

# بازشناسی مقاوم و توأم گفتار مستقیم و تلفنی با استخراج مناسب بردارهای بازنمایی و اصلاح آنها توسط معکوس سازی شبکه‌های عصبی

منصور ولی و سید علی سید صالحی

می‌روند. این تنوعات ناشی از منابعی مثل نویز اضافه شونده، انواع کانال صوتی، مشخصات میکروفن، لهجه‌ها و شیوه‌های مختلف گویش و غیره می‌باشند. این عوامل باعث عدم تطبیق بین داده‌های تعلیم و دادگان آزمون شده، منجر به کاهش کیفیت بازشناسی خواهند شد [۱].

روشهای متعددی برای غلبه بر این عدم تطابق پیشنهاد شده است. این روشها شامل تنظیم پارامترهای مدل با یادگیری بدون سرپرستی یا با سرپرستی و یا روشهای پردازش سیگنال به منظور از بین بردن اثر نامطلوب نویز روی سیگنال می‌باشند [۲] و [۳]. گفتار تلفنی در مقایسه با گفتار مستقیم (تهیه شده در شرایط عاری از هرگونه نویز محیط) بسیاری از این تنوعات گفتار را در بردارد. در طراحی بسیاری از سیستم‌های بازشناسی گفتار تلفنی، برای بدست آوردن صحت بازشناسی مطلوب، مدل بازشناسی را با دادگان تلفنی تعلیم مجدد می‌دهند و به این ترتیب یک مدل جدید بازشناسی برای دادگان تلفنی نسبت به دادگان گفتار مستقیم بدست می‌آورند. اما تعلیم مجدد مدل برای دادگان تلفنی چه از نظر تهیه دادگان تلفنی و چه از نظر زمانی بسیار هزینه‌بر است. علاوه بر این از آنجایی که مدل روی دادگان تلفنی با پهنای باند محدود تعلیم می‌بیند بنابراین کیفیت بازشناسی نسبت به گفتار مستقیم پایین‌تر خواهد آمد [۴]. یک راه حل دیگر برای طراحی سیستم بازشناسی گفتار تلفنی، نگاشت غیرخطی از بردارهای بازنمایی گفتار تلفنی به بردارهای بازنمایی گفتار مستقیم است که آن را با استفاده از شبکه‌های عصبی MLP می‌توان محقق نمود اما در این روش اگرچه نیازی به طراحی مدل جدید بازشناسی گفتار تلفنی نمی‌باشد اما دقت بازشناسی نسبت به روشهای قبلی بهبود چندانی ندارد. علاوه بر این به دلیل ضرورت همزمانی در سیگنال‌های گفتار تلفنی و مستقیم، این دو دسته دادگان باید همزمان ضبط شوند بنابراین نمی‌توان از مجموعه دادگان گفتاری مجزای تلفنی و مستقیم خود استفاده نمود [۵].

مطالعات روانشناسی نشان می‌دهد که انسان از دانش<sup>۳</sup> قبلی خود برای یادگیری بهتر فعلی خود بهره می‌گیرد. علاوه بر این چندین وظیفه مرتبط را می‌تواند بطور توأمان یاد بگیرد [۶]. نقطه ضعف سیستم‌های بازشناسی گفتار در این است که آنها براساس دادگان تعلیمی خود یاد می‌گیرند و از قابلیت‌های انسان در یادگیری بهره‌ای نمی‌برند. با دخیل کردن دانش قبلی در روند یادگیری، چندین خاصیت مطلوب مانند کاهش زمان یادگیری و تعمیم بهتر حاصل خواهد شد. انسان در مواجهه شدن با یک مورد یادگیری جدید از دانشهای فراگرفته شده قبلی خود بهره می‌گیرد [۷]. به عنوان مثال به نظر می‌رسد کیفیت بالای بازشناسی گفتار تلفنی در انسان تحت

چکیده: در حال حاضر تلاش فراگیری برای طراحی سیستم‌های بازشناسی گفتار مقاوم نسبت به تنوعات گفتار صورت می‌گیرد. یکی از این تنوعات، گفتار تلفنی نسبت به گفتار مستقیم (تهیه شده در شرایط عاری از هرگونه نویز محیط) می‌باشد. در مقاله حاضر با بهره‌گیری از پارامترهای طیفی LHCB<sup>۱</sup> و طراحی یک سری آزمایشهای عملی مشخص می‌گردد که این نوع بازنمایی برای طراحی سیستم‌های بازشناسی گفتار تلفنی و سیستم‌های بازشناسی توأم گفتار مستقیم و تلفنی که مبتنی بر شبکه‌های عصبی باشد نسبت به روش متداول MFCC مناسب‌تر است. سپس با استخراج بردارهای بازنمایی LHCB از گفتار مستقیم و تلفنی و طراحی مدل بازشناسی گفتار مبتنی بر شبکه عصبی MLP، یک سیستم بازشناسی توأم گفتار مستقیم و تلفنی ساخته می‌شود. آنگاه با استفاده از معکوس سازی شبکه‌های عصبی به روش گرادین بردارهای بازنمایی گفتار تلفنی به سمت بردارهای بازنمایی گفتار مستقیم اصلاح می‌گردد و با تعلیم شبکه دیگری روی دادگان اصلاح شده تلفنی و دادگان مستقیم دست نخورده، افزایش ۱/۴٪ در صحت بازشناسی گفتار تلفنی حاصل شده است. در مرحله بعد با استفاده از معکوس سازی عمومی شبکه‌های عصبی هر دو دسته بردارهای بازنمایی گفتار مستقیم و تلفنی به گونه‌ای اصلاح می‌شوند که بیشتر حاوی اطلاعات آوایی گفتار باشند و سایر تنوعات تا جای ممکن حذف شوند. با تعلیم شبکه دیگری روی این دادگان اصلاح شده افزایش ۲/۹۸٪ در صحت بازشناسی گفتار تلفنی و ۱/۶۸٪ در صحت بازشناسی گفتار مستقیم بدست آمده است.

کلید واژه: بازشناسی مقاوم گفتار، بازنمایی، شبکه عصبی، معکوس سازی.

## ۱- مقدمه

در سالهای اخیر کاربرد سیستم‌های بازشناسی گفتار روز به روز رو به افزایش است. بسیاری از سیستم‌های بازشناسی گفتار در کاربردهای محدود دقت بازشناسی بالایی دارند اما همین سیستمها در کاربردهای واقعی و فراگیر دچار مشکل خواهند شد. یکی از این مشکلات وجود تنوعات<sup>۲</sup> گفتار در شرایط متنوع صوتی است که این سیستمها در آنها بکار

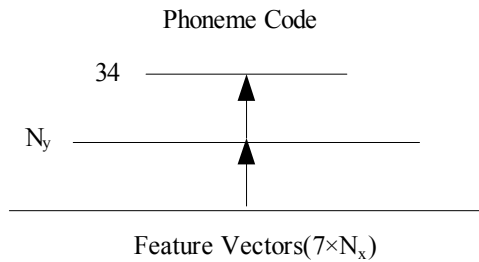
این مقاله در تاریخ ۹ خرداد ماه ۱۳۸۴ دریافت و در تاریخ ۳ آذر ماه ۱۳۸۴ بازنگری شد. این تحقیق توسط مرکز تحقیقات مخابرات براساس قرارداد شماره ۵۰۰/۱۴۵۲/ت پشتیبانی شده است.

منصور ولی، دانشکده مهندسی پزشکی، دانشگاه صنعتی امیرکبیر، خیابان حافظ، تهران، ایران (email: m\_vali@bme.aut.ac.ir).

سید علی سید صالحی، دانشکده مهندسی پزشکی، دانشگاه صنعتی امیرکبیر، خیابان حافظ، تهران، ایران (email: ssalehi@aut.ac.ir).

1. Logarithm of Energies in Hanning Critical Band Filter Banks

2. Variations



شکل ۱: مدل بازشناسی.

## ۲- ساختار مدل بازشناسی

مدل‌های بازشناسی بکار رفته در این تحقیق، مبتنی بر شبکه‌های عصبی MLP بوده و ساختار کلی آنها مشابه شکل ۱ می‌باشد. ساختار نورونی این شبکه بصورت  $34 - N_y - 7 \times N_x$  می‌باشد. که در آن  $N_x$  طول بردار بازنمایی است. بنابراین در ورودی شبکه بردار بازنمایی فریم جاری به همراه بردارهای بازنمایی ۳ فریم مجاور چپ و راست فریم جاری (در مجموع ۷ فریم) قرار می‌گیرند. در خروجی شبکه نیز به تعداد آواهای موجود در دادگان تعلیم یعنی ۳۴ نورون در نظر گرفته شده است.  $N_y$  تعداد نورون‌های لایه مخفی شبکه است. انتخاب مقدار  $N_y$  بر اساس این قاعده است که تعداد الگوهای تعلیمی به شبکه باید ۴ تا ۱۰ برابر تعداد وزن‌های مجهول شبکه انتخاب شوند [۱۲]. کمتر از ۴ برابر خاصیت تعمیم شبکه را ضعیف کرده و اعتبار مدل را پایین می‌آورد و بزرگتر از ۱۰ برابر نیز منجر به بزرگی شبکه و طولانی شدن زمان تعلیم آن می‌گردد درحالیکه تأثیر ناچیزی در بهبود دقت بازشناسی می‌گذارد.

