

# مقایسه و ارزیابی کارآیی انواع روش‌های استخراج پارامترهای بازنمایی و هنجارسازی در بازشناسی مستقل از گوینده گفتار

سید علی سیدصالحی  
استادیار

مهدی رحیمی‌نژاد  
کارشناسی ارشد

دانشکده مهندسی پزشکی، دانشگاه صنعتی امیرکبیر

## چکیده

فرآیند بازشناخت گفتار در اولین گام خود نیازمند یک نمایش مناسب از اطلاعات سیگنال گفتاری است. این بازنمایی باید شامل ویژگی‌هایی از سیگنال گفتار باشد که حاوی اطلاعات مفیدی جهت متمایز کردن آواهای مختلف باشند. برای استخراج این بیان مناسب اطلاعات سیگنال یا پارامترهای بازنمایی، روش‌های متفاوتی ارائه شده‌است که در سیستم‌های بازشناخت مختلفی مورد استفاده قرار گرفته‌اند. از رایج‌ترین این روش‌ها، روش‌های طیفی می‌باشند که در آنها تحلیل طیفی از سه طریق تبدیل فوریه زمان کوتاه، مدل پیشگویی خطی و تبدیل وولت صورت می‌گیرد. در این مقاله متداول‌ترین روش‌های استخراج پارامترهای بازنمایی مورد مقایسه قرار گرفته‌اند و با اعمال تغییراتی در الگوریتم برخی از آنها و نیز استفاده از روش‌های نرمالیزه متفاوت، مناسب‌ترین آنها در بکارگیری برای بازشناسی مستقل از گوینده گفتار عاری از نویز، به کمک یک شبکه عصبی جلوسو با تأخیر زمانی، معرفی گردیده‌اند. در ارزیابی‌های انجام شده روش‌های طیفی الهام گرفته از سیستم شنوایی انسان و مبتنی بر تبدیل فوریه بالاترین صحت بازشناخت را ارائه داده‌اند.

## کلمات کلیدی

پارامترهای بازنمایی، هنجارسازی به گوینده، هنجارسازی به بلندی صدا، LHCB، MFCC، PLP، LPCC، وولت، بازشناخت گفتار مستقل از گوینده.

## A Comparative Study of Representation Parameters Extraction And Normalization Methods for Speaker Independent Recognition of Speech

M. Rahiminejad  
M.Sc.

S. A. Seyyed Salehi  
Assistant Professor

Biomedical Engineering Department,  
Amirkabir University of Technology

## Abstract

*An automatic speech recognition system at The first step for a proper representation of speech input which contains all useful information of the signal. These representations must be the voice features, which may distinguish different phonemes of a language. There are several methods for feature extraction that are used in ASR systems. The most common methods are spectral ones. The spectral analysis is based on Short Time Fourier Transform (STFT), Linear Prediction (LP), or Discrete Wavelet Transform (DWT). In this paper, these feature extraction methods, their modified versions, and some normalization algorithms were studied and the most suitable feature extraction methods for time delay neural network classifier were identified. It was shown that the perceptual*

*STFT based features are the best, among the features studied in this paper.*

## Keywords

*Representation parameters, speaker normalization, loudness normalization, LHC, MFCC, PLP, LPCC, wavelet, speaker independent*

## مقدمه

فرآیند بازشناخت گفتار نیازمند به بیان مناسب سیگنال گفتار ورودی است که اطلاعات مفید آنرا دربر داشته باشد. این اطلاعات گفتار شامل ویژگی‌هایی از سیگنال گفتار می‌باشد که دارای اطلاعات مفیدی جهت متمایز کردن آواهای مختلف می‌باشند. از آنجا که در بازشناخت گفتار، اولین گام، استخراج پارامترهای بازنمایی می‌باشد و چنانچه اطلاعات مفیدی در این مرحله از دست برود دیگر در مراحل بعدی ممکن است قابل جبران نباشد، استخراج پارامترهای بازنمایی را می‌توان یک مرحله زیربنایی در سیستم بازشناخت گفتار به حساب آورد.

تاکنون روشهای مختلفی جهت استخراج پارامترهای بازنمایی از سیگنال گفتار معرفی و مورد استفاده قرار گرفته‌اند. اکثر روشهای متداول استخراج ویژگی، روشهای طیفی می‌باشند یعنی پارامترهای بازنمایی در حوزه فرکانس و از روی طیف سیگنال بدست آورده می‌شوند. برخی نتایج حاکی از این است که برای بازشناخت مستقل از گوینده آواها به کمک شبکه عصبی بازگشتی<sup>۱</sup>، بازنمایی در دامنه طیفی مستقل از اینها را به دو دسته کلی تقسیم کرد؛ روشهایی که از آنالیزهای خطی استفاده می‌کنند و روشهایی که از سیستم شنوایی انسان الهام می‌گیرند. در سالهای اخیر گرایش به سمت ایده گرفتن از سیستم شنوایی انسان برای استخراج پارامترهای بازنمایی بوجود آمده است. چراکه سیستم شنوایی انسان یک سیستم بازشناخت با عملکردی بسیار مطلوب می‌باشد. دنگ<sup>۲</sup> معتقد است که موفقیت نهایی در ساخت سیستمهای بازشناخت با تواناییهای مشابه مغز انسان به فهم کامل پردازشهای موجود در سیستم طبیعی شنوایی انسان وابسته است. بدین منظور وی یک مدل غیر خطی برای گوش داخلی پیشنهاد داده و با آزمایش آن به این نتیجه رسیده است که پردازشهای غیرخطی و الهام گرفته از مدل گوش داخلی، تواناییهای قوی تری در پردازش گفتار نسبت به روشهای آنالیز خطی فراهم نموده‌اند[۴]. همچنین هالتون<sup>۳</sup> با انجام یک سری آزمایشات جهت آشکارسازی و طبقه بندی ویژگیهای گفتاری آواهای واکدار به این نتیجه رسیده است که الگوریتمهای بازشناخت گفتار برمبنای مدل شنوایی جایگزینهای مؤثری نسبت به روشهای طیفی معمولی جهت استخراج فورمنتها هستند [۵].

تحلیل‌های طیفی که در روشهای طیفی استخراج پارامترهای بازنمایی مورد استفاده قرار می‌گیرند عبارتند از تبدیل فوریه زمان کوتاه، مدل پیشگویی خطی و تبدیل ویولت. تبدیل فوریه گسسته زمان کوتاه یک روش کلاسیک در تحلیل طیفی سیگنال گفتار می‌باشد. با استفاده از رابطه تبدیل فوریه از روی نمونه‌های زمانی سیگنال، نمونه‌های فرکانسی آن ساخته می‌شود. در مدل پیشگویی خطی، یک مدل ریاضی تمام قطب برای سیستم تولید گفتار فرض می‌گردد و ضرایب این تابع تمام قطب که حاوی اطلاعات مربوط به فورمنتها هستند تخمین زده می‌شوند. تبدیل ویولت نیز چون برای سیگنال‌های غیرایستاد مناسب است و خاصیت چند رزولوشنی دارد، در سالهای اخیر توجه محققان را جهت استفاده از آن در تحلیل طیفی سیگنال گفتار و برای استخراج پارامترهای بازنمایی به خود جلب کرده است.

با توجه به این مطلب که در طبقه بندی آواها، میزان صحت طبقه بندی به این نکته بستگی دارد که روش طبقه بندی تا چه حد با روش بازنمایی سیگنال انطباق داشته باشد [۶]، بررسی روشهای بازنمایی مختلف برای یک طبقه بندی کننده خاص مورد نیاز است. در این مقاله هدف مقایسه ای بین روشهای مختلف استخراج پارامترهای بازنمایی در شرایطی یکسان در بکارگیری برای بازشناسی مستقل از گوینده گفتار عاری از نویز می‌باشد. در این سیستم از یک شبکه عصبی جلوسو با تأخیر زمانی به عنوان طبقه بندی کننده استفاده می‌شود. علاوه بر روشهای بازنمایی متداول، با بهبود الگوریتم برخی از آنها و نیز استفاده از روشهای هنجارسازی متفاوت جهت هنجارساختن پارامترهای بازنمایی، مناسبترین روشهای بازنمایی جهت استفاده در سیستم فوق معرفی می‌گردند.

در ادامه مقاله، دادگان مورد آزمایش و شبکه عصبی که به عنوان طبقه‌بندی کننده در آزمایشات مورد استفاده قرار گرفته است در بخش ۲ معرفی می‌گردند و در بخش ۳ به معرفی روش‌های هنجارسازی پارامترهای بازنمایی که در [۱] برای پارامترهای مورد استفاده توسط شبکه عصبی مناسب تشخیص داده شده‌اند پرداخته می‌شود. در بخش ۴ روش‌های رایج استخراج پارامترهای بازنمایی که با توجه به گزارشات مقالات و کارهای قبلی انتخاب شده‌اند، معرفی می‌شوند. در بخش ۵ نتایج حاصل از پارامترهای مختلف بیان شده با هم مقایسه می‌گردند و جمع‌بندی و نتیجه‌گیری نهایی در بخش ۶ آورده شده است.

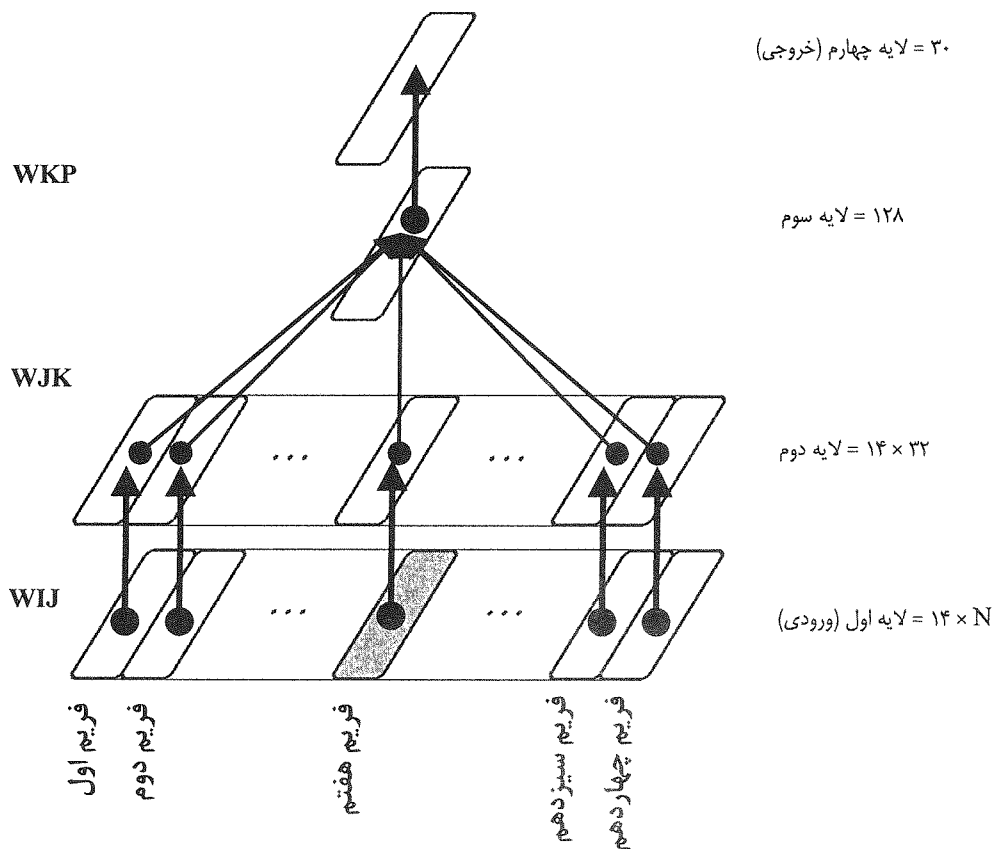
## ۱- مدل شبکه عصبی بازشناس و دادگان آزمایش

