

تنظیم تفريقي طيفي به منظور بهبود كارايي سيستم‌هاي بازشناسي گفتار

مهران صفابيانی^{*}، حسين ثامتى، باقر باباعلى، محمدتقى منظوري شلماني

چكيده

در اين مقاله روشی جديد برای تنظیم يک فیلتر تفريقي طيفي ارائه می‌گردد، به گونه‌ای که اين فیلتر بيشترین تأثير را بر بهبود نتایج يک سيستم بازشناسي گفتار داشته باشد. در حال حاضر واحدهای بازشناسي و بهسازی گفتار به صورت دو واحد مستقل عمل می‌کنند، بدین صورت که در ابتدا الگوريتم‌هاي بهسازی گفتار بروي سیگنال گفتار اعمال می‌شوند و سپس سیگنال بهسازی شده به واحد بازشناسي گفتار وارد می‌شود. در اين سيستم‌ها فرض می‌شود که بهترین الگوريتم بهسازی گفتار مطمئناً منجر به بهترین نتیجه بازشناسي خواهد شد، ولی بازشناسي گفتار يک مسئله بندی الگوهاست که از بردارهای ويژگی استخراج شده از سیگنال گفتار برای دسته بندی استفاده می‌کند. بنابراین تنها در صورتی نتایج بازشناسي افزایش خواهد یافت که اين ويژگی‌ها درست‌نمایي[†] دنباله آوایي صحیح را نسبت به سایر دنباله‌های آوایي رقیب افزایش دهد. در روش جدید ارائه شده، فیلتر تفريقي طيفي به نحوی تنظیم می‌شود که درست‌نمایي ويژگی‌های استخراج شده از سیگنال خروجی اين فیلتر بيشينه شود. با بكارگيري اين روش دقت سيستم‌هاي بازشناسي گفتار بروي دادگان فارس‌داد نويزي شده به ميزان ۱۸٪ افزایش می‌يابد.

كلمات کليدي

بهسازی گفتار، بازشناسي گفتار، تفريقي طيفي، بيشينه کردن درست‌نمایي.

Calibration of Spectral Subtraction for Improving Performance of Speech Recognition Systems

Mehran Safayani, H. Sameti, B. Babaali, M.T. Manzuri Shalmani
Computer Engineering Dept. Sharif University of technology

Abstract

In this paper, we present a novel approach for adjusting a spectral subtraction filter which has the most effect on improving speech recognition performance. Currently speech enhancement and recognition units work as two independent sets. At first speech enhancement stage process the speech data then enhanced speech signal is processed by speech recognition stage. These systems make the assumption that the best speech enhancement algorithm, will result in the best recognition performance. However speech recognition is a pattern classification problem that uses the extracted feature vectors from the speech signal, so if these feature vectors increase the likelihood of the correct transcription with respect to other competing incorrect hypothesis, speech recognition performance will increase surely. In our new method, parameters of the spectral subtraction are optimized to maximize the likelihood of the recognition features extracted from the resultant output signal. By incorporating the speech recognition system into adjusting filter parameter, speech recognition performance increase up to 18% on FARSDAT noisy database.

Keywords

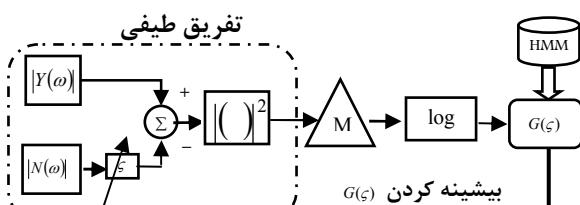
Speech Enhancement, Speech Recognition, Spectral Subtraction, Maximizing Likelihood.

