

پیش‌گویی قابلیت فهم همخوان‌ها در افراد دارای شنوایی عادی با استفاده از مدل‌های میکروسکوپی دارای معیار فاصله متفاوت در بازشناساگر خودکار گفتار

مسعود گراوانچی‌زاده^۱، علی فلاح^۲ و میرعلی اعتراف اسکویی^۳

^۱ دانشکده مهندسی برق و کامپیوتر - دانشگاه تبریز - تبریز - ایران

^۳ دانشکده توانبخشی - دانشگاه علوم پزشکی تبریز - تبریز - ایران

چکیده

در این مطالعه، نرخ تشخیص همخوان‌های موجود در ساختار هجایی «واکه - همخوان - واکه»، در آزمون‌های شنوایی و دو مدل میکروسکوپی ادراک گفتار مورد بررسی قرار می‌گیرد. چنین ساختار هجایی در زبان فارسی و ترکی آذری وجود ندارد؛ با وجود این، نتایج آزمون‌های شنوایی نشان می‌دهد که شنونده آذری یا فارسی‌زبان در شرایط بدون نوفه، قادر به تشخیص صحیح همخوان‌ها هستند. برای این پژوهش که در آن هدف، تشخیص صحیح آواها و نه کلمات بامعنی است، استفاده از این دادگان صوتی فاقد معنی مناسب است، چون با استفاده از این دادگان، دانش زبانی شنوندگان در پیش‌بینی کلمات نادیده گرفته می‌شود. نتایج آزمون‌های شنوایی با نتایج دو مدل میکروسکوپی که بر پایه دستگاه شنوایی انسان است، مقایسه می‌شود. تفاوت دو مدل در مرحله نهایی استخراج ویژگی به‌منظور استفاده در شناساگر خودکار گفتار DTW است. در مدل میکروسکوپی اول، در مرحله پایانی استخراج ویژگی، از فیلتر ۸ هرتز و در مدل دوم، از فیلتربانک مدولاسیون استفاده می‌شود. در ادامه، نرخ تشخیص صحیح آواها در مقادیر مختلف سیگنال به نوفه با استفاده از معیارهای فاصله اقلیدسی و لگاریتمی با یکدیگر مقایسه می‌شود. در این تحقیق، نرخ تشخیص همخوان‌ها برای شنونده آذری‌زبان مورد بررسی قرار گرفته است. در کنار جنبه تجربی این مطالعه، نوآوری این مقاله در بررسی دو معیار فاصله مختلف برای مدل هلوب و نیز مقایسه مستقیم دو مدل میکروسکوپی در پیش‌بینی میانگین نرخ تشخیص و نیز نرخ تشخیص تک‌تک همخوان‌ها است.

واژگان کلیدی: قابلیت فهم، ادراک گفتار، مدل میکروسکوپی، بردار ویژگی، نرخ تشخیص آوا، معیار فاصله، شناساگر خودکار گفتار.

۱- مقدمه

آزمون‌های شنوایی به‌عنوان دقیق‌ترین معیار ارزیابی در الگوریتم‌های بهبود کیفیت و جداسازی سیگنال‌های گفتار مطرح هستند. در کنار این آزمون‌ها، روش‌های ارزیابی بدون استفاده از شنونده انسانی وجود دارند که متخصصان شنوایی‌سنجی و مهندسان را از انجام آزمون‌های وقت‌گیر شنوایی بی‌نیاز می‌کند. دو روش ساده برای ارزیابی میزان بهبود کیفیت گفتار و جداسازی وجود دارند که روش اول، مقایسه نسبت سیگنال به نوفه^۱ ورودی با نسبت سیگنال به

نوفه خروجی و روش دوم، مقایسه ماسک باینری طراحی‌شده با ماسک باینری ایده‌آل است (Wang et al., 2005). مطالعاتی که بر روی معیارهای ارزیابی کیفیت گفتار صورت گرفته، نشان می‌دهد که بهبود در مقدار سیگنال به نوفه به‌حتم به معنی بهبود کیفیت گفتار نیست (Loizou, P., 2007)؛ از این رو مدل‌هایی مطرح شده‌اند که برای پیش‌گویی قابلیت ادراک گفتار مورد استفاده قرار گرفته‌اند. مطالعه میزان تأثیر نوفه بر قابلیت فهم گفتار منجر به ابداع معیارهای استانداردشده‌ای مانند شاخص بیان^۲ و

² Articulation Index (AI)

¹ Signal-to-Noise Ratio (SNR)

شاخص قابلیت فهم^۱ و شاخص انتقال گفتار^۲ شده است. در شاخص بیان و شاخص قابلیت فهم، مقدار سیگنال به نوفه در باندهای فرکانسی مختلف محاسبه می‌شود و این مقادیر سیگنال به نوفه پس از وزن دهی براساس اهمیت باندهای فرکانسی، با هم جمع می‌شوند. این مقدار به دست آمده، با استفاده از تبدیل غیرخطی وابسته به داده‌های گفتار در آزمون‌های شنوایی، به شاخص قابلیت فهم نگاشت می‌شود. در شاخص انتقال گفتار، تابع تبدیل مدولاسیون^۳ برای بررسی میزان مدولاسیون‌های تخریب‌شده ناشی از اکوستیک اتاق برای سیگنال گفتار، مورد استفاده قرار می‌گیرد. از شاخص انتقال گفتار برای پیش‌بینی قابلیت فهم گفتار در محیط‌های نوفه‌ای و پژواک‌دار استفاده شده است. به‌تازگی، روشی جدیدی در مرجع (Taghia et al 2014) از نقطه‌نظر تئوری اطلاعات برای پیش‌بینی قابلیت فهم ارائه شده است. در این روش، اطلاعات متقابل^۴ بین پوش زمانی سیگنال تمیز و سیگنال آلوده در باندهای فرکانسی مختلف تخمین زده می‌شود. تمامی این مدل‌های نامبرده شده، مدل‌های ماکروسکوپی نامیده می‌شوند. علت این نام‌گذاری، پیش‌بینی نرخ متوسط تشخیص جملات و عدم بررسی نرخ‌های تشخیص واحدهای کوچک‌تر گفتار مانند آوا، سیلاب و کلمه است. در کنار مدل‌های ماکروسکوپی، مدل‌های میکروسکوپی پیشنهاد شده‌اند که در آن‌ها سعی شده است از پردازشی مشابه دستگاه شنوایی انسان استفاده شود. همچنین در این مدل‌ها، برخلاف مدل‌های ماکروسکوپی، نرخ تشخیص برای واحدهای کوچک‌تر گفتار نیز بررسی می‌شود. ضعف مدل‌های میکروسکوپی در پیچیدگی و زمان لازم برای پیاده‌سازی است (Bronkhorst et al., 2000). دو مدل میکروسکوپی پیش‌بینی ادراک گفتار بر مبنای مدل اولیه پردازش شنوایی گوش انسان، پیشنهاد شده‌اند. این مدل اولیه پیش‌تر توسط داو (Dau et al., 1996) ارائه شده است. الگوی استخراج‌شده از مراحل پیش‌پردازش به‌عنوان بردار ویژگی در شناساگر خودکار گفتار^۵ (Sakoe et al., 1978) استفاده می‌شود. همچنین، از این بردار ویژگی در دستگاه بازشناسی گفتار با استفاده از مدل مخفی مارکوف^۶ نیز استفاده شده است (Tchortz et al., 1999). هلوب