## ۳- دادگان گفتار

دادگان گفتاری مورد نیاز برای تعلیم شبکه‌ها از مجموعه دادگان فارس‌دات و فارس‌دات تلفنی انتخاب شده‌اند [۱۳]. این مجموعه دادگان شامل دو جمله مشترک بیان شده توسط ۶۴ گویشور تلفنی و ۲۰۰ گویشور مستقیم (غیرتلفنی) می‌باشد. این دو جمله تمام آواهای فارسی را در بر دارند اما طبیعی است که تمام ترکیب‌های آوایی را نخواهند داشت. بنابراین یک سیستم بازشناسی گفتار پیوسته مستقل از گوینده ولی وابسته به متن خواهد بود. علت این انتخاب به هدف ما برمی‌گردد که درحقیقت پیدا کردن روش‌های بهبود بردارهای بازنمایی با استفاده از قابلیت‌های شبکه‌های عصبی است که وابسته به متن بودن دادگان به قدرت تعمیم نتایج این تحقیق لطمه‌ای وارد نمی‌کند. به این ترتیب از این مجموعه دادگان ۷۵٪ برای تعلیم شبکه‌ها و ۲۵٪ دیگر آن برای آزمون شبکه‌ها تخصیص داده شده‌اند. دو جمله مشترک انتخاب شده به صورت زیر می‌باشد.

"با روشن شدن هوا تظاهر کنندگان به سوی مجلس شورای ملی شروع به راهپیمایی کردند"  
 "مگر مژده اول چراغ قوه را خاموش نکرد"  
 نرخ نمونه برداری سیگنال گفتار برای گویشهای تلفنی ۸ کیلوهرتز و برای گویشهای غیرتلفنی برابر ۱۶ کیلوهرتز می‌باشد.

تأثیر دانش قبلی او از گفتار مستقیم است که قبلاً در محیط عاری از نویز فراگرفته شده است.

برای استفاده از دانش قبلی در یادگیری بهتر فعلی، توجه ویژه‌ای به تکنیک‌های شبکه‌های عصبی صورت گرفته است [۸] و [۹]. با استفاده از شبکه‌های عصبی می‌توان با چند روش مختلف از دانش فراگرفته قبلی استفاده نمود؛ از قبیل برنامه‌ریزی وزنهای شبکه [۸]، ایجاد مثالهای مجازی [۱۰] و استفاده از هدفهای خروجی اضافه که این تکنیک با عنوان مدل بازشناسی چند منظوره شناخته می‌شود [۹] و [۱۱].

مقاله حاضر گامی در جهت بکارگیری دانش حاصل از تعلیم مدل روی گفتار مستقیم در بازشناسی بهتر گفتار مستقیم و تلفنی است. ما معتقدیم طراحی یک مدل واحد برای بازشناسی گفتار مستقیم و تلفنی باعث خواهد شد که در بازشناسی گفتار تخریب شده تلفنی بتوان از دانش نهفته در شبکه، ناشی از یادگرفتن گفتار مستقیم بهره گرفته و بردارهای بازنمایی گفتار تلفنی را بهبود بخشید. برای این منظور ابتدا لازم است بردارهای بازنمایی به گونه‌ای از گفتار استخراج شوند که برای دو نوع گفتار تلفنی و مستقیم سازگار باشند و امکان تعلیم توأمان آنها به یک مدل شبکه‌عصبی فراهم شود. در این تحقیق مشخص خواهد شد که بردارهای بازنمایی LHCB برای این منظور بهتر از MFCC عمل می‌نمایند. برای استفاده از دانش گفتار مستقیم در بهبود بازشناسی گفتار تلفنی از تکنیک بازشناسی چند منظوره در شبکه‌های عصبی استفاده شده است. به این ترتیب که بازشناسی آواهای گفتار، توأمان با تشخیص نوع کانال گفتار (تلفنی یا مستقیم) در یک شبکه MLP صورت می‌گیرد. سپس از دو روش معکوس‌سازی شبکه (معکوس‌سازی به روش گرادیان و معکوس‌سازی عمومی) استفاده شده و بردارهای بازنمایی در جهت حذف تنوعات از آنها اصلاح می‌شوند. آنگاه با تعلیم یک مدل شبکه عصبی دیگر روی این دادگان جدید، بازشناسی آواهای گفتار مستقیم و تلفنی با صحت بالاتری صورت می‌پذیرد.

در ادامه مقاله ابتدا در بخشهای ۲ و ۳ به ترتیب ساختار شبکه‌های عصبی MLP و دادگان گفتاری استفاده شده در این تحقیق معرفی می‌شوند. سپس در بخش ۴ یک بحث تحلیلی روی دو نوع پارامتر بازنمایی گفتار LHCB و MFCC خواهد شد و نحوه استخراج بهینه بردارهای بازنمایی برای دادگان تلفنی بر اساس طراحی چند آزمون روی مدل‌های شبکه عصبی به دست می‌آید. در بخش ۵ بر اساس یک سری آزمایش ثابت می‌گردد که بردارهای بازنمایی LHCB برای بازشناسی گفتار تلفنی و نیز بازشناسی‌های توأمان تلفنی و مستقیم نسبت به روش MFCC در مدل‌های بازشناسی مبتنی بر شبکه‌عصبی ترجیح دارند. در بخش ۶ مقاله مبانی نظری معکوس‌سازی به روش گرادیان و معکوس‌سازی عمومی شبکه‌های عصبی بیان می‌گردد سپس در بخش ۷ با تعلیم توأمان شبکه روی بردارهای بازنمایی LHCB گفتار مستقیم و تلفنی، با استفاده از دو شیوه معکوس‌سازی به روش گرادیان و معکوس‌سازی عمومی، بردارهای بازنمایی در جهت افزایش صحت بازشناسی گفتار اصلاح می‌شوند. در پایان پس از بحث و تحلیل در بخش ۸، نتایج این تحقیق در بخش ۹ جمع‌بندی می‌شوند. کلیه برنامه‌ها در محیط نرم‌افزار MATLAB پیاده‌سازی شده‌اند.

## ۴-۲- طراحی آزمایش برای تعیین محدوده فرکانسی مناسب در استخراج بردارهای بازنمایی از گفتار تلفنی

محدوده فرکانسی مناسب برای استخراج بازنمایی از گفتار تلفنی (دارای محدوده فرکانسی بین صفر تا ۴ کیلوهرتز) را معمولاً بین ۱۲۵ تا ۳۷۰۰ هرتز انتخاب می کنند [۱۶]. زیرا که اثرات کانال های تلفنی مختلف در محدوده فرکانسهای پایین و بالای طیف ایجاد تخریب می کند. و این بخش از طیف گفتار تلفنی، اطلاعات مفیدی برای بازشناسی ندارد. اما از آنجایی که مدل بازشناسی در این تحقیق از نوع شبکه های عصبی MLP می باشد و شبکه این توانایی را دارد که مؤلفه های غیرمفید ورودی را با تخصیص وزنهای ضعیف به آنها کنار گذاشته و از مؤلفه های مفید ورودی استفاده کند؛ این سؤال مطرح می شود که آیا برای بازشناسی گفتار تلفنی در مدل های بازشناسی مبتنی بر شبکه های عصبی نیز لازم است که قسمتهای ابتدایی طیف گفتار کنار گذاشته شود؟ به همین منظور آزمایش هایی به صورت زیر طراحی می شوند.

مدل بازشناسی بکار رفته در این آزمایش ها، مبتنی بر شبکه های عصبی MLP بوده و ساختار نورونی آنها مطابق شکل ۱ به صورت  $7 \times N - 70 - 34$  می باشد.

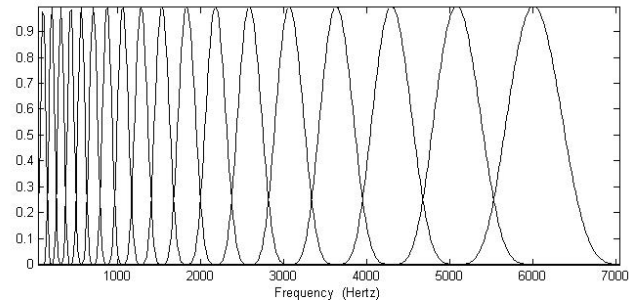
دادگان مورد نیاز برای تعلیم شبکه ها نیز مطابق بخش ۳ دو جمله مشترک بیان شده توسط تمام گویندگان تلفنی (۶۴ گویشور) استفاده شده است. از این مجموعه دادگان ۴۸ گوینده برای تعلیم و ۱۶ گوینده برای آزمون استفاده شده است.