* دانشگاه صنعتی شریف، دانشکده مهندسی کامپیوتو، {Sameti, Manzuri}@sharif.edu , {Safayani , Babaali}@ce.sharif.edu

[†] Likelihood

است، استفاده می‌شود. در این مرحله ضریب تفیریق طیفی به گونه‌ای پیدا می‌شود که فاصله مابین درستنمایی دنباله آوای صحیح و درستنمایی فرضیه نادرست با بیشترین امتیاز، بیشینه شود. در مرحله دوم با استفاده از فیلتر استخراج شده در مرحله اول بازناسی بروی همه مجموعه آزمون انجام می‌شود. در شکل (۱) نمودار بلوکی این روش نشان داده شده است. سیستم بازناسی بکاررفته بر مبنای مدل مخفی مارکوف بوده و بردارهای ویژگی بر اساس MFCC استخراج شده‌اند. این روش را می‌توان به صورت تفیریق طیفی چند کاناله نیز تعمیم داد. بطوریکه یک ضریب مجزا بروی هر باند فرکانسی اعمال گردد.

باقی‌مانده این مقاله در ۵ بخش ارائه می‌شود. در بخش دوم روش‌های تفیریق طیفی مرور می‌شود. در بخش سوم روش جدید ارائه شده به منظور تنظیم فیلتر تفیریق طیفی بر مبنای نتایج واحد بازناسی گفتار بیان می‌شود. در بخش چهارم پایگاه داده و سیستم بازناسی گفتار بکار رفته توضیح داده می‌شود و آزمایش‌های انجام شده و تحلیل نتایج در بخش پنجم توضیح داده می‌شود و سرانجام در بخش ششم نتیجه مقاله ارائه شده است.



شکل (۱): نمودار بلوکی الگوریتم تنظیم کردن تفیریق طیفی بر اساس نتایج بازناسی

۲- تفیریق طیفی

یکی از معروف‌ترین روش‌های بهسازی گفتار، روش تفیریق طیفی است. در این روش اندازه طیف سیگنال اصلی توسط کم کردن اندازه طیف نویز از اندازه طیف سیگنال نویزی بدست می‌آید. در این روش فرض می‌شود که نویز ناهمبسته با سیگنال صحبت بوده و با سیگنال صحبت می‌شود. اگر این فرض برقرار باشد که در مورد نویز پس‌زمینه برقرار است، می‌توان رابطه سیگنال با نویز را به صورت (۱) نوشت و تبدیل فوریه این رابطه به صورت (۲) بدست می‌آید. همچنین توان طیف سیگنال نویزی به صورت (۳) تخمین زده می‌شود.

$$y(k) = s(k) + n(k) \quad (1)$$

$$Y(w) = S(w) + N(w) \quad (2)$$

$$|Y(w)|^2 \approx |S(w)|^2 + |N(w)|^2 \quad (3)$$

که در آن $|S(w)|$ ، $|Y(w)|$ به ترتیب اندازه طیف سیگنال نویزی و تمیز است. از آنجاییکه طیف نویز نمی‌تواند به طور مستقیم