(Holube et al., 1996)، با استفاده از مدل اولیه پردازش شنوایی انسان و شناساگر گفتار DTW، اولین مدل میکروسکوپی را برای پیش‌بینی قابلیت فهم آواها، ارائه کرده است. یورگنز (Jürgen et al., 2009)، با حفظ ساختار کلی مدل هلوب و با ایجاد تغییراتی در بردار ویژگی، مدل میکروسکوپی جدیدی ارائه نموده است. موضوع جدیدی که در مدل یورگنز مطرح است، استفاده از معیار فاصله متفاوت در DTW است. بر حسب این معیار، نرخ تشخیص تغییر می‌کند و در نتیجه یکی از مدل‌ها برای برآزش داده‌های آزمون شنوایی مناسب‌تر خواهد بود. تاکنون مقایسه‌ای در مورد نرخ تشخیص این دو مدل و نیز نوع معیار فاصله مناسب در آنها، صورت نگرفته است. در این مقاله، نرخ تشخیص آواها برای افراد دارای شنوایی عادی مورد بررسی قرار می‌گیرد. در ادامه، نرخ تشخیص آواها در دو مدل میکروسکوپی مورد مطالعه قرار خواهد گرفت و مدل و معیار فاصله بهینه تعیین می‌شود. همچنین، نرخ تشخیص تک‌آواها در دو مدل مقایسه خواهد شد.

بخش‌بندی ادامه مقاله به این ترتیب است که در بخش دوم، ساختار مدل‌های شنوایی، بازشناساگر DTW، دادگان مورد استفاده و شرایط انجام تست شنوایی با جزئیات شرح داده می‌شود. در ادامه و در بخش سوم، نتایج تست شنوایی و شبیه‌سازی با رویکرد مقایسه مدل‌ها و تأثیر معیار فاصله انتخابی بر میزان پیش‌بینی قابلیت فهم، مورد تحلیل قرار می‌گیرد. در پایان و در بخش چهارم نیز بحث و نتیجه‌گیری‌های لازم از نتایج بالا لحاظ شده است.

۲- ساختار مدل‌ها

در شکل (۱-۲)، مدل ادراک گفتار پیشنهادشده توسط هلوب نشان داده شده است (Holube et al., 1996). در قسمت بالای این مدل، مرحله آموزش و در قسمت پایین، مرحله آزمون داده‌های صوتی آمده است. سیگنال‌های گفتار آلوده به نوفه، به‌عنوان ورودی‌های آموزش و آزمون، به مدل اعمال می‌شوند.

در اولین مرحله پیش‌پردازش، سیگنال از فیلتربانک گاماتون عبور داده می‌شود که این فیلتربانک مدلی برای فیلترینگ سیگنال در حلقون گوش انسان است. فیلتربانک گاماتون با استفاده از روش هومن (Hohmann, 2002) پیاده‌سازی شده است. در فیلتربانک گاماتون پاسخ ضربه فیلتر به‌صورت زیر است:

¹ Speech Intelligibility Index (SII)

² Speech Transmission Index (STI)

³ Modulation Transfer Function (MTF)

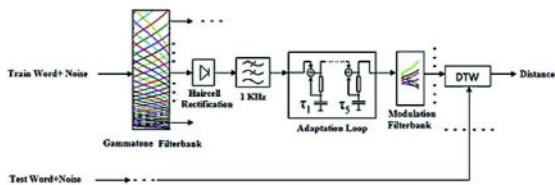
⁴ Mutual Information

⁵ Automatic Speech Recognizer (ASR)

⁶ Dynamic Time Warping (DTW)

⁷ Hidden Markov Model (HMM)

سیگنال‌های غیریابستان، مانند سیگنال گفتار، شروع^۴ و پایان^۵های سیگنال در خروجی حلقه‌ها، به‌طور برجسته‌تر تقویت می‌شود. هر مسیر بازگشتی در این حلقه‌ها شامل یک مقسم و یک فیلتر پایین‌گذر با ثابت زمانی τ_i است. سیگنال ورودی به حلقه بر خروجی فیلتر پایین‌گذر تقسیم می‌شود. برای سیگنال با مقدار ثابت C ، سیگنال خروجی برای هر حلقه انطباق به‌صورت $O = \sqrt{C}$ است. بنابراین خروجی پنج حلقه برابر $O = C^{1/32}$ خواهد بود که به این صورت عملیات فشرده‌سازی لگاریتمی سیگنال ورودی تخمین زده می‌شود. با تغییر ثابت‌های زمانی با استفاده از ظرفیت خازن‌ها، می‌توان پاسخ سامانه به نوسانات سیگنال ورودی را تغییر داد. اگر نوسانات سیگنال ورودی در مقایسه با ثابت‌های زمانی حلقه‌های انطباق سریع‌تر باشد، سیگنال به‌طور تقریبی بدون تغییر از حلقه‌ها عبور می‌کند و در غیر این صورت، سیگنال فشرده می‌شود؛ حتی هنگامی که سیگنال تحریک در ورودی حلقه‌ها قطع می‌شود، به دلیل پردازش غیرخطی تا مدت زمانی در خروجی، تحریک وجود دارد. بنابراین، فشرده‌سازی سیگنال و انطباق زمانی^۱ به‌طور هم‌زمان انجام می‌گیرد. در این پردازش غیرخطی، اطلاعات مربوط به پوش سیگنال حفظ می‌شود. در مدل هلوب، مرحله قبل از DTW، عبور سیگنال از فیلتر پایین‌گذر با ثابت زمانی بیست میلی ثانیه (که به‌طور تقریبی معادل فیلتر هشت هرتز با شیب شش دسیبل در هر اکتاو می‌باشد) است. شماتیک مدل دوم که مربوط به کار یورگنز است، در شکل (۲-۲) نشان داده شده است (Jürgen et al., 2009). تفاوت آن با مدل هلوب، در مرحله پیش پردازش قبل از DTW است.



