایین‌تر آن نسبت به روش متداول MFCC است.

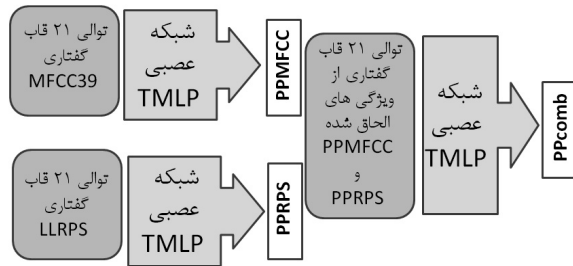
از طرف دیگر نتیجه سطر آخر جدول (۱) نشان می‌دهد که مدل‌سازی صوتی HMM با ویژگی‌های مبتنی بر RPS (PPRPS) محاسبه شده از روش TANDEM (که همان مقادیر احتمالاتی پسین پس‌پردازش شده از خروجی شبکه عصبی TMLP هستند)، توانسته است در حالت افزوده شدن به ویژگی‌های متداول MFCC، دقت بازشناسی واج ۷۵/۸۴٪ را کسب کند که بهبود کارایی حدود ۰/۱۶٪ نسبت به سامانه پایه شامل روش مدل‌سازی HMM و بردار ویژگی MFCC نشان می‌دهد.

۸- روش‌های ترکیب خروجی سامانه‌های بازشناسی

استفاده هم‌زمان از بردارهای ویژگی متمایز که حاوی اطلاعات متفاوتی از یک سیگنال هستند، می‌تواند منجر به افزایش کارایی عمل‌کرد سامانه‌های بازشناسی شود (Kittler, 1998; Morgan, 2004; Nejadgholi, 2009). در این بخش

سال ۱۳۹۲ شماره ۱ پیاپی ۱۹

PPRPS) استفاده می‌شود. در شکل (۶) نمای کلی این روش ترکیب غیرخطی نشان داده شده است.



(شکل ۶) - روندنمای کلی روش ترکیب غیرخطی مقادیر احتمال پسین بوسیله ساختار شبکه عصبی TMLP.

پرواضح است که مشابه با مباحث مطرح شده در بخش ۴، خروجی این شبکه نیز تخمینی از مقدارهای احتمال پسین را تولید خواهد کرد. با توجه به شکل (۶)، این روش شامل دو طبقه شبکه عصبی است که در یک ساختار کلی سلسله‌مراتبی^۱ گنجانده شده‌اند. در این ساختار، شبکه‌های عصبی موجود در طبقه اول وظیفه استخراج تخمین اولیه مقادیر احتمال پسین (PPRPS و PPMFCC) را برعهده دارند و تنها شبکه عصبی موجود در طبقه دوم مسؤول ترکیب غیرخطی تخمین‌های اولیه و تولید مقدار احتمال پسین نهایی (PPComb) است.

۹- آزمایش‌ها بر روی روش‌های ترکیب و بحث و بررسی نتایج به‌دست آمده

این بخش مشتمل بر آزمایش‌های انجام گرفته برای بررسی عمل کرد هر یک از روش‌های ترکیب اشاره شده در بخش (۸) است. در نمودار شکل (۷)، درصد دقت بازشناسی قلابی برای روش‌های ترکیب جمع و ضرب وزن‌دار برحسب تغییرات وزن w نشان داده شده است. پرواضح است که برای مقدار $w=0$ ، خروجی ترکیب معادل خروجی روش MFCC39 و برای $w=1$ ، خروجی ترکیب معادل خروجی روش LLRPS است.

نشان خواهیم داد که چگونه با استفاده از روش‌های مناسب ترکیب، تخمین به‌دست آمده از مقادیر احتمالات پسین در خروجی شبکه‌های حاوی بردار ویژگی LLRPS و MFCC39 (که بالاترین نرخ بازشناسی را در برداشت) قابلیت بهبود یافتن دارد. شواهد تجربی نشان می‌دهند که اگر ماهیت اطلاعات کسب‌شده در روش‌های استخراج ویژگی با یکدیگر متفاوت باشند، کارایی سامانه‌های بازشناسی به‌وسیله روش‌های ترکیب مناسب بهبود خواهند یافت (Johnson, 2005; Valente, 2010).