۱- مقدمه

دقت سیستم‌های بازناسی گفتار در حضور نویزهای جمع‌شونده و انعکاسی به شدت افت می‌کند. برای حذف اثر نویز روشهای بسیاری پیشنهاد شده است که می‌توان از تفیریق طیفی، فیلتر وینر، فیلتر ورقی، فیلتر کالمون و سایر روشها نام برد [۱]. تفیریق طیفی یکی از عمومی‌ترین روشهای، در بهسازی گفتار است [۳]. ایده اصلی آن بسیار ساده بوده و به راحتی قابل پیاده‌سازی است. این روش طیف سیگنال اصلی را توسط تفیریق کردن طیف نویز از طیف سیگنال نویزی بدست می‌آورد. در این روش فرض می‌شود که نویز ناهمبسته بوده و با سیگنال جمع می‌شود. به منظور افزایش کارایی این روش مطالعات بسیاری صورت گرفته است که بیشتر آنها کیفیت سیگنال گفتار را بر اساس معیارهایی نظری افزایش توان سیگنال به نویز یا تجربیات شنیداری انسان افزایش می‌دهند. هنگامی که این روشهای قبل از سیستم‌های بازناسی بکار گرفته شود، فرض می‌شود که با افزایش کیفیت سیگنال صحبت، دقت سیستم‌های بازناسی گفتار نیز افزایش خواهد یافت. البته اگر این روشهای دنباله‌ای از ویژگی‌هایی تولید کنند که درست-نمایی دنباله آوایی صحیح را نسبت به سایر فرضیه‌ها افزایش دهد، نتیجه بازناسی قطعاً بهبود خواهد یافت ولی معیار افزایش توان سیگنال به نویز یا سایر معیارهای مبتنی بر شکل موج لزوماً باعث بهبود دقت بازناسی نمی‌شوند زیرا مسئله بازناسی گفتار یک مسئله دسته‌بندی الگو است در حالیکه طراحی فیلتر یک مسئله پردازش سیگنال است. برای مثال یکی از مشکلات رایج تفیریق طیفی نویز موزیکال است. این نویز آزاردهنده باعث کاهش دقت بازناسی گفتار می‌شود در حالی که ممکن است توان سیگنال به نویز افزایش یافته باشد. سلتزر^۱ در [۴] از این فرضیه استفاده کرده و نشان می‌دهد که نتیجه بازناسی گفتار با تنظیم پارامترهای فیلتر بر اساس نتایج بازناسی افزایش چشم‌گیری می‌یابد. در کار ایشان که بازناسی گفتار از راه دور به کمک آرایه میکروفون است، یک فیلتر و- جمعی^۲ ارائه می‌شود که سیستم بازناسی گفتار در تنظیم پارامترهای آن دخیل است و ضرایب شکل دادن پرتو^۳ بر اساس بیشینه کردن درستنمایی دنباله آوایی صحیح تطبیق داده می‌شود. در این روش فرض می‌شود که با بیشینه کردن و یا حتی افزایش درستنمایی دنباله آوایی صحیح، دقت بازناسی افزایش خواهد یافت. در [۵] یک نسخه زیرباند از این الگوریتم ارائه شده است و ضرایب فیلتر در حوزه فرکانس بهینه می‌شود.

در این مقاله یک روش جدید به منظور تنظیم یک فیلتر تفیریق طیفی بر اساس نتایج بازناسی گفتار ارائه می‌شود. در این روش واحد بازناسی گفتار در تخمین پارامتر فیلتر تخمین طیفی دخیل می‌شود و نشان داده می‌شود با انجام این کار نتایج بازناسی گفتار افزایش خواهد یافت. این روش شامل دو مرحله است. مرحله اول، بهینه‌سازی است که در این مرحله از گفتاری که دنباله آوایی اش از قبل معلوم

در این رابطه $O(\zeta)$ بیانگر دنباله مشاهدات است که تابعی از پارامتر فیلتر تفريقي طيفی است و λ_k نشان دهنده مدل k ام است. از روی $L(k, \zeta)$ ها، می‌توان متغير مفيid $H(k, \zeta)$ را بصورت (۹) تعريف نمود:

$$H(k, \zeta) = \log \frac{L(k, \zeta)}{\sum_{j=1}^K L(j, \zeta)} \quad (9)$$

که در آن K تعداد همه مدل‌هاست. اگر C نشان دهنده اندیس مدل صحیح باشد، می‌توان ضریب فیلتر تفريقي طيفی را به صورتی بهینه کرد که درست نمایی $H(c, \zeta)$ بیشینه گردد و در نتیجه دقت سیستم بازناسی افزایش یابد که به صورت (۱۰) بیان می‌گردد.

$$H(c, \hat{\zeta}) = \arg \max_{\zeta} \left\{ \log \frac{L(c, \zeta)}{\sum_{j=1}^K L(j, \zeta)} \right\} \quad (10)$$

در این مقاله به منظور کاهش پیچیدگی محاسباتی، در مخرج عبارت (۱۰) فقط مدل برزنه (مدل با بیشترین امتیاز) قرار می‌گیرد و تابع هدفی به صورت (۱۱) بدست می‌آید.

$$H(c, \hat{\zeta}) = \arg \max_{\zeta} \{ \log(L(c, \zeta)) - \log(L(r, \zeta)) \} \quad (11)$$