### ۴-۲-۱- بررسی تأثیر انرژی فیلترهای فرکانس پایین بانک در دقت بازشناسی گفتار تلفنی

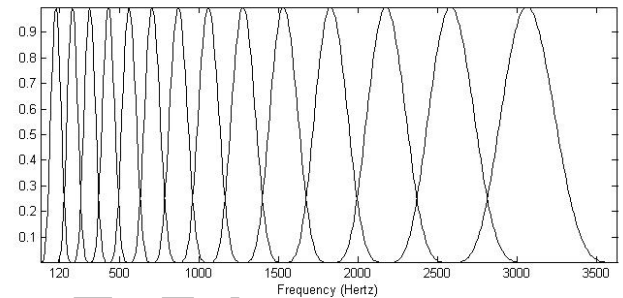
مطابق شکل های ۲ و ۳ از ۱۸ بانک فیلتر طراحی شده برای استخراج پارامترهای LHCب یا MFCC از گفتار مستقیم، ۱۴ تا ابتدایی آن در محدوده طیف گفتار تلفنی قرار می گیرند. هدف این است که از میان بردار ۱۴ مؤلفه ای انرژی فیلترها، تأثیر چند مؤلفه اول آن در صحت بازشناسی آواهای گفتار تلفنی به دست آید. برای این منظور در مرحله اول یک مدل بازشناسی گفتار مبتنی بر شبکه عصبی توسط بردارهای بازنمایی ۱۴ مؤلفه ای LHCب تعلیم داده می شود (انتخاب  $N=14$ ). مرحله دوم با حذف مؤلفه اول این بردارهای بازنمایی، شبکه عصبی دوم تعلیم داده می شود (انتخاب  $N=13$ ) و در مرحله سوم با حذف دو مؤلفه اول بردارهای بازنمایی، شبکه عصبی سوم تعلیم داده می شود (انتخاب  $N=12$ ) و بالاخره در مرحله چهارم با حذف سه مؤلفه اول بردارهای بازنمایی، شبکه چهارم تعلیم داده می شود (انتخاب  $N=11$ ). نتایج آزمون این چهار شبکه در جدول ۱ آمده است. از روی نتایج جدول ۱ که جهت درک بهتر در شکل ۴ نیز ترسیم شده اند؛ مشخص است که حذف انرژی فیلتر اول از بردارهای بازنمایی تلفنی تأثیر مثبت در صحت بازشناسی آوا داشته است ولی انرژی فیلترهای دوم و سوم بردارهای بازنمایی نباید کنار گذاشته شوند و حاوی اطلاعات مفیدی برای بازشناسی آواها می باشند. طبیعتاً نتیجه این آزمون برای استخراج پارامترهای MFCC از گفتار تلفنی نیز قابل تعمیم است. چرا که پارامترهای MFCC هم تبدیل کسینوسی همین بردارهای LHCب هستند.

### ۴-۲-۲- تحلیل حساسیت شبکه بازشناسی گفتار تلفنی نسبت به انرژی فیلترهای بانک

برای حذف انرژی فیلتر اول از بردارهای بازنمایی علاوه بر روش فوق، بر اساس روش تحلیلی زیر نیز می توان آن را اثبات کرد. در این روش یک



شکل ۲: بانک فیلتر برای استخراج پارامترهای بازنمایی از گفتار مستقیم.



شکل ۳: بانک فیلتر برای استخراج پارامترهای بازنمایی از گفتار تلفنی.

## ۴- استخراج بازنمایی های MFCC و LHCب

### ۴-۱- معرفی پارامترهای MFCC و LHCب

یکی از بهترین و رایجترین روشهای استخراج بازنمایی از گفتار مستقیم و تلفنی استخراج پارامترهای MFCC به همراه مشتقات اول و دوم آنها و اعمال هنجارسازی<sup>۱</sup> طولی (کم کردن میانگین بردارهای بازنمایی هر گویش از هر بردار بازنمایی و تقسیم بر انحراف معیار آن بردارها) روی آنها می باشد [۱۴].

شکل ۲ یک بانک فیلتر ۱۸ تایی نوع هنینگ که در محدوده فرکانسی ۸ کیلوهرتز توزیع شده است را نشان می دهد. محدوده صفر تا ۸ کیلوهرتز محدوده فرکانسی گفتار مستقیم است. پارامترهای LHCب لگاریتم انرژی گفتار در هر یک از این فیلترها می باشند که مجموعاً ۱۸ پارامتر انرژی می شود [۱۵] و پارامترهای MFCC با اعمال تبدیل کسینوسی روی پارامترهای LHCب در تعداد دلخواه بدست می آیند (معمولاً ۱۲ پارامتر در بازشناسی گفتار استخراج می گردد که با اضافه کردن لگاریتم انرژی کل فریم جمعاً ۱۳ پارامتر می گردد). بنابراین پارامترهای LHCب در حوزه طیف سیگنال گفتار هستند در حالی که پارامترهای MFCC در حوزه کپستروم می باشند.

برای گفتار با پهنای باند ۴ کیلوهرتز (گفتار تلفنی) تنها ۱۴ فیلتر از این ۱۸ فیلتر در محدوده فرکانسی آن واقع می شوند. این محدوده از بانک فیلتر در شکل ۳ نشان داده شده است. پارامترهای LHCب لگاریتم انرژی این ۱۴ فیلتر هستند. و پارامترهای MFCC با محاسبه تبدیل کسینوسی روی این ۱۴ پارامتر LHCب بدست می آیند. به این ترتیب ۱۳ پارامتر MFCC بدست آمده از گفتار با پهنای باند ۴ کیلوهرتز به دلیل همین تبدیل کسینوسی هیچ شباهتی به ۱۳ پارامتر بدست آمده از گفتار با پهنای باند ۸ کیلوهرتز ندارند. درحالی که انرژی ۱۴ فیلتر اول بانک (پارامترهای LHCب) برای هر دو گفتار مشابه همدیگر هستند.

1. Normalization

علت کوچک شدن مقدار حساسیت برای این مؤلفه از بردار بازنمایی را می‌توان اینگونه تحلیل کرد که شبکه در طول تعلیم به دلیل مؤثر نبودن این ورودی در ایجاد تمایز بین آواها، اثر آن را با وزنهای ضعیف‌تری که به آن تخصیص می‌دهد کاهش خواهد داد. بنابراین لازم است این مؤلفه از بردار بازنمایی حذف شود تا این بار اضافی از دوش شبکه برداشته شود. بنابراین در ادامه تحقیق برای استخراج بازنماییهای MFCC و LHCب از گفتار تلفنی، انرژی فیلتر اول بانک تأثیر داده نمی‌شود و فیلتر چهاردهم نیز تا محدوده ۳۵۵۰ هرتز را پوشش داده است.

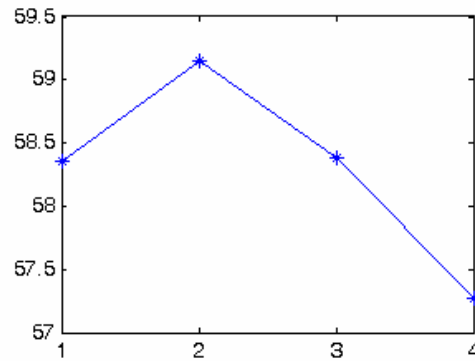
### ۳-۴- ساختن بردارهای بازنمایی با اضافه کردن مشتقات اول و دوم پارامترها به آنها

از آنجایی که در انتخاب بهینه بردارهای بازنمایی گفتار از مشتقات اول و دوم پارامترها در کنار خود پارامترها نیز استفاده می‌شود بنابراین بردارهای بازنمایی MFCC به همراه مشتقات اول و دوم آنها برای هر دو نوع گفتار مستقیم و تلفنی شامل ۳۹ مؤلفه می‌گردد ( $13+13+13=39$  مؤلفه). اضافه کردن مشتقات اول و دوم پارامترهای LHCب به آنها به صورت زیر انجام می‌شود. برای گفتار مستقیم یک بردار بازنمایی شامل ۱۸ پارامتر LHCب و مشتقات اول و دوم آنها که جمعاً شامل ۵۴ پارامتر می‌باشد ساخته می‌شود. اما از آنجایی که هدف نهایی تعلیم یک شبکه واحد بازشناسی گفتار مستقیم و تلفنی است باید طول بردارهای بازنمایی گفتار مستقیم و تلفنی با هم برابر باشند. به همین جهت در ساختن بردار بازنمایی گفتار تلفنی علاوه بر اینکه اولین مؤلفه از بردار ۱۴ تایی LHCب تلفنی صفر قرار داده می‌شود (دلیل آن در بخش ۴-۲ آمده است) ۴ مؤلفه صفر که مربوط به فیلترهای محدوده ۳/۶ تا ۸ کیلوهرتز است نیز به انتهای آن اضافه شده و یک بردار ۱۸ تایی از پارامترهای LHCب تلفنی ساخته می‌شود. سپس با اضافه کردن مشتقات اول و دوم این بردار به آن، بردار بازنمایی به طول ۵۴ مشابه بردار بازنمایی گفتار مستقیم بدست می‌آید. حال نه تنها طول بردار بازنمایی گفتار مستقیم و تلفنی برابر همدیگر هستند بلکه مؤلفه‌های نظیر هم از بردارهای بازنمایی گفتار مستقیم و تلفنی، مربوط به فیلترهای نظیر هم در بانک فیلتر هستند و همین باعث تشابه این دو بردار بازنمایی شده که امکان تعلیم توأمان به یک شبکه را فراهم می‌نماید.