اکنون فرض کنیم که بردارهای $PPRPS(i)$ و $PPMFCC(i)$ به‌ترتیب بردارهای شامل تخمین مقادیر احتمال پسین به‌دست آمده از خروجی شبکه‌های عصبی شامل ویژگی‌های LLRPS و MFCC39 برای قاب گفتاری i ام باشند، در این‌صورت خروجی ترکیب شده مقدارهای احتمال پسین، $PPComb(i)$ ، می‌تواند به‌وسیله یکی از هفت روش ترکیب معرفی شده در ذیل حاصل شود:

۱- روش جمع وزن‌دار:

$$PPComb(i) = w \times PPRPS(i) + (1-w) \times PPMFCC(i) \\ \text{where } 0 \leq w \leq 1$$

۲- روش ضرب وزن‌دار:

$$PPComb(i) = PPRPS(i)^w \times PPMFCC(i)^{(1-w)} \\ \text{where } 0 \leq w \leq 1$$

۳- روش میانگین: حالت خاص روش ۱ با مقدار $w=0.5$.

۴- روش ضرب: حالت خاص روش ۲ با مقدار $w=0.5$.

۵- روش جمع وزن‌دار متناسب با معکوس آنتروپی:

$$PPComb(i) = w_{pps}(i) \times PPRPS(i) + w_{mfcc}(i) \times PPMFCC(i) \\ \text{where } w_j(i) = \frac{h_j^{-1}(i)}{h_{pps}(i)^{-1} + h_{mfcc}(i)^{-1}}$$

که در آن $h(i)$ مقدار آنتروپی به‌دست آمده از بردار مقادیر احتمال پسین از قاب گفتاری i ام است (Misra, 2003; Kazemi, 2011).

۶- روش بیشینه (Max): در این روش برداری که حاوی بزرگ‌ترین مقدار احتمال پسین باشد، انتخاب خواهد شد.

۷- روش ترکیب غیرخطی با شبکه عصبی TMLP:

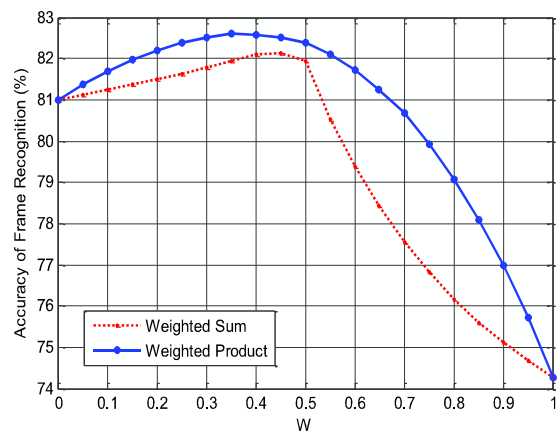
$$PPComb(i) = TMLP\{PPRPS(i), PPMFCC(i)\}$$

در این حالت، دوباره از یک ساختار شبکه عصبی TMLP برای ترکیب غیرخطی احتمالات پسین به‌دست آمده از خروجی دو شبکه MFCC39 (PPMFCC) و LLRPS

¹ Hierarchical

(جدول ۲) - مقایسه نتایج آزمون دقت بازشناسی قاب و واج برای روش‌های متنوع ترکیب بردارهای PPMFCC و PPRPS

روش ترکیب خروجی شبکه‌ها	% دقت قاب	% دقت واج	شماره
جمع وزن دار ($w=0.45$)	۸۲/۱۳	۷۳/۵۷	۱
ضرب وزن دار ($w=0.35$)	۸۲/۵۹	۷۴/۴۵	۲
میانگین ($w=0.5$)	۸۱/۹۵	۷۳/۲۶	۳
ضرب ($w=0.5$)	۸۲/۳۷	۷۴/۲۲	۴
جمع وزن دار متناسب با معکوس آنتروپی	۸۱/۷۸	۷۳/۴۲	۵
بیشینه (Max)	۸۱/۶۱	۷۳/۱۰	۶
غیرخطی $TMLP$	۸۴/۹۲	۷۶/۲۷	۷



(شکل ۷) - نمودار درصد دقت بازشناسی قاب بر حسب وزن W در حالت روش‌های ترکیب جمع (Sum) و ضرب (Product) وزن دار.

با توجه به نتایج ارائه شده در جدول (۲)، استفاده از روش‌های ترکیب، همواره منجر به بهبود نتایج دقت بازشناسی قاب و واج نسبت به نتایج به دست آمده از هر یک از سامانه‌ها در حالت بدون ترکیب جدول (۱) شده است. بیشترین بهبود کارایی با روش ترکیب غیرخطی $TMLP$ (سطر آخر جدول) به دست آمده است که به ترتیب با درصد دقت بازشناسی قاب و واج $۸۴/۹۲\%$ و $۷۶/۲۷\%$ ، توانسته است نسبت به روش پایه MFCC39 (با مدل بازشناس $TMLP$) به ترتیب بهبود $۳/۹۴\%$ و $۴/۰۲\%$ داشته باشد. همچنین با مقایسه درصد دقت بازشناسی واج این روش (ترکیب غیرخطی خروجی‌های PPMFCC و PPRPS با مدل شبکه عصبی $TMLP$) با نتیجه به دست آمده از سامانه پایه مبتنی بر HMM (ارائه شده در ردیف ۴ جدول (۱))، دیده می‌شود که روش ترکیب غیرخطی پیشنهادی (با مدل شبکه عصبی $TMLP$) توانسته است بهبود عمل کردی در حدود یک درصد در نتیجه دقت بازشناسی واج نشان دهد.

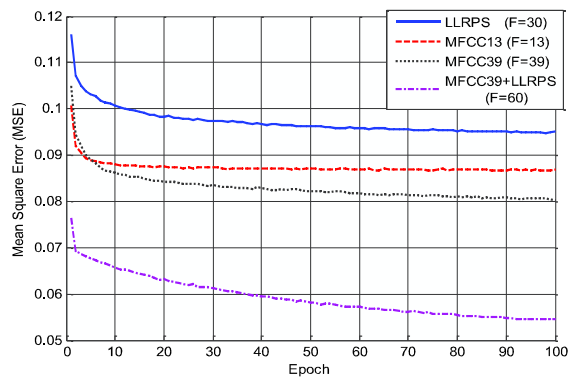
با توجه به نتایج جداول (۱) و (۲) موارد ذیل قابل نتیجه‌گیری است:

- در شرایط یکسان، عمل کرد بازشناسی گفتار با مدل HMM بهتر از مدل شبکه عصبی $TMLP$ است.
- ترکیب غیرخطی اطلاعات در سطح خروجی با روش شبکه عصبی $TMLP$ نتیجه بازشناسی بهتری نسبت به ترکیب غیرخطی اطلاعات در سطح ویژگی با ساختار TANDEM و روش HMM داشته است.

- در مقایسه با نتایج ارائه شده در (Shekofteh, 2013) که در آن از اعمال نگاشت خطی LDA بر روی ویژگی‌های اولیه

نتایج به دست آمده از نمودار شکل (۷) نشان می‌دهد که استفاده از روش‌های ترکیب جمع وزن دار و ضرب وزن دار می‌توانند منجر به بهبود نتایج سیستم بازشناسی ترکیبی نسبت به سیستم بازشناس پایه با ویژگی‌های متداول MFCC39 شود. همچنین کارایی روش ترکیبی ضرب وزن دار از روش جمع وزن دار بالاتر است.

با اعمال تابع لگاریتم در رابطه روش ترکیب ضرب وزن دار، این روش تبدیل به روش جمع وزن دار بر روی مقادیر لگاریتم احتمال پسین می‌شود. از آنجا که تابع لگاریتم، مقادیر کمتر از یک (محدوده مقادیر احتمالات پسین) را به صورت غیرخطی و با وضوح مناسبی بازنمایی می‌کند، نتایج به دست آمده از شکل (۷) مؤید این مطلب است که استفاده از تابع لگاریتم در آنالیز مقادیر احتمال پسین برای کاربردهای بازشناسی بسیار مؤثرتر است. در ادامه و در جدول (۲) نتایج آزمون دقت بازشناسی قاب و واج برای تمامی روش‌های ترکیب بردارهای PPRPS و PPMFCC آورده شده است. در اینجا تنها از اعمال الگوریتم ویتربی بر روی مقدار ترکیبی احتمال پسین (PPcomb) برای تعیین خروجی بازشناسی در سطح آوایی واج استفاده شده است.



(شکل ۸) - نمودار خطای تعلیم شبکه TMLP در صد تکرار اول برای بردار ویژگی‌های ورودی مختلف با بُعد ورودی F.

با توجه به نمودارهای خطای رسم شده در شکل (۸)، شبکه TMLP دربرگیرنده ویژگی‌های PPMFCC+PPRPS، خطای تعلیم کمتری را نسبت به دیگر سامانه‌های مورد استفاده کسب کرده است که این نیز می‌تواند ملاک مناسبی برای توجیه عمل کرد بهتر آن در فرآیند یادگیری و نتایج ارائه شده در جدول‌های (۱) و (۲) باشد.

همان‌طور که در بخش مقدمه اشاره شد، در روش‌های متداول استخراج ویژگی از سیگنال گفتار به‌طور معمول از روش‌های مبتنی بر طیف و کپستروم (مانند MFCC، LFBE و PLP) استفاده می‌شود که در الگوریتم آنها اطلاعات فاز سیگنال گفتار به‌طور کامل حذف شده است. نتایج به‌دست آمده در آزمایش‌های این بخش مؤید این مطلب است که ترکیب مناسب ویژگی‌های کپستروم (MFCC) که فاقد اطلاعات فاز سیگنال است با ویژگی‌های پیش‌نهادهی (LLRPS) استخراج شده از RPS که حاوی اطلاعات حذف نشده فاز سیگنال است، توانسته است به‌علت هم‌افزودن اطلاعات آنها، موجب بهبود عمل کرد نهایی سامانه بازشناسی در حالت ترکیبی شود.

۱۰- نتیجه‌گیری

در این مقاله روشی مبتنی بر استفاده از اطلاعات فضای بازسازی شده فاز (RPS) سیگنال گفتار برای بهبود عمل کرد یک سامانه بازشناسی گفتار پیوسته واج، مطرح و بررسی شد. در این روش علاوه بر استفاده از بردار ویژگی‌های متداول گفتاری، از اطلاعات به‌دست آمده از یک روش نوین استخراج ویژگی مبتنی بر فضای بازسازی شده فاز استفاده شده است. این ویژگی‌ها (LLRPS) بر مبنای ساختار هندسی توزیع سیگنال گفتار جاسازی شده در حوزه RPS حاصل شده‌اند.

LLRPS اعمال شده است، استفاده از نگاشت غیرخطی ویژگی‌ها با روش شبکه عصبی و ساختار TANDER، منجر به بهبود نتایج دقت بازشناسی واج از ۵۶/۶۹ به ۶۰/۰۵ در بازشناسی با HMM شده است.

همان‌طور که از نتایج جدول (۲) مشاهده شد، عمل کرد روش ترکیب غیرخطی TMLP، با اعمال مجدد ویژگی‌های PPRPS و PPMFCC به شبکه عصبی TMLP در یک ساختار سلسله‌مراتبی، توانست بالاترین نتیجه دقت بازشناسی را کسب کند. یک توجیه مناسب برای بهبود عمل کرد این روش، بررسی نوع ویژگی‌های ورودی به شبکه‌های عصبی در هر یک از طبقات ساختار سلسله‌مراتبی شکل (۶) است. به‌عنوان مثال ویژگی متداول MFCC از نوع ویژگی‌های صوتی مبتنی بر مدل منبع-فیلتر^۱ است که درجه بالایی از تغییرات غیرزبانی^۲ مانند مشخصه‌های گوینده و محیط (مانند نوفه و کانال) را شامل می‌شود. این در حالی است که ویژگی‌های مبتنی بر احتمالات پسین دارای توزیع تُنک^۳ بوده و به‌دلیل خاصیت تفکیک‌پذیری بالاتر آن نسبت به ویژگی‌های صوتی، کمتر تحت تأثیر مشخصه‌های هم‌تولیدی^۴ گفتاری قرار می‌گیرند (Ellis, 2001; Sivasdas, 2002). بنابراین یادگیری الگوهای گفتاری توسط شبکه‌های عصبی موجود در طبقه اول ساختار شکل (۶)، دارای پیچیدگی فراوان‌تری نسبت به طبقه دوم آن است که باید آموزش الگوهای گفتاری را از روی ویژگی‌های مبتنی بر احتمالات پسین بر روی توالی قاب‌های گفتاری انجام دهد.

برای بررسی بهتر این موضوع، در شکل (۸) نمودار خطای تعلیم (MSE) شبکه عصبی TMLP در حالت استفاده از ویژگی‌های ورودی LLRPS با بُعد ویژگی F=30، MFCC13 با F=13، MFCC39 با F=39 و شبکه عصبی TMLP در طبقه دوم ساختار شکل (۶) با بردار ترکیب PPMFCC+PPRPS و F=60، برای صد تکرار اول مرحله آموزش آنها آورده شده است.