که در آن c اندیس مدلی است که واقعاً بایستی بازناسی شود و r اندیس مدلی است که هم‌اکنون بازناسی شده است و ممکن است که مدل صحیح نباشد. فرض می‌شود که بیشینه کردن (۱۱) نسبت به ζ منجر به افزایش فاصله، بین درست‌نمایی مدل صحیح و سایر مدل‌های رقیب می‌شود. در این فرض یک پارادوکس وجود دارد. اگر دنباله واجی صحیح و به عبارتی دنباله مدل‌های صحیح معلوم باشند، دیگر نیازی به بازناسی گفتار نیست. این مسئله بدین صورت توجیه می‌گردد که در الگوریتم جدید یک مرحله بهینه‌سازی وجود دارد و کاربر بایستی یک جمله‌ای که دنباله آوایی آن از قبل معلوم است بیان کند. سپس سیستم با استفاده از این دانش، بهینه‌سازی را انجام می‌دهد و یک فیلتر مناسب تنظیم می‌کند و بعد از آن همه مجموعه آزمون از این فیلتر به منظور بهینه‌سازی استفاده می‌کنند. $(L(c, \zeta)$ و $L(r, \zeta)$ به ترتیب به صورت (۱۲) و (۱۳) بیان می‌شود.

$$L(c, \zeta) = \sum_{i=1}^T \log P(X_i(\zeta) | s_i) + \log(s_1, s_2, \dots, s_T) \quad (12)$$

$$L(r, \zeta) = \sum_{i=1}^T \log P(X_i(\zeta) | \hat{s}_i) + \log(\hat{s}_1, \hat{s}_2, \dots, \hat{s}_T) \quad (13)$$

تخمين زده شود، در این رابطه نويز توسط ميانگيري بروي M فريمي که تصور می‌شود که نويز است به صورت (۴) بدست می‌آيد:

$$\left| \hat{N}(w) \right|^2 = \frac{1}{M} \sum_{i=0}^M \left| Y_i(w) \right|^2 \quad (4)$$

يکی از مشکلات تفريقي طيفی وجود نويز موزیکال است که باعث افت کیفیت سیگنال گفتار می‌شود. بروتی^۵ [۶] با انجام تغیيراتی در نسخه اصلی این الگوریتم روشی برای کاهش اثر این نويز ارائه داد که در آن تخمين بهتری از نويز بدست می‌آورد و از اینکه سیگنال حاصل از یک حد خاصی کمتر شود ممانعت می‌کند. این الگوریتم به صورت (۵) ارائه شده است.

(۵)

$$\left| \hat{Y}(w) \right|^2 = \begin{cases} \left| \hat{S}(w) \right|^2 - \alpha \left| \hat{N}(w) \right|^2 & \left| \hat{S}(w) \right|^2 > \beta \left| \hat{N}(w) \right|^2 \\ \beta \left| \hat{N}(w) \right|^2 & \text{else} \end{cases} \quad \text{for all } w$$

که در آن α پارامتر تفريقي طيفی است و β حداقل اندازه طيف را تعیین می‌کند. مقادیر زياد α ممکن است باعث اعوجاج شود. برای ممانعت از این مشکل مقدار α بصورت فريمی به فريمی تطبیق داده می‌شود. تفريقي طيفی ای که در این مقاله استفاده می‌شود در رابطه (۶) و (۷) نشان داده شده است. در این فرمول‌ها ضریب ζ یا بصورت تجربی بدست می‌آید و یا با توجه به توان سیگنال به نويز شکل موج خروجی تنظیم می‌شود که اگر مشخصات محیط تغیير کند این ضریب بایستی دوباره تنظیم شود. در این مقاله روشی برای تنظیم این پارامتر ارائه شده است بطوریکه بیشترین تأثیر را بر نتایج بازناسی داشته باشد. البته روش ارائه شده قابل تعمیم برای سایر فرمول‌های تفريقي طيفی نیز هست.