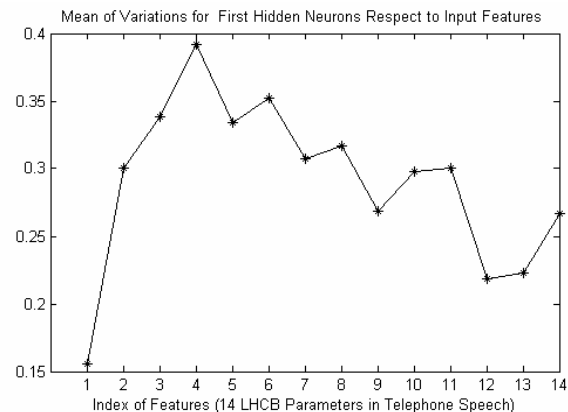
### ۵- مقایسه بردارهای بازنمایی MFCC و LHCب در بازشناسی گفتار مبتنی بر شبکه‌های عصبی MLP

اکنون دو دسته بردار بازنمایی MFCC (۳۹ مؤلفه‌ای) و LHCب (۵۴ مؤلفه‌ای) روی مدل‌های بازشناسی گفتار تلفنی، گفتار مستقیم و بازشناسی توأمان گفتار مستقیم و تلفنی بررسی شده و در عمل کاربرد این دو دسته بازنمایی ارزیابی می‌شوند.

توجه به این نکته ضروری است که اگرچه طول بردارهای بازنمایی MFCC کمتر از LHCب می‌باشد و ظاهراً مقایسه بین دو نوع بردار بازنمایی با طول غیریکسان منطقی به نظر نمی‌رسد ولی همانگونه که در بخش ۴-۱ مطرح شد پارامترهای ۱۳ تایی MFCC از تبدیل کسینوس روی پارامترهای ۱۸ تایی LHCب به دست می‌آیند. بنابراین تعداد فیلترها در بانک فیلتر اعمال شده روی طیف سیگنال، برای هر دو نوع بازنمایی یکسان است. به عبارتی دیگر می‌توان گفت محتوای اطلاعاتی این دو نوع بازنمایی یکسان، اما نحوه بیان آنها متفاوت هستند. به این ترتیب مقایسه بین کارایی آنها منطقی است.



شکل ۴: صحت بازشناسی گفتار تلفنی مربوط به چهار شبکه فوق.



شکل ۵: میزان حساسیت نورونهای لایه مخفی شبکه به ۱۴ مؤلفه بردار LHCب در ورودی شبکه.

جدول ۱: نتایج حذف پارامترهای فرکانس پایین از پارامترهای LHCب.

بردار بازنمایی	صحت بازشناسی آواها
۱۴ پارامتر LHCب	۵۸,۳۵
حذف پارامتر اول	۵۹,۱۵
حذف دو پارامتر اول	۵۸,۳۸
حذف سه پارامتر اول	۵۷,۲۷

دید فراگیر نسبت به میزان تأثیر همه مؤلفه‌های بردار ورودی در عملکرد شبکه می‌توان پیدا کرد.

اگر در شبکه اول تعلیم داده شده روی بردارهای ۱۴ تایی LHCب، برای ۷۰ نورون لایه مخفی شبکه، میزان حساسیت خروجی این نورونها را نسبت به ورودیهای شبکه محاسبه کنیم یعنی  $|\partial y_j / \partial x_i|$  و این مقدار را روی تمامی دادگان تعلیم شبکه میانگین گیری کنیم خواهیم داشت

$$S_{ji} = \sum_{\forall x_i} \left| \frac{\partial y_j}{\partial x_i} \right| \quad j = 1, 2, \dots, 70 \quad i = 1, 2, \dots, 14 \quad (1)$$

اکنون اگر بخواهیم میزان حساسیت نورونهای لایه مخفی شبکه را نسبت به هر یک از این ۱۴ مؤلفه ورودی محاسبه کنیم کافی است  $S_{ji}$ ها را روی ۷۰ نورون لایه مخفی میانگین گیری کنیم یعنی

$$S_i = \frac{1}{70} \sum_{j=1}^{70} S_{ji} \quad i = 1, 2, \dots, 14 \quad (2)$$

بردار ۱۴ تایی  $S$  نمایانگر حساسیت نورونهای لایه مخفی شبکه به ۱۴ مؤلفه بردار بازنمایی ورودی است که در شکل ۵ نشان داده شده است. بر طبق این نمودار واضح است که مؤلفه اول بردار بازنمایی تلفنی تأثیر بسیار کمی در تعیین وضعیت نورونهای شبکه دارد و می‌توان آنرا حذف کرد.

از آوای گفتار مستقیم و تلفنی با همدیگر، امکان طراحی سیستم های بازشناسی توأم گفتار مستقیم و تلفنی را فراهم می نماید. علاوه بر این از روی نتایج جدول ۲ می توان مشاهده کرد که شبکه بازشناسی توأم گفتار مستقیم و تلفنی برای بازنمایی های LHCب تلفنی، حدود ۱/۲ درصد صحت بالاتری نسبت به شبکه بازشناسی خالص تلفنی (آزمایش دوم) دارد. زیرا دانش نهفته از تعلیم دادگان گفتار مستقیم در شبکه، باعث رشد بازشناسی گفتار تلفنی شده است. البته برای گفتار مستقیم نتیجه بازشناسی نسبت به آزمایش اول کمی پایین تر بدست آمده است اما در هر حال به رجحان داشتن بازنمایی های LHCب نسبت به MFCC لطمه ای نمی زند.

### ۲-۵- نتایج حاصل از آزمونها

بر اساس بحثهای استدلالی و آزمایش های مطرح شده در مجموع نتایج زیر بدست می آید.

الف. استخراج بازنمایی MFCC به همراه هنجارسازی مناسب آنها یکی از بهترین روشهای استخراج بازنمایی برای گفتار مستقیم است.

ب. برای مدلهای بازشناسی گفتار تلفنی مبتنی بر شبکه های عصبی به دلایلی که در متن بحث شد استخراج بازنمایی های طیفی مانند LHCب نتایج بهتری نسبت به بازنمایی های حوزه کپستروم (MFCC) دارد.

ج. در مدلهای بازشناسی تلفنی-مستقیم مبتنی بر شبکه های عصبی نیز به دلیل ضرورت سازگاری بین مؤلفه های بردار بازنمایی در گفتار تلفنی و مستقیم باز هم بازنمایی های طیفی مانند LHCب بر بازنمایی های حوزه کپستروم (MFCC) ترجیح پیدا می کنند.

نتیجه سوم برای ما حاوی نکته ارزشمندی است که می توان با طراحی مدلهای بازشناسی توأم گفتار مستقیم و تلفنی نه تنها کمبود دادگان جمع آوری شده تلفنی را از طریق دادگان مستقیم تأمین نمود بلکه با بکارگیری روشهای تطبیقی و نیز شیوه های معکوس سازی شبکه های عصبی بردارهای بازنمایی تلفنی را به سمت بردارهای بازنمایی مستقیم اصلاح کرد و نهایتاً به سیستم های بازشناسی گفتار مستقیم-تلفنی که نتایج بهتری در بازشناسی گفتار تلفنی و هم گفتار مستقیم در مقایسه با سیستم های مجزای بازشناسی گفتار تلفنی یا مستقیم دارند دسترسی پیدا کرد. در ادامه ایده فوق مورد بررسی و آزمون قرار می گیرد.