(شکل ۲-۲): شمای مدل پیش‌بینی ادراک گفتار ارائه شده توسط یورگنز. مراحل پیش‌پردازش، به جز مرحله آخر که از فیلتربانک مدولاسیون به جای فیلتر ۸ هرتز استفاده شده است، مشابه مدل هلوب است. IR در این مدل ماتریس دوبعدی است: بعد اول برای فیلتربانک گاماتون و بعد دوم برای فیلتربانک مدولاسیون تعریف می‌شود (Jürgen, et al., 2009).

⁴ Onset
⁵ Offset
⁶ Adaptation

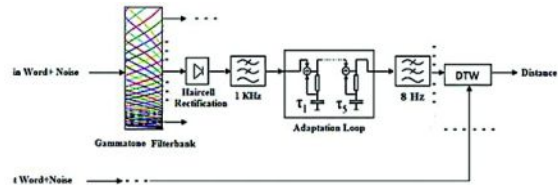
$$g_r[n] = n^{-1} \tilde{a}^n \quad n \geq 0, \quad (1)$$

$$\tilde{a} = \lambda \cdot \exp(j\beta),$$

که در آن γ مرتبه فیلتر، β فرکانس نوسان و λ پارامتر میرایی یا پهنای باند است. این پاسخ ضربه، با سری قراردادن چهار فیلتر مرتبه اول میان‌گذر با ضرایب مختلط پیاده‌سازی می‌شود. پاسخ فرکانسی یک فیلتر گاماتون از رابطه زیر به‌دست می‌آید:

$$H(z) = \frac{1}{(1 - \tilde{a}z^{-1})^4}. \quad (2)$$

برای استخراج بردارهای ویژگی، سیگنال ابتدا از ۲۷ فیلتر گاماتون عبور داده می‌شود. این فیلترها، طبق مقیاس پهنای باند مستطیلی معادل^۱ و یک فیلتر به‌ازای هر پهنای باند مستطیلی معادل در محور فرکانس، توزیع شده‌اند. این فیلتربانک، بازه فرکانسی ۲۳۴ هرتز تا ۸ کیلوهرتز را پوشش می‌دهد. در ادامه، خروجی هر فیلتر، از یکسوساز نیم‌موج عبور داده می‌شود. پس از عبور از یکسوساز نیم‌موج، سیگنال توسط فیلتر پایین‌گذر مرتبه اول با فرکانس قطع یک کیلوهرتز، فیلتر می‌شود. یکسوساز نیم‌موج و فیلتر یک کیلوهرتز، برای مدل‌سازی سلول‌هایی مویی^۲ گوش استفاده شده‌اند. در مرحله بعدی، خروجی مدل سلول‌های مویی از پنج حلقه انطباق^۳ متوالی با ثابت‌های زمانی $\tau_1 = 5 \text{ ms}$, $\tau_2 = 50 \text{ ms}$, $\tau_3 = 129 \text{ ms}$, $\tau_4 = 253 \text{ ms}$ و $\tau_5 = 500 \text{ ms}$ ، که برای مدل‌سازی پردازش‌های غیر خطی گوش استفاده شده است، مطابق با (Dau et al., 1996)، عبور داده می‌شود. حلقه‌های انطباق، سیگنال‌های زمانی ایستان را به‌طور تقریبی به‌صورت لگاریتمی فشرده می‌سازند. در مورد



(شکل ۲-۱): شمای مدل پیش‌بینی قابلیت فهم هلوب.

پیش‌پردازش شامل عبور از فیلتربانک گاماتون، یکسوسازی نیم‌موج و حلقه‌های انطباق است که برای مدل‌سازی برخی از آزمون‌های شنوایی مورد استفاده قرار گرفته است. برای آزمایش‌های مربوط به ادراک گفتار از شناساگر DTW استفاده می‌شود که برای سنجش فاصله بین کلمه الگو با کلمه آزمون مورد استفاده قرار می‌گیرد (Holube et al., 1996).

¹ Equivalent Rectangular Bandwidth (ERB)
² Haircells
³ Adaptation Loop

در مدل یورگنز، از فیلتربانک مدولاسیون^۱ مطابق کار داو (Dau et al., 1997) استفاده شده است. این فیلتربانک، شامل چهار کانال مدولاسیون به‌ازای هر کانال فرکانسی گامتون است که یک فیلتر پایین‌گذر با فرکانس قطع ۲/۵ هرتز و سه فیلتر میان‌گذر با فرکانس‌های مرکزی ۵، ۱۰ و ۱۶/۷ هرتز را در بر می‌گیرد. پهنای باند برای فیلتر پایین‌گذر و فیلترهای با فرکانس‌های مرکزی ۵ و ۱۰ هرتز، برابر ۵ هرتز و برای فیلتر ۱۶/۷ برابر ۸/۳ هرتز است. فیلترهای مدولاسیون با استفاده از رابطه بازگشتی زیر پیاده‌سازی می‌شوند:

$$y[n] = e^{-\pi\beta\Delta} e^{j2\pi f_0\Delta} y[n-1] + (1 - e^{-\pi\beta\Delta}) x[n], \quad (3)$$

که در آن β پهنای باند فیلتر، Δ پریود نمونه‌برداری و f_0 فرکانس مرکزی فیلتر است. تابع تبدیل رابطه (۳) برابر است با:

$$H(z) = \frac{1 - e^{-\pi\beta\Delta}}{1 - e^{-\pi\beta\Delta} e^{j2\pi f_0\Delta} z^{-1}}, \quad (4)$$

که تابع تبدیل فیلتر میان‌گذر مرتبه اول است.

پس از کاهش نرخ نمونه‌برداری به مقدار صد هرتز، نمایش داخلی سیگنال^۲ به‌دست می‌آید. بنابراین، نمایش‌های داخلی سیگنال گفتار برای مدل‌های هلوب و یورگنز، به‌ترتیب، به‌صورت بردار یک‌بعدی و ماتریس دوبعدی در هر ده میلی‌ثانیه خواهد بود. واحد عناصر این بردار و ماتریس، واحد مدل^۳ نامیده می‌شود. یک واحد مدل معادل با یک دسیبل شدت صوت است (Dau et al., 1997).

۲-۱- بازشناساگر خودکار DTW

در شکل‌های (۱-۲) و (۲-۲)، مسیر پردازشی بالا برای سیگنال مرحله یادگیری و مسیر پایین برای سیگنال آزمون در نظر گرفته شده است. الگوهای به‌دست‌آمده از این دو مسیر، توسط شناساگر خودکار گفتار DTW (Sakoe et al., 1978) مقایسه می‌شوند. در مدل‌های میکروسکوپی، از DTW برای حالتی که تفاوت سیگنال الگو و آزمون تنها در نوبه افزوده‌شده به آنهاست، به‌منظور پیش‌بینی قابلیت فهم گفتار استفاده شده است. یک تبدیل زمانی^۴ خاص برای انطباق بهینه نمایش‌های داخلی سیگنال الگو و سیگنال آزمون، با محاسبه ماتریس فاصله D به‌دست می‌آید. هر عنصر $D(i, j)$ در این ماتریس، از محاسبه فاصله بین ماتریس

ویژگی مربوط به نمایش داخلی الگو در شاخص زمانی i و ماتریس ویژگی مربوط به نمایش داخلی سیگنال آزمون در شاخص زمانی j به‌دست می‌آید. سپس، یک مسیر به‌هم پیوسته از روی این ماتریس فاصله به‌دست می‌آید. این مسیر به‌گونه‌ای انتخاب می‌شود که مقدار نهایی فاصله بین دو الگوی آموزش و آزمون در نقطه انتهایی مسیر حداقل شود. با مقایسه تمامی سیگنال‌های آزمون با سیگنال الگو، سیگنال دارای کمترین فاصله، به‌عنوان الگوی تشخیص داده‌شده تعیین می‌شود.

۲-۲- الگوی فاصله

ساده‌ترین رویکرد در DTW استفاده از فاصله اقلیدسی بین بردارهای ویژگی الگو ($IR_{template}$) و آزمون (IR_{test}) است:

$$D_{Euclidean} = \sqrt{\sum_f \sum_{f_{mod}} (IR_{templ}(i, f, f_{mod}) - IR_{test}(i, f, f_{mod}))^2}, \quad (5)$$

که در آن، f بیان‌گر کانال فرکانسی مربوط به فیلتربانک گامتون و f_{mod} مربوط به باند فرکانسی مدولاسیون در مدل یورگنز است. عمل جمع‌بندی در رابطه (۵) در مدل هلوب، تنها برای خروجی تمامی فیلترهای گامتون صورت می‌گیرد؛ و شاخصی برای مدولاسیون در محاسبه فاصله وجود ندارد. در حالتی که اختلاف عناصر ماتریس الگو و ماتریس آزمون دارای توزیع گوسی باشد، فاصله اقلیدسی بهینه خواهد بود؛ اما هنگامی که این توزیع گوسی نباشد، فاصله اقلیدسی دیگر بهترین معیار، و به‌عبارتی ساده‌تر، کم‌ترین مقدار فاصله نخواهد بود. در مدل یورگنز، برای ماتریس‌های ویژگی نشان داده شده است، هنگامی که تفاوت سیگنال الگو و سیگنال آزمون تنها در نوبه افزوده‌شده به آنهاست، فاصله بین این ماتریس‌های ویژگی از توزیع کوشی تبعیت می‌نماید. برای توزیع کوشی، معیار فاصله لگاریتمی بهینه است (Jürgens et al., 2009). فاصله لگاریتمی از رابطه زیر به‌دست می‌آید:

$$D_{Log} = \sum_f \sum_{f_{mod}} \log \left(1 + \frac{1}{2} (IR_{templ}(i, f, f_{mod}) - IR_{test}(i, f, f_{mod}))^2 \right). \quad (6)$$

با استفاده از دو معیار فاصله و دو مدل میکروسکوپی می‌توان چهار حالت را برای بررسی نتایج مربوط به نرخ‌های تشخیص برای پیش‌بینی قابلیت فهم گفتار در نظر گرفت.

¹ Modulation Filter Bank (MFB)

² Internal Representation (IR)

³ Model Unit (MU)

⁴ Time-Transformation

ده فرد دارای شنوایی طبیعی (نه مرد و یک زن) که بین ۲۲ تا ۲۹ سال دارند، در آزمون‌های شنوایی شرکت کرده‌اند.

۲-۴- شرایط آزمون شنوایی و آزمون مدل

دادگان آزمون به دسته‌هایی تقسیم‌بندی می‌شوند که تنها در آوای وسط با یکدیگر متفاوت باشند. برای آزمون‌های شنوایی، فرد آزمون‌شونده پس از شنیدن فایل صوتی با استفاده از موش‌واره، یکی از گزینه‌ها را از روی صفحه نمایش انتخاب می‌کند. این انتخاب از بین یکی از سیزده گزینه موجود صورت می‌گیرد. از این رو شرایط آزمون به صورت ارزیابی بسته^۳ است. شنونده امکان بازپخش مجدد یک داده صوتی را برای دفعات دلخواه دارد. از آنجایی که دانش زبانی افراد در آزمون شنوایی با توجه به بی‌معنابودن مواد گفتار نادیده گرفته می‌شود، امکان مقایسه دقیق‌تر مدل‌ها با آزمون شنوایی وجود دارد. همچنین، علاوه بر مقایسه نرخ تشخیص آواها، می‌توان ماتریس ابهام^۴ را نیز در مورد آنها مورد بررسی قرار داد.

مقدار جذر میانگین مربعات^۵ سیگنال گفتار با حذف بازه‌های سکوت در ابتدا و انتهای فایل‌های صوتی در شصت دسیبل تنظیم می‌شود. نوفه ایستان که دارای طیف زمان-طولانی مشابه سیگنال گفتار^۶ است (Dreschler et al., 2001)، در شدت خاصی که بر مبنای نسبت سیگنال به نوفه^۷ مورد نظر تنظیم شده است، چهارصد میلی‌ثانیه قبل از شروع پخش سیگنال گفتار به داده‌های صوتی افزوده می‌شود. زمان چهارصد میلی‌ثانیه به‌منظور مقداردهی اولیه به خروجی حلقه‌های انطباق است. برای سیگنال آلوده به نوفه، در ابتدا و انتها از شیب هن^۷ دارای بازه زمانی صد میلی‌ثانیه استفاده می‌شود. بعد از محاسبه نمایش داخلی برای سیگنال تولیدشده، بردارهای ویژگی مربوط به چهارصد میلی‌ثانیه پخش نوفه پیش از شروع گفتار، حذف می‌شوند. این عمل به‌منظور در نظر گرفتن اطلاعات لازم از زمان پخش آواها صورت می‌گیرد.

برای ارزیابی مدل، نمونه‌های از نوفه در نسبت سیگنال به نوفه^۷ مورد نظر به هر کدام از دادگانی که تفاوت آنها تنها در هجای وسط است، افزوده می‌شود و نمایش

برای پیاده‌سازی فیلتربانگ گاماتون، از کدهای موجود در پایگاه اینترنتی دانشگاه Oldenburg استفاده شده که در آن روش هومن (Hohmann, 2002) پیاده‌سازی شده است؛ همچنین، مدل سلول‌های مویی و حلقه‌های انطباق، با استفاده از فایل‌های MEX موجود در تارنما همین دانشگاه شبیه‌سازی شده است. کدهای مربوط به فیلتر هشت هرتز، فیلتربانگ مدولاسیون و DTW توسط نویسندگان مقاله در محیط C و در قالب MEX تهیه شده‌اند.

۲-۳- داده‌های گفتاری و افراد شرکت‌کننده

در آزمون شنوایی

دادگان گفتاری، دارای ساختار «واکه-همخوان-واکه»^۱ هستند که از پایگاه دادگان OLLO (Wesker et al., 2005) موجود در تارنمای دانشگاه الدنورگ^۲ انتخاب شده‌اند. ساختار هجایی «واکه-همخوان-واکه» که توسط گویندگان آلمانی ادا شده است، به‌طور معمول در زبان‌های فارسی و ترکی آذری وجود ندارد؛ اما از آنجایی که هدف از استفاده از این دادگان، تنها تشخیص همخوان است، نتایج آزمون‌های شنوایی که در ابتدا برای سیگنال تمیز بدون نوفه انجام گرفت، نشان داد که شنوندگان آذری زبان در تشخیص همخوان‌ها در این ساختار هجایی مشکلی ندارند و استفاده از این دادگان در آزمون‌های شنوایی برای سیگنال‌های نوفه‌ای، بلامانع است.

دادگان OLLO از نوع «واکه-همخوان-واکه» دارای واژه‌های یکسان در ابتدا و انتها هستند. برای مقایسه مستقیم نتایج مدل‌ها با آزمون‌های شنوایی، از سیگنال‌های صوتی مشابهی در مدل و آزمون استفاده می‌شود. فایل‌های صوتی انتخاب‌شده، مربوط به گوینده مرد با شناسه S10M_NO و دارای سرعت بیان طبیعی آواها است. همخوان‌های مورد استفاده یکی از سیزده آوای /د/، /ت/، /گ/، /ک/، /ف/، /س/، /ب/، /پ/، /ل/، /م/، /ن/، /ش/ و /ل/ است که در بین یکی از پنج واژه /آه/، /اه/، /ای/، /او/، /ا/ قرار می‌گیرد. در نتیجه ۶۵ فایل صوتی متفاوت برای آزمون شنوایی و آزمون مدل‌ها خواهیم داشت. از همخوان /ts/ پایگاه داده OLLO که دارای صدایی مابین صدای دو همخوان /ت/ و /س/ می‌باشد، به‌علت ناآشنایی شنوندگان با آن استفاده نشده است.

³ Closed Test

⁴ Confusion Matrix

⁵ Root-Mean-Square

⁶ ICRA noise

⁷ Hann

¹ Vowel-Consonant-Vowel (VCV)

² Oldenburg

۳- شبیه‌سازی و نتایج آزمون‌های شنوایی

۱-۳- میانگین نرخ تشخیص

شکل (۱-۳) میانگین نرخ تشخیص آواها به درصد برحسب مقدار سیگنال به نوفه برای تمامی آواها را نشان می‌دهد. میله‌های خطا^۳ در شکل (۱-۳)، انحراف معیار برای ده شنونده را نشان می‌دهد. تابع روان‌سنجی برازش شده برای تشخیص صحیح همخوان‌ها، دارای شیب ۷٪/dB و آستانه ادراک گفتار برابر ۹/۹۲- دسیبل است. نتایج آزمون‌های شنوایی با نتایج به‌دست‌آمده از اندازه‌گیری‌های یورگنز مطابقت خوبی دارد. مقدار میانگین قدر مطلق خطا بین اندازه‌گیری‌ها و مقاله بالا کمتر از ۵٪ است. بیشترین مقدار خطا در نسبت سیگنال به نوفه ۲۰- دسیبل و در حدود ۱۰٪ است که بیشتر ناشی از شرایط آزمایشگاهی و افراد استفاده‌شده در آزمون شنوایی است.

در شکل (۲-۳)، نتایج مربوط به شبیه‌سازی‌های دو مدل که در هر کدام از مدل‌ها نیز از دو معیار فاصله استفاده شده، رسم شده است. با توجه به شکل (۲-۳)، استفاده از معیار فاصله متفاوت در مدل هلوب، تأثیر چندانی در تغییر نرخ تشخیص ندارد؛ اما در مدل یورگنز، معیار فاصله، تأثیر قابل ملاحظه‌ای در نرخ تشخیص دارد. استفاده از فاصله لگاریتمی، نرخ تشخیص را به‌ویژه در سیگنال به نوفه‌های میانی نسبت به فاصله اقلیدسی بهبود می‌بخشد. البته این بهبود به معنی مناسب‌بودن مدل نیست؛ چون ممکن است مدل دارای نرخ تشخیص پایین‌تر، نتایج آزمون شنوایی را بهتر پیش‌بینی کند.

برای مقایسه مدل‌ها با نتایج آزمون شنوایی، پارامترهای تابع روان‌سنجی برازش شده برای هر کدام از منحنی‌های شکل (۲-۳) با پارامترهای تابع روان‌سنجی مربوط به آزمون شنوایی، مقایسه شده است. پارامترهای حاصل از برازش تابع روان‌سنجی در آزمون‌های شنوایی و مدل‌ها با دو معیار فاصله مختلف در جدول (۱-۳) نشان داده شده است. کمترین تفاوت در آستانه ادراک گفتار در مدل‌ها نسبت به آستانه ادراک گفتار در آزمون شنوایی، مربوط به مدل یورگنز با استفاده از فاصله لگاریتمی است. استفاده از فاصله اقلیدسی در مدل یورگنز، نتیجه بدتری در پیش‌بینی مقدار آستانه ادراک گفتار در مقایسه با مدل هلوب دارد.

داخلی سیگنال نوفه‌ای به‌عنوان الگویی که می‌بایست تشخیص داده شود، ذخیره می‌شود. نوفه‌ای متفاوت از نوفه اضافه‌شده به الگو، این بار به هر کدام از سیگنال‌های گفتار آزمون افزوده می‌شود که یکی از این سیگنال‌های آزمون همان سیگنال استفاده شده در الگوی ذخیره‌شده است. با محاسبه فاصله هر سیگنال آزمون با الگوهای ذخیره‌شده و محاسبه الگوی دارای کم‌ترین فاصله در DTW، نرخ تشخیص صحیح آواها در سیگنال به نوفه‌های مختلف محاسبه می‌شود. برای افزایش دقت نتایج شبیه‌سازی در مدل‌ها، برای هر کدام از دادگان الگو، ده بار شبیه‌سازی تکرار می‌شود. در نتیجه با داشتن ۶۵ فایل صوتی، در مجموع ۶۵۰ بار مقایسه الگو با سیگنال‌های آزمون صورت می‌گیرد. حداکثر مدت شبیه‌سازی مربوط به مدل یورگنز با استفاده از فاصله لگاریتمی است که این زمان در حدود پنج ساعت است.

۲-۵- آزمون‌های شنوایی

نرخ تشخیص سیزده هجا با استفاده از دستگاه شنوایی‌سنجی AC40 که قابلیت پخش سیگنال گفتار در شدت دسیبل دلخواه را دارد، در آزمایشگاه گفتار دانشکده توانبخشی دانشگاه تبریز که مجهز به یک اتاقک عایق در برابر صداست، انجام گرفته است. از نسبت سیگنال به نوفه‌های صفر، -۵، -۱۰، -۱۵ و -۲۰ دسیبل در آزمون‌های شنوایی و آزمون مدل‌ها استفاده شده است. ترتیب پخش هجاها به‌صورت تصادفی است. شنونده می‌تواند یکی از سیزده گزینه موجود را انتخاب کند. قبل از انجام آزمون شنوایی، تمامی افراد با داده‌های صوتی آشنا شده‌اند. برای توصیف میانگین قابلیت فهم برای تمامی هجاها، تابع مدل یا تابع روان‌سنجی^۱ که با رابطه زیر توصیف می‌شود:

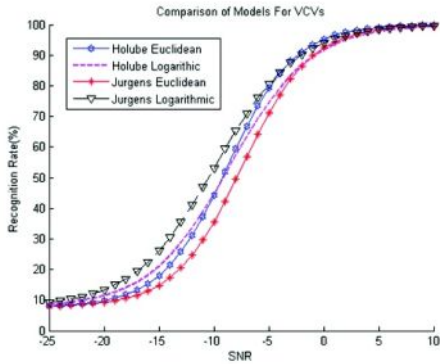
$$P(L, L_{50}, S_{50}) = \frac{1-g}{1 + \exp(4 \cdot S_{50} \cdot (L_{50} - L))} + g, \quad (7)$$

به میانگین نرخ تشخیص شنوندگان در تمامی سیگنال به نوفه‌ها، برازش می‌شود که در آن S_{50} شیب تابع روان‌سنجی، L_{50} آستانه ادراک گفتار^۲، L مقدار سیگنال به نوفه و g احتمال تشخیص صحیح یک آوا به‌صورت تصادفی است که برابر ۰/۰۷ است. عملیات برازش با استفاده از cftool نرم‌افزار متلب انجام می‌شود.

³ Error Bars

¹ Psychometric

² Speech Reception Threshold (SRT)



شکل ۳-۲: نتایج پیش‌بینی نرخ تشخیص در مدل‌ها برای پیش‌بینی مقدار قابلیت فهم گفتار با استفاده از تابع روان‌سنجی برازش‌شده بر شبیه‌سازی‌ها.

جدول ۳-۱: فهرست پارامترهای تخمین‌زده‌شده تابع روان‌سنجی برازش‌شده بر نرخ تشخیص آواها در اندازه‌گیری‌ها و دو مدل میکروسکوپی با دو معیار فاصله.

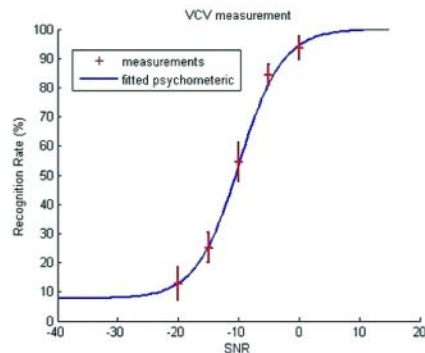
	SRT (dB)	اختلاف با SRT آزمون شنوایی	شیب (%/dB)	پیرسون (r^2)
اندازه‌گیری‌ها	-۹/۹۲۸۲	صفر	۷	۱
مدل هلوب با فاصله اقلیدسی	-۸/۷۴۰۷	۱/۱۹	۸/۳	۰/۹۷۸۵
مدل هلوب با فاصله لگاریتمی	-۸/۵۶۴۱	۱/۳۶	۶/۸۴	۰/۹۸۵۵
مدل یورگنز با فاصله اقلیدسی	-۷/۴۰۹۶	۲/۵۱	۸/۱۵	۰/۹۵۹۰
مدل یورگنز با فاصله لگاریتمی	-۹/۸۸۹۳	۰/۰۳۸۹	۶/۷۷	۰/۹۹۵۳

۳-۲- نرخ تشخیص هر کدام از آواها

نمودار میله‌ای شکل (۳-۳)، نتایج مربوط به نرخ تشخیص هر کدام از همخوان‌ها در آزمون‌های شنوایی را نشان می‌دهد. حداکثر نرخ تشخیص مربوط به آواهای /ات، /اس/ و /اش/ است. نرخ تشخیص برای آواهای /ام، /او، /اب/ و /اک/ به‌طور متوسط در تمامی سیگنال به نوفه‌ها، پایین‌تر از بقیه است. در شکل (۳-۴)، نتایج مربوط به نرخ تشخیص دو مدل رسم شده است. نمودار میله‌ای اول در شکل (۳-۴)، مربوط به مدل هلوب با استفاده از فاصله اقلیدسی و نمودار دوم در این شکل، نتایج مدل یورگنز با استفاده از فاصله لگاریتمی است. این نمودارها نشان می‌دهند که نرخ تشخیص دو مدل برای سه آوای /ات، /اس/ و /اش/ مشابه آزمون‌های شنوایی به‌طور مشخص بیشتر از بقیه آواهاست.

مدل هلوب با استفاده از فاصله اقلیدسی دومین پیش‌بینی نزدیک به آستانه ادراک گفتار اندازه‌گیری شده را دارد. همچنین، در پیش‌بینی شیب تابع روان‌سنجی، مدل هلوب با فاصله لگاریتمی و مدل یورگنز با فاصله لگاریتمی، نزدیک‌ترین پیش‌بینی را نسبت به آزمون‌های شنوایی دارد. مربع ضرایب همبستگی پیرسون^۱ (r^2) در ستون آخر جدول (۳-۱) آمده است. ضریب همبستگی پیرسون (r) بین دو جمعیت آماری از میانگین‌گیری حاصل ضرب مشاهدات هنجارسازی‌شده دو جمعیت به‌دست می‌آید. مشاهده هنجارسازی‌شده برای هر جمعیت از اختلاف مقدار مشاهده از میانگین جمعیت تقسیم بر انحراف معیار جمعیت به‌دست می‌آید. مقدار «یک» برای مربع این ضریب نشان‌دهنده همبستگی کامل و مقدار «صفر» نشان‌دهنده عدم همبستگی بین داده‌های دو جمعیت است. بیشترین همبستگی با مقادیر آزمون شنوایی در مدل یورگنز با فاصله لگاریتمی مشاهده می‌شود.

استفاده از معیار فاصله متفاوت، در مدل هلوب تفاوت چندانی را در نتایج نشان نمی‌دهد (اختلاف کمتر از ۰/۲ دسیبل است)؛ اما در مدل یورگنز، تفاوت مقادیر آستانه ادراک گفتار پیش‌بینی‌شده با استفاده از دو فاصله مختلف، در حدود ۲/۵ دسیبل است.



شکل ۳-۱: نتایج آزمون شنوایی برای اندازه‌گیری متوسط مقدار قابلیت فهم گفتار بر حسب مقدار سیگنال به نوفه. در آزمون شنوایی از نوفه دارای طیف مشابه سیگنال گفتار (ICRA noise) استفاده شده است. میله‌های خطای مربوط به انحراف معیار نرخ تشخیص بین ده شنونده و نیز تابع روان‌سنجی برازش‌شده بر اندازه‌گیری‌ها در شکل نشان داده شده است.

^۱ Pearson



(جدول ۳-۲): مربع ضرایب همبستگی پیرسون بین ستون‌های ماتریس ابهام آزمون شنوایی با ماتریس‌های ابهام مدل هلوب و یورگنز در $SNR = -15dB$

مدل یورگنز با فاصله	مدل هلوب با فاصله	همخوان
لگاریتمی	اقلیدسی	اد/
۰/۰۰۷	۰/۰۲۲	اد/
۰/۸۷	۰/۱۵	ات/
۰/۱۰۴	۰/۰۰۶	اک/
۰/۰۵۶	۰/۵۰	اک/
۰/۰۵۳	۰/۰۵۴	اف/
۰/۹۶	۰/۸۶	اس/
۰/۰	۰/۲۲	اب/
۰/۱۳	۰/۰۰۹	اپ/
۰/۱۵	۰/۰۲۳	او/
۰/۰۷۶	۰/۰۳	ام/
۰/۰۰۰۱	۰/۰۰۲	ان/
۰/۹۲	۰/۹۵	اش/
۰/۰۶۹	۰/۰۰۵	ال/

بدتر شدن پیش‌بینی مدل با تغییر فاصله اقلیدسی به فاصله لگاریتمی است، در حالی که در مدل یورگنز عکس این روند مشاهده می‌شود.

یکی از مزیت‌های مدل‌های میکروسکوپی در مقایسه با مدل‌های ماکروسکوپی امکان ارزیابی نرخ تشخیص تک‌تک آواها با استفاده از نمودار میله‌ای نرخ تشخیص صحیح و نیز ماتریس ابهام است. نرخ تشخیص آواهای ات، اس، ا و اش/ در آزمون‌های شنوایی دارای مقادیر بیشتری نسبت به بقیه آواها به‌ویژه در سیگنال به نوفه‌های پایین هستند. مدل‌های میکروسکوپی شبیه‌سازی‌شده در این پژوهش نیز نتایج مشابهی را در تشخیص صحیح این سه همخوان نسبت به بقیه نشان می‌دهند. بیشتر بودن نرخ تشخیص صحیح مدل یورگنز در مقایسه با مدل هلوب در سیگنال به نوفه‌های کم منجر به افزایش شیب تابع روان‌سنجی در مدل هلوب می‌شود. مدل یورگنز علاوه بر تخمین دقیق‌تر مقدار آستانه ادراک گفتار آزمون شنوایی، شیب تابع روان‌سنجی را نیز با اختلاف کمتری در مقایسه با مدل هلوب تخمین می‌زند.

پژوهش‌هایی که در ادامه کار پیشنهاد می‌شود، مطالعه پیش‌بینی قابلیت فهم در حضور پژواک و گویندگان دیگر است. پژواک می‌تواند قابلیت فهم سیگنال گفتار را تا حد زیادی کاهش دهد. ارائه مدل میکروسکوپی که بتواند علاوه بر پیش‌بینی قابلیت فهم در شرایط ایده‌ال، برای شرایط پژواک‌دار نیز تخمین مناسبی از آستانه ادراک گفتار داشته باشد، می‌تواند قدم بعدی در تکمیل مدل‌های میکروسکوپی ادراک گفتار باشد.

تقدیر و تشکر

نویسندگان این مقاله از هم‌فکری خانم نگین صالحی در دانشکده توانبخشی دانشگاه علوم پزشکی تبریز در استفاده از آزمایشگاه صوت و انجام آزمون‌های شنوایی کمال تشکر را دارند.

۵- مراجع

Bronkhorst, A. W., 2000. The Cocktail Party Phenomenon: A Review of Research on Speech Intelligibility in Multiple Talker Conditions, Acta Acoustical United With Acustica, vol. 86, pp. 117–128.