¹ Source-Filter Model
² Nonlinguistic
³ Sparse
⁴ Coarticulation

Alsteris, L.D., Paliwal, K.K., 2007. Short-time phase spectrum in speech processing: A review and some experimental results. *Digital Signal Processing*, vol. 17, pp. 578–616.

Awrejcewicz J., 1990. Bifurcation portrait of the human vocal cord oscillation. *Journal of Sound Vibrations*. vol. 136, pp. 151–156.

Banbrook, M., McLaughlin, S., 1996. Dynamical modelling of vowel sounds as a synthesis tool. In *Proc. ICSLP*, pp. 1981-1984.

Berry, D.A., Herzel, H., Titze, I.R., Krischer K., 1994. Interpretation of biomechanical simulations of normal and chaotic vocal fold oscillations with empirical eigenfunctions. *The Journal of the Acoustical Society of America*, vol. 95, pp. 3595–3604.

Bijankhan M., Sheykhzadegan J., Roohani M.R., Zarrintare R., Ghasemi S.Z. and Ghasedi M.E., 2003. TFarsDat - The Telephone Farsi Speech Database. In *Proc. EuroSpeech*, Geneva, Switzerland.

Bourlard, H., Morgan, N., 1994. *Connectionist speech recognition - a hybrid approach*. Norwell, MA: Kluwer.

Chen, B., Zhu, Q., Morgan, N., 2005. Tonotopic multi-layer perceptron a neural network for Learning long-term temporal features for speech recognition. In *Proc. ICASSP*, pp. 945-948.

Davis, S.B., Mermelstein, P., 1980. Comparison of parametric representations for monosyllable word recognition in continuously spoken sentences. *IEEE Trans. Speech and Audio Processing*, vol. 28(4), pp. 357-366.

Ellis, D., Singh, R., Sivasdas, S., 2001. Tandem acoustic modeling in large vocabulary recognition. In *Proc. ICASSP*, pp. 517–520.

Furui, S., 1986. Speaker-independent isolated word recognition using dynamic features of speech spectrum. *IEEE Trans. Acoustics, Speech and Signal Processing*, vol. 34 (1), pp. 52-59.

FARSDAT, Persian speech database: <http://catalog.elra.info/product_info.php?products_id=18>.

Hagmuller, M., Kubin, G., 2006. Poincare pitch marks. *Speech Communication*, vol. 48, pp. 1650–1665.

Hegde, R. M., Murthy, H.A., Gadde, V.R.R., 2007. Significance of the modified group delay feature in speech recognition, *IEEE Trans. Audio, Speech and Language Processing*, vol. 15 (1), pp. 190–202.

برای این منظور باید در مرحلهٔ تعلیم، یکسری مدل‌های مانیفولد واجی با مدل‌سازی GMM تعلیم داده شوند. ویژگی‌های نهایی، از امتیاز شباهت (مبتنی بر لگاریتم درست‌نمایی) تراژکتوری سیگنال گفتار جاسازی شده در RPS با GMM‌های گفتاری تعلیم‌یافته در RPS محاسبه شوند. در اینجا نشان دادیم که با استفاده از مدل بازشناس غیرخطی TMLP، عمل کرد ویژگی پیشنهادی LLRPS از ویژگی‌های متداول MFCC پایین‌تر است. اما با ترکیب مناسب خروجی‌های به‌دست آمده از دو سیستم بازشناس حاوی اطلاعات متمایز نسبت به یکدیگر هستند، می‌توانیم نتایج دقت قاب و واج بازشناسی شده را به‌ترتیب به مقدار $3/94\%$ و $4/02\%$ نسبت به بهترین سامانه بازشناس پایه مبتنی بر مدل TMLP افزایش دهیم. همچنین با استفاده از ساختار TANDEM، نشان داده شد که الحاق ویژگی‌های مبتنی بر RPS می‌تواند کارایی ویژگی‌های MFCC را حدود $0/6\%$ بهبود دهد.

مراجع

رحیمی نژاد، مهدی؛ سیدصالحی، سیدعلی؛ ۱۳۸۲. مقایسه و ارزیابی کارایی انواع روش‌های استخراج پارامترهای بازنمایی و هنجارسازی در بازشناسی مستقل از گویندهٔ گفتار، نشریه امیرکبیر، ش. ۵۵.

جعفری، ایوب؛ الماس‌گنج، فرشاد؛ نبی بیدهندی، مریم؛ ۱۳۸۹. مدل‌سازی غیرخطی قطع پوانکاره سیگنال گفتار در ترکیب با تحلیل حوزهٔ فرکانس به‌منظور افزایش صحت عمل کرد سیستم‌های بازشناسی گفتار، نشریه فنی و مهندسی مدرس، دوره ۱۰، ش. ۳، ص. ۵۵ - ۷۰.

ولی، منصور؛ سیدصالحی، سیدعلی؛ ۱۳۸۵. بازشناسی مقاوم و توأم گفتار مستقیم و تلفنی با استخراج مناسب بردارهای بازنمایی و اصلاح آنها توسط معکوس‌سازی شبکه‌های عصبی، نشریه مهندسی برق و مهندسی کامپیوتر ایران، سال ۴، ش. ۱، ص. ۲۱ - ۲۹.

شکفته، یاسر؛ الماس‌گنج، فرشاد؛ ۱۳۸۶. بهبود بازشناسی گفتار با استفاده از شبکه‌های عصبی دربرگیرنده الگوهای زمانی، مجموعه مقالات سومین کنفرانس بین‌المللی فناوری اطلاعات و دانش (IKT2007)، دانشگاه فردوسی مشهد، آذر ۱۳۸۶.

Abarbanel, H.D.I., 1996. *Analysis of observed chaotic data*. Springer, New York.

- Kennel, M.B., Brown, R., Abarbanel, H.D.I., 1992. Determining embedding dimension for phase-space reconstruction using a geometrical construction. *Physical review A*, vol. 45 (6), pp. 3403–3411.
- Kittler, J., Hatef, M., Duin, R.P.W., Matas, J., 1998. On combining classifiers. *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 20(4), pp. 226–239.
- Kokkinos, I., Maragos, P., 2005. Nonlinear speech analysis using models for chaotic systems. *IEEE Trans. Speech Audio Processing*, vol. 13, pp. 1098–1109.
- Lindgren, A.C., Johnson, M.T., Pavinelli, R.J., 2003. Speech recognition using reconstructed phase space features. In *Proc. ICASSP, China*, pp. 61–63.
- Lindgren, A.C., Johnson, M.T., Pavinelli, R.J., 2004. Joint frequency domain and reconstructed phase space features for speech recognition. In *Proc. ICASSP, Montreal, Canada*, pp. I-533–I-536.
- Liu, X., Pavinelli, R.J., Johnson, M.T., 2003. Vowel classification by global dynamic modeling. In *Proc. NOLISP*, pp. 111–114.
- Misra, H., Bourlard, H., Tyagi, V., 2003. New entropy based combination rules In HMM/ANN multi-stream ASR. In *Proc. ICASSP*, pp. 741–744.
- Morgan, N., Chen, B.Y., Zhu, Q., Stolcke, A., 2004. TRA-Pping conversational speech: Extending TRAP/Tandem approaches to conversational telephone speech recognition. In *Proc. ICASSP*, pp. 536–539.
- Narayanan, N.K., Thasleema, T.M., Prajith, P., 2012. Reconstructed state space model for recognition of consonant - vowel utterances using support vector machines. *International Journal of Artificial Intelligence and Applications*, vol. 3(2), pp. 101-119.
- Narayanan, S.S., Alwan, A.A., 1995. A nonlinear dynamical systems analysis of fricative consonants. *Acoustical Society of America Journal*, vol. 97, pp. 2511-2524.
- Nejadgholi, I., Seyyedsalehi, S.A., 2009. Nonlinear normalization of input patterns to speaker variability in speech recognition neural networks. *Neural Computing and Applications*, vol. 18, pp. 45–55.
- Paliwal, K., Alsteris, L., 2005. On the usefulness of STFT phase spectrum in human listening tests. *Speech Communication*, vol. 45, pp. 153–170.
- Park, J., Diehl, F., Gales, M.J.F., Tomalin, M., Woodland, P.C., 2011. The efficient incorporation of MLP features into automatic speech recognition
- Hermansky, H., 1990. Perceptual linear predictive (PLP) analysis of speech. *Journal of the Acoustic Society of America*, vol. 87(4), pp. 1738-1752.
- Hermansky, H., Ellis, D.P., Sharma, S., 2000. Tandem connectionist feature extraction for conventional HMM systems. In *Proc. ICASSP*, pp. 1635-1638.
- Herzel, H., Berry, D., Titze, I., Steinecke, I., 1995. Nonlinear dynamics of the voice: signal analysis and biomechanical modeling. *Chaos*, vol. 5, pp. 30–34.
- HTK (v.3.4), Hidden Markov Model Toolkit: <<http://htk.eng.cam.ac.uk/>>
- Indrebo, K.M., Pavinelli, R.J., Johnson, M.T., 2006. Sub-banded reconstructed phase spaces for speech recognition. *Speech Communication*, vol. 48, pp. 760-774.
- Ikbal, S., Misra, H., Hermansky, H., Magimai, M., 2012. Phase autocorrelation (PAC) features for noise robust speech recognition, *Speech Communication*, vol. 54, pp. 867–880.
- Jafari, A., Almasganj, F., 2012. Using nonlinear modeling of reconstructed phase space and frequency domain analysis to improve automatic speech recognition performance. *International Journal of Bifurcation and Chaos*, vol. 22(3).
- Jafari, A., Almasganj, F., NabiBidhendi, M., 2010. Statistical modeling of speech Poincaré sections in combination of frequency analysis to improve speech recognition performance. *Chaos*, vol. 20, pp. 033106:1-11.
- Jiang, J.J., Zhang, Y., 2002. Chaotic vibration induced by turbulent noise in a two-mass model of vocal folds. *The Journal of the Acoustical Society of America*, vol. 112, pp. 2127–2133.
- Jiang, J.J., Zhang, Y., McGilligan, C., 2006. Chaos in voice, from modeling to measurement. *Journal of Voice*, vol. 20(1), pp. 2-17.
- Johnson, M.T., Pavinelli, R.J., Lindgren, A.C., Ye, J., Liu, X., Indrebo, K.M., 2005. Time-domain isolated phoneme classification using reconstructed phase spaces. *IEEE Trans. Speech Audio Processing*, vol. 13(4), pp. 458–466.
- Kantz, H., Schreiber, T., 1997. *Nonlinear Time Series Analysis* Cambridge University Press, Cambridge, England.
- Kazemi, A.R., Sobhanmanesh, F., 2011. MLP refined posterior features for noise robust phoneme recognition. *Scientia Iranica, Trans. D: Computer Science & Engineering and Electrical Engineering*, vol. 18, pp. 1443–1449.