### ۶- معکوس سازی شبکه های عصبی

الگوریتم های معکوس سازی شبکه های عصبی را می توان به دو دسته اصلی تقسیم کرد:

(۱) جستجوی کامل<sup>۱</sup>

(۲) جستجوی تک عنصری<sup>۲</sup>

روشهای جستجوی کامل معمولاً زمانی که ابعاد و محدوده تغییرات هر یک از متغیرهای ورودی محدود باشد کاربرد دارد. این روش به دنبال پیدا کردن مجموعه جوابهای ممکن در ورودی شبکه است که منجر به یک خروجی مطلوب می گردند. در روشهای جستجوی تک عنصری در هر بار جستجو یک جواب ممکن یا یک جواب خاص از این مجموعه جوابها

جدول ۲: صحت بازشناسی آوای غیرسکوت، برای دادگان آزمون دو نوع بازنمایی

مدل بازشناسی	MFCC	LHCب
گفتار مستقیم	۸۲٫۴۸	۸۲٫۶۷
گفتار تلفنی	۶۴٫۱۸	۶۷٫۱۲
توأم مستقیم	۷۹٫۷۸	۸۱٫۹۲
گفتار تلفنی	۶۱٫۱۱	۶۸٫۳۹

### ۵-۱- مقایسه بردارهای بازنمایی MFCC و LHCب در بازشناسی گفتار مستقیم و تلفنی

برای این منظور سه دسته آزمایش طراحی شده است. مدل بازشناسی بکار رفته در تمام آزمایشها، مبتنی بر شبکه های عصبی MLP مشابه شکل ۱ می باشد. ساختار نورونی این شبکه برای بردارهای بازنمایی LHCب بصورت ۳۴-۷۰-۷×۵۴ و برای بردارهای بازنمایی MFCC بصورت ۳۴-۷۰-۷×۳۹ می باشد. برای این قسمت از تحقیق از دو جمله مشترک بیان شده توسط تمام گویندگان تلفنی (۶۴ گویشور) و گویندگان غیرتلفنی (۱۰۰ گویشور) استفاده شده است. از هر مجموعه دادگان ۷۵٪ برای تعلیم شبکه ها و ۲۵٪ دیگر آن برای آزمون شبکه ها تخصیص داده شده است. دو جمله مشترک انتخاب شده همان جملاتی هستند که در بخش ۳ معرفی شدند.

آزمایش اول: تعلیم بردارهای بازنمایی MFCC و LHCب گفتار مستقیم به دو شبکه مجزا.

آزمایش دوم: تعلیم بردارهای بازنمایی MFCC و LHCب گفتار تلفنی به دو شبکه مجزا.

آزمایش سوم: تعلیم بردارهای بازنمایی MFCC مستقیم و تلفنی به یک شبکه و بردارهای بازنمایی LHCب مستقیم و تلفنی به یک شبکه دیگر.

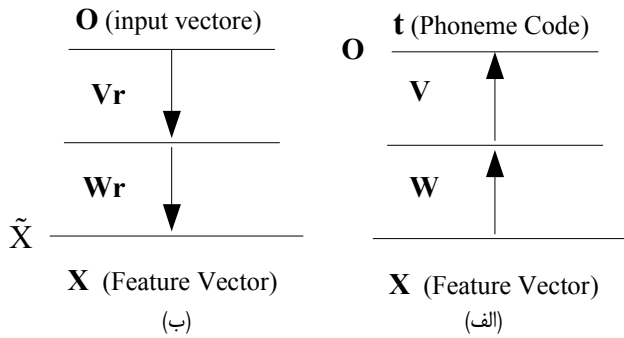
بر اساس آزمایشهای مذکور نتایج صحت بازشناسی آواها برای دادگان آزمون این دو نوع بازنمایی در جدول ۲ آمده است. صحت بازشناسی بر اساس درصد صحت فریم های غیرسکوت بیان شده است. زیرا که تعداد فریم های سکوت گفتار از سایر آواها بیشتر بوده و صحت بازشناسی آنها برای هر نوع بازنمایی بسیار بالا است. و مانع از مشهود شدن تفاوت روش های مختلف بازشناسی گفتار و انواع بهبود در بازنمایی که مد نظر این مقاله است می گردد.

نتایج آزمایش اول نشان می دهد که بازشناسی گفتار مستقیم برای این دو دسته بازنمایی تفاوت چندانی با هم ندارند و از آنجایی که طول بردار بازنمایی MFCC از طول بردار بازنمایی LHCب کوتاه تر است برای بازشناسی گفتار مستقیم، پارامترهای MFCC ترجیح خواهند داشت.

در آزمایش دوم نتایج بهتر بازنمایی LHCب نسبت به MFCC برای بازشناسی گفتار تلفنی، گویای آن است که بازنمایی LHCب برای دادگان تخریب شده و نویزی مثل گفتار تلفنی مناسب تر است.

مقایسه نتایج حاصل از آزمایش سوم در مقایسه با نتایج آزمایش های اول و دوم نشان می دهد که بکارگیری بردارهای بازنمایی MFCC در شبکه بازشناسی توأم گفتار مستقیم و تلفنی، نسبت به شبکه های بازشناسی مجزای گفتار مستقیم و تلفنی نتایج ضعیف تری می دهد. زیرا که بازنمایی های MFCC برای گفتار مستقیم و تلفنی تطابق منطقی با هم ندارند. بنابراین در طراحی سیستم های بازشناسی گفتار توأم تلفنی و مستقیم نمی تواند کاربرد داشته باشد. اما استفاده از بازنمایی های طیفی گفتار مانند پارامترهای LHCب به دلیل تطابق منطقی بازنمایی های نظیر

1. Exhaustive Search  
2. Single Element Search



شکل ۶: (الف) شبکه MLP مستقیم و (ب) شبکه MLP معکوس.

خروجی واقعی این شبکه نیز بردار  $\tilde{X}$  خواهد شد و خطای خروجی برای آن نیز برابر  $|X - \tilde{X}|$  می‌گردد. نحوه محاسبه خروجی شبکه و خطای تعلیم در روابط (۵) و (۶) آمده است. روش تعلیم شبکه معکوس شکل ۶-ب مشابه روش تعلیم شبکه مستقیم (شکل ۶-الف) است. به این ترتیب که خطای خروجی شبکه نسبت به اهداف خروجی برخلاف جهت گرادیان پس‌انتشار شده، وزنهای شبکه اصلاح می‌گردند. این روند تا آنجا که صحت بازشناسی روی دادگان آزمون به ماکزیمم مقدار خود یا خطای خروجی پایین‌تر از یک حد آستانه برسد ادامه پیدا می‌کند.

$$O = f(V \times f(W \times X)) \quad \tilde{X} = f(W_r \times f(V_r \times O)) \quad (۵)$$

$$E = |t - O| \quad E = |X - \tilde{X}| \quad (۶)$$

اگرچه خطای نهایی تعلیم در این شبکه نسبت به شبکه مستقیم بزرگتر است اما به میانگین تقریبی مجموعه جوابهای ممکن منتهی می‌گردد [۱۱]. به همین جهت هم در ادامه از این روش با نام روش معکوس سازی عمومی یاد خواهد شد. در ادامه از این دو روش معکوس سازی شبکه‌های عصبی برای اصلاح بردارهای بازنمایی گفتار استفاده می‌شود.

## ۷- اصلاح بردارهای بازنمایی با استفاده از معکوس سازی شبکه‌های عصبی

### ۷-۱- انتخاب بردار بازنمایی و دادگان

همانگونه که در بخش ۵ نشان داده شد بردارهای بازنمایی طیفی مانند LHCB برای گفتار مستقیم و تلفنی و نیز برای انواع گفتار تلفنی تطابق بهتری دارند. بنابراین در این بخش از تحقیق که هدف ما ارتقاء کیفیت بازشناسی گفتار در سیستم بازشناسی توأم گفتار مستقیم و تلفنی است بردارهای بازنمایی LHCB به همراه مشتقات اول و دوم آنها که در مجموع طولی برابر ۵۴ مؤلفه پیدا می‌کنند مورد استفاده قرار می‌گیرند و قبلاً در بخش ۴-۳ بیان شد که نه تنها طول این بردارهای بازنمایی، برای گفتار مستقیم و تلفنی برابر هم‌دیگر هستند بلکه مؤلفه‌های نظیر هم از این دو نوع بردار بازنمایی، مربوط به فیلترهای نظیر هم در بانک فیلتر هستند و تشابه این دو بردار بازنمایی در مؤلفه‌های غیرصفر، امکان تعلیم توأمان به یک شبکه را فراهم می‌نماید.

دادگان استفاده شده برای این بخش از تحقیق، دادگانی است که در بخش ۳ مقاله به صورت کامل معرفی شده است.

بدست می‌آید که در برخی از روشهای آن وابسته به مقادیر اولیه شروع الگوریتم است [۱۷].

### ۶-۱- جستجوی تک‌عنصری به روش گرادیان

توجه به معکوس‌سازی تک‌عنصری شبکه‌های عصبی از طریق گرادیان<sup>۱</sup> اولین بار توسط ویلیامز [۱۸] مطرح شد. این روش‌ها با استفاده از روش استاندارد بهینه‌سازی با پس‌انتشار خطا برای اصلاح یک ورودی اولیه کار می‌کنند. جستجو با یک بردار ورودی  $x^*$  شروع می‌شود. اگر  $x_k^i$ ،  $k$  امین مؤلفه بردار  $x^i$  باشد، با هدف کاهش گرادیان، رابطه بازگشتی زیر به دست می‌آید

$$x_k^{i+1} = x_k^i - \eta \frac{\partial E}{\partial x_k^i} \quad (۳)$$

که  $\eta$  طول گام و  $t$  اندیس تکرار است. با فرض یک توپولوژی جلوسوی عمومی، تکرار برای معکوس‌سازی در (۳) می‌تواند به صورتی که در ادامه می‌آید حل شود: با فرض  $\delta_k = \partial E / \partial x_k$  و  $k \in I$  که برای هر نرون

$$\delta_j = \begin{cases} \phi_j'(o_j)(o_j - t_j) & j \in O \\ \phi_j'(o_j) \sum_{m \in H, O} \delta_j w_{jm} & j \in I, H \end{cases} \quad (۴)$$

$I$  و  $H$  و  $O$ : به ترتیب مجموعه‌های نرون‌های ورودی، لایه مخفی و خروجی؛

$w_{jm}$ : مقدار وزن از نرون  $j$  به نرون  $m$ ؛

$\phi_j'$ : مشتق تابع غیرخطی نرون  $j$ ؛

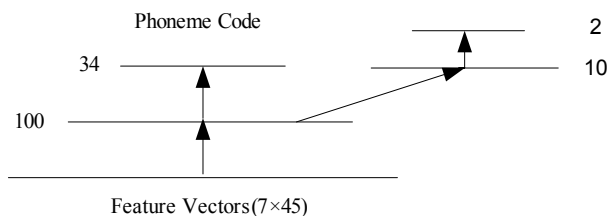
$o_j$ : میزان فعالیت نرون  $j$ ؛

$t_j$ : خروجی مطلوب برای نرون  $j$ .

