

بازشناسی مقاوم گفتار با دادگان مفقود با استفاده از روشی مبتنی بر کوراریانس

حجت محمدنژاد^۱، منصور ولی^۲ ایمان اسمعیلی^۳

^۱دانشکده فنی دانشگاه شاهد، mohammadnejad@shahed.ac.ir

^۲دانشکده فنی دانشگاه شاهد، vali@shahed.ac.ir

^۳دانشکده فنی دانشگاه شاهد، iesmaili@shahed.ac.ir

چکیده - سیستم‌های بازشناسی گفتار در مواجهه با گفتارهایی که در اثر نویز تخریب شده‌اند عملکرد ضعیفی از خود نشان می‌دهند. روش‌های دادگان مفقود سعی دارند که تاثیر نویز را با حذف یا تخمین مؤلفه‌های بازنمایی نامعتبر تخریب شده توسط نویز، جبران کنند و بازشناسی را با استفاده از مؤلفه‌های معتبر باقیمانده و تخمین زده شده انجام دهند. این روش‌ها، بردارهای ویژگی ورودی به مدل بازشناسی را بدون اینکه در سیستم بازشناسی تغییری ایجاد کنند، بهبود می‌بخشند. در این تحقیق ما یک روش جبران ویژگی را پیشنهاد خواهیم داد که در آن، مؤلفه‌های نامعتبر در نمایش اسپکتروگرام (نمایش زمانی-فرکانسی) با مؤلفه‌های تخمینی بدست آمده از مؤلفه‌های معتبر و اطلاعات آماری که نسبت به دادگان تمیز وجود دارد، جایگزین می‌شوند. در حالت کلی، روش‌های جبران سازی در اثر استخراج ویژگی‌های کپستروم از اسپکتروگرام بازسازی شده، در مواجهه با نویز غیرایستاد، نرخ بهبود زیادی در صحت بازشناسی دارند. روش‌های بازسازی دادگان مفقود از لحاظ محاسباتی کم‌هزینه بوده، چرا که نیاز به آموزش چندباره مدل بازشناسی برای دادگان نویزی و اصلاح مدل را از بین می‌برد.

کلید واژه- بازشناسی مقاوم گفتار، جبران سازی مبتنی بر ویژگی، دادگان مفقود

نسبت سیگنال به نویز (SNR) بالا و به تناسب، برخی نواحی

دیگر با معیار سیگنال به نویز پایین را نمایش می‌دهد. [۲]

در این روش‌ها نواحی تخریب شده توسط نویز با معیار SNR پایین در نمایش اسپکتروگرام گفتار، به عنوان نواحی مفقود شناسایی شده و نامعتبر فرض می‌شوند و با نمایش ماسک از اسپکتروگرام پاک می‌شوند. سپس از نواحی معتبر با معیار SNR بالا، برای جبران نواحی مفقود استفاده می‌شود.

بیشتر تکنیک‌های بازشناسی گفتار که از روش ویژگی‌های مفقود بهره می‌گیرند مبتنی بر اصلاح شیوه‌ای است که در آن، سیستم بازشناسی، احتمال کلاس‌ها و یا حالت‌ها را برای دستیابی به ویژگی‌های نامعتبر محاسبه می‌کند [۳]. این روش‌ها مدل بازشناسی را محدود می‌کنند که در حوزه لگاریتم طیف، جائیکه نواحی مفقود مشخص شده‌اند، عمل بازشناسی را صورت دهد و این یک اشکال جدی است زیرا اگر چه روش‌های ویژگی‌های مفقود، رویکردهای بسیار موثری هستند اما صحت بازشناسی، زمانی که از ویژگی‌های لگاریتم طیفی برای بازشناسی استفاده می‌شود ضعیف تر از سایر ویژگی‌های مورد استفاده نظیر ضرایب کپستروم می‌باشد.

۱- مقدمه

صحت بازشناسی، برای گفتاری که توسط نویز تخریب شده است کاهش چشمگیری می‌یابد مخصوصاً اگر سیستم بازشناسی بر روی دادگان تمیز آموزش داده شده باشد. روش‌های متعددی برای جبران اثرات نویز بر روی سیستم‌های بازشناسی گفتار مطرح شده‌اند. تقریباً بیشتر آنها نویز را ایستاد فرض می‌کنند بنابراین در بازشناسی گفتارهای نویزی غیر ایستاد، کارایی نخواهند داشت [۱].

روش‌های مبتنی بر ویژگی‌های مفقود بطور خاص این نوید را می‌دهند که قابلیت جبران اثرات نویز غیر ایستاد را دارا هستند. روش‌های دادگان مفقود، الگوریتم جبران سازی هر گونه از نویز را شامل می‌شوند و قابلیت خوبی در بازشناسی مقاوم گفتار در حضور سطوح بالایی از نویز، از خود نشان داده‌اند.

این روش بر اساس این حقیقت است که چون انرژی سیگنال گفتار و نویز در باندهای فرکانسی متمایز، مختلف می‌باشند، نویز در نواحی متمایز اسپکتروگرام تاثیر متفاوتی می‌گذارد. به عبارتی دیگر اسپکتروگرام گفتار نویزی، برخی نواحی با

قابل اجرا باشد نباید هیچ دانش اولیه‌ای برای شناسایی مولفه های نامعتبر بکار گرفته شود و به همین دلیل یکی از مشکل ترین قسمت های روش های مبتنی بر ویژگی های مفقود، شناسایی مولفه های معتبر و نامعتبر در نمایش اسپکتروگرام است. برای این منظور می بایست طیف نویز موجود در تک تک مولفه ها را تخمین زد و از این معیار برای شناسایی مولفه های نامعتبر استفاده کرد که به این رویکرد، روش مبتنی بر SNR گفته می شود.

برای تخمین SNR مولفه های طیفی، تخمینی از طیف توان نویز مخرب نیاز است. برای این منظور چند فریم اول هر گویش بعنوان سکوت فرض شده و میانگین طیف توان آنها، بعنوان طیف توان نویز اولیه در نظر گرفته می شود [۴]. در ادامه برای بدست آوردن تخمینی از تغییرات آهسته نویز برای هر مولفه در نمایش اسپکتروگرام، می توان از رابطه بازگشتی (۱) استفاده کرد [۵].

$$\hat{N}_p(m, k) = \begin{cases} (1-\lambda)\hat{N}_p(m-1, k) + \lambda Y_p(m, k) \\ \text{if } Y_p(m, k) < \beta \hat{N}_p(m, k) \\ \hat{N}_p(m-1) \\ \text{otherwise} \end{cases} \quad (1)$$

در رابطه فوق، λ و β بترتیب برابر ۰/۹۵ و ۲ می باشند $Y_p(m, k)$ و $N_p(m, k)$ بترتیب نشان دهنده طیف k امین باند فرکانسی سیگنال گفتار نویزی و سیگنال نویز در m امین فریم می باشند. λ معیار دنبال کننده تغییرات کند و تند نویز می باشد. معیارهای مختلفی برای شناسایی مولفه های نامعتبر از روی طیف تخمینی نویز پیشنهاد شده اند. یکی از این معیارها که توسط Al-Maleki ارائه شده، معیار انرژی منفی می باشد که تشخیص مولفه های نامعتبر، مبتنی بر این اساس است که انرژی در آن مولفه کمتر از انرژی تخمینی نویز باشد و یا عبارت دیگر $Y_p(m, k)$ نامعتبر خواهد بود اگر رابطه (۲) برقرار باشد [۶]:

$$|Y_p(m, k)| \leq |\hat{N}_p(m, k)| \quad (2)$$

از طرفی دیگر اگر SNR هر مولفه ای زیر صفر dB باشد معیار SNR ، باندهای طیفی را نامعتبر تشخیص می دهد. برای تخمین SNR ، تخمین طیف سیگنال تمیز نیز مورد نیاز است. این امر بوسیله کسر طیف نویز از طیف سیگنال نویزی میسر می شود که رابطه (۳) آن را نمایش می دهد [۷].

روش های مبتنی بر ویژگی های مفقود، در عمل از یک دانش اولیه از خواص آماری گفتار تمیز که توسط سیستم بازشناسی گفتار، مدل شده اند استفاده می کنند. [۳]

در این تحقیق ما یک روش بازسازی نواحی تخریب شده از نمایش اسپکتروگرام توسط نویز را به عنوان یک مرحله پیش پردازش قبل از بازشناسی مطرح می کنیم که از آن به عنوان بازسازی مبتنی بر همبستگی (correlation-based reconstruction) و یا بازسازی مبتنی بر کواریانس (covariance-based reconstruction) یاد می شود و در نهایت نمایش اسپکتروگرام بازسازی شده می تواند به یک مجموعه ویژگی انتخابی، نظیر کپستروم تبدیل شده و توسط مدل بازشناسی استاندارد مورد استفاده قرار گیرد. این روش مزایای دیگری نیز دارد که نیاز به انواع مدل بازشناسی را از بین برده و امکان استفاده از اطلاعاتی که بطور آشکارا در مدل بازشناسی بکار گرفته نمی شوند، نظیر ارتباط زمانی بین مؤلفه های بردارهای لگاریتم طیفی، فراهم می شود.

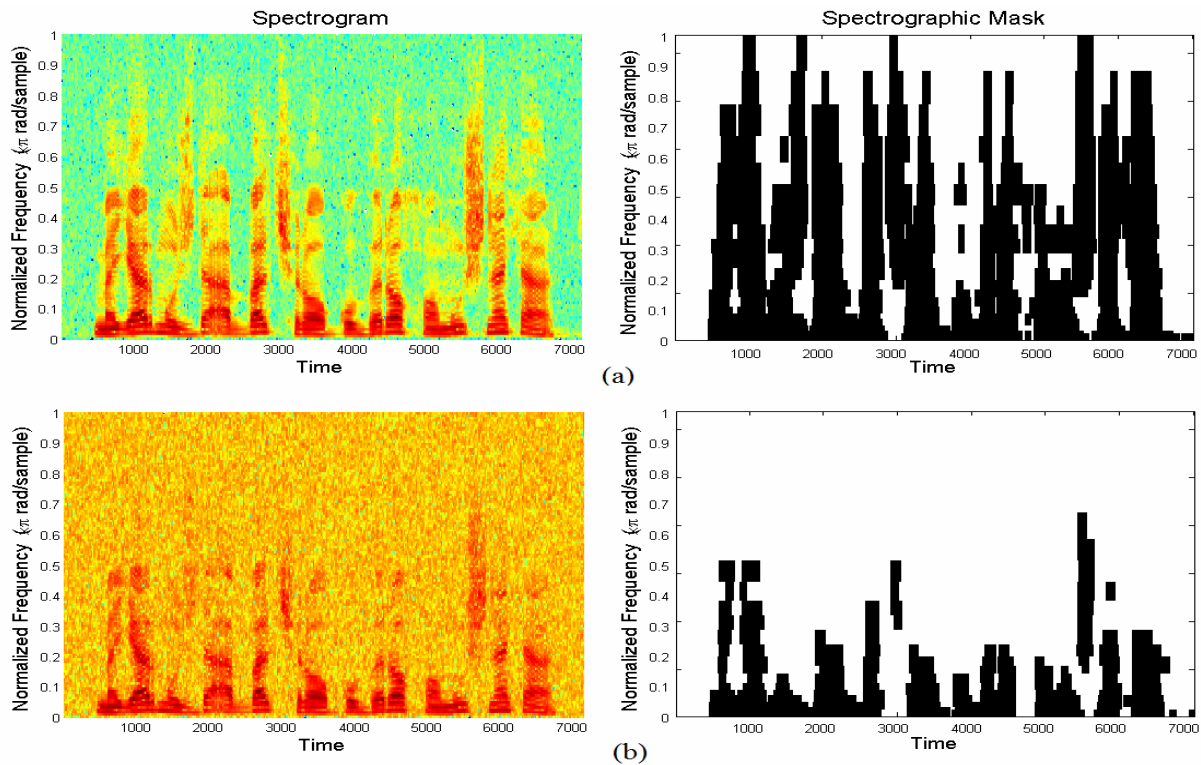
الگوریتمی که در این تحقیق مطرح شده است، از ساده ترین مشخصات آماری توزیع لگاریتم طیف گفتار تمیز به منظور بازسازی ناحیه های تخریب شده از اسپکتروگرام استفاده می کند. نتایج ارزیابی نشان می دهد که این روشها زمانیکه نواحی تخریب شده از اسپکتروگرام، بطور کامل تشخیص داده شوند، می توانند بطور قابل توجهی در بازشناسی، موثر باشند.

در روش های بازسازی خودکار مبتنی بر ویژگی های مفقود، لازم است که نواحی تخریب شده از اسپکتروگرام هم بطور خودکار تشخیص داده شوند. در این تحقیق از دو معیار متفاوت انرژی منفی و همچنین SNR برای تشخیص نواحی مفقود استفاده شده است. ارزیابی ها نشان می دهد که استفاده از این دو معیار در ترکیب با روشهای مبتنی بر ویژگی های مفقود مطرح شده در این تحقیق، بهبود قابل توجهی در صحت بازشناسی گفتار آلوده به نویز سفید خواهد داشت.

در بخش ۲ نحوه شناسایی مولفه های مفقود از روی اسپکتروگرام سیگنال مطرح می شود. در بخش ۳ روش بازسازی مبتنی بر کواریانس، شرح داده می شود. در بخش ۴ عملکرد این روش برای SNR های مختلف ارزیابی شده و در نهایت نتیجه گیری از مقاله در بخش ۵ صورت خواهد گرفت.

۲- شناسایی مولفه های مفقود

برای اینکه روش های مبتنی بر ویژگی های مفقود در عمل



شکل ۱: (a) نمایش اسپکتروگرام و ماسک اسپکتروگرام یک گویش تمیز (b) نمایش اسپکتروگرام و ماسک اسپکتروگرام همان گویش تخریب شده توسط نویز سفید با SNR معادل ۱۰ dB

باشد. برای نویزهای غیر ایستا و زودگذر این روش منجر به نتایج ضعیفتری در شناسایی نواحی مفقود اسپکتروگرام خواهد شد [۴].

۳- روش بازسازی مبتنی بر کواریانس

یکی دیگر از روش‌های مبتنی بر ویژگی‌های مفقود روش بازسازی مبتنی بر کواریانس می‌باشد. در این روش بردارهای طیفی لگاریتمی استخراج شده از گفتار، نمونه‌های یک فرایند تصادفی گوسی ایستا فرض می‌شوند. دانش اولیه درباره سیگنال‌های گفتار تمیز بصورت پارامترهای آماری نظیر مقدار میانگین بردارها و کواریانس بین مولف‌های آنها، نمایش داده می‌شوند.

$\mu(k)$ نشان دهنده میانگین k امین باند فرکانسی از m امین فریم بردار طیفی $X(m, k)$ می‌باشد. کواریانس بین k_1 امین باند فرکانسی از m امین بردار طیفی $X(m, k_1)$ و k_2 امین باند فرکانسی از $m + \xi$ امین بردار طیفی $X(m + \xi, k_2)$ ، نیز با $c(\xi, k_1, k_2)$ نمایش داده می‌شود که مقدار نرمالیزه شده نظیر آن $r(\xi, k_1, k_2)$ می‌باشد که از طریق روابط زیر بدست می‌آیند:

$$\hat{X}_p(m, k) = \begin{cases} Y_p(m, k) - \hat{N}_p(m, k) & \text{if } Y_p(m, k) - \hat{N}_p(m, k) > \gamma Y_p(m, k) \\ \gamma Y_p(m, k) & \text{otherwise} \end{cases} \quad (3)$$

که در آن $X_p(m, k)$ تخمین k امین مولفه طیفی سیگنال تمیز را در m امین فریم نمایش می‌دهد. طبق معیار SNR ، اگر طیف تخمینی سیگنال تمیز از طیف تخمینی نویز کمتر باشد مولفه طیفی $Y_p(m, k)$ نامعتبر فرض می‌شود که در رابطه (۴) نشان داده شده است.

$$\hat{X}_p(m, k) < \hat{N}_p(m, k) \quad (4)$$

در عمل، بهترین شناسایی ویژگی‌های مفقود زمانی حاصل می‌شود که دو معیار انرژی منفی و معیار SNR بصورت ترکیبی در شناسایی مولفه‌های نامعتبر بکار گرفته شوند. شکل ۱ نمونه‌ای از نمایش اسپکتروگرام و ماسک اسپکتروگرامی را که بوسیله تخمین نویز بدست آمده نشان می‌دهد.

بطور معمول، شناسایی ویژگی‌های مفقود براساس تخمین نویز، زمانی عملی و موثر است که نویز، ایستا و یا شبه ایستا

۴- ارزیابی

در این تحقیق روش بازسازی ویژگی های مفقود مبتنی بر کواریانس بر روی بخشی از مجموعه دادگان فارس دات پیاده سازی شده است که شامل ۲۰۰ گویش از ۱۰۰ گوینده مختلف می باشد.

در این تحقیق، از ضرایب کپستروم فرکانس مل استخراج شده از ویژگی های لگاریتم طیفی و مشتقات آن که از مجموعه دادگان تمیز بدست آمده اند، برای آموزش مدل بازشناسی HMM استفاده شده است.

برای در نظر گرفتن محیط های نویزی مختلف، سیگنال گفتار تمیز را بصورت مجازی با SNR های ۲۰ الی صفر dB توسط نویز سفید گوسی تخریب کرده و روش جبران سازی مبتنی بر کواریانس روی آن اعمال شده است. دادگان مفقود در سیگنال های گفتار نویزی توسط ماسک اسپکتروگرام، شناسایی شده و توسط الگوریتم Bounded MAP بازسازی شده اند. بردارهای طیفی بازسازی شده که بیانگر ویژگی های لگاریتم طیفی فرکانس مل هستند به ضرایب کپستروم فرکانس مل (MFCC) تبدیل شده و با افزودن مشتقات اول و دوم آن ها، بردار بازنمایی برای بازشناسی دادگان نویزی فراهم می شود.

مدل بازشناسی HMM توسط ۱۴۰ گویش برچسب گذاری شده، برای هر واج از مجموعه دادگان فارس دات آموزش داده می شود. ۳۴ مدل مخفی مارکوف برابر با ۳۴ واج فارسی، برای هر دسته از بردارهای بازنمایی آموزش داده شده است. مدل های مخفی مارکوف شامل سه حالت چپ به راست به همراه شانزده تابع چگالی گوسی برای هر حالت هستند. پارامترهای آماری نظیر میانگین، کواریانس متقابل و کواریانس نرمالیزه شده شرح داده شده در بخش ۳، برای بازسازی ویژگی های مفقود، از مجموعه دادگان تمیز استخراج می شوند.

مجموعه دادگان تست شامل ۶۰ گویش نویزی آلوده به نویز سفید گوسی با SNR های ۲۰ الی صفر dB می باشد. جدول ۱ و شکل ۲، صحت بازشناسی را برای دادگان نویزی و همچنین دادگان نویزی اصلاح شده توسط روش دادگان مفقود مبتنی بر کواریانس نمایش می دهد.

با توجه به جدول ۱ مشهود است که جبران سازی با روش مبتنی بر ویژگی های مفقود عملکرد قابل توجهی در محیط های نویزی با SNR های مختلف بر روی صحت بازشناسی مدل داشته است.

$$\mu(k) = E[X(m,k)]$$

$$c(\xi, k_1, k_2) = E(X(m, k_1) - \mu_{k_1})(X(m + \xi, k_2) - \mu_{k_2}) \quad (5)$$

$$r(\xi, k_1, k_2) = \frac{c(\xi, k_1, k_2)}{\sqrt{c(\xi, k_1, k_1)c(\xi, k_2, k_2)}}$$

در رابطه فوق [۱] عملگر امید ریاضی را نشان می دهد. این پارامترها از بردارهای طیفی استخراج شده از مجموعه دادگان تمیز بدست می آیند.

برای تخمین مولفه های طیفی نامعتبر در m امین بردار طیفی $X(m)$ ، همه ی آن ها را در بردار $X_u(m)$ گرد می آوریم سپس همه ی مولفه های طیفی معتبر در نمایش اسپکتروگرام را که دارای کواریانس نرمالیزه شده ی حداقل ۰/۵ با حداقل یکی از مولفه های بردار $X_u(m)$ باشند را شناسایی و در بردار $X_r^{(n)}(m)$ مرتب می کنیم. مقدار میانگین و کواریانس بین مولفه های $X_u(m)$ و $X_r^{(n)}(m)$ و همچنین کواریانس متقابل بین آنها بوسیله پارامترهای میانگین و کواریانس حاصل از دانش اولیه ی فرایند گوسی مورد نظر ساخته می شوند. نظر به اینکه فرایند گوسی فرض شده است، مقدار صحیح مولفه های نامعتبر $X_u(m)$ را می توان بوسیله الگوریتم MAP تخمین زد که در زیر رابطه نهایی آن آورده شده است.

$$\hat{X}_u(m) = \mu_u + C_{ru} \cdot C_{rr}^{-1} \cdot (Y_r(m) - \mu_r) \quad (6)$$

که در آن μ_u و μ_r بترتیب نشان دهنده مقدار میانگین مولفه های $X_u(m)$ و $X_r^{(n)}(m)$ و همچنین C_{ru} و C_{rr} بترتیب کواریانس و کواریانس متقابل بین مولفه های $X_u(m)$ و $X_r^{(n)}(m)$ می باشند. $\hat{X}_u(m)$ نیز معرف مقادیر بازسازی شده بردار ویژگی های مفقود نظیر $X_u(m)$ می باشد.

برای بازشناسی بهتر لازم است مقادیر بازسازی شده ی ویژگی های مفقود در یک حیطه خاص محدود شوند که آن را در ترکیب با رابطه فوق اصطلاحاً Bounded MAP می نامند. رابطه (۷) نحوه اعمال این محدودیت را نشان می دهد.

$$\hat{X}_u(m,k) = \begin{cases} \hat{X}_u(m,k) & \text{if } \hat{X}_u(m,k) \leq X_u(m,k) \\ X_u(m,k) & \text{otherwise} \end{cases} \quad (7)$$

از آنجایی که سیگنال نویزی در اثر اضافه شدن نویز مخرب به گفتار تمیز حاصل شده است در نتیجه مقادیر لگاریتم طیفی سیگنال گفتار نویزی را می توان به عنوان یک کران بالا برای مقادیر بازسازی شده در نظر گرفت.

جدول ۱: صحت بازشناسی برای دادگان تمیز و نویزی اصلاح نشده و نویزی بازیابی شده در ازای SNR های ۲۰ الی ۰ dB در حیطه لگاریتم طیفی

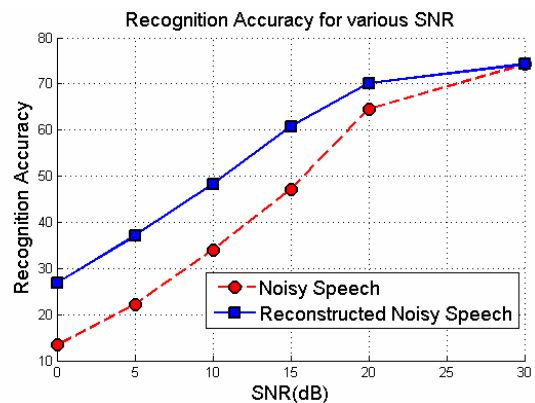
Recognition Accuracy	Clean	SNR 20	SNR 15	SNR 10	SNR 5	SNR 0
Noisy Speech	74.22	64.50	47.22	33.99	22.23	13.51
Noisy Reconstructed Speech	74.22	70.14	60.74	48.29	37.07	26.75

در صورتی که (۱) خاصیت کران بالایی برای جایگزینی مقادیر بازسازی شده در نمایش اسپکتروگرام سیگنال گفتار نویزی در نظر گرفته شود و (۲) روند بازسازی، مبتنی بر جداسازی طیفی باشد، بازسازی بهتر و به تبع آن بازشناسی بهتری بدست خواهد آمد.

ما نشان دادیم که روش مبتنی بر ویژگی‌های مفقود، در بازسازی مقادیر مفقود در مواجهه با نویز ایستا، در حد بسیار بالایی موثر است. بطوریکه برای سیگنال با $SNR=0$ dB به میزان ۱۳/۲۵٪ افزایش در صحت بازشناسی گفتار حاصل شده است.

مراجع

- [1] X. Xiong, "Speech Enhancement with Applications in Speech Recognition" A First Year Report Submitted to the School of Computer Engineering of the Nanyang Technological University. 2006
- [2] R. P. Lippman, "Using Missing Feature Theory to Actively Select Features for Robust Speech Recognition with Interruptions, Filtering and Noise", *Proc. Eurospeech*, 1997
- [3] M. P. Cooke, A. Morris, and P. D. Green, "Recognizing Occluded Speech", *ESCA Tutorial and Workshop on Auditory Basis of Speech Perception*, Keele University, 1996.
- [4] B. Raj, R. Stern "Improving recognition accuracy in noise by using partial spectrographic information: Missing Feature in ASR" *IEEE SIGNAL PROCESSING MAGAZINE*, pp. 101-116, 2005.
- [5] H.G. Hirsch and C. Ehrlicher, "Noise estimation techniques for robust speech recognition," in *Proc. IEEE Conf. Acoustics, Speech Signal Processing, Detroit, Michigan*, pp. 153-156, 1995.
- [6] M. El-Maliki and A. Drygajlo, "Missing features detection and handling for robust speaker verification," in *Proc. Eurospeech, Budapest, Hungary*, pp. 975-978, 1999.
- [7] S. Boll, "Suppression of acoustic noise in speech using spectral subtraction," *IEEE Trans. Acoust. Speech Signal Processing*, vol. 27, no. 2, pp. 113-120, 1979.



شکل ۲: صحت بازشناسی بدست آمده از روش جبران سازی مبتنی بر ویژگی‌های مفقود در SNR های بین صفر الی ۲۰ dB در حیطه لگاریتم طیفی.

۵- نتیجه گیری

در این مقاله عملکرد روش جبران سازی مبتنی بر کواریانس که ویژگی‌های مفقود یا نامعتبر را از روی نمایش لگاریتم طیفی سیگنال گفتار نویزی بازیابی می‌کند، شرح داده شد. این روش از یک سری اطلاعات آماری بدست آمده از بردارهای لگاریتم طیفی دادگان تمیز برای بازسازی نواحی تخریب شده توسط نویز استفاده می‌کند.