consonant classification. *International Journal of Speech Technology*, vol. 15(2), pp. 227-239.

Tohidypour, H.R., Seyyedsalehi, S.A., Behbood, H., Roshandel, H., 2012. A new representation for speech frame recognition based on redundant wavelet filter banks. *Speech Communication*, vol. 54(2), pp. 256-271.

Valente, F., 2010. Multi-stream speech recognition based on Dempster-Shafer combination rule. *Speech Communication*, vol. 52(3), pp. 213-222.

Vaziri, G., Almasganj, F., Behroozmand, R., 2010. Pathological assessment of patients' speech signals using nonlinear dynamical analysis. *Computers in Biology and Medicine*, vol. 40(1), pp. 54-63.

White, H., 1989. Learning in artificial neural networks: A statistical perspective, *Neural Computation*, pp. 425-464.

Ye, J., Johnson, M.T. M.T., Povinelli, R.J., 2003. Phoneme classification over reconstructed phase space using principal component analysis. In Proc. NOLISP, Le Croisic, France, pp. 11-16.

Ye, J., Povinelli, R.J., Johnson, M.T., 2002. Phoneme Classification Using naïve Bayes Classifier in Reconstructed Phase Space. In Proc. IEEE Digital Signal Processing Workshop, Atlanta, Georgia.

Yu, S., Zheng, D., Feng, X., 2006. A new time domain feature parameter for phoneme classification. In Proc. WESPAC IX 2006, Seoul, Korea.

Zavaliagos, G., Zhao, Y., Schwartz, R., Makhoul, J., 1994. A hybrid seg-mental neural net/hidden Markov model system for continuous speech recognition, *IEEE Trans. Speech Audio Processing*, vol. 2 (1) pp. 151-160.

Zhu, Q., Chen, B., Morgan, N., Stolcke, A., 2004. On using MLP features in LVCSR. In Proc. ICSLP.

systems. *Computer Speech and Language*, vol. 25, pp. 519-534.

Pinto, J., Garimella, S., Magimai-Doss, M., Hermansky, H., Bourlard, H., 2011. Analysis of MLP-Based hierarchical phoneme posterior probability estimator. *IEEE Trans. Audio Speech Language Processing*, vol. 19(1), pp. 225-241.

Pitsikalis, V., Maragos, P., 2002. Speech analysis and feature extraction using chaotic models. In Proc. ICASSP, Orlando, Florida, pp. 533-536.

Pitsikalis, V., Maragos, P., 2009. Analysis and classification of speech signals by generalized fractal dimension features. *Speech Communication*, vol. 51(12), pp. 1206-1223.

Povinelli, R.J., Johnson, M.T., Lindgren, A.C., Roberts, F.M., Ye, J., 2006. Statistical models of reconstructed phase spaces for signal classification. *IEEE Trans. Signal Processing*, vol. 54, pp. 2178-2186.

Povinelli, R.J., Johnson, M.T., Lindgren, A.C., Ye, J., 2004. Time series classification using Gaussian mixture models of reconstructed phase spaces. *IEEE Trans. Knowledge and Data Engineering*, vol. 16, pp. 779-783.

Sameti, H., Veisi, H., Bahrani, M., Babaali, B., Hosseinzadeh, K., 2011. A large vocabulary continuous speech recognition system for Persian language. *EURASIP Journal on Audio, Speech, and Music Processing*, vol. (2011-1), pp. 1-12.

Sauer, T., Yorke, J.A., Casdagli, M., 1991. Embedology. *Journal of Statistical Physics*, vol. 65, pp. 579-616.

Shekofteh, Y., Almasganj, F., 2013. Feature extraction based on speech attractors in the reconstructed phase space for automatic speech recognition systems, *ETRI Journal*, vol. 35(1), pp. 100-108.

Sivadas, S., Hermansky, H., 2002. Hierarchical tandem feature extraction, In Proc. ICASSP, pp. 809-812.

Sun, J., Zheng, N., Wang, X., 2007. Enhancement of Chinese speech based on nonlinear dynamics. *Signal Processing*, vol. 87, pp. 2431-2445.

Takens, F., 1980. Detecting strange attractors in turbulence. In Proc. Dynamical System Turbulence, pp. 366-381.

Thasleema, T.M., Prajith, P., Narayanan, N.K., 2012. Time-domain non-linear feature parameter for

یاسر شکفته مدارک کارشناسی،
 کارشناسی ارشد و دکترای خود را
 در رشته مهندسی پزشکی (گرایش
 بیوالکترونیک) به ترتیب در سال‌های
 ۱۳۸۴، ۱۳۸۷ و ۱۳۹۲ از دانشکده
 مهندسی پزشکی دانشگاه صنعتی امیرکبیر اخذ نموده است.
 زمینه‌های تحقیقاتی مورد علاقه ایشان شامل پردازش

سیگنال‌های حیاتی، شناسایی الگو، مدل‌سازی سیستم‌های بیولوژیکی و سیستم‌های بازشناسی گفتار است.

نشانی رایانامه ایشان عبارتست از:

y_shekofteh@{aut.ac.ir,yahoo.com}



فرشاد الماس گنج مدارک

کارشناسی و کارشناسی ارشد خود را در رشته مهندسی برق (گرایش الکترونیک) به ترتیب در سال‌های ۱۳۶۳ و ۱۳۶۷ از دانشگاه صنعتی

امیرکبیر اخذ نموده است. وی سپس در سال ۱۳۷۷ مدرک دکترای خود را در گرایش مهندسی پزشکی از دانشکده مهندسی برق دانشگاه تربیت مدرس اخذ نمود. ایشان هم‌اکنون عضو هیئت علمی دانشکده مهندسی پزشکی دانشگاه صنعتی امیرکبیر با سمت و درجه علمی دانشیار هستند. زمینه‌های تحقیقاتی مورد علاقه ایشان شامل پردازش سیگنال‌های رقمی، شناسایی الگو و مدل‌سازی صوتی و زبانی در سیستم‌های بازشناسی گفتار است.

نشانی رایانامه ایشان عبارتست از:

almas@aut.ac.ir

Archive of S.I.P.

فصلنامه



سال ۱۳۹۲ شماره ۱ پیاپی ۱۹

www.SIP.ir