مشابه روش استاندارد پس‌انتشار خطا، مشتقات نرون،  $\delta_j$ ها باید با ترتیب برعکس از خروجی به ورودی حل شوند. به این ترتیب مشتق خطای خروجی نسبت به بردارهای ورودی شبکه، مطابق رابطه بازگشتی (۴) محاسبه شده و این روند آنقدر تکرار می‌شود تا مقدار خطای خروجی یا شیب آن از یک حد آستانه کمتر شود. عدم وجود فیدبک، تنها فرض لازم برای اتصالات نرون‌ها در (۴) است. در ادامه از این روش معکوس‌سازی شبکه‌های عصبی، با نام روش معکوس‌سازی گرادیان یاد خواهد شد.

### ۶-۲- جستجوی تک‌عنصری به روش تعلیم شبکه معکوس

در این روش پس از تعلیم شبکه مستقیم، مطابق شکل ۶-الف به‌ازای بردار بازنمایی ورودی  $X$  مقدار هدف، بردار  $t$  (کد آوایی نظیر  $X$ ) می‌باشد. اما مقدار خروجی واقعی شبکه پس از تعلیم، بردار  $O$  می‌شود و خطای شبکه برای این ورودی،  $|t - O|$  است. حال شبکه MLP دیگری مطابق شکل ۶-ب به عنوان شبکه معکوس تعلیم داده می‌شود که ورودیهای این شبکه خروجیهای شبکه مستقیم (بردار  $O$ ) و اهداف خروجی این شبکه بردارهای بازنمایی است که در ورودی شبکه شکل ۶-الف مورد استفاده قرار گرفته‌اند (بردار بازنمایی  $X$ ). در حالی که



شکل ۸: ساختار مدل بازشناسی دو منظوره.

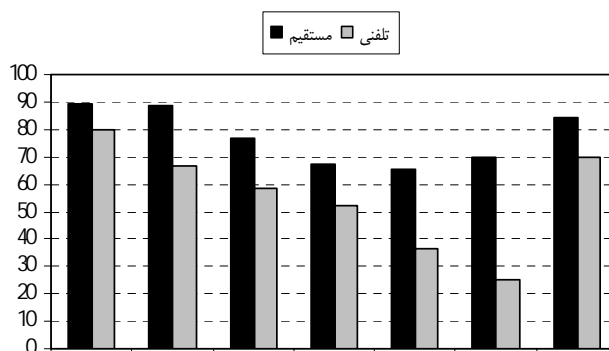
شامل یک لایه پنهان ۱۰ نورونی و دو نورون در خروجی است که برای تشخیص نوع کانال (مستقیم یا تلفنی) دو بیت تخصیص داده شده است. این یک مدل بازشناسی دو منظوره است؛ بازشناسی آوا و بازشناسی نوع کانال که ساختار آن در شکل ۸ نشان داده شده است. نتایج آزمون این شبکه به همراه نتایج آزمون شبکه مرجع در جدول ۳ آورده شده است.

نتایج نشان می دهد که این شبکه با صحت بسیار بالایی نوع کانال تلفنی را از کانال گفتار مستقیم تمایز می دهد.  $100\% - 99/95\%$  و افت ناچیز در بازشناسی آواهای گفتار مستقیم  $0/18\%$  و تلفنی  $0/27\%$  نسبت به شبکه مرجع کاملاً قابل توجه است. زیرا شاخه اضافه شده برای شناسایی نوع کانال گفتار، در این شبکه نسبت به شبکه مرجع یک بار اضافه ای است که بر آن تحمیل شده است. به عبارتی وزنه های لایه اول این شبکه باید برای برآورده شدن دو هدف (شناسایی درست آواها و تشخیص درست نوع کانال) بهینه شوند.

حال بردارهای بازنمایی تلفنی به ورودی شبکه داده شده و کد دوبیتی خروجی شبکه (مربوط به تشخیص کانال گفتار) بصورت تحمیلی برابر کد کانال گفتار مستقیم قرار داده می شود و خطای خروجی مطابق روابط (۳) و (۴) به منظور اصلاح ورودی پس انتشار می شود. نکاتی که در اعمال این روش پس انتشار بر روی ورودی شبکه در نظر گرفته شده است به صورت زیر می باشد:

۱. خروجی مطلوب شبکه برای ۳۴ نورون شناسایی کننده آواها همان کد باینری آوای شناسایی شده توسط شبکه فرض می شود و خطای ناشی از این کد و خروجی واقعی شبکه نیز برای اصلاح ورودی پس انتشار می شود. به این ترتیب ورودی شبکه نه تنها در جهت تشخیص نوع کانال مستقیم اصلاح می گردد بلکه اطلاعات آوایی قبلی اش را نیز حفظ می کند.
۲. از میان هفت فریم متوالی که در ورودی شبکه قرار می گیرند تنها فریم میانی اصلاح می گردد. ولی تمام فریمهای دادگان تلفنی تعلیم و آزمون به این شیوه اصلاح می شوند.
۳. پس از اتمام الگوریتم اصلاح دادگان تلفنی، این دادگان اصلاح شده تلفنی به ورودی مدل بازشناسی دو منظوره داده شد و  $97/8\%$  از این دادگان از نوع گفتار مستقیم شناسایی شد و صحت بازشناسی آوا نیز تقریباً ثابت ماند. این نتیجه حاکی از موفقیت در اصلاح دادگان تلفنی در جهت نزدیک شدن به دادگان گفتار مستقیم است.

پس از اتمام مراحل فوق با رسم بردارهای بازنمایی تلفنی اصلاح شده بر روی بردارهای اولیه نظیر آنها متوجه این نکته شدیم که بیشترین اصلاح و تغییر در مؤلفه های صفر بردارهای بازنمایی اتفاق افتاده است. به عبارتی دیگر، مؤلفه هایی از بردارهای بازنمایی تلفنی که نسبت به مؤلفه های نظیر از بردارهای بازنمایی گفتار مستقیم وجود ندارند به نوعی از روی دانش نهفته در مدل بازشناسی دو منظوره تخمین زده شده اند.



شکل ۷: درصد بازشناسی دسته های آوایی مختلف در دادگان گفتار تلفنی و مستقیم.

## ۲-۷- مدل بازشناسی مرجع

در اولین مرحله یک مدل مرجع برای بازشناسی توأم گفتار مستقیم و تلفنی مبتنی بر شبکه های عصبی MLP مطابق شکل ۱ طراحی می گردد. ساختار نورونی این شبکه بصورت  $34-100-7 \times 54$  می باشد. خروجی شبکه دارای توابع غیرخطی از نوع تانژانت هایپربولیک<sup>۱</sup> هستند. انتخاب این نوع تابع غیرخطی بجای تابع غیرخطی سیگموئید جهت تسریع در همگرایی شبکه است. مقادیر وزنه های اولیه شبکه بصورت تصادفی و در محدوده ۱- تا ۱ در نظر گرفته شده اند و ضریب یادگیری در شروع تعلیم برابر  $0/01$  و در حین تعلیم در هر تکرار با ضریب  $0/95$  کاهش می یابد که باعث تعلیم بهینه شبکه می گردد. ضریب ممنتوم<sup>۲</sup> نیز برابر  $0/2$  در نظر گرفته شده است.

شبکه مذکور چندین بار با مقادیر وزنه های تصادفی متفاوت اولیه تعلیم داده می شود و نهایتاً صحت بازشناسی گفتار تلفنی برابر  $69/76 \pm 0/1$  درصد و صحت بازشناسی گفتار مستقیم برابر  $84/1 \pm 0/2$  درصد بدست آمده است. صحت بازشناسی بر حسب درصد صحت فریم های غیرسکوت بیان شده است. به این ترتیب انحراف معیار در بازشناسی دادگان تلفنی برابر  $0/1\%$  و برای دادگان گفتار مستقیم برابر  $0/2\%$  می باشد. این مقادیر میزان اعتبار مدل بکار رفته در این تحقیق را نشان می دهد.

نمودار میله ای شکل ۷ درصد صحت بازشناسی آواها برای شش دسته مختلف آواهای گفتار یعنی واکه ها، سایشی ها، شبه واکه ها، بست ها، انفجاری ها و سایشی-انفجاری ها را به تفکیک نشان می دهد. نکات قابل توجه این نمودار بصورت زیر است:

الف: بیشترین صحت بازشناسی برای هر دو دسته دادگان در واکه های گفتار حاصل شده است.

ب: بیشترین اثر تخریبی کانال تلفنی بر روی آواهای انفجاری و سایشی-انفجاری واقع شده است.

## ۳-۷- اصلاح بردارهای بازنمایی با استفاده از معکوس سازی گرادیان

در این روش ابتدا یک شبکه عصبی واحد برای بازشناسی گفتار مستقیم و تلفنی مشابه شبکه عصبی مدل مرجع تعلیم داده می شود. با این تفاوت که یک شاخه اضافه به لایه پنهان شبکه اضافه می گردد که خود

1. Hyperbolic
2. Momentum



جدول ۳: نتایج بازشناسی آواهای گفتار مستقیم و تلفنی بر حسب درصد.

شبکه عصبی	صحت بازشناسی مستقیم	صحت بازشناسی تلفنی	تشخیص مستقیم یا تلفنی
مرجع	۸۴/۱±۰/۲	۶۹/۷۶±۰/۱	---
مدل دو منظوره	۸۳/۹۲	۶۹/۳۹	۹۹/۸۵- ۱۰۰
معکوس گردایان	۸۳/۹۴	۷۱/۳۴	---
معکوس عمومی	۸۵/۷۸	۷۲/۶۵	---

## ۸- بحث و تحلیل نتایج

گفتار تلفنی و مستقیم محتوی دو دسته اطلاعات هستند؛ بخش اول، اطلاعات آوایی گفتار است که در هر دو دسته دادگان مشترک است با این تفاوت که در گفتار تلفنی نسبت به گفتار مستقیم قسمتی از این اطلاعات از بین رفته است. بخش دوم، اطلاعات غیرآوایی گفتار یا همان تنوعات گفتار است که ناشی از اثرات گوینده، میکروفون، کانال و غیره است که برای این دو دسته دادگان متفاوت می‌باشد.

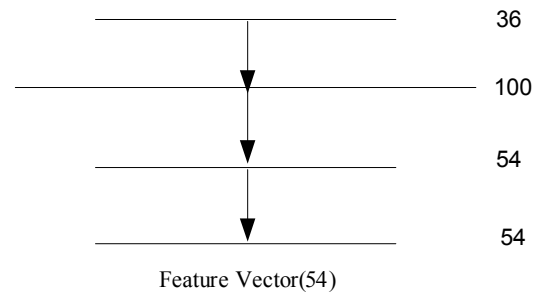
تعلیم یک شبکه عصبی واحد روی این دو دسته دادگان به عنوان مدل مرجع در اولین گام کمک می‌کند تا شبکه به‌عنوان یک واحد محاسباتی موازی و گسترده، بخش مشترک اطلاعات این دو دسته دادگان (اطلاعات آوایی) را بیرون بکشد. و برای گفتار تلفنی از دانش گفتار مستقیم بهره برده، صحت بازشناسی بهتری نسبت به یک سیستم بازشناسی تلفنی به تنهایی داشته باشد (نتایج جدول ۲ برای بازنمایی‌های LHCB). به این ترتیب در اولین مرحله یک گام به سمت مقاوم‌سازی سیستم بازشناسی گفتار برداشته‌ایم.

مدل بازشناسی دو منظوره از روی بردارهای بازنمایی گفتار که دربردارنده تنوعات مختلف از قبیل گوینده، آواها، لهجه‌ها و نوع کانال است، اطلاعات مشترک آوایی و نوع کانال تلفنی یا مستقیم (غیرمشترک) را شناسایی می‌کند. در معکوس سازی به روش گردایان با شروع از مقدار اولیه که همان بردار بازنمایی تلفنی است تدریجاً به سمت بردار بازنمایی جدیدی حرکت می‌کند که ضمن حفظ اطلاعات آوایی قبلی‌اش اثر کانال از روی بردارهای بازنمایی تلفنی تا حدی جبران گردد. بنابراین تنها یکی از تنوعات گفتار (نوع کانال) تا حدودی حذف شده است.

در روش معکوس‌سازی عمومی از روی اطلاعات آوایی و نوع کانال گفتار که در خروجی مدل دو منظوره بدست آمده بود، توسط شبکه معکوس، بردارهای بازنمایی دوباره بازنمایی می‌شوند. از آنجایی که در ورودی این شبکه غیر از اطلاعات کانال بسیاری از تنوعات گفتار تا حدودی حذف شده‌اند و برای اصلاح بردارهای بازنمایی نیز نوع کانال روی کانال گفتار مستقیم تنظیم می‌شود. بنابراین بردارهای بازنمایی جدید تلفنی و مستقیم بیشتر حاوی اطلاعات آوایی گفتار هستند و اثر کانال تلفنی نیز بصورت خاص از روی بردارهای بازنمایی تلفنی حذف شده است. به این ترتیب در روش معکوس‌سازی عمومی برای هر دو نوع داده اصلاح شده مستقیم و تلفنی، صحت بازشناسی بالاتری نسبت به شبکه مرجع بدست آمده است.

برای نشان دادن موفقیت این دو روش معکوس‌سازی در بهبود بردارهای بازنمایی گفتار تلفنی، شکل ۱۰ میزان افزایش صحت بازشناسی گفتار تلفنی برای شش دسته آوایی مختلف نسبت به شبکه مرجع را

Phoneme Recognitions + Channel Code (Tel/Direct)



شکل ۹: ساختار شبکه عصبی معکوس.

حال که بردارهای بازنمایی گفتار مستقیم و تلفنی به هم نزدیک شده‌اند. یک شبکه جدید بازشناسی آوا مشابه شبکه مرجع روی دو دسته از بردارهای بازنمایی گفتار مستقیم و تلفنی اصلاح شده تعلیم داده می‌شود. نتایج آزمون این شبکه در جدول (۳) با عنوان معکوس بروش گردایان آمده است. به این ترتیب نزدیک به ۱/۴٪ افزایش در صحت بازشناسی گفتار تلفنی بدست آمده است. در بازشناسی گفتار مستقیم با توجه به انحراف معیار ۰/۲٪ که در نتیجه بازشناسی شبکه مرجع داشتیم، رشدی بدست نیامده که این نتیجه با توجه به اینکه هیچ گونه اصلاحی روی بردارهای بازنمایی گفتار مستقیم ایجاد نشده بود طبیعی و قابل انتظار است.

## ۷-۴- اصلاح بردارهای بازنمایی با استفاده از معکوس

### سازای عمومی

در این روش، هدف این است که معکوس شبکه بازشناسی دو منظوره با یک مدل شبکه عصبی MLP دیگر تخمین زده شود. ساختار این شبکه در شکل ۹ رسم شده است و مبانی نظری مربوط به معکوس‌سازی عمومی نیز در بخش ۶-۲ تبیین شده است. شبکه دارای دو لایه پنهان با توابع غیرخطی تانژانت هایپربولیک و لایه خروجی با نورونهای خطی می‌باشد. سایر پارامترهای شبکه مشابه شبکه مرجع در نظر گرفته شده است. برای تعلیم این شبکه همانگونه که ذکر شد ورودیهای آن را خروجیهای مدل دو منظوره شکل ۸ و اهداف خروجی آن را بردارهای بازنمایی نظیر آنها قرار می‌دهیم.

پس از تعلیم شبکه معکوس، نتایج بازشناسی آواهای تلفنی شبکه مدل دو منظوره را هم برای دادگان تعلیم و هم برای دادگان آزمون به ورودی شبکه معکوس می‌دهیم و کد دوبیتی مربوط به نوع کانال گفتار را بجای "تلفنی" بودن "مستقیم" قرار می‌دهیم خروجیهای بدست آمده از شبکه معکوس یک مجموعه بردارهای بازنمایی جدید است که می‌توان گفت مجموعه بردارهای بازنمایی مستقیمی است که از روی دادگان تلفنی تخمین زده شده است. به همین ترتیب نتایج بازشناسی گفتار مستقیم در خروجی شبکه مدل دو منظوره را برای هم دادگان تعلیم و هم دادگان آزمون به شبکه معکوس می‌دهیم تا تخمین شبکه معکوس از بردارهای بازنمایی گفتار مستقیم نیز بدست آید.

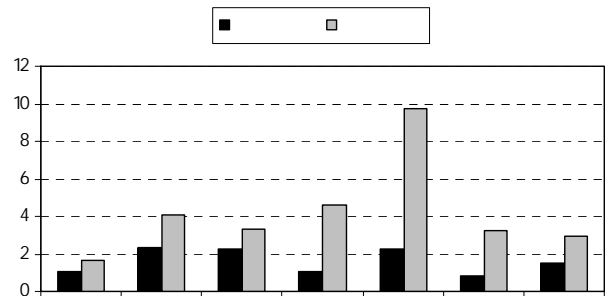
حال دادگان تعلیم جدید بدست آمده را به شبکه MLP دیگری مشابه شبکه مرجع، تعلیم می‌دهیم و صحت بازشناسی آواهای گفتار را برای دادگان آزمون جدید بررسی می‌نماییم. نتایج آزمون این شبکه جدید در جدول (۳) آمده است. افزایش صحت بازشناسی دادگان تلفنی در این روش نسبت به شبکه مرجع ۲/۸۹٪ و افزایش صحت در بازشناسی دادگان مستقیم برابر ۱/۶۸٪ می‌باشد. به این ترتیب روش معکوس‌سازی عمومی در اصلاح دادگان موفق‌تر از معکوس‌سازی به روش گردایان بوده است.



- [4] A. Martin, J. Fiscus, B. Fisher, D. Pallet, and M. Przybocki, "System Descriptions and Performance Summary," presented at the *Conversational Speech Recognition Workshop: DARPA Hub-5E Evaluation*, Baltimore, Maryland, US, May 1997.
- [5] D. Yuk and J. Flanagan, "Telephone speech recognition using neural networks and hidden Markov models," in *Proc. ICASSP*, pp. 157-160, 1999.
- [6] S. Thrun, "Is learning the n-th thing any easier than learning the first?" *Advances in Neural Information Processing Systems*, MIT Press, 1996.
- [7] S. Ben-David and R. Schuller, "Exploiting task relatedness for multiple task learning," *Lecture Notes in Computer Science*, vol. 2777, pp. 567 - 580, 2003.
- [8] C. W. Omlin and C. L. Giles, "Training second-order recurrent neural networks using hints," in *Proc. of the Ninth International Conference on Machine Learning.*, pp. 363-368, 1992.
- [9] S. Parveen and P. Green, "Multitask learning in connectionist robust ASR using recurrent neural networks," in *Proc. Eurospeech*, pp. 1813-1816, Geneva, Switzerland, Sep. 2003.
- [10] P., Niyogi and *et al.* "Incorporating prior information in machine learning by creating virtual examples," in *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2196-2209, Nov. 1998.
- [۱۱] الف. نژادقلی، *بازشناخت مقاوم گفتار نسبت به تنوعات مختلف گوینده در شبکه های عصبی بازشناخت گفتار*، پایان نامه کارشناسی ارشد، دانشگاه صنعتی امیرکبیر، دانشکده مهندسی پزشکی، ۱۳۸۲.
- [12] A. Blumer, A. Ehrenfeucht, D. Haussler, and M. Warmuth, "Learnability and the Vapnik-Chervon-Nenkis dimation," *J. Ass. Comput. Match.*, vol.36, no.4, pp. 929-965, 1989.
- [13] M. Bijankhan, J. Seikhzadeghan, M. R. Roohani, Y. Samareh, K. Lucas, M. Tebyani., "FARSDAT: the speech database of Farsi spoken language," in *Proc. SST-94*, pp. 826-831, Perth, Australia, 1994.
- [14] S. B. Davis and P. Mermelstein, "Comparison of parametric representations of monosyllabic word recognition in continuously spoken sentences," *IEEE Trans. ASSP*, vol. 28, no. 4, pp. 357-366, Aug. 1980.
- [۱۵] م. رحیمی نژاد، *توسعه و بهبود کیفیت روشهای استخراج پارامترهای بازنمایی در سیستم های بازشناخت گفتار*، پایان نامه کارشناسی ارشد، دانشگاه صنعتی امیرکبیر، دانشکده مهندسی پزشکی، ۱۳۸۱.
- [16] J. Han and W. Gao, "Robust telephone speech recognition based on channel compensation," *Journal of Pattern Recognition Society*, vol. 32, no.6, pp. 1061-1067, Jun. 1999.
- [17] C. A. Jensen, *et al.*, "Inversion of feedforward neural networks: algorithms and applications," *Proceedings of the IEEE*, vol. 87, no. 9, pp. 1536-1549, Sep. 1999.
- [18] R. J., Williams, "Inverting a connectionist network mapping by backpropagation of error," in *Proc 8th Annu. Conf. Cognitive Science Society*, pp. 859-865, 1986.

**منصور ولی** مدرک کارشناسی خود را در مهندسی برق- الکترونیک از دانشگاه صنعتی اصفهان در سال ۱۳۷۶ و مدرک کارشناسی ارشد را از دانشگاه صنعتی شریف در مهندسی برق- بیوالکترونیک در سال ۱۳۷۸ دریافت نموده است. وی در حال حاضر دانشجوی دوره دکتری مهندسی پزشکی- بیوالکترونیک در دانشگاه صنعتی امیرکبیر می باشد. زمینه های علمی مورد علاقه ایشان پردازش سیگنال های حیاتی، پردازش گفتار و شبکه های عصبی مصنوعی و زیستی می باشد.

**سید علی سید صالحی** مدرک کارشناسی خود را در مهندسی برق از دانشگاه صنعتی شریف در سال ۱۳۶۱، کارشناسی ارشد را در مهندسی برق از دانشگاه صنعتی امیرکبیر در سال ۱۳۶۷ و دکتری خود را در مهندسی برق- بیوالکترونیک از دانشگاه تربیت مدرس در سال ۱۳۷۴ دریافت نموده است. وی در حال حاضر استادیار دانشکده مهندسی پزشکی دانشگاه صنعتی امیرکبیر می باشد. زمینه های پژوهشی مورد علاقه ایشان پردازش و بازشناسی گفتار، شبکه های عصبی مصنوعی و زیستی، مدل سازی عملکرد مغز و پردازش خطی و غیرخطی سیگنال می باشد.



شکل ۱۰: درصد افزایش بازشناسی دسته های آوایی مختلف گفتار تلفنی.

نشان می دهد و بهترین بهبود در هر دو روش برای آوای انفجاری گفتار بدست آمده است. آن چه که مطابق شکل ۷ یکی از بیشترین دسته های آوایی تخریب شده توسط کانال تلفنی است.

نکته مهم و لازم به ذکر این است که اگرچه درصد حضور آوای مختلف در دادگان تعلیم و آزمون شبکه ها متفاوت هستند (به عنوان مثال بیشترین تعداد را واکه ها و کمترین تعداد را سایشی- انفجاری ها تشکیل می دهند) اما از آنجایی که در تمام مراحل این تحقیق، دادگان تعلیم و آزمون شبکه ها تغییر نکرده اند بلکه همان دادگان از روی دانش آموخته شده در شبکه اصلاح شده و دوباره به مدل جدید تعلیم داده شده اند بنابراین تفاوت درصد حضور آوای مختلف به توفیق این روش ها در بهبود بردارهای بازنمایی خدش های وارد نمی کنند.

## ۹- نتیجه گیری

در این تحقیق به منظور طراحی یک سیستم واحد بازشناسی گفتار مستقیم و تلفنی سعی گردید بردارهای بازنمایی گفتار بصورت بهینه انتخاب گردد. نشان داده شد که بردارهای بازنمایی طیفی LHCB برای بازشناسی گفتار تلفنی و بازشناسی توأم گفتار مستقیم و تلفنی در مدل های بازشناسی مبتنی بر شبکه های عصبی مناسب تر از بردارهای بازنمایی رایج MFCC هستند. آنگاه با استفاده از معکوس سازی شبکه های عصبی به روش گرادیان بردارهای بازنمایی گفتار تلفنی به سمت بردارهای بازنمایی گفتار مستقیم اصلاح گردید و با تعلیم شبکه دیگری روی دادگان اصلاح شده تلفنی و دادگان مستقیم دست نخورده، افزایش ۱/۴٪ در صحت بازشناسی گفتار تلفنی حاصل گردید. در مرحله بعد با استفاده از معکوس سازی عمومی شبکه های عصبی هر دو دسته بردارهای بازنمایی گفتار مستقیم و تلفنی به گونه ای اصلاح شدند که بیشتر حاوی اطلاعات آوایی گفتار باشند و سایر تنوعات تا جای ممکن حذف شوند. سپس با تعلیم شبکه دیگری روی این دادگان اصلاح شده، افزایش ۲/۹۸٪ در صحت بازشناسی گفتار تلفنی و ۱/۶۸٪ در صحت بازشناسی گفتار مستقیم بدست آمد.

## مراجع

- [1] S. Fouri, "Robust methods in automatic speech recognition and understanding," in *Proc. Eurospeech*, pp. 1993-1997, Geneva, Switzerland, 2003.
- [2] Y. Gong, "Speech recognition in noisy environments: A survey," *Speech Communication*, vol. 16, no. 3, pp. 261-291, Apr. 1995.
- [3] C. H. Lee and Q. Huo, "On adaptive decision rules and decision parameter adaptation for automatic speech recognition," in *Proceedings of the IEEE*, vol. 88, pp. 1241- 1269, Aug. 2000.