برای طبقه‌بندی آواها در این مقاله از یک شبکه عصبی جلوسوی زمانی استفاده شده است. این شبکه براساس نگاشت مسیر بردار بازنمایی به فضای آواها طراحی شده است و در واقع توسط این شبکه درصد صحت بازشناخت فریم‌ها یا به بیان دیگر نمونه‌های مسیر بردارهای بازنمایی، استخراج شده توسط هر کدام از روش‌ها، تعیین می‌گردد. توسط تکنیک‌هایی می‌توان از بازشناخت فریم‌ها به بازشناخت آواها رسید اما چون هدف مقایسه عملکرد روش‌های مختلف استخراج پارامترهای بازنمایی است و مهم یکسان بودن شرایط آزمایش برای تمام پارامترها است، لذا مقایسه صحت بازشناخت فریم‌ها نیز خود یک معیار مناسب جهت مقایسه عملکرد روش‌های استخراج پارامترهای بازنمایی می‌باشد. همانطور که در شکل ۱ نشان داده شده است، شبکه عصبی مورد استفاده دارای ۱۴ تأخیر زمانی در لایه ورودی می‌باشد، یعنی در ورودی به ۱۴ فریم کنار هم نگاه می‌کند و دارای دو لایه پنهان به ابعاد  $32 \times 14$  و ۱۲۸ می‌باشد، لایه خروجی نیز دارای ۳۰ نورون، به تعداد آواهایی که باید بازشناسی شوند، می‌باشد. تابع تصمیم‌گیری هر نورون یک تابع غیر خطی از نوع سیگموئید است که باعث تصمیم‌گیری نرم در خروجی هر نورون می‌گردد. این شبکه بر اساس طی مسیر بازنمایی بر روی دادگان آزمون لغزیده و همواره یکی از خروجیها که بیانگر یکی از آواها است فعال می‌شود. در تعلیم شبکه و برای بدست آوردن وزنه‌های مطلوب بین لایه‌ها از روش پس انتشار خطا استفاده می‌گردد و وزنه‌ها در جهتی اصلاح می‌شوند که مجموع مربعات خطای خروجی (اختلاف بین خروجی شبکه و خروجی مطلوب) به حداقل برسد. در آموزش شبکه ضریب یادگیری شبکه و ضریب ممنوم به ترتیب  $0.1$  و  $0.7$  انتخاب شده‌اند و برای اینکه شرایط آموزش و شرایط کار برای تمام حالت‌ها و برای تمام ویژگیهای مختلف یکسان باشد، ضریب یادگیری شبکه در طول آموزش ثابت می‌ماند.

با توجه به اینکه سیستم مورد نظر برای بازشناخت گفتار پیوسته طراحی شده است، لذا دادگانی که برای آموزش و آزمون سیستم انتخاب می‌شوند، باید جملات پیوسته باشند. یعنی گوینده یک جمله را به طور معمولی و بدون ایجاد وقفه‌های عمدی در بین کلمات ادا می‌کند. از طرفی چون سیستم بازشناخت مستقل از گوینده است لذا برای آموزش مدل بازشناخت تا آنجا که ممکن است باید از تنوع گویندگان مختلف استفاده کرد تا مدل بازشناخت در هنگام آموزش فقط به یک سری جملات بیان شده توسط یک یا چند گوینده خاص وابسته نشود و تا آنجا که ممکن است گوینده‌های متفاوت را ببیند. برای این منظور از دو جمله ۴۰۰ و ۵۰۰ از ۱۰۱ گوینده از مجموعه دادگان فارسی‌زبان که در محیطی عاری از نویز ضبط شده و با فرکانس ۴۴۱۰۰ هرتز نمونه برداری شده‌اند، استفاده می‌شود [۱۳]. جملات ۷۱ گوینده برای تعلیم شبکه و جملات ۳۰ گوینده دیگر به عنوان دادگان آزمون مورد استفاده قرار می‌گیرند. در این دادگان آواها تقطیع نشده‌اند و صرفاً برحسب شماره نمونه‌ها برچسب‌دهی گردیده‌اند. لذا سیگنال‌های گفتاری و برچسب‌دهی آن بدون تقطیع در دسترس می‌باشند.

## ۲- روش‌های هنجارسازی پارامترهای بازنمایی

در بسیاری از سیستمهای بازشناخت گفتار، استفاده از روشهای استخراج پارامترهای بازنمایی به تنهایی پارامترهای مناسبی در اختیار قرار نمی‌دهند و پس از حصول بردارهای ویژگی، نیاز است که این پارامترها نسبت به عوامل مختلفی به هنجار<sup>۴</sup> شوند [۷]-[۱۰]. با توجه به اینکه مدل بازشناخت مورد استفاده یک شبکه عصبی جلو سو با تأخیر زمانی می‌باشد، در این قسمت به معرفی برخی از روشهای نرمالیزه پارامترهای بازنمایی که طی یک سری آزمایشات برای پارامترهای مورد استفاده برای شبکه عصبی مفید شناخته شده‌اند [۱]، می‌پردازیم.



شکل (۱) ساختار شبکه عصبی جلوسوی زمانی.

## الف - هنجارسازی عرضی پارامترهای بازنمایی

هنگامیکه یک گوینده، یک آوای خاص را چندبار با بلندی صدای متفاوت بیان کند، پارامترهایی که از این آوا بدست می‌آیند، با هم یکسان نیستند. علت این است که بلندی صدای گوینده در انرژی کل طیف تأثیر مستقیم می‌گذارد و لذا هنگامیکه انرژی طیفی در باندهای فرکانسی مختلف به عنوان پارامترهای بازنمایی محاسبه می‌شوند، این انرژی کل در این مقادیر پارامترها نمود پیدا می‌کند و لذا صحت بازشناسی سیستم بازشناخت گفتار کاهش پیدا می‌کند. لذا کاهش حساسیت پارامترهای بازنمایی نسبت به بلندی صدا می‌تواند باعث بهبود عملکرد سیستم بازشناخت گفتار شود.

در این روش هنجارسازی، میانگین  $N$  پارامتر بردار بازنمایی بدست آمده از هر فریم حساب شده و سپس این مقدار میانگین از هر یک از پارامترها کاسته می‌شود [۱]. چون این روش بر روی پارامترهای حاصله از بانک فیلترها اعمال می‌شود، و پارامترها لگاریتم انرژی داخل هر فیلتر هستند، میانگین پارامترها با انرژی کل فریم متناسب می‌باشد و کم کردن این مقدار از پارامترهای بردار بازنمایی مانند آن است که پارامترهای بازنمایی به تقریبی از انرژی کل فریم به هنجار شده باشند. اگر  $\bar{x}$  بردار بازنمایی یک فریم باشد و  $x_i$  پارامتر  $i$  ام (لگاریتم انرژی فیلتر  $i$  ام) این بردار باشد، می‌توان این نحوه هنجارسازی را به صورت زیر بیان کرد.

$$\bar{x} = [x_1, x_2, K, x_N]$$

$$m = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\bar{x}_{\text{new}} = \bar{x} - m = [x_1 - m, x_2 - m, K, x_N - m] \quad (1)$$

که  $m$  میانگین پارامترهای بردار بازنمایی یک فریم و  $\bar{x}_{new}$  بردار به هنجار شده می‌باشند.

## ب - هنجارسازی طولی پارامترهای بازنمایی

در این روش، هر پارامتر بازنمایی به میانگین آن در کل دادگان به هنجار می‌شود. این روش هنجارسازی باعث می‌شود که دامنه تغییرات همه پارامترها تقریباً به محدوده بین  $-1$  تا  $+1$  منتقل گردد و در دو طرف صفر تقریباً متوازن باشند [۱]. در این روش، ابتدا بردار میانگین پارامترها روی کل دادگان محاسبه می‌شود.

$$\bar{k}_1 = \frac{1}{M} \sum_{i=1}^M \bar{x}_i \quad (2)$$

که در این رابطه  $M$  برابر تعداد بردارهای بازنمایی در کل دادگان می‌باشد و  $\bar{x}_i$ ،  $i$  امین بردار بازنمایی است. پس از بدست آمدن بردار میانگین  $\bar{k}_1$ ، این بردار از تمام بردارهای بازنمایی کم می‌گردد.

$$\bar{q}_i = \bar{x}_i - \bar{k}_1 \quad (3)$$

سپس بردار میانگین نرم اول (میانگین قدر مطلق) [۱] و یا بردار میانگین نرم دوم (میانگین مجذور) [۱۲] پارامترهای  $\bar{q}_i$  روی کل دادگان محاسبه می‌گردد. در استفاده از میانگین نرم اول داریم

$$\bar{k}_2 = \frac{1}{M} \sum_{i=1}^M |\bar{q}_i| \quad (4)$$

و برای بدست آوردن پارامترهای به هنجار شده، تمامی بردارهای  $\bar{q}_i$  عنصر به عنصر به بردار  $\bar{k}_2$  تقسیم می‌شوند.

$$\bar{x}_i^{new} = \left\{ \frac{q_{in}}{k_{2n}} \mid n = 1, 2, K, N \right\} \quad (5)$$

در رابطه فوق،  $\bar{x}_i^{new}$  بردار بازنمایی به هنجار شده و  $N$  تعداد پارامترهای هر بردار بازنمایی است.  $q_{in}$  و  $k_{2n}$  به ترتیب  $n$  امین مؤلفه از بردارهای  $\bar{q}_i$  و  $\bar{k}_2$  هستند. اما در استفاده از میانگین نرم دوم داریم

$$\bar{k}_2 = \frac{1}{M} \sum_{i=1}^M \bar{q}_i^2 \quad (6)$$

و برای بدست آوردن پارامترهای به هنجار شده، تمامی مؤلفه‌های بردارهای  $\bar{q}_i$  به جذر مؤلفه‌های بردار  $\bar{k}_2$  تقسیم می‌شوند.

$$\bar{x}_i^{new} = \left\{ \frac{q_{in}}{\sqrt{\alpha k_{2n}}} \mid n = 1, 2, K, N \right\} \quad (7)$$

در رابطه فوق  $\alpha$  ضریبی است که بواسطه آن واریانس پارامترهای به هنجار شده تعیین می‌گردد. و مقدار  $\alpha$  برابر عکس واریانس پارامترهای به هنجار شده می‌باشد. به عنوان مثال اگر می‌خواهیم واریانس پارامترهای به هنجار شده برابر  $0.5$  باشد، باید  $\alpha=2$  انتخاب شود.

این نحوه هنجارسازی موجب می‌شود که مقادیر مؤلفه‌های مختلف دادگان در فضای  $N$  بعدی ورودی شبکه عصبی در

محللهای مناسبتری برای شکل‌گیری مدل نگاشت فضای بازنمایی به آواها قرار گیرند و از طرفی با هم دامنه شدن پارامترهای بازنمایی نسبت به هم، ارزش همه مؤلفه‌ها از دید شبکه عصبی، تقریباً یکسان گردد و لذا شبکه بهتر می‌تواند عمل طبقه‌بندی را روی پارامترهای به‌هنجار شده انجام دهد [۱].

### ب - هنجارسازی طولی پارامترهای بازنمایی، جداگانه برای هر گوینده

این روش نیز مانند روش هنجارسازی قبلی است با این تفاوت که در روش قبلی، بردارهای  $\bar{k}_1$  و  $\bar{k}_2$  روی تمام دادگان آموزش محاسبه می‌شدند اما در این روش بردارهای  $\bar{k}_1$  و  $\bar{k}_2$  روی دادگان مربوط به هر گوینده و به صورت جدا محاسبه می‌شوند و دادگان هر گوینده دارای یک بردار  $\bar{k}_1$  و یک بردار  $\bar{k}_2$  مختص به خود می‌باشد [۱]. مزیت این روش به روش قبل در این است که بردارهای  $\bar{k}_1$  و  $\bar{k}_2$  تخمینی از مشخصات گوینده مربوطه را دربر دارند و پس از هنجارسازی بردارهای بازنمایی مربوط به دادگان یک گوینده، این پارامترها تا حدودی نسبت به گوینده به‌هنجار می‌شوند. اما در روش قبل دو بردار  $\bar{k}_1$  و  $\bar{k}_2$  میانگینی از اطلاعات همه گویندگان را تخمین می‌زنند.

### ۳- روش‌های استخراج پارامترهای بازنمایی

روش‌های طیفی استخراج پارامترهای بازنمایی از نظر نوع تحلیل طیفی به سه دسته تقسیم می‌شوند که در این قسمت متداول‌ترین روش‌های هر دسته معرفی می‌شوند.

#### الف - روش‌های طیفی مبتنی بر تبدیل فوریه زمان کوتاه

در روش‌های طیفی استخراج پارامترهای بازنمایی بخصوص در روش‌هایی که از تحلیل فوریه جهت بدست آوردن طیف سیگنال گفتار استفاده می‌شود، عموماً از بانک فیلتر جهت محاسبه انرژی طیف حول فرکانس‌هایی مشخص به عنوان پارامترهای بازنمایی استفاده می‌گردد. تعداد این فیلترهای میانگذر در سیستم‌های مختلف متفاوت است ولی معمولاً بین ۶ الی ۲۰ فیلتر مورد استفاده قرار می‌گیرند. افزایش فیلترها معمولاً موجب بهبود کیفیت بازشناخت می‌شود ولی در صورت کاهش پهنای باند فیلترها تا حد کمتر از فرکانس گام<sup>۵</sup> صدای گوینده، کیفیت بازشناخت افت می‌کند [۲]. تنظیم فواصل مابین فیلترها به صورت غیرخطی و در مقیاس مل<sup>۶</sup> یا بارک<sup>۷</sup> که مقیاس‌هایی الهام گرفته از سیستم شنوایی انسان می‌باشند، انجام می‌گیرد. روابط تبدیل مقیاس هرتز به این دو مقیاس عبارتند از [۳۰]:

$$f_{\text{mel}} = 2595 \log \left( 1 + \frac{f_{\text{Hz}}}{700} \right) \quad (8)$$

$$f_{\text{Bark}} = 6 \ln \left( \frac{f_{\text{Hz}}}{600} + \sqrt{\left( \frac{f_{\text{Hz}}}{600} \right)^2 + 1} \right) \quad (9)$$

هر دو مقیاس تقریباً شبیه به هم بوده و تا فرکانس یک کیلوهرتز بصورت تقریباً خطی و بالاتر از این فرکانس به صورت لگاریتمی می‌باشند. پارامترهای طیفی مبتنی بر تبدیل فوریه که در این مقاله مورد مطالعه قرار گرفته اند عبارتند از LHCBC<sup>۸</sup>، LHCBC<sup>۹</sup>، MFCC<sup>۱۰</sup> و PLP<sup>۱۱</sup> که در ادامه به معرفی آنها می‌پردازیم.

پارامترهای LHCBC، لگاریتم انرژی خروجی فیلتر بانک‌های مجذور هنینگ و متساوی الفاصله در مقیاس بارک با پهنای یک و فاصله یک نسبت به هم می‌باشند [۲]. الگوریتم استخراج پارامترهای LHCBC بصورت زیر است:

گام اول - انتخاب یک فریم از سیگنال به طول N نمونه (در پیاده‌سازی تمام الگوریتم‌ها در این مقاله طول فریم ۲۳ میلی‌ثانیه در نظر گرفته شده است که با توجه به فرکانس نمونه‌برداری ۴۴/۸ کیلوهرتز از سیگنال گفتار، طول هر فریم N=۱۰۲۴ نمونه خواهد بود).

گام دوم - حذف مقدار dc فریم  
 گام سوم - ضرب فریم در پنجره زمانی همینگ با رابطه زیر

$$w(n) = \begin{cases} 0.54 - 0.46 \cos\left(\frac{2\pi n}{N-1}\right) & ; 0 \leq n \leq N-1 \\ 0 & ; \text{otherwise} \end{cases} \quad (10)$$

گام چهارم - محاسبه تبدیل فوریه زمان کوتاه N نقطه ای از فریم پنجره گذاری شده  $X(k)$   
 گام پنجم - محاسبه طیف توان  $|X(k)|^2$   
 گام ششم - اعمال بانک فیلترهای مجذور هنینگ روی طیف توان. پاسخ تبدیل گسسته فوریه یک فیلتر هنینگ به صورت زیر است:

$$\psi(k) = \begin{cases} 0.5 - 0.5 \cos\left(\frac{2\pi k}{M}\right) & ; 0 \leq k \leq M \\ 0 & ; \text{Otherwise} \end{cases} \quad (11)$$

گام هفتم - محاسبه لگاریتم خروجی هر فیلتر به عنوان ویژگی

$$E_j = \sum_{k=1}^N |\psi_j(k)|^2 |X(k)|^2 \quad ; j = 0, 1, K, M \quad (12)$$

$$C_j = \log(1 + E_j) \quad (13)$$

که M تعداد فیلترهای بکار رفته در بانک فیلتر می باشد.  
 برای بدست آمدن پارامترهای LHCBC کافیست که پارامترهای LHCB را به حوزه کپستروم انتقال داد. برای این منظور باید از پارامترهای  $C_j$  عکس تبدیل فوریه گسسته گرفت. اما چون  $C_j$  مقادیری حقیقی و متقارن می باشد بجای عکس تبدیل فوریه گسسته می توان از تبدیل کسینوسی گسسته استفاده کرد:

$$c_m = \sum_{j=1}^M \ln(1 + E_j) \cos\left(m \left(\frac{2j-1}{2}\right) \frac{\pi}{M}\right) \quad ; 1 \leq m \leq L \quad (14)$$

در رابطه فوق  $E_j$  مجموع انرژی فیلتر  $i$  ام است و M تعداد کل فیلترها می باشد. L نیز تعداد ضرایب مورد نظر در دامنه کپستروم است. مقدار L همواره کوچک تر یا مساوی M می باشد. اگر  $L=M$  اختیار شود، از تمام ضرایب کپستروم به عنوان پارامترهای باز نمایی استفاده کرده ایم. اما به خاطر اینکه آنالیز کپستروال کانال دو سیگنال مربوط به تحریک و فیلتر لوله صوتی را به صورت جمع در می آورد و چون لگاریتم طیف فیلتر لوله صوتی، تغییرات آرام دارد ولی سیگنال تحریک باعث تغییرات شدید در سیگنال گفتار می شود، لذا ضرایب ابتدایی کپستروم بیانگر طیف فیلتر لوله صوتی می باشند و ضرایب بالایی کپستروم مربوط به سیگنال تحریک هستند [۱۴]. به همین خاطر فقط چند ضریب اول حفظ می شوند و بقیه ضرایب دور ریخته می شوند. با توجه به مطالعاتی که در [۲] بر روی تعداد فیلترهای فرکانسی برای بدست آوردن پارامترهای LHCB و LHCBC انجام گرفته است، تعداد آنها برابر  $M=18$  اختیار شده است. در سیستم های باز شناخت گفتار تعداد ضرایب مورد استفاده کپستروم (L) کمتر یا مساوی ۱۵ اختیار می شود که در این مقاله نیز  $L=15$  اختیار شده است. همانطور که گفته شد آنالیز کپستروال دو سیگنال مربوط به تحریک و ناحیه صوتی را به صورت خطی تفکیک پذیر می کند اما این تفکیک ایده آل نیست و این دو سیگنال بر همدیگر اثر می گذارند و چند ضریب اول کپستروم علاوه بر اطلاعات ناحیه صوتی و اطلاعات گفتاری، حاوی اطلاعات سیگنال تحریک و اطلاعات گوینده نیز می باشد. جهت کاهش این اثرات در چند ضریب اول، از اعمال تابع وزن بر ضرایب

کپستروم استفاده می‌شود. در این تابع وزن، ارزش کمتری به ضرایبی که متأثر از اطلاعات گوینده می‌باشند، داده می‌شود و ضرایبی که فقط حاوی اطلاعات ناحیه صوتی و اطلاعات گفتاری می‌باشند، ارزش بیشتری پیدا می‌کنند. در مراجع متفاوت توابع وزن متفاوتی جهت اعمال بر ضرایب کپستروم تعریف شده‌اند. معمول ترین تابع وزن مناسب برای ضرایب کپسترومی که براساس خروجی بانک فیلترها بدست می‌آید به صورت زیر می‌باشد [۱۵] و [۳۱]:

$$g(m) = 1 + \frac{L}{2} \sin\left(\frac{m\pi}{L}\right) \quad ; \quad 1 \leq m \leq L \quad (15)$$

که در رابطه فوق  $L$  برابر تعداد ضریب کپستروم انتخاب شده است.

ضرایب کپستروم در مقیاس مل (MFCC)، در بسیاری از سیستم‌های بازشناخت گفتار به عنوان پارامترهای بازنمایی مطرح است و این روش استخراج ویژگی بیشتر از سایر روش‌ها مورد استفاده قرار می‌گیرد. الگوریتم این روش همانند الگوریتم روش LHCBC است و همانند آن روش پس از پنجره گذاری و تبدیل فوریه زمان کوتاه، یک بانک فیلترها که نحوه توزیع آن به صورت غیر خطی است، بر طیف توان اعمال می‌گردد و پس از محاسبه لگاریتم مجموع خروجی هر فیلتر با اعمال تبدیل کسینوسی گسسته، ضرایب کپستروم بدست می‌آیند و در آخر با اعمال تابع وزن بر روی ضرایب کپستروم، پارامترهای بازنمایی MFCC حاصل می‌شوند.

یکی از تفاوت اساسی بین پارامترهای MFCC و LHCBC در مقیاس غیرخطی است که برای محور فرکانس و توزیع فیلترها، برگزیده شده است. همانطور که دیده شد، برای پارامترهای LHCBC و LHCBC از مقیاس بارک استفاده شده است و برای پارامترهای MFCC از مقیاس مل استفاده می‌شود. تفاوت مهم دیگر در شکل فیلترهای استفاده شده در بانک فیلترها می‌باشد. در پارامترهای MFCC از فیلترهای مثلثی استفاده می‌شود، حال آنکه در پارامترهای LHCBC از فیلترهای مجذور هنینگ استفاده می‌شود.

یکی دیگر از روش‌های استخراج پارامترهای بازنمایی که بر اساس سیستم شنوایی انسان و الهام گرفته شده از آن می‌باشد، پارامترهای پیشگویی خطی ادراکی (PLP) است که در سال ۱۹۹۰ و توسط هرمانسکی معرفی شده است. در این روش سه ویژگی از فیزیولوژی شنوایی به منظور رسیدن به تخمینی از طیف شنوایی مورد استفاده قرار گرفته است:

۱- رزولوشن طیفی بر مبنای باندهای بحرانی

۲- منحنی بلندی یکنواخت<sup>۱۲</sup> صدا

۳- قانون توان شدت-بلندی<sup>۱۳</sup> صدا

ویژگی اول تبدیل محور فرکانسی از مقیاس هرتز به مقیاس بارک است. این تبدیل که برگرفته از سیستم شنوایی انسان است در روش LHCBC نیز مورد استفاده قرار گرفته است و دو ویژگی بعدی ویژگی‌هایی است که هرمانسکی در این روش مورد توجه قرار داده است و بیان نموده است که باعث بهبود عملکرد سیستم‌های بازشناخت گفتار گشته است [۱۶] و [۱۷]. با در نظر گرفتن سه ویژگی فوق و اعمال آنها بر روی سیگنال گفتار، طیف شنوایی به کمک یک مدل تمام قطب تخمین زده می‌شود. البته برای اعمال این سه ویژگی فیزیولوژی، سیگنال گفتار به کمک تبدیل فوریه زمان کوتاه از حوزه زمان به حوزه فرکانس انتقال می‌یابد و از این جهت این روش جزو روش‌های استخراج پارامترهای بازنمایی است که براساس تبدیل فوریه و در حوزه فرکانس انجام می‌گیرند. الگوریتم استخراج پارامترهای PLP را می‌توان به صورت زیر خلاصه کرد:

گام اول - انتخاب یک فریم از سیگنال به طول  $N$  نمونه

گام دوم - حذف مقدار dc فریم

گام سوم - ضرب فریم در پنجره زمانی همینگ به طول  $N$  نمونه

گام چهارم - محاسبه تبدیل فوریه  $N$  نقطه از فریم پنجره گذاری شده

گام پنجم - محاسبه طیف توان

گام ششم - اعمال بانک فیلتر به شکل دوزنقه نامتقارن بر روی طیف توان به همراه اعمال منحنی بلندی یکنواخت صدا. رابطه

هر فیلتر دوزنقه‌ای نامتقارن در مقیاس بارک به صورت زیر است:



$$\psi_j(z) = \begin{cases} 10^{(z-z_j+0.5)} & ; z \leq z_j - 0.5 \\ 1 & ; z_j - 0.5 < z < z_j + 0.5 \\ 10^{-2.5(z-z_j-0.5)} & ; z \geq z_j + 0.5 \end{cases} \quad (16)$$

که در رابطه فوق  $z$  فرکانس در مقیاس بارک است و  $z_j = 0.9994j$  مرکز فیلتر  $z$ ام می‌باشد. و منحنی بلندی یکنواخت صدا به صورت زیر تقریب زده شده است:

$$E(\omega) = 1.151 \sqrt{\frac{(\omega^2 + 144 \times 10^4) \omega^2}{(\omega^2 + 16 \times 10^4)(\omega^2 + 961 \times 10^4)}} \quad (17)$$

که  $\omega$  فرکانس زاویه‌ای می‌باشد. اگر خروجی فیلتر  $z$ ام پس از اعمال منحنی بلندی یکنواخت را  $F_j$  بنامیم، این مقدار توسط رابطه زیر حساب می‌گردد.

$$F_j = E(\omega_j) \sum_{k=1}^N \psi_j(k) |X(k)|^2 \quad ; j = 0, 1, K, 18 \quad (18)$$

گام هفتم - اعمال تبدیل شدت - بلندی صدا

$$\Phi_j = \sqrt[3]{F_j} \quad (19)$$

گام هشتم - اعمال مدل تمام قطب از مرتبه  $L=5$   
گام نهم - تبدیل ضرایب مدل تمام قطب به ضرایب کپستروم حقیقی

$$\hat{X}(n) = \begin{cases} 0 & ; n < 0 \\ \ln(G) & ; n = 0 \\ a_n + \sum_{k=1}^n \left(\frac{k}{n}\right) x(n) a_{n-k} & ; 0 < n \leq L \end{cases} \quad (20)$$

که  $a_n$  ها ضرایب مدل تمام قطب،  $G$  بهره مدل تمام قطب و  $x(n)$  نیز ضرایب کپستروم مختلط می‌باشند.  
گام دهم - اعمال وزن بر ضرایب کپستروم

$$\gamma(n) = n \quad ; 1 \leq n \leq L \quad (21)$$

## ب - روش‌های طیفی مبتنی بر مدل پیشگویی خطی

در این روش‌ها طیف توان سیگنال بر اساس مدل پیشگویی خطی بدست می‌آید. یک مدل ریاضی مؤثر ولی در عین حال ساده برای فرایند گفتار در نظر گرفته می‌شود که شامل یک منبع تحریک و یک فیلتر مدل کننده لوله صوتی است. این فیلتر لوله صوتی به کمک یک تابع تمام قطب مدل می‌گردد و به کمک ضرایب این مدل که ضرایب پیشگویی خطی<sup>۱۴</sup> (LPC) می‌باشند، فورمنتهای طیف گفتار قابل استخراج هستند. بطور کلی بر اساس این روش طیفی دو نوع پارامتر بازنمایی حاصل می‌گردد. پارامترهای L<sub>PCC</sub><sup>۱۵</sup> که بطور مستقیم ضرایب پیشگویی خطی به ضرایب کپستروم تبدیل می‌شوند و پارامترهای Mel-L<sub>PCC</sub> که از بانک فیلتر جهت استخراج ویژگی استفاده می‌شود [۲۰].

الگوریتم استخراج پارامترهای LPCC به صورت زیر است:  
 گام اول - انتخاب یک فریم از سیگنال به طول N نمونه  
 گام دوم - حذف مقدار dc فریم  
 گام سوم - ضرب فریم در پنجره زمانی همینگ به طول N نمونه  
 گام چهارم - محاسبه ضرایب پیشگویی خطی (LPC) به کمک روش لوینسن - دربین  
 گام پنجم - تبدیل ضرایب پیشگویی خطی (LPC) به ضرایب کپستروم پیشگویی خطی (LPCC)

$$\hat{h}(n) = \begin{cases} 0 & ; n < 0 \\ \ln(\xi) & ; n = 0 \\ a(n) + \sum_{k=1}^{n-1} \left(\frac{k}{n}\right) h(n)h(k)a(n-k) & ; 0 < n \leq M \end{cases} \quad (22)$$

در رابطه فوق، M مرتبه مدل تمام قطب می‌باشد.  $\hat{a}(n)$  ضرایب مدل تمام قطب و  $\xi$  خطای ناشی شده از تخمین گر مرتبه M ام در الگوریتم لوینسن - دربین می‌باشد که معادل بهره مدل تمام قطب می‌باشد. پارامترهای  $\hat{h}(n)$  ضرایب کپستروم مختلط پیشگویی خطی می‌باشند.

گام ششم - انتخاب L ضریب اول کپستروم به عنوان پارامترهای LPCC  
 گام هفتم - ضرب تابع وزن در پارامترهای LPCC

در الگوریتم فوق مرتبه مدل تمام قطب متناسب با فرکانس نمونه برداری انتخاب می‌شود [۱۸] و برای سیگنالی با فرکانس نمونه برداری ۱۶ کیلوهرتز این مقدار معمولاً برابر ۱۴ یا ۱۶ اختیار می‌گردد و در این حالت از ۱۲ ضریب اول کپستروم به عنوان پارامترهای بازنمایی استفاده می‌نماییم. تابع وزنی که برای پارامترهای LPCC استفاده می‌شود مانند پارامترهای MFCC و LHCBC از رابطه (۱۵) بدست می‌آید.

برای بدست آوردن پارامترهای Mel-LPCC از الگوریتمی مانند LHCBC استفاده می‌شود با این تفاوت که در پارامترهای Mel-LPCC برای بدست آوردن طیف سیگنال از مدل پیشگویی خطی استفاده می‌گردد نه از تبدیل فوریه. و پس از بدست آمدن طیف از بانک فیلتر مجذور هنینگ در مقیاس مل استفاده می‌شود و انرژی هر باند به دامنه کپستروم انتقال می‌یابد. اما برای آنکه بتوان مقایسه بهتری بین پارامترهای LHCBC و Mel-LPCC انجام داد، در پارامترهای Mel-LPCC بجای مقیاس مل از مقیاس بارک استفاده می‌کنیم و در ادامه از آن با عنوان پارامترهای Bark-LPCC یاد خواهیم کرد.

## پ - روش‌های طیفی مبتنی بر تبدیل ویولت

برخلاف تبدیل فوریه، تبدیل ویولت گسسته<sup>۱۶</sup> (DWT) برای سیگنالهای غیرایستاد نیز تعریف شده است و همچنین دارای خاصیت چند رزولوشنی<sup>۱۷</sup> می‌باشد [۲۱]. براساس دو خاصیت فوق، در چند سال اخیر گرایش به سمت استفاده از تبدیل ویولت افزایش یافته است [۲۳]-[۲۵].

یکی از روشهای استفاده از تبدیل ویولت، بکار بردن لگاریتم انرژی خروجی فیلترهایی است که تبدیل ویولت ایجاد می‌کند [۲۲]، [۲۳]، [۲۷]. تبدیل ویولت را می‌توان با بانک فیلتر و عمل کاهش نرخ نمونه برداری معادل سازی کرد [۲۸]. لذا تبدیل ویولت یک سیگنال، مشابه اعمال یک بانک فیلتر با باند فرکانسهای متفاوت بر روی سیگنال می‌باشد و ضرایب ویولت حاصله، خروجی هر کدام از فیلترها می‌باشند که کاهش نرخ نمونه برداری نیز بر روی آنها اعمال شده است و به کمک قضیه پارسوال می‌توان انرژی خروجی هر فیلتر را حساب کرد. در این حالت پس از انتخاب یک فریم پنجره گذاری شده از سیگنال با توجه به اینکه در رزولوشن‌های به اندازه کافی بالا نمونه‌های سیگنال خیلی به ضرایب مقیاس نزدیک هستند، می‌توان فرض کرد که نمونه‌های فریم پنجره گذاری شده، ضرایب مقیاس در مقیاس بالا می‌باشند و از روی آنها ضرایب مقیاس و ویولت در مقیاس پایین‌تر را حساب کرد [۲۱]. معمولاً در سیگنالهای گفتار، پس از ۶ مرحله تقسیم فرکانسی و بدست آوردن ضرایب مقیاس در ۶

مرحله، ضرایب مقیاس مرحله بعدی آنقدر کوچک می‌شوند که قابل چشم پوشی می‌باشند. لذا شش باند فرکانسی خواهیم داشت که به کمک قضیه پارسوال انرژی این شش باند حساب شده و لگاریتم آنها به عنوان پارامترهای  $DWTBF^{18}$  مورد استفاده قرار می‌گیرد.

روش دیگر جهت استفاده بهینه از انرژی فیلترهای مختلف، استفاده از تجزیه بسته ویولت<sup>19</sup> (WPD) بجای DWT است. WPD مانند DWT است با این تفاوت که در انتخاب باندهای فرکانسی تنوع وجود دارد. با استفاده از این خاصیت، بانک فیلترهای معادل مقیاس بارک یا مل که الهام گرفته از شنوایی انسان می‌باشند، توسط تبدیل فوریه قابل پیاده‌سازی هستند [24]، [26]، [29]. الگوریتم استخراج پارامترهای بازنمایی پارامترهای  $MFWP^{20}$  را می‌توان بصورت زیر خلاصه کرد:

گام اول - انتخاب یک فریم از سیگنال به طول N نمونه

گام دوم - حذف مقدار dc

گام سوم - ضرب فریم در پنجره زمانی همینگ

گام چهارم - اعمال WPD بر فریم پنجره گذاری شده

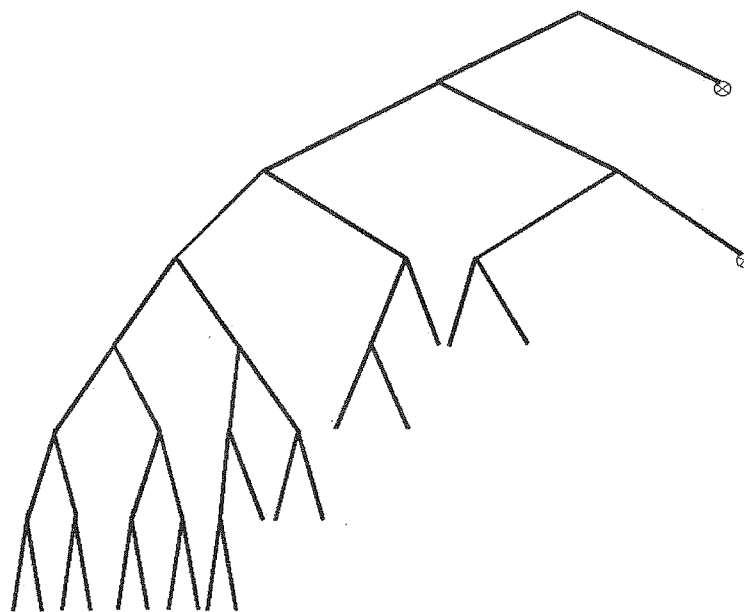
گام پنجم - انتخاب باندهای فرکانسی مشابه بانکهای فیلتر مل، مطابق شکل ۲

گام ششم - محاسبه انرژی داخل هر باند فرکانس با استفاده از قضیه پارسوال

گام هفتم - اعمال تبدیل کسینوس گسسته بر لگاریتم انرژی داخل هر باند فرکانس

گام هشتم - انتخاب پانزده ضریب اول کیستروم به عنوان پارامترهای بازنمایی

گام نهم - اعمال تابع وزن بر ضرایب کیستروم



شکل (۲) درخت انتخابی در WPD جهت استخراج پارامترهای MFWP.

جدول (۱) دسته بندی ۲۹ آوایی زبان فارسی<sup>۲۱</sup>.

/@/, /a/, /e/, /o/, /u/, /i/	واکه‌ها
/y/, /l/, /m/, /n/, /r/	شبه واکه‌ها
/b/, /d/, /q/, /g/, /ʔ/, /p/, /t/, /k/	انفجاری‌ها
/f/, /v/, /s/, /z/, /ʃ/, /*/, /h/, /x/	سایشی‌ها
/j/, /#/	سایشی - انفجاری‌ها

در ادامه و در بخش ۵، نتایج حاصل از ارزیابی عملکرد پارامترهای بازنمایی برای بازشناخت مستقل از گوینده گفتار مورد

## ۴- نتایج آزمایشات

در این بخش به نتایج بدست آمده از پارامترهای رایج بازنمایی که در بخش ۴ به معرفی آنها پرداخته شد و نیز پارامترهای بازنمایی که از بهبود آنها و اعمال هنجارسازی های مختلف حاصل گردید می‌پردازیم. همانطور که در بخش ۲ اشاره شد از دو جمله ۴۰۰ و ۵۰۰ دادگان فارسی‌دات به عنوان دادگان آزمایش استفاده شده است. این دو جمله تقریباً ۲۹ آوای فارسی را دربردارند و عملکرد شبکه هم بر روی کل آواها و هم بطور جداگانه برای دسته آواهای مختلف طبق جدول ۱ بررسی می‌شود. در اولین آزمایشات پارامترهای LHCB با هنجارسازی‌های مختلف مورد بررسی قرار گرفت تا بهترین نوع هنجارسازی مشخص گردد. همانطور که در جدول ۲ مشاهده می‌شود در مورد هنجارسازی عرضی پارامترهای LHCB یکبار نیز از حضور میانگین مولفه‌های یک بردار به عنوان یک پارامتر جدا (پارامتر نوزدهم) استفاده شده است. این پارامتر که تقریبی از انرژی می‌باشد باعث بهبود عملکرد بازشناسی شده است. هنجارسازی‌های طولی به خصوص هنجارسازی طولی نرم دوم با واریانس ۰/۵ برای پارامترهای به هنجار شده نسبت به پارامترهای عرضی از عملکرد بهتری برخوردار هستند. این برتری بطور محسوس در واکه‌ها، شبه‌واکه‌ها و سایشی‌ها مشاهده می‌گردد. نتایج هنجارسازی طولی پارامترها بطور جداگانه برای هر گوینده نشان می‌دهد که هنجارسازی طولی پارامترها بطور جداگانه برای هر فرد نسبت به روش هنجارسازی طولی یکجا برای کل افراد عملکرد بهتری دارد. چرا که در این حالت اطلاعات مربوط به هر گوینده از پارامترهای بازنمایی همان گوینده تا حدی حذف

جدول (۲) مقایسه کارایی روشهای مختلف هنجارسازی در بهبود کیفیت بازشناسی آواها که بر روی پارامترهای بازنمایی LHCB اعمال گردیده‌اند.

درصد صحت بازشناخت فریم روی دادگان آزمون						
نوع بازشناسی		کل	واکه	شبه‌واکه	انفجاری	سایشی - انفجاری
۱	هنجارسازی عرضی پارامترها + LHCB	۷۹/۶۲	۸۳/۴۲	۶۷/۷۹	۶۴/۹۳	۸۲/۳۸
۲	پارامتر نوزدهم + LHCB + هنجارسازی عرضی پارامترها	۸۰/۵۴	۸۴/۳۶	۶۸/۵۱	۶۷/۱۹	۸۳/۰۷
۳	LHCB + هنجارسازی طولی نرم ۱ پارامترها	۸۱/۰۶	۸۴/۱۹	۷۰/۶۱	۶۴/۵۲	۸۵/۲۸
۴	LHCB + هنجارسازی طولی نرم ۲ پارامترها (واریانس = ۱)	۸۲/۱۱	۸۵/۵۸	۷۱/۶۷	۶۶/۹۹	۸۳/۸۴
۵	LHCB + هنجارسازی طولی نرم ۲ پارامترها (واریانس = ۰/۵)	۸۲/۷۶	۸۶/۶۰	۷۱/۳۳	۶۷/۷۴	۸۵/۰۳
۶	هنجارسازی طولی پارامترها + LHCB برای هر فرد جدا	۸۳/۹۱	۸۸/۴۹	۷۲/۳۸	۷۰/۳۵	۸۴/۰۴

می‌شود ولی در دو روش هنجارسازی طولی نرم اول و نرم دوم پارامترهای بازنمایی، اطلاعاتی که به عنوان اطلاعات گوینده از پارامترهای بازنمایی کم می‌شود میانگینی از اطلاعات تمام گویندگان دادگان آزمایش می‌باشد. انتقال انرژی بانک فیلترها به حوزه کپستروم و محاسبه ضرایب کپستروم، یکی دیگر از مراحل پردازش سیگنال گفتار است که در برخی بازنمایی‌ها از جمله MFCC، LHCB و PLP مورد استفاده قرار می‌گیرد. همانطور که اشاره شد در پارامترهای MFCC از مقیاس مل و در پارامترهای LHCB از مقیاس بارک استفاده می‌شود و تفاوت دیگری که در روش استخراج این دو نوع پارامتر بازنمایی وجود دارد شکل بانک فیلتر می‌باشد. در MFCC از فیلترهای مثلثی و در LHCB مانند پارامترهای LHCB از فیلترهای مجذور هنینگ استفاده می‌شود. برای اینکه بتوان یک مقایسه بین اثر مقیاس‌های مل و بارک انجام داد

یکبار نیز پارامترهای MFCC را با استفاده از بانک‌های فیلتر مجذور هنینگ استخراج می‌نماییم. در تعلیم شبکه عصبی با استفاده از هر کدام از سه پارامتر فوق الذکر، شبکه عصبی در کمینه موضعی<sup>۲۲</sup> افتاده بطور کامل تعلیم نمی‌بیند و صحت بازشناسی فریم برای کل آواها در حدود ۷۳٪ بدست می‌آید. شبکه عصبی به دلایل مختلفی در کمینه موضعی قرار می‌گیرد که یکی از آنها می‌تواند تغییرات زیاد دامنه برخی از مؤلفه‌های بردار بازنمایی باشد. با اعمال هنجارسازی طولی نرم دوم پارامترهای بازنمایی مشکل تعلیم شبکه عصبی برای هر سه نوع پارامتر بازنمایی فوق حل می‌شود. در این حالت میزان صحت بازشناسی فریم‌ها روی کل آواها برای پارامترهای به هنجار شده MFCC با بانک فیلترهای مثلثی و مجذور هنینگ به ترتیب برابر ۸۰/۲۶٪ و ۸۰/۹۱٪ بدست آمد. همین صحت بازشناسی برای پارامترهای LHCBC برابر ۸۱/۴۰٪ حاصل شد. این مقادیر نشان می‌دهند که عملکرد پارامترهای LHCBC از پارامترهای MFCC اندکی بهتر می‌باشد. علت این مطلب را تا حدی می‌توان در عملکرد بهتر نوع پنجره‌های طیفی انتخابی در LHCBC یعنی پنجره‌های مجذور هنینگ دانست. همچنین اثر مقیاس‌های مل و بارک تقریباً شبیه به هم می‌باشد. در جدول ۳ صحت بازشناخت ناشی از این پارامترها نشان داده شده است. با مقایسه این جدول و جدول ۲ مشاهده می‌گردد که پارامترهای LHCBC با هنجارسازی طولی نرم دوم نسبت به پارامترهای کپستروم MFCC و LHCBC و با همین هنجارسازی عملکرد بهتری برای شبکه عصبی ایجاد کرده است.

پارامترهای طیفی دیگری که براساس تبدیل فوریه زمان کوتاه بدست می‌آیند پارامترهای PLP می‌باشند. پارامترهای PLP یکی دیگر از پرکاربردترین پارامترهای بازنمایی در سیستم‌های بازشناخت گفتار می‌باشد. هنگامیکه در کنار پارامترهای PLP، تقریبی از انرژی فریم نیز به عنوان پارامتری جدا استفاده می‌شود، صحت بازشناسی فریم روی کل آواها ۳٪ رشد کرده و از مقدار ۷۸/۴۰٪ به ۸۱/۴۶٪ می‌رسد. در جدول ۴ صحت بازشناخت فریم روی دسته آواهای مختلف ناشی از پارامترهای PLP در مقایسه با پارامترهای LHCBC با هنجارسازی های عرضی و طولی نرم دوم آورده شده است. مشاهده می‌شود که پارامترهای PLP با احتساب پارامتری که تقریبی از انرژی فریم می‌باشد، عملکرد مناسبی برای شبکه عصبی مهیا ساخته است.

جدول (۳) نتایج حاصل از پارامترهای LHCBC و MFCC

درصد صحت بازشناخت فریم روی دادگان آزمون						
نوع بازشناسی	کل	واکه	شبه‌واکه	انفجاری	سایشی	سایشی- انفجاری
MFCC + هنجارسازی طولی نرم ۲ پارامترها (واریانس = ۰/۵)، فیلتر مثلثی	۸۰/۲۶	۸۵/۰۶	۶۵/۴۹	۶۲/۳۵	۸۲/۵۳	۶۴/۱۲
MFCC + هنجارسازی طولی نرم ۲ پارامترها (واریانس = ۰/۵)، فیلتر مجذور هنینگ	۸۰/۹۱	۸۵/۴۶	۶۶/۸۸	۶۴/۳۴	۸۳/۰۷	۵۷/۰۶
LHCBC + هنجارسازی طولی نرم ۲ پارامترها (واریانس = ۰/۵)	۸۱/۴۰	۸۵/۵۴	۶۸/۷۸	۶۳/۸۳	۸۴/۵۱	۶۰/۱۷

جدول (۴) مقایسه نتایج حاصل از پارامترهای LHCBC و PLP

درصد صحت بازشناخت فریم روی دادگان آزمون						
نوع بازشناسی	کل	واکه	شبه‌واکه	انفجاری	سایشی	سایشی- انفجاری
PLP + پارامتر انرژی +	۸۱/۴۶	۸۲/۹۶	۷۲/۴۷	۶۸/۰۲	۸۶/۲۳	۷۰/۹۰
PLP	۷۸/۴۰	۸۰/۸۴	۷۰/۲۳	۵۸/۹۶	۸۲/۵۵	۶۲/۷۱
LHCBC + پارامتر نوزدهم + هنجارسازی عرضی پارامترها	۸۰/۵۴	۸۴/۳۶	۶۸/۵۱	۶۷/۱۹	۸۳/۰۷	۶۲/۷۱
LHCBC + هنجارسازی طولی نرم ۲ پارامترها (واریانس = ۰/۵)	۸۲/۷۶	۸۶/۶۰	۷۱/۳۳	۶۷/۲۴	۸۵/۰۳	۵۸/۴۷

در بررسی پارامترهای طیفی مبتنی بر مدل تمام قطب و پیشگویی خطی باید به این نکته دقت داشت که مرتبه مدل متناسب با فرکانس نمونه برداری سیگنال گفتار باید انتخاب گردد. لذا چون دادگان فارس‌دات با فرکانس ۴۴/۱ کیلوهرتز نمونه‌برداری شده‌است و مدل تمام قطب مرتبه بالا تخمین مناسبی از طیف بدست نمی‌دهد، در ابتدا فرکانس نمونه برداری به ۱۶ کیلوهرتز کاهش داده می‌شود. مرتبه مدل تمام قطب برابر ۱۴ و یکبار نیز برابر ۱۶ انتخاب شده است و برای اینکه تغییری در شرایط آزمایش بوجود نیاید، برای این روش‌ها نیز مانند روش‌های قبلی از فیلتر پیش تأکید استفاده نشده است. البته برای اینکه اثر این فیلتر بر روی عملکرد پارامترهای بازنمایی در بازشناخت گفتار مشخص گردد، یکبار نیز برای استخراج پارامترهای LPCC سیگنال گفتار را از یک فیلتر پیش تأکید  $P(z) = 1 - 0.97z^{-1}$  عبور داده‌ایم. در این روش‌ها پس از حصول ضرایب مدل تمام قطب (ضرایب پیشگویی خطی) به کمک الگوریتم لوینسن - دربین، یا بطور مستقیم این ضرایب به ضرایب کپستروم تبدیل می‌شوند (LPCC) و یا از روی مدل تمام قطب طیف توان تخمین زده شده و با استفاده از بانک فیلتر انرژی باندهای فرکانسی مختلف بدست آمده و به حوزه کپستروم انتقال می‌یابند (Bark-LPCC). همانطور که در جدول ۵ مشاهده می‌شود بجز در سایشی - انفجاری‌ها در بقیه موارد وجود فیلتر پیش تأکید باعث بهبود عملکرد پارامترهای LPCC شده است. نتایج جدول ۵ نشان می‌دهد که پارامترهای LPCC عملکرد مناسبی در جهت بازشناخت گفتار نداشته است. برای بهبود این پارامترها از بهره مدل تمام قطب به عنوان تقریبی از انرژی فریم در کنار دیگر پارامترهای LPCC استفاده شد. بهره مدل تمام قطب برای واکدارها با پرلود گام متناسب است. نتایج مندرج در جدول ۵ مؤید این نکته است که اضافه شدن بهره مدل تمام قطب باعث بهبود عملکرد پارامترهای LPCC شده است.

الگوریتم استخراج پارامترهای Bark-LPCC فقط در نوع تحلیل طیفی با روش استخراج پارامترهای LHCBC متفاوت است و همانطور که در جدول ۶ نشان داده شده است، نتایج نسبتاً ضعیفی از این پارامترها نسبت به پارامترهای LHCBC حاصل شده است. این نتایج بیان می‌کند که مدل تمام قطب استفاده شده نمی‌تواند به خوبی تبدیل فوریه زمان کوتاه طیف توان سیگنال را تقریب بزند. صحت بازشناخت فریم روی کل آواها برای پارامترهای Bark-LPCC با هنجارسازی طولی نرم دوم پارامترها برابر ۷۶/۲۰٪ می‌باشد حال آنکه همین صحت بازشناخت برای پارامترهای LHCBC با همین هنجارسازی برابر ۸۱/۴۰٪ بوده است.

جدول (۵) نتایج حاصل از پارامترهای LPCC.

درصد صحت بازشناخت فریم روی دادگان آزمون							
نوع بازشناسی		کل	واکه	شبه‌واکه	انفجاری	سایشی - انفجاری	
۱	LPCC (مرتبه مدل = ۱۴)	۷۵/۶۶	۷۸/۸۶	۶۱/۶۴	۶۱/۶۰	۷۶/۵۹	۶۱/۸۶
۲	فیلتر پیش تأکید + LPCC (مرتبه مدل = ۱۴)	۷۷/۰۰	۸۰/۸۵	۶۵/۸۲	۶۳/۸۰	۷۹/۰۵	۴۷/۷۴
۳	بهره مدل تمام قطب + LPCC (مرتبه مدل = ۱۴)	۷۷/۷۱	۷۹/۷۳	۶۸/۱۵	۶۰/۷۴	۸۳/۳۵	۵۳/۱۱
۴	LPCC (مرتبه مدل = ۱۶)	۷۵/۵۹	۷۷/۷۹	۶۲/۵۶	۶۰/۲۳	۸۱/۳۶	۴۲/۶۶

جدول (۶) نتایج حاصل از پارامترهای Bark-LPCC در مقایسه با LHCBC.

درصد صحت بازشناخت فریم روی دادگان آزمون							
نوع بازشناسی		کل	واکه	شبه‌واکه	انفجاری	سایشی - انفجاری	
۱	بهره مدل تمام قطب + Bark-LPCC + هنجارسازی طولی نرم ۲ پارامترها	۷۶/۲۰	۸۰/۸۰	۶۵/۰۶	۵۶/۲۵	۷۴/۴۱	۴۷/۷۴
۲	LHCBC + هنجارسازی طولی نرم ۲ پارامترها (واریانس = ۰/۵)	۸۱/۴۰	۸۵/۵۴	۶۸/۷۸	۶۳/۸۳	۸۴/۵۱	۶۰/۱۷

پارامترهای مبتنی بر تبدیل گسسته ویولت<sup>۲۳</sup> (DWTBF) و همچنین پارامترهای بازنمایی به کمک تجزیه بسته ویولت و بانک‌های مشابه مقیاس مل<sup>۲۴</sup> (MFWP) پارامترهایی هستند که در سال‌های اخیر برای بازشناخت گفتار مورد توجه قرار گرفته‌اند. موارد استفاده این نوع پارامترها بیشتر در بیان ویژگی‌ها و طبقه‌بندی گروهی خاصی از آواها مانند انفجاری‌ها بوده است و یا برای بازشناسی‌های محدود مانند بازشناسی ۹ رقم از یکدیگر و یا بازشناسی آواهای منقطع صورت گرفته است. در این مقاله عملکرد این پارامترها را برای بازشناسی گفتار پیوسته بررسی کرده‌ایم و با اصلاح الگوریتم استخراج این پارامترها باعث بهبود عملکرد آن شده‌ایم.

در تجزیه بسته ویولت برخلاف تبدیل ویولت گسسته انتخاب باندهای فرکانسی در اختیار ما می‌باشد و معمولاً باندهایی انتخاب می‌شوند که موجب تقسیم‌بندی محور فرکانسی مشابه مقیاس مل یا بارک گردند. در روش MFWP که در این مقاله مورد بررسی قرار گرفته است، تجزیه بسته ویولت تا هفت مرحله انجام می‌گیرد و سپس باندهای فرکانسی مشابه مقیاس مل انتخاب می‌گردد و سپس مانند روش استخراج پارامترهای MFCC، انرژی هر باند محاسبه شده و به ضرایب کپسترم تبدیل می‌گردند. صحت بازشناخت حاصل از پارامترهای MFWP با پارامترهای MFCC در جدول ۷ مقایسه شده است. نتایج نشان می‌دهد که عملکرد پارامترهای MFWP بسیار ضعیفتر می‌باشد. در پارامترهای DWTBF نیز تبدیل ویولت تا مقیاس شش انجام گرفته و سپس لگاریتم انرژی هر باند فرکانسی به عنوان پارامتر بازنمایی مورد استفاده قرار گرفته است. نتایج حاصل از این پارامترها نیز چندان قابل توجه نیست. در روش اصلاح شده که در این مقاله برای استخراج پارامترهای بازنمایی به کمک تبدیل ویولت گسسته پیشنهاد شده است، بجای آنکه تبدیل ویولت گسسته را روی هر فریم اعمال کنیم، از کل سیگنال گفتار تبدیل ویولت گسسته می‌گیریم و سپس ضرایب ویولت متناسب با هر فریم از سیگنال را جدا می‌کنیم. در واقع فریم‌بندی را روی ضرایب ویولت انجام می‌دهیم نه روی سیگنال گفتار. سپس لگاریتم انرژی هر باند فرکانسی را به عنوان پارامترهای بازنمایی مورد استفاده قرار می‌دهیم. نتایج حاصل از آزمایشات نشان می‌دهد که این پارامترهای بازنمایی اصلاح شده (MDWTBF<sup>۲۵</sup>) نسبت به پارامترهای DWTBF عملکرد بهتری دارند و صحت بازشناخت فریم روی کل آواها از ۷۳/۷۹٪ به ۷۷/۸۶٪ رسیده است. در جدول ۷ نتایج حاصل از این دو نوع پارامتر نیز مورد مقایسه قرار گرفته‌اند. علت این برتری رفع اثر

جدول (۷) نتایج حاصل از پارامترهای DWTBF و MFWP

درصد صحت بازشناخت فریم روی دادگان آزمون							
نوع بازشناسی		کل	واکه	شبه‌واکه	انفجاری	سایشی - انفجاری	
۱	هنجارسازی طولی نرم ۲ + DWTBF (واریانس = ۰/۵)	۷۳/۷۹	۷۹/۲۹	۵۹/۵۴	۵۵/۶۸	۷۴/۸۷	۵۷/۹۴
۲	هنجارسازی طولی نرم ۲ + MDWTBF (واریانس = ۰/۵)	۷۷/۸۶	۷۹/۶۸	۶۵/۹۱	۶۱/۸۷	۸۲/۸۵	۷۱/۴۷
۳	هنجارسازی طولی نرم ۲ + MFWP (واریانس = ۰/۵)	۷۱/۵۰	۷۶/۳۸	۵۵/۱۰	۴۸/۵۹	۷۱/۵۳	۵۶/۵۰
۴	MFCC + هنجارسازی طولی نرم ۲ پارامترها (واریانس = ۰/۵ ، فیلتر مجذور هنینگ)	۸۰/۹۱	۸۵/۴۶	۶۶/۸۸	۶۴/۳۴	۸۳/۰۷	۵۷/۰۶

جدول (۸) نتایج حاصل از ترکیب پارامترهای LHCB با DWTBF و MDWTBF.

درصد صحت بازشناخت فریم روی دادگان آزمون							
نوع بازشناسی		کل	واکه	شبه‌واکه	انفجاری	سایشی - انفجاری	
۱	هنجارسازی طولی پارامترها + LHCB	۸۱/۰۶	۸۴/۱۹	۷۰/۶۱	۶۴/۵۲	۸۵/۲۸	۵۴/۸۰
۲	LHCB_DWTBF + هنجارسازی طولی پارامترها	۸۲/۵۲	۸۶/۴۴	۷۱/۶۹	۶۶/۵۴	۸۴/۱۹	۶۷/۲۳
۳	LHCB_MDWTBF + هنجارسازی طولی پارامترها	۸۳/۲۱	۸۵/۷۹	۷۱/۶۹	۷۰/۳۵	۸۷/۹۶	۷۴/۰۱

مخرب پنجره‌گذاری بر طیف توان می‌باشد.

با ترکیب پارامترهای LHCب و MDWTBF یا DWTBF و در کنارهم قراردادن آنها در یک بردار بازنمایی جدید، صحت بازشناخت نسبت به پارامترهای LHCب به تنهایی، رشد نشان داده است و صحت بازشناخت فریم روی کل آواها برای پارامترهای LHCب با هنجارسازی طولی نرم اول که برابر با  $0.81/0.6$  می‌باشد در ترکیب با پارامترهای MDWTBF و DWTBF به ترتیب به مقادیر  $0.83/0.21$  و  $0.82/0.52$  رسیده است. پارامترهای بازنمایی مبتنی بر تبدیل ویولت دارای اطلاعات مفیدی از تغییرات سریع سیگنال گفتار می‌باشند و ترکیب این اطلاعات با اطلاعات پارامترهای LHCب باعث می‌شود تا شبکه عصبی بهتر بتواند آواهای مختلف را طبقه‌بندی کند. نتایج حاصل از این پارامترها برای دسته آواهای مختلف در جدول ۸ مورد مقایسه قرار گرفته‌اند.

پارامترهای دیگری که در این مقاله مورد بررسی قرار گرفتند ادغام برخی از پارامترها با ضرایب دلتای مربوطه می‌باشد. در سیستم‌های بازشناخت گفتار که از مدل‌های استاتیک بازشناخت گفتار استفاده می‌کنند، از ضرایب دلتا جهت استفاده از اطلاعات دینامیکی سیگنال گفتار استفاده می‌شود. اما شبکه عصبی جلوسو با تأخیر زمانی به علت اینکه در ورودی به ۱۴ فریم متوالی نگاه می‌کند، مدل دینامیک‌های ورودی را می‌آموزد و اطلاعات دینامیکی گسترده‌تری نسبت به ضرایب دلتا از سیگنال گفتار استخراج می‌کند. به عبارت دیگر ضرایب دلتا از نظر اطلاعات دینامیکی نمی‌تواند اطلاعات جدیدی در اختیار شبکه عصبی قرار دهد. اما چون اطلاعات گوینده در فریم‌های نزدیک به هم به آرامی تغییر می‌کند، در ضرایب دلتا که پارامترهای دو فریم نزدیک به هم از یکدیگر کم می‌شوند، تاحدودی اطلاعات گوینده حذف می‌گردد. و لذا با اضافه شدن

جدول (۹) نتایج حاصل از اضافه کردن ضرایب دلتا.

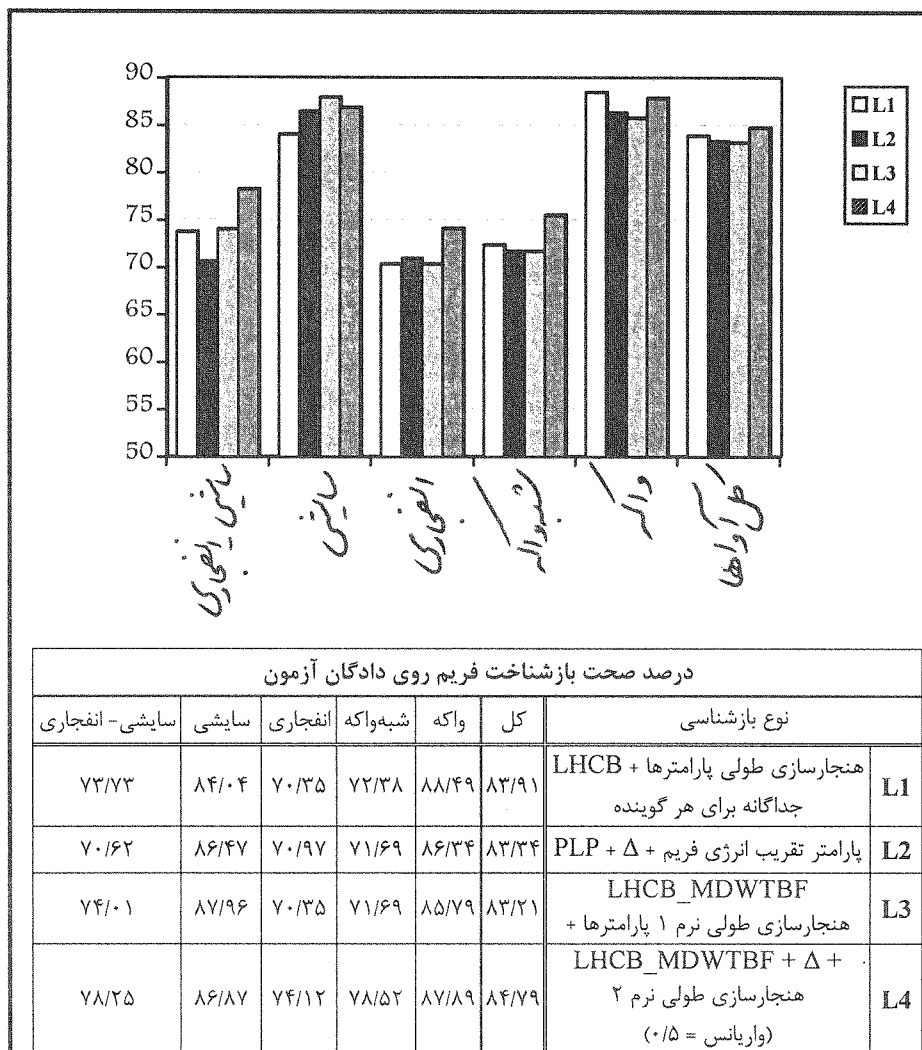
درصد صحت بازشناخت فریم روی دادگان آزمون						
نوع بازشناسی	کل	واکه	شبه‌واکه	انفجاری	سایشی	سایشی- انفجاری
۱ هنجارسازی عرضی پارامترها + LHCب + $\Delta$	۸۲/۹۸	۸۵/۴۹	۷۲/۳۴	۶۹/۶۶	۸۶/۷۵	۷۴/۰۱
۲ هنجارسازی عرضی پارامترها + LHCب	۸۰/۵۴	۸۴/۳۶	۶۸/۵۱	۶۷/۱۹	۸۳/۰۷	۶۲/۷۱
۳ پارامتر تقریب انرژی فریم + $\Delta$ + PLP	۸۲/۳۴	۸۶/۳۴	۷۱/۶۹	۷۰/۹۷	۸۶/۴۷	۷۰/۶۲
۴ پارامتر تقریب انرژی فریم + PLP	۸۱/۴۵	۸۲/۹۶	۷۲/۴۷	۶۸/۰۲	۸۶/۲۳	۷۰/۹۰
۵ DWTBF + $\Delta$	۷۴/۹۱	۷۸/۳۴	۶۰/۰۵	۵۷/۴۱	۷۹/۸۵	۵۹/۰۴
۶ DWTBF	۷۲/۳۱	۷۷/۱۰	۵۶/۹۶	۵۱/۳۷	۷۹/۵۰	۴۴/۰۷
۷ MDWTBF + $\Delta$	۷۸/۵۵	۸۰/۹۸	۶۸/۳۴	۶۲/۳۵	۸۲/۸۰	۶۱/۳۰
۸ MDWTBF	۷۷/۸۶	۷۹/۶۸	۶۵/۹۱	۶۱/۸۷	۸۲/۸۵	۷۱/۴۷
۹ LHCب_ MDWTBF + $\Delta$ + هنجارسازی طولی نرم ۲ (واریانس = ۰/۵)	۸۴/۷۹	۸۷/۸۹	۷۸/۵۲	۷۴/۱۲	۸۶/۸۷	۷۸/۲۵

ضرایب دلتا به پارامترهای بازنمایی، شبکه عصبی قادر خواهد بود که با دقت بیشتری آواها را طبقه‌بندی کند. آزمایش بر روی چهار نوع پارامتر بازنمایی LHCب، PLP، DWTBF و MDWTBF نشان داد که اضافه شدن ضرایب دلتا به پارامترهای بازنمایی حدود ۲٪ صحت بازشناسی فریم روی کل آواها را افزایش می‌دهد که نتایج آزمایشات مربوطه در جدول ۹ درج شده است. همانطور که در ردیف آخر این جدول دیده می‌شود در کنار هم قراردادن پارامترهای LHCب و MDWTBF در یک بردار بازنمایی و استفاده از هنجارسازی طولی نرم ۲ (واریانس = ۰/۵) و همچنین اضافه کردن ضرایب دلتا باعث رشد قابل توجه درصد صحت بازشناخت کل فریمها تا میزان ۸۴/۷۹٪ شده است.

در بین پارامترهای بازنمایی مختلف که در این مقاله مورد بررسی قرار گرفتند، پارامترهای LHCب با هنجارسازی طولی



پارامترها بطور جداگانه برای هر گوینده با صحت بازشناخت فریم ۸۳/۹۱٪، پارامترهای PLP به همراه ضرایب دلتا با صحت بازشناخت فریم برابر ۸۳/۳۴٪، پارامترهای بازنمایی تشکیل شده از پارامترهای LHCب و MDWTBF با هنجارسازی طولی نرم اول با صحت بازشناخت فریم برابر ۸۳/۲۱٪ و پارامترهای بازنمایی تشکیل شده از پارامترهای LHCب و MDWTBF با هنجارسازی طولی نرم دوم به همراه ضرایب دلتا با صحت بازشناخت فریم برابر ۸۴/۷۹٪ باعث بهترین عملکرد شبکه عصبی جلوسو با تأخیر زمانی در بازشناخت فریم شده‌اند. نتایج حاصل از این چهار بازنمایی برتر در شکل ۳ درج شده است.



شکل (۳) نتایج حاصل از چهار بازنمایی برتر.

## ۵- جمع‌بندی و نتیجه‌گیری

در این مقاله روش‌های رایج استخراج پارامترهای بازنمایی مورد بررسی و مقایسه قرار گرفت و با اعمال تغییراتی در الگوریتم برخی از آنها و بهبود پارامترهای بازنمایی و همچنین با استفاده از روش‌های هنجارسازی متفاوت، مناسب‌ترین پارامترهای بازنمایی در بکارگیری برای بازشناسی مستقل از گوینده گفتار عاری از نویز با استفاده از شبکه عصبی جلوسو با تأخیر زمانی معرفی گردید. در بین پارامترهای مورد بررسی، پارامترهای LHCب با هنجارسازی طولی پارامترها بطور جداگانه برای هر گوینده با صحت بازشناخت فریم ۸۳/۹۱٪، پارامترهای PLP به همراه ضرایب دلتا با صحت بازشناخت فریم برابر ۸۳/۳۴٪ و پارامترهای بازنمایی تشکیل شده از پارامترهای LHCب و MDWTBF با هنجارسازی طولی نرم اول با صحت بازشناخت فریم برابر ۸۳/۲۱٪ باعث بهترین عملکرد شبکه عصبی جلوسو با تأخیر زمانی در بازشناخت فریم‌ها شده‌اند.

در آزمایشات انجام شده مشاهده شد که استفاده از تقریبی از انرژی فریم به عنوان یک پارامتر جدا و در کنار سایر پارامترها

باعث بهبود عملکرد شبکه جلوسوی زمانی در طبقه‌بندی آواها می‌گردد. علت این رشد به این خاطر است که وقتی تقریبی از انرژی فریم بطور مجزا توسط شبکه عصبی مورد توجه قرار می‌گیرد به عنوان یک عامل مهم برای تفکیک فریم‌های کم انرژی و پر انرژی عمل می‌کند.

نتیجه دیگر اثر مثبت فیلتر پیش تأکید در عملکرد پارامترهای بازنمایی می‌باشد. آزمایش انجام شده بر روی پارامترهای LPCC مؤید این نکته است که وجود فیلتر پیش تأکید رشدی حدود ۱/۵٪ در بازشناخت فریم روی کل آواها دارد. از بین هنجارسازی‌های بررسی شده هنجارسازی طولی پارامترها با پارامترهای به هنجار کننده جداگانه برای هر گوینده و نیز هنجارسازی طولی نرم دوم پارامترها با واریانس ۰/۵ برای پارامترهای به هنجار شده، بهترین عملکرد را برای شبکه عصبی به همراه داشته‌اند.

## ۶- قدردانی

این تحقیق با حمایت و همکاری پژوهشکده پردازش هوشمند علائم به انجام رسیده است.

## زیر نویسها

- 1- Recurrent Neural Networks
- 2- Deng
- 3- Holton
- 4- Normal
- 5- Pitch Frequency
- 6- Mel Scale
- 7- Bark Scale
- 8- Logarithm of square Hanning Critical Band filter banks
- 9- Logarithm of square Hanning Critical Band filter banks Cepstral
- 10- Mel Frequency Cepstral Coefficients
- 11- Perceptual Linear Prediction
- 12- Equal Loudness
- 13- Intensity-Loudness
- 14- Linear Prediction Coefficients
- 15- Linear Prediction Coefficients Cepstral
- 16- Discrete Wavelet Transform
- 17- Multi Resolution
- 18- Discrete Wavelet Transform Based Feature
- 19- Wavelet Packet Decompose
- 20- Mel Frequency Wavelet Packet

۲۱- برخی از این علائم که استاندارد نمی‌باشند عبارتند از:

/\* برای صدای "ژ"، /#/ برای صدای "چ"، /\$/ برای صدای "ش"، /X/ برای صدای "خ"، /?/ برای پست چاکنای(حمزه)، /q/ برای صدای "ق"، /@/ برای صدای "آ"

- 22- Local minimum
- 23- Discret Wavelet Transform Based Features
- 24- Mel Frequency Wavelet Packet
- 25- Modified DWTBF

## مراجع

- [۱] س.ع. سیدصالحی، هنجارسازی خطی و غیرخطی پارامترهای بازنمایی در بازشناسی گفتار، دلدست تدوین، ۱۳۸۲
- [۲] س.ع. سیدصالحی، بازشناخت گفتار پیوسته فارسی با استفاده از مدل عملکردی مغز انسان در درک گفتار، پایان‌نامه دکتری، دانشگاه تربیت مدرس، دانشکده فنی و مهندسی، بهمن ۱۳۷۴

- [3] T. Robinson, J. Holdsworth, R. Patterson, F. Fallside, A Comparison of Preprocessors for The Cambridge Recurrent Error Propagation Network Speech Recognition System., Proceedings of ICSLP, pp. 1033-1036, (1990)
- [4] L.I. Deng, Processing of Acoustic Signals in a Cochlear Model Incorporating Laterally Coupled Suppressive Elements., Journal of Neural Networks, Vol. 5, (1992)
- [5] T. Holton, S.D. Love, S.P. Gill, Formant and Pitch-Pulse Detection Using Models of Auditory Signal Processing., Proceeding of ICSLP'92, pp. 81-84, (1992)
- [6] B. Chigier, H.C. Leung, The Effects of Signal Representations, Phonetic Classification Techniques, and The Telephone Network., Proceedings of ICSLP (1992)
- [7] R. Chengalvarayan, Robust Energy Normalization Using Speech/Nonspeech Discriminant for German Connected Digital Recognition., Proceedings of Eurospeech'99, Budapest, Hungary, (1999)
- [8] M. Naito, L. Deng, Y. Sagisaka, Model-Based Speaker Normalization Methods for Speech Recognition., Proceedings of Eurospeech'99, Budapest, Hungary, (1999)
- [9] L. Welling, R. Heab-Umbach, X. Aubert, N. Haberland, A Study on Speaker Normalization Using Vocal Tract Normalization and Speaker Adaptive Training., Proceedings of ICASSP, pp. 797-800, (1998)
- [10] C. Furlanello, D. Ggiuliani, E. Trentin, D. Falavigna, Application of Generalized Radial Basis Functions in Speaker Normalization and Identification., IEEE International Symposium on Circuits and Systems, Vol. 3, pp. 1704-1707, (1995)
- [11] S.A. Seyyed Salehi, A Neural Network Speech Recognition Based on the Both Acoustic Steady Portions and Transitions., Proceedings of ICASP, Beijing (2000)
- [12] D.E. Olsan, J.C. Harris, M.H. Capps, Computerized Polygraph Scoring System., Journal of Forensic Sciences (1996)
- [13] M. Bijankhan and et. al., FARSDAT – The Speech Database of Farsi Spoken Language., Proc. SST-94, Perth, pp. 826-831, (1994)
- [14] C. Becchetti, L.P. Ricotti, Speech Recognition, John Wiley & Sons., LTD, (2000)
- [15] J.H.L. Hansen, S. Bou-Ghazale, A Computer Study of Traditional and Newly Proposed Features for Recognition of Speech Under Stress., IEEE Transaction on Speech and Audio Processing. Vol. 8, No. 4, pp. 429-442, July (2000)
- [16] H. Hermansky, Perceptual Linear Predictive (PLP) Analysis of Speech., Journal of Acoustic Society of America, Vol. 87, No. 4, pp. 1738-1752, April (1990)
- [17] J.C. Junqua, H. Wakita, H. Hermansky, Evaluation and Optimization of Perceptually Based ASR Front-End., IEEE Transactions on Speech and Audio Processing, Vol. 1, No. 1, pp. 39-48, January (1993)
- [18] J.R. Deller, J.G. Proakis, J.H.L. Hansen, Discrete-Time Processing of Speech Signals., Mcmillan Publishing Company, (1993)
- [19] Q. Li, F.K. Soong, O. Siohan, An Auditory System-Based Feature for Robust Speech Recognition., Proceedings of Eurospeech (2001)
- [20] H. Matsumoto, M. Moroto, Evaluation of MEL-LPC Cepstrum in a Large Vocabulary Continuous Speech Recognition., Proceedings of ICSLP (2000)
- [21] C.S. Burrus, R.A. Gopinath, H. Guo, Introduction to Wavelets and Wavelet Transforms., Prentice Hall Inc., (1998)
- [22] O. Farooq, S. Datta, Speaker Independent Phoneme Recognition by MLP Using Wavelet Features., Proceedings of ICSLP (2000)
- [23] O. Farooq, S. Datta, Dynamic Feature Extraction by Wavelet Analysis., Proceedings of ICSLP (2000)
- [24] O. Farooq, S. Datta, Mel Filter-Like Admissible Wavelet Packet Structure for Speech Recognition., IEEE Signal Processing Letters, Vol. 8, No. 7, pp. 196-198, July (2001)
- [25] J.N. Gowdy, Z. Tufekci, Mel-Scaled Discrete Wavelet Coefficients for Speech Recognition., Proceedings of ICASSP (2000)
- [26] E. Lukasik, Wavelet Packets Based Features Selection for Voiceless Plosives Classification., IEEE Transactions on Speech and Audio Processing, pp. 689-692, (1999)
- [27] C.J. Long, S. Datta, Wavelet Based Feature Extracion for Phoneme Recognition., Proceedings of ICSLP (1996)
- [28] P.P. Vaidanathan, Multirate Systems and Filter Banks, Prentice Hall Inc., (1993)
- [29] I. Cohen, Enhancement of Speech Using Bark-Scaled Wavelet Packet Decomposition., Proceedings of Eurospeech (2001)
- [30] B. Davis, P. Mermelstein, Comparison of Parametric Representations for Monosyllabic Word Recognition in Continuously Spoken Sentences., IEEE Transactions on Audio and Speech and Signal Processing, pp. 357-366, July (1987)
- [31] B. H. Juang, L. Rabiner, J. Wilpon, On the Use of Bandpass Lifterin in Speech Recognition., IEEE Transactions on Audio and Speech and Signal Processing, pp. 947-954, August (1980)