Dau, T., Püschel, D. and Kohlrausch, A., 1996. A quantitative model of the “effective” signals proce-

۴- بحث و نتیجه‌گیری

در این پژوهش، نرخ تشخیص همخوان‌ها برای افراد دارای شنوایی طبیعی با استفاده از آزمون شنوایی در سیگنال به نوفه‌های مختلف، اندازه‌گیری شد. در ادامه، نتایج دو مدل میکروسکوپی در پیش‌بینی نرخ تشخیص صحیح آواها مورد بررسی قرار گرفت. با مقایسه توابع روان‌سنجی برازش‌شده بر آزمون شنوایی با توابع روان‌سنجی برازش‌شده بر مدل‌ها، نشان داده شد که هر دو مدل قابلیت پیش‌بینی صحیح مقادیر آستانه ادراک گفتار در آزمون‌های شنوایی را دارند. در مورد متوسط نرخ تشخیص برای تمامی آواها، دقیق‌ترین پیش‌بینی مربوط به مدل یورگنز با فاصله لگاریتمی است. در این مدل با تغییر معیار فاصله لگاریتمی به اقلیدسی، اختلاف آستانه ادراک گفتار مدل با آستانه ادراک گفتار آزمون شنوایی از مقدار ناچیز 0.0389 دسیبل به مقدار $2/51$ دسیبل افزایش می‌یابد. از این رو انتخاب معیار فاصله در مدل یورگنز اهمیت زیادی دارد. همچنین، با تغییر معیار فاصله در مدل هلوب، اختلاف زیادی در مقادیر آستانه ادراک گفتار به‌وجود می‌آید. تفاوت آستانه ادراک گفتار مدل هلوب با استفاده از معیار اقلیدسی با آستانه ادراک گفتار آزمون شنوایی برابر $1/19$ دسیبل است. با تغییر معیار فاصله اقلیدسی به فاصله لگاریتمی در این مدل، این اختلاف به $1/36$ دسیبل می‌رسد. نکته قابل توجه در مورد مدل هلوب،

فصل ۵



Wesker, T., Meyer, B., Wagener, K., Anemüller, J., Mertins, A., and Kollmeier, B., 2005. Oldenburg log-atom speech corpus (OLLO) for speech recognition experiments with humans and machines, in Proceedings of Interspeech, pp. 1273-1276, Lisbon, Portugal.



مسعود گراوانچی‌زاده مدرک

کارشناسی را در رشته مهندسی الکترونیک در سال ۱۳۶۵ از دانشگاه تبریز اخذ کرده است. سپس، مدارک کارشناسی ارشد

و دکتری در رشته پردازش سیگنال را به ترتیب، در سال‌های ۱۳۷۴ و ۱۳۸۰ از دانشگاه Ruhr-University Bochum کسب کرده است. ایشان هم‌اکنون عضو هیأت علمی دانشکده مهندسی برق و کامپیوتر دانشگاه تبریز با مرتبه علمی دانشیاری هستند. زمینه‌های تحقیقاتی مورد علاقه ایشان بهبود کیفیت گفتار، مکان‌یابی و جداسازی سیگنال گفتار، پردازش سیگنال‌های تصادفی و پردازش سیگنال دوگوشی است.

نشانی رایانامه ایشان عبارت است از:

geravanchizadeh@tabrizu.ac.ir



علی فلاح مدارک کارشناسی و

کارشناسی ارشد در رشته برق گرایش مخابرات را به ترتیب در سال‌های ۱۳۸۵ و ۱۳۸۸ از دانشگاه‌های صنعتی امیرکبیر و دانشگاه شاهد اخذ کرده است. وی

هم‌اکنون دانشجوی دکتری رشته مهندسی برق گرایش مخابرات سیستم در دانشگاه تبریز است. زمینه‌های پژوهشی مورد علاقه ایشان پردازش سیگنال آرایه‌ای، پردازش سیگنال گفتار، پردازش سیگنال دوگوشی و مدل‌های ادراک شنوایی است.

نشانی رایانامه ایشان عبارت است از:

ali.fallah@tabrizu.ac.ir

ssing in the auditory system: I. Model structure, J. Acoust. Soc. Am, vol. 99, pp. 3615-3622.

Dau, T. and Kohlrausch, A., 1997. Modeling auditory processing of amplitude modulation I. Detection and masking with narrowband-carriers, J. Acoust. Soc. Am, vol. 102, pp. 2893-2905.

Dreschler, W. A., Verschuure, H., Ludvigsen, C. and Westermann, S., 2000. ICRA noises: Artificial noise signals with speech-like spectral and temporal properties for hearing instrument assessment, Audiology, vol. 40, pp. 148-157.

Hohmann, V., 2002. Frequency analysis and synthesis using a gammatone filterbank, Acta Acoustical United With Acustica, vol. 88, pp. 433-442.

Holube, I. and Kollmeier, B., 1996. Speech Intelligibility Prediction in Hearing-Impaired Listeners Based on a Psychoacoustically Motivated Perception Model, J. Acoust. Soc. Am, vol. 100, pp. 1703-1716.

Jürgens, T. and Brand, T., 2009. Microscopic prediction of speech recognition for listeners with normal hearing in noise using an auditory model, J. Acoust. Soc. Am, vol. 126, pp. 2635-2648.

Loizou, P., 2007. Speech Enhancement: Theory and Practice, CRC Press, Boca Raton: FL.

Press, W., Teukolsky, S., Vetterling, W. T. and Flannery, B. P., 1992. Numerical Recipes in C, Cambridge University, Press.

Sakoe, H. and Chiba, S., 1978. Dynamic programming algorithm optimization for spoken word recognition, IEEE Trans. Acoust, Speech, Signal Process, vol 26, pp. 43-49.

Tchorz, J. and Kollmeier, B., 1999. A model of an auditory perception as front end for automatic speech recognition, J. Acoust. Soc. Am, vol. 106, pp. 2040-2050.

Taghia, J. and Martin, R., 2014. Objective intelligibility measures based on mutual information for speech subjected to speech enhancement processing, IEEE Trans. Audio, Speech, Signal Process, vol 22, pp. 6-16.

Wang, D. L. and Brown, G., 2005. Computational Auditory Scene Analysis: Principals, Algorithms and Applications, New York: IEEE Press, Wiley-Interscience.





میر علی اسکوئی مدارک

کارشناسی و کارشناسی ارشد را در رشته فیزیوتراپی، به ترتیب، در سال‌های ۱۳۷۰ و ۱۳۷۳ از دانشگاه علوم پزشکی ایران و تهران اخذ کرده است. سپس، با اعطای بورسیه

تحصیلی خارج از کشور در سال ۱۳۸۰ عازم کشور کانادا شده و در سال ۱۳۸۵ دکترای تخصصی و Post Doctorate Fellowship خود را در زمینه بیومکانیک عضله از University of Calgary کشور کانادا اخذ کرد. ایشان هم‌اکنون دانشیار دانشکده توانبخشی دانشگاه علوم پزشکی تبریز بوده و زمینه‌های مورد علاقه ایشان بیومکانیک عضله، بیماری‌های عضلانی اسکلتی، ارگونومی و پردازش سیگنال‌های الکتریکی عضله است.

نشانی رایانامه ایشان عبارت است از:

eterafoskouei@tbzmed.ac.ir