$$\left| \hat{S}(k) \right| = \left| Y(k) \right| - \zeta \left| \hat{N}(k) \right| \quad (6)$$

$$\left| \hat{S}(k) \right| = \begin{cases} \left| \hat{S}(k) \right| & \left| Y(k) \right| > \zeta \left| \bar{N}(k) \right| \\ \left| \bar{N}(k) \right| & \text{otherwise} \end{cases} \quad (7)$$

۳- بهینه سازی تفريقي طيفی به منظور بازناسی مقاوم گفتار

در این قسمت نحوه بدست آوردن یک فیلتر تفريقي طيفی مناسب بر مبنای افزایش فاصله بین درست‌نمایی دنباله آوایی صحیح نسبت به سایر دنباله‌های آوایی بیان می‌شود. فرض کنید که $L(k, \zeta)$ امتیاز مدل واحد k ام (یا بطور خلاصه مدل k ام) نسبت به پارامتر فیلتر تفريقي طيفی ζ باشد که به صورت (۸) بیان می‌گردد.

$$L(k, \zeta) = P(O(\zeta) | \lambda_k) \quad (8)$$

$$\nabla_{\zeta} G(X) = \left[\begin{array}{l} \left(-\sum_t \sum_{k=1}^K \left(\frac{\exp(\rho_c)}{\sum_{j=1}^K \exp(\rho_c)} \right) \Sigma_{S_c(t),k}^{-1} \vartheta_c \frac{\partial X_t^c(\zeta)}{\partial \zeta} \right) + \\ \left(\sum_t \sum_{k=1}^K \left(\frac{\exp(\rho_r)}{\sum_{j=1}^K \exp(\rho_r)} \right) \Sigma_{S_r(t),k}^{-1} \vartheta_r \frac{\partial X_t^r(\zeta)}{\partial \zeta} \right) \end{array} \right] \quad (19)$$

که در آن ρ_c به صورت (۲۰) و ρ_r به صورت (۲۱) بیان می-

شوند و $\frac{\partial X_t^c(\zeta)}{\partial \zeta}$ یک ماتریس ژاکوبین است که از مشتق جزئی هر مؤلفه بردار ویژگی نسبت به ضریب فیلتر تفیری طیفی بدست می‌آید.

$$\rho_c = -\frac{1}{2} \vartheta_c^H \Sigma_{S_c(t),k}^{-1} \vartheta_c \quad (20)$$

$$\rho_r = -\frac{1}{2} \vartheta_r^H \Sigma_{S_r(t),k}^{-1} \vartheta_r \quad (21)$$

به طور کلی الگوریتم این روش در ادامه بیان می‌شود:

۱. یک مقدار پیش فرض برای ضریب فیلتر در نظر گرفته می-شود برای مثال ضریب برابر با صفر.
۲. کاربر یک جمله‌ای را که دنباله آوابی آن از قبل معلوم است بیان می‌کند.
۳. سیستم بازشناسی بهترین دنباله حالت را با استفاده از الگوریتم پیتری، دنباله آوابی معلوم و ضریب فیلتر تفیری طیفی بدست می‌آورد.
۴. با استفاده از روش گرادیان، (۱۶) نسبت به ضریب فیلتر تفیری طیفی بیشینه می‌شود.
۵. اگر فیلتر تولید شده همگرا شد به مرحله ۶ برو و گرنه به مرحله ۳ برگرد.
۶. ضریب بیشینه بدست آمده به عنوان پارامترهای فیلتر نهایی بکار می‌رود.

۴- پایگاه دادگان و سیستم بازشناسی گفتار

به منظور آموزش و آزمایش این سیستم از دادگان فارس دات [۷] استفاده گردیده است. این پایگاه داده شامل ۶۰۸۰ گفتار فارسی است که توسط ۳۰۴ گوینده بیان شده است. این گویندگان از ۱۰ منطقه جغرافیایی مختلف در ایران انتخاب شده‌اند بنابراین ۱۰ لهجه رایج فارسی در این پایگاه وجود دارد. نسبت مرد به زن در این پایگاه دو به یک است و ۴۰۵ جمله در آن وجود دارد که ۲۰ جمله بازی هر

که در آنها X_i بردار ویژگی i امین فریم است و فرض می‌شود که درست‌نمایی دنباله واجی صحیح و همچنین دنباله واجی غیر صحیح با بیشترین امتیاز، با یک دنباله حالت که بیشترین امتیاز را دارد بیان می‌گردد که به ترتیب به صورت $\hat{S} = \{\hat{s}_1, \hat{s}_2, \dots, \hat{s}_T\}$ و $S_C = \{s_1, s_2, \dots, s_T\}$ بیان می‌گردد. همچنین بردارهای ویژگی i به روش MFCC استخراج می‌گردد. که می‌توان این ویژگی‌ها را به صورت رابطه (۱۴) بیان نمود.

$$X_i = DCT(\log(M | S_i|^2)) \quad (14)$$

که در آن X_i بیانگر ضرایب MFCC فریم i ام گفتار، S_i خروجی فیلتر تفیری طیفی و M بیانگر ماتریس ضرایب وزن دار فیلترهای مل هستند. توزیع‌های درون حالت‌های مدل مخفی مارکوف توسط چند تلفیق گوسی مدل می‌شوند، بنابراین می‌توان رابطه (۱۱) را به صورت (۱۵) بدست آورد.

$$(15)$$

$$H(c, \hat{\zeta}) = \arg \max_{\zeta} \left\{ \left(\sum_{t=1}^{|S_C|} \log \sum_{k=1}^K N(X(t, \zeta), \mu_{S_c(t), k}, \Sigma_{S_c(t), k}) \right) - \left(\sum_{t=1}^{|S_r|} \log \sum_{k=1}^K N(X(t, \zeta), \mu_{S_r(t), k}, \Sigma_{S_r(t), k}) \right) \right\}$$

که در آن N توزیع نرمال، $|S_C|$ تعداد کل فریمهای تلفیق‌های گوسی و $\mu_{S_c(t), k}$ ، $\mu_{S_r(t), k}$ به ترتیب بردار میانگین و ماتریس کوواریانس توزیع نرمال حالت k ام هستند. برای بیشینه کردن عبارت (۱۵) تابع (ζ) به صورت (۱۶) تعریف می‌شود.

$$(16)$$

$$G(\zeta) = \left[\left(\sum_{t=1}^{|S_c|} \log \sum_{k=1}^K \exp \left(-\frac{1}{2} \vartheta_c^H \Sigma_{S_c(t), k}^{-1} \vartheta_c \right) \right) - \left(\sum_{t=1}^{|S_r|} \log \sum_{k=1}^K \exp \left(-\frac{1}{2} \vartheta_r^H \Sigma_{S_r(t), k}^{-1} \vartheta_r \right) \right) \right]$$

که در آن ϑ_c به صورت (۱۷) و ϑ_r به صورت (۱۸) تعریف می‌شوند.

$$\vartheta_c = X(t, \zeta) - \mu_{S_c(t), k} \quad (17)$$

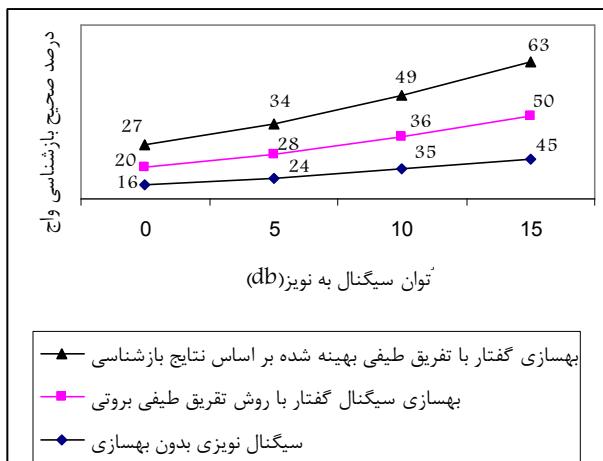
$$\vartheta_r = X(t, \zeta) - \mu_{S_r(t), k} \quad (18)$$

که در آنها $X(t, \zeta)$ بردار ویژگی در فریم t ام نسبت به پارامتر فیلتر ζ است، $\mu_{S_c(t), k}$ و $\mu_{S_r(t), k}$ به ترتیب برابر بردار میانگین مربوط به مدل صحیح و مدل بازشناسی شده است. گرادیان (۱۶) نسبت به ζ به صورت (۱۹) بدست می‌آید. با استفاده از (۱۶) و گرادیان آن می‌توان (۱۵) را با استفاده از روش گرادیان مزدوج [۲] بیشینه کرد و ضریب بیشینه فیلتر تفیری طیفی را بدست آورد.

جدول (۱): درصد صحیح بازناسی واج بروی ۹۹ جمله فارسی

بهینه‌سازی به روش بیشینه کردن درست نمائی	بدون بهینه‌سازی	توان سیگنال به نویز (db)
۲۷	۱۶	۰
۳۴	۲۴	۵
۴۹	۳۵	۱۰
۶۳	۴۵	۱۵

در این قسمت الگوریتم بهینه سازی تفیریق طیفی ارائه شده با یک روش تفیریق طیفی مبتنی بر توان سیگنال به نویز مقایسه می‌شود. یکی از متدالوں ترین الگوریتم‌های تفیریق طیفی، الگوریتمی است که توسط بروتی [۷] ارائه گردید که رابطه این روش در (۶) ارائه شده است. در این قسمت در ابتدا با استفاده از الگوریتم بروتی دادگان نویزی که شامل ۱۰۰ جمله فارسی دات است، بهینه سازی می‌شود و سپس این سیگنال‌های بهینه سازی شده به سیستم بازناسی وارد می‌گردد. نتایج حاصل از اعمال این روش در شکل (۲) با روش بهینه سازی تفیریق طیفی به روش بیشینه کردن درست نمائی مقایسه شده است. همانگونه که از شکل (۲) مشاهده می‌شود، روش تفیریق طیفی با استفاده از بیشینه کردن درست نمائی در همه موارد نسبت به تفیریق طیفی بروتی بهبود داشته است و این بهبود در دادگان با توان سیگنال به نویز ۱۵ دسی‌بل به طور واضح مشاهده می‌گردد که روش بروتی تنها ۵ درصد نتایج بازناسی را افزایش داده است درحالیکه روش جدید ۱۸ درصد دقت بازناسی را افزایش داده است. با بررسی سیگنال‌های خروجی روش تفیریق طیفی بروتی مشاهده شد، که این روش مقدار زیادی نویز موزیکال به سیگنال اضافه می‌کند که باعث می‌شود که بسیاری از واج‌ها به درستی بازناسی نشود و در نتیجه اعمال این روش در ورودی سیستم بازناسی گفتار نتایج بازناسی را به میزان قابل توجهی افزایش نمی‌دهد.



شکل (۲): مقایسه الگوریتم تفیریق طیفی بهینه شده با تفیریق طیفی بروتی

گوینده است. هر گوینده ۱۸ جمله را به صورت تصادفی بیان کرده است باضافه دو جمله که در بین همه گویندگان مشترک است. این جمله‌ها بر اساس ۱۰۰۰ کلمه فارسی شکل می‌گیرند. پایگاه داده در محیط با نویز کم ضبط شده است و توان سیگنال به نویز آن به طور متوسط ۳۱ دسی‌بل است. به منظور آزمودن سیستم از ۱۰۰ جمله آن استفاده شده است و باقی جمله‌ها برای آموزش استفاده گردیده است. به منظور نشان دادن اثر الگوریتم جدید بر نتایج بازناسی گفتار از یک سیستم بازناسی واج‌های گیسته مبتنی بر مدل مخفی مارکوف استفاده می‌شود. HMM‌های به کار برده شده از نوع چگالی پیوسته و با تلفیق‌های گوسی بوده و پرش مجاز بین حالت‌های آن‌ها فقط به صورت چپ به راست است. در آزمایش‌های ما مدل سازی گفتار براساس واج و با استفاده از دادگان فارسی دات انجام گردید و سعی شد از بهترین پارامترهای ممکن استفاده شود. بردارهای ویژگی استخراج شده از فریم‌های گفتار، شامل ضرب مل-C0-(C11) همراه با مشتقات زمانی اول و دوم آن‌ها می‌باشد. مدل‌های مخفی مارکوف با استفاده از الگوریتم k-means segmental می‌باشند و توپولوژی آن‌ها برای تمام واج‌ها یکسان در نظر گرفته شده است. لازم به ذکر است که در دادگان فارسی دات کل آواهای زبان فارسی به ۴۳ واج تقسیم‌بندی شده است. از این ۴۳ واج، ۱۲ واج مربوط به قسمت گفتار واج‌های انسدادی و بقیه واج‌های استاندارد زبان فارسی می‌باشند. قسمت‌های غیر گفتاری سیگنال (مانند سکوت) نیز به عنوان یک واج در نظر گرفته شده‌اند؛ بنابراین در مجموع ۴۴ واج می‌تواند مدل گردد. ولی در این پژوهه با توجه به عملکرد بهتر بازناسی، با ادغام قسمت گفتار واج‌های انسدادی با قسمت رهش و همچنین ترکیب بعضی از واج‌های شبیه با یکدیگر، تعداد مدل‌ها به ۳۰ مدل واج کاهش یافته است.

۵- آزمایش‌ها

در ابتدا دادگان آزمون شامل ۱۰۰ جمله فارسی دات با نویز سفید جمع گردید و در توان سیگنال به نویزهای مختلف، دادگان آزمون آماده شدند. سپس یک جمله آن به منظور بهینه‌سازی استفاده گردید و با استفاده از روش توضیح داده شده در بخش (۳) فیلتر مناسب تنظیم گردید و از آن در مرحله بازناسی بر روی ۹۹ جمله ۱۸ افزایش داده شد. درصد بازناسی واج‌های صحیح در جدول (۱) نشان داده شده است. همانگونه که در این جدول مشاهده می‌شود نتایج بازناسی در توان سیگنال به نویزهای مختلف افزایش داشته است که برای نمونه می‌توان به افزایش ۱۸ درصدی دقت بازناسی واج در توان سیگنال به نویز ۱۵ دسی‌بل اشاره کرد.

۶- نتیجه

در این مقاله روشی جدید به منظور بهینه‌سازی یک فیلتر تفريق طيفی به منظور استفاده در ورودی سیستم‌های بازناسی گفتار ارائه گردید و نشان داده شد که بکارگیری اين روش نتایج بازناسی را نسبت به روش‌های متداول تفريق طيفی افزایش می‌دهد. با بکارگیری روش جدید ارائه شده مطمئن هستیم که اين فیلتر ویژگی‌هایی از سیگنال صحبت را که برای سیستم بازناسی مهم تر است تقویت می‌کند و در نتیجه نتایج سیستم‌های بازناسی گفتار افزایش می‌یابد. این روش را می‌توان به صورت تفريق طيفی چند کاناله نیز گسترش داد و همچنین بر روی سایر الگوريتم‌های بهینه‌سازی گفتار نیز اعمال کرد.

مراجع

- [1] J. R. Deller, J. G. Proakis, H. L. Hansen, "Discrete Time Processing of Speech Signals", Macmillan, New York, NY, 1993.
- [2] E. Polak, "Computational methods in optimizations", New York: Academic Press, 1971.
- [3] S.F. Boll, "Suppression of Acoustic Noise in Speech using Spectral Subtraction", IEEE Trans Acoustics, Speech & Signal Processing, April 1979.
- [4] M. L. Seltzer, B. Raj, R. M. Stern, "Likelihood-Maximizing Beamforming for Robust Hands-Free Speech Recognition", IEEