

# بهبود بازشناسی گفتار پیوسته فارسی با استفاده از خصوصیات پروزودیک گفتار در سطح واج‌ها

کارولوکس  
استاد

دانشکده فنی و مهندسی، دانشگاه تهران

سیدمحمدرضا هاشمی گلپایگانی  
استاد

دانشکده مهندسی پزشکی، دانشگاه صنعتی امیرکبیر

فرشاد الماس گنج  
استادیار

دانشکده مهندسی پزشکی، دانشگاه صنعتی امیرکبیر

محمود بی جن خان  
استادیار

دانشکده ادبیات و علوم انسانی، دانشگاه تهران

سیدعلی سیدصالحی  
استادیار

دانشکده مهندسی پزشکی، دانشگاه صنعتی امیرکبیر

## چکیده

در بازشناسی خودکار واج‌های گفتار پیوسته فارسی مشکلاتی مشاهده می‌شوند که بسیاری از آنها از ماهیت سیگنال گفتار ناشی می‌شوند و وابستگی زیادی به روش بازشناسی واج ندارند. در این گزارش به شرح روش‌هایی پرداخته ایم که برای رفع تعدادی از این مشکلات فراهم آمده‌اند. ویژگی‌هایی که برای این منظور به کار رفته‌اند از نوع ویژگی‌های پروزودیک گفتار در سطح واج‌ها می‌باشند. افزودن این روش‌ها به روش عادی ما در بازشناسی واج‌ها که از ترکیب شبکه‌های عصبی و قواعد آوایی زبان فارسی شکل گرفته است، باعث گردید که افزایشی بیش از ۱/۵ درصد در صحت بازشناسی واج‌ها برای ۱۴۰ جمله آزمایشی که از داده‌گان «فارس دات» برداشته شده‌اند، داشته باشیم و در مجموع به درصد صحت بازشناسی واج ۸/۷۴ درصد برای این ۱۴۰ جمله برسیم.

## *Better Phoneme Recognition by Using Phoneme's Prosodic Features*

F. Almasganj  
Assistant Professor  
Amirkabir University

C. Lucus  
Professor  
Tehran University

S. M. R. Hashemi Golpaygani  
Professor  
Amirkabir University

M. Bijankhan  
Assistant Professor  
Tehran University

S. A. S. Salehi  
Assistant Professor  
Amirkabir University

### Abstract

*Many of the errors that occur in automatic phoneme recognition of continuous speech recognition are arisen from the complicated nature of speech signal such that the known phoneme recognition approaches can not handle them. Phoneme prosodic features can be used to alter some of these errors. In this paper, we report 1.5% decreasing of phoneme recognition error rate by using some prosodic features on the level of phonemes, such as energy, aspiration and lips protrusion. This approach has evaluated on a set of 140 test sentences that are taken from "FARSDAT" database, yeilding 74.8% phoneme recongnition accuracy rate.*

### مقدمه

در بازشناسی واج‌های گفتار پیوسته فارسی مشکلاتی رخ می‌دهند که روش‌های عمده بازشناسی در حل برخی از آنها دچار مشکل می‌شوند. این نوع مشکلات که از ماهیت سیگنال گفتار ناشی می‌شوند، خوب است به طور جداگانه و به عنوان مسائل خاص بازشناسی واج‌ها مورد مطالعه قرار گیرند و راه حل مخصوص هر کدام مهیا گردد.

آنچه که در این گزارش به آن پرداخته‌ایم، معرفی چند مورد مهم از این نوع مشکلات و ارائه راه‌حل‌های مناسب برای آنهاست. موارد مورد بحث با استفاده از ویژگی‌هایی که از انواع ویژگی‌های پروژودی در سطح واج‌ها و خصوصیات اسپکتروگرام<sup>۱</sup> سیگنال تشکیل شده‌اند، مورد بررسی قرار گرفته و راه‌حل‌هایی جهت بازشناسی دقیق‌تر آنها عنوان گردیده‌اند.

ویژگی‌های پروژودیک استفاده شده از دو نوع ویژگی‌های پروژودیک کلاسیک و غیر کلاسیک می‌باشند. منحنی فرکانس پایه شدت سیگنال صوتی در کنار اثراتی که مربوط به دو عامل سرعت و ریتم هستند، در مجموع به عنوان خصوصیات پروژودیک کلاسیک گفتار [۱] شناخته می‌شوند. غیر از این تعریف، نوع گسترده‌تری از خصوصیات پروژودیک در سطح واج‌ها تعریف می‌شوند که آنها را با عنوان ویژگی‌های پروژودیک غیر کلاسیک در سطح واج‌ها نام می‌بریم. این تعریف فراگیر پروژودی در سطح واج‌ها [۲] به صورت زیر عنوان می‌شود:

پروژودی در سطح واج‌ها مجموعه‌ای از وضعیت‌های دستگاه تولید گفتار است که می‌تواند به صورت لایه‌ای مستقل بر روی لایه واج‌ها اثر گذارد.

هر واج از تعدادی وضعیت تولیدی دستگاه گفتار تشکیل شده است که اینها می‌توانند بر روی واج‌های مجاور خود اثر گذارند. با توجه به این تعریف از پروژودی، هر کدام از این وضعیت‌های تولیدی که بر روی واج‌های قبلی یا بعدی خود اثر می‌گذارند، به عنوان یک عامل پروژودیک تلقی می‌شوند. به این ترتیب لایه پروژودی شامل ویژگی‌هایی است که از لحاظ بازشناسی واج‌ها حشو هستند. برای مثال واج /p/ دارای ویژگی‌هایی است که در هر حالتی از تولید /p/ موجود

هستند و نشان می‌دهند که واج تولیدی /p/ است. ولی ویژگی‌های دیگری هستند که بنا بر موقعیت قرار گرفتن /p/ در گفتار تولید شده و یا حذف می‌گردند، مثل دمش<sup>۲</sup>. این ویژگی‌ها تأثیری در /p/ بودن این واج ندارند، بلکه حاوی اطلاعات اضافه‌ای می‌باشند. این نوع ویژگی‌ها را می‌توان در لایه پروژودی در سطح واج‌ها قرار داد. در این گزارش روش‌هایی ارائه شده‌اند که از دو نوع ویژگی دمش و گردی لب<sup>۴</sup> در کنار ویژگی پروژودیک طول زمانی و خصوصیات اسپکتروگرام سیگنال گفتار در جهت افزایش درصد صحت بازشناسی واج‌ها استفاده می‌کنند.

برای ارزیابی روش‌های عنوان شده از یک داده گان تست ۱۴۰ جمله‌ای از ۱۴ گوینده مختلف، برداشته شده از داده گان «فارس دات» [۳] استفاده نموده‌ایم. داده‌های یادگیری مورد لزوم برای بخش‌هایی که نیاز به اطلاعات یادگیری دارند، نیز از «فارس دات» برداشته شده‌اند. برای اندازه‌گیری درصد صحت بازشناسی واجها از برنامه استاندارد "NIST" استفاده نموده‌ایم که مخفف عبارت "National Institute of Standard and Technology" می‌باشد.

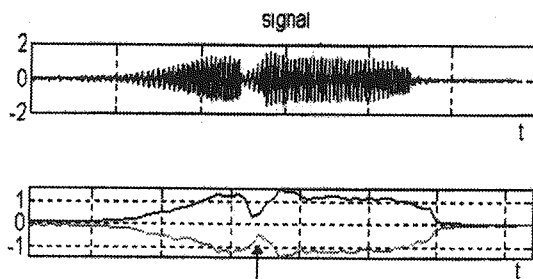
## ۲- روش کار

### ۲-۱- جداسازی واژه‌های به هم چسبیده

در فرایند بازشناسی واج‌های گفتار پیوسته فارسی به مواردی برخورد می‌کنیم که نشانگر واژه‌های طولانی می‌باشند. تعدادی از این موارد مربوط به حالتی می‌شوند که همخوان بین دو واژه (معمولاً /r/ یا /h/) به علت کوتاه یا ناسوده تلفظ شدن، حذف شده است و دو واژه به صورت یک واژه طولانی بازشناسی شده‌اند. این موارد را می‌توان با مقایسه کشش واژه بازشناسی شده با میانگین کشش آن واژه مشخص نمود.

واحدی که برای اندازه‌گیری سرعت تولید گفتار پیوسته بکار برده‌ایم، برحسب تعداد هجاهای تولیدی در یک ثانیه می‌باشد. متوسط سرعت تولید جملات در دادگان «فارس دات» برابر ۴/۶۳۷ هجا بر ثانیه می‌باشد، [۴] که آن را جهت نرمالیزه کردن سرعت جملات مختلف به کار می‌بریم. برای مثال اگر جمله‌ای سریعتر از این حد بیان شده باشد، کشش واج‌های آن را به نسبتی افزایش می‌دهیم که سرعت جمله مورد نظر به میانگین ذکر شده برسد.

مرز اینها را در محل کاهش دامنه سیگنال قرار می دهیم.



شکل (۱) سیگنال صوتی کلمه «ورود» به همراه مرزهای بالا و پایین دامنه آن.

## ۲-۲- دممش

یکی از تفاوت های مهم واج های /k/, /t/, /p/ به ترتیب با واج های /g/, /d/, /b/ در وجود ویژگی «دمش» در گروه اول است. وقوع دممش در گروه اول که همخوان هایی بی واکه<sup>۵</sup> هستند، یک عامل متمایز کننده این دو گروه همخوان از یکدیگر است و می تواند در فرایند اصلاح واج ها (در مواقعی که بین همخوانی از گروه اول یا دوم، تردید وجود دارد) مؤثر واقع شود.

برای بازشناسی «دممش» در این واج ها از یک شبکه عصبی چهار لایه پس انتشار خطا [۵] استفاده کرده ایم. باتوجه به اینکه «دممش» در بخش انتهایی همخوان های انفجاری بی واک رخ می دهد، بازنمایی های مربوط به ۷ فریم اخر واج های مورد نظر را در ورودی شبکه عصبی قرار داده ایم.

هر فریم با ۱۸ پارامتر معرفی می شود. این پارامترها توسط مجموعه ای از فیلترها که با مقیاس «بارک»<sup>۶</sup> تنظیم شده اند، تولید می گردند [۶]. برای استخراج این بازنمایی، از فریم هایی ۲۳ میلی ثانیه ای سیگنال بعد از اعمال پنجره همینگ<sup>۷</sup> تبدیل فوریه گسسته می گیریم. سپس مجموعه ای از ۱۸ فیلتر که در مقیاس «بارک» [۷] مراکز آنها در  $Z_k = k$  قرار گرفته است به تبدیل فوریه هر فریم اعمال می نمایم تا ۱۸ پارامتر مورد نظر به دست آید. تبدیل محور فرکانس از هرتز «بارک» به کمک رابطه زیر انجام می شود:

$$Z = 6 \ln [f/600 + \sqrt{(f/600)^2 + 1}] \quad (1)$$

شکل فیلترها در حوزه فرکانس با مقیاس «بارک» به صورت مجذور پنجره «هنینگ»<sup>۸</sup> در فواصل مشخص

بعد از نرمالیزه کردن سرعت تولید جملات مختلف «فارس دات»، میانگین کشش واکه های زبان فارسی به صورت جدول (۱) به دست می آید [۴].

جدول (۱) میانگین کشش واکه های زبان فارسی در گفتار پیوسته (نرمالیزه شده)

نوع واکه	میانگین کشش واکه (msec)
A	۱۴۱/۳۶
a	۱۱۹/۷۷
e	۸۰/۹۴
o	۱۰۱/۶۰
u	۱۱۶/۸۰
i	۱۰۶/۱۰

برای محاسبه میانگین کشش واکه های مختلف از دادگان «فارس دات» استفاده نموده ایم. باید توجه کرد که در مقایسه کشش واکه های بازشناسی شده با این مقادیر متوسط، سرعت تولید گفتار نقش مهمی را ایفا می کند و کشش واکه ها قبل از مقایسه با متوسط موجود باید ابتدا نسبت به سرعت تولید گفتار نرمالیزه شوند و سپس مقایسه صورت گیرد.

اگر کشش واکه بازشناسی شده بیش از ۱/۵ برابر میانگین کشش آن واکه بود، امکان اینکه این واکه از دو واکه تشکیل شده باشد، افزایش می یابد. در چنین مواردی به بررسی تغییرات لحظه ای دامنه سیگنال صوتی این واکه طولانی می پردازیم. همین عمل را در مواقعی که روش بازشناسی واج های مورد استفاده ما امکان حضور یک همخوان در وسط واکه بازشناسی شده را با عدم قطعیت اعلام داشته است، می توانیم مورد استفاده قرار دهیم. شکل ۱ سیگنال صوتی کلمه «ورود» را به همراه مرزهای پایین و بالای انرژی آن و نشانه ای که حد فاصل بین واکه /o/ و /u/ است را نشان می دهد. مطابق این شکل در حدفاصل بین دو واکه /o/ و /u/ یک کاهش قابل ملاحظه در دامنه سیگنال دیده می شود که مشخص کننده مرز بین دو واکه است. در مواردی که این کاهش دامنه به بیش از ۳۰٪ حداکثر دامنه (این آستانه به طور تجربی تعیین شده است) برسد، واکه طولانی بازشناسی شده را به دو واکه مجزا تبدیل می کنیم و

شده زیر است:

$$C_k = (0.5 + 0.5 \cos(\pi(Z - Z_k)))^2 \quad Z_k + 1 > Z > Z_k - 1 \quad \text{در فاصله} \quad (2)$$

$C_k = 0$  و در خارج این فاصله

مشخصات شبکه عصبی استفاده شده در جدول (۱) آمده است:

جدول (۱) مشخصات شبکه عصبی بازشناسی کننده «دمش».

لایه ورودی	لایه پنهان اول	لایه پنهان دوم	لایه خروجی
$7 \times 19 = 133$	۲۲۴	۷۲	۲

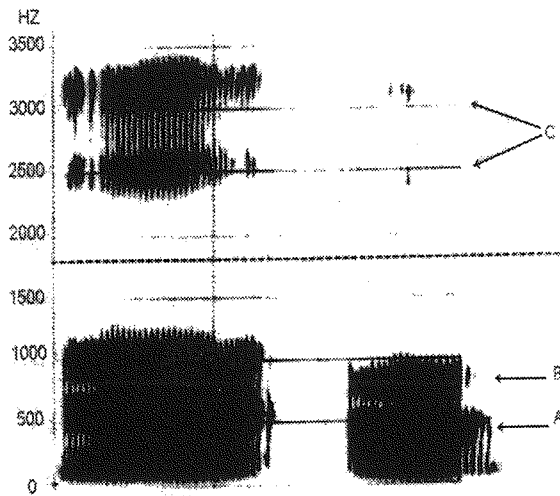
فعال شدن یکی از نرون های خروجی شبکه به معنی وجود «دمش» و فعال شدن دیگری به معنی عدم وجود «دمش» است. در مواردی که هر دو نرون فعال شوند، تصمیم نهایی براساس نرونی گرفته می شود که فعال تر است و خروجی بالاتری دارد.

اطلاعات تعلیم شبکه از دادگان «فارس دات» انتخاب شده اند و شامل تعداد زیادی واج انفجاری از نوع «واکدار» و «بی واک» از گویندگان مختلف می شوند.

بعد از اتمام مراحل یادگیری شبکه عصبی، آن را بر روی واج های انفجاری «داده آزمون» که از ۱۴۰ جمله از دادگان «فارس دات» تشکیل می شود، آزمایش نمودیم. درصد بازشناسی «دمش» برابر با ۹۳/۵ درصد به دست آمد.

### ۲-۳- گردی لب ها

در موضوع گردی لب ها «لب ها» مورد اشتباه شدن واکه های /o/ و /A/ را که مورد بسیار مهمی است و در حین بازشناسی واج ها به تعداد قابل توجهی رخ می دهد، انتخاب نمودیم. از تفاوت های مهم این دو واکه، گردی بیشتر لب در واکه /o/ را می توان مورد توجه قرار داد. ویژگی گردی لب در واکه /o/ بسیار قویتر از واکه /A/ ظاهر می شود. این ویژگی باعث کم شدن فاصله فورمنت های اول و دوم می شود [۷]. تفاوت دیگری هم بین این واکه ها وجود دارد و آن ویژگی «ارتفاع زبان» است. ارتفاع زبان و فورمنت اول نسبت به هم حالت معکوس دارند [۸]. بنابر این فورمنت اول در واکه /o/ که ارتفاع زبان در آن بالاتر است، پایین ترمی باشد. شکل ۲ اسپکتروگرام دو واکه /o/ و /A/ را نشان می دهد.



شکل (۲) از سمت چپ به ترتیب اسپکتروگرام دو واکه /o/ و /A/ - نشانه های A و B به ترتیب فورمنت های اول و دوم واکه ها را نشان می دهند.

همانگونه که انتظار می رود فورمنت های اول و دوم واکه /o/ پایین تر از /A/ هستند. نکته دیگری که در این شکل دیده می شود، تضعیف فرکانس های بالاتر در /o/ می باشد که اثر دیگر گردی لب است (در محل نشانه C). باتوجه به تفاوت هایی که در اسپکتروگرام این دو واکه دیده می شوند، ویژگی هایی را از روی طیف انرژی سیگنال صوتی دو واکه انتخاب نموده ایم که باتوجه به آنها می توان تعیین نمود، با کدامیک از این واکه ها مواجه هستیم. این ویژگی ها در جدول ۲ آورده شده اند.

جدول (۲) مشخصات ویژگی هایی که برای جداسازی

واکه های /o/ و /A/ به کار رفته اند.

ویژگی اول	«انرژی بالای ۱۱۰۰ هرتز» بر «انرژی زیر ۱۱۰۰ هرتز»
ویژگی دوم	«انرژی زیر ۵۰۰ هرتز» بر «انرژی بالای ۵۰۰ هرتز»
ویژگی سوم	«انرژی بین ۹۰۰ و ۱۳۲۰ هرتز» بر «انرژی زیر ۴۵۰ هرتز»

متوسط و انحراف معیار این نسبت ها را برای این دو واکه در ۱۴۰ جمله «داده تست» محاسبه نمودیم. جدول ۳ این مقادیر را نشان می دهد.

جدول (۳) میانگین و انحراف معیار سه ویژگی از

طیف انرژی واکه های /o/ و /A/

انحراف معیار	متوسط	نوع واکه	نوع ویژگی
۱۴۴	۲۰۵	A	ویژگی اول
۲۴۹	۸۰	O	ویژگی اول
۲/۴	۳/۹	O	ویژگی دوم
۱/۲	۰/۵	A	ویژگی دوم
۵۴۶	۸۵۷	A	ویژگی سوم
۲۷۶	۱۰۰	O	ویژگی سوم

۱۰۰، برابر ۱ فرض کرده ایم.

بخش باقیمانده تابع امکان دو واکه را با استفاده از رابطه توزیع نرمال محاسبه نموده ایم. در محاسبه توابع امکان /o/ و /A/ برای ویژگی های دیگر هم از همین روش پیروی شده است. برای بازشناسی /o/ و /A/، ابتدا ویژگی های اول تا سوم را برای سیگنال واکه محاسبه می کنیم و سپس امکان /A/ یا /o/ بودن را با قرار دادن مقادیر این ویژگی ها در توابع امکان مربوطه، برای هر ویژگی به طور جداگانه محاسبه می نماییم. سپس به ترتیب وزن های ۰/۵، ۱، ۱ را به مقادیر امکان محاسبه شده بر اساس ویژگی های اول تا سوم اعمال می کنیم و نتایج حاصل شده را با یکدیگر جمع می کنیم. تصمیم گیری نهایی در مورد اینکه این واکه /a یا /A/ است، باتوجه به مجموع مقادیر امکان محاسبه شده برای آنها گرفته می شود.

در بازشناسی /A/ و /b/ در ۱۴۰ جمله «داده تست»، برداشته شده از دادگان «فارس دات» ۵۹ مورد اشتباه بین این واکه ها مشاهده می گردید که با کمک گرفتن از روش فوق به ۲۶ مورد کاهش یافت.

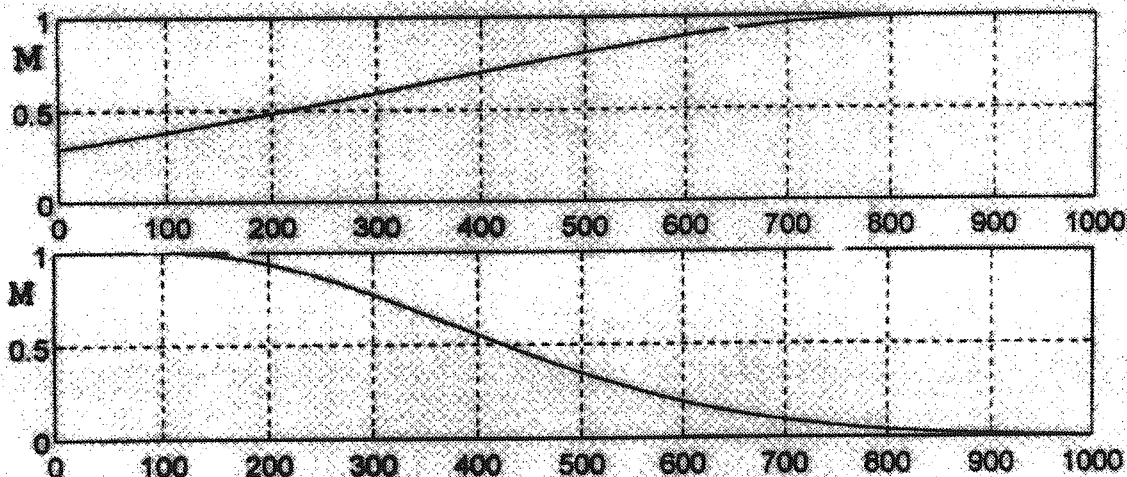
### ۳- نتیجه گیری

در این مقاله به روش هایی اشاره شد که با استفاده از ویژگی های پررودیک و خصوصیات اسپکتروگرام سیگنال گفتار در سطح واج ها به اصلاح درصد صحت بازشناسی آنها می پرداختند. کارایی این روش ها با

در تعریف ویژگی های اول تا سوم به دنبال این بوده ایم که ویژگی هایی را بیابیم که بتوانند حداکثر فاصله را بین واکه های /o/ و /A/ ظاهر نمایند. بعد از یافتن میانگین و انحراف معیار این ویژگی ها برای واکه های /o/ و /A/ از تابع توزیع نرمال [۹] برای ساختن تابع امکان آنها سود می جوییم:

$$m(x) = e^{-(x-x_m)^2 / 2\sigma^2}$$

نمونه ای از توابع امکان /o/ و /A/ برای ویژگی سوم در شکل ۳ دیده می شوند. همانگونه که در شکل ۳ دیده می شود، امکان واکه /A/ بودن را برای مقادیر بالاتر از ۸۵۷، و امکان واکه /o/ بودن را برای مقادیر کمتر از



شکل (۳) توابع امکان واکه های /o/ و /A/ باتوجه به ویژگی سوم. به ترتیب از بالا به پایین.

## زیرنویس‌ها

- 1- Spectrogram
- 2- Rhythm
- 3- Aspiration
- 4- Protrusion
- 5- Unvoiced
- 6- Bark
- 7- Hamming
- 8- Hanning

افزایش ۱/۵ درصدی صحت بازشناسی واج‌ها و رساندن آن به مقدار ۷۴/۸ درصد که مقدار قابل توجهی است، نشان می‌دهد که ویژگی‌های پروزودیک باید در بازشناسی گفتار فارسی مورد توجه بیشتری قرار گیرند و با توسعه این رویکرد و استفاده از ویژگی‌های پروزودیک بیشتر در سطح واج‌ها می‌توان به نتایج قابل توجهی در زمینه اصلاح روش‌های بازشناسی واج‌ها رسید. به خصوص زمانی که روش‌های کلاسیک بازشناسی گفتار به اشباع می‌رسند، این خصوصیات حاوی اطلاعات فراوانی هستند که با استفاده از آنها می‌توان بر بسیاری از مشکلات بازشناسی گفتار غلبه نمود.

## مراجع

- [1] Crystal, D. "Encyclopedia of Language and Linguistic", London, pp. 169-173, 1992.
- [2] Lass, R. "Phonology and Introduction to Basic Concepts", Cambridge University Press, London, 1983.
- [3] Bijankhan, M., "FARSDAT-The Speech Data base of FARSI Spoken Language", SST'94, pp. 836-831, 1994.
- [4] الماس گنج، ف. «تجزیه و تحلیل ساختار گفته‌های زبان فارسی با استفاده از اطلاعات پروزودیک سیگنال گفتار»، رساله دکتری، دانشگاه تربیت مدرس، ۱۳۷۷.
- [5] Lippman, R. P. "An Introduction to Computing with Neural Nets", IEEE ASSP Mag., April, pp. 4-22, 1987.
- [6] سیدصالحی، س، ع، «بازشناخت گفتار پیوسته فارسی با استفاده از مدل عملکردی مغز انسان در درک گفتار» رساله دکتری، دانشگاه تربیت مدرس، ۱۳۷۴.
- [7] Hermansky, H., "Perceptual Linear Predictive (PLP) Analysis of Speech", Journal of Acoustical Society of America, 87 (94) pp. 1738-1752, 1990.
- [8] Beddor, P.S., "Predicting the Structure of Phonological Systems", Journal of Phonetics, 48, pp. 83-107, 1991.
- [9] Papoulis, A., "Probability, Random Variables and Stochastic Processes", Mc Graw-Hill Inc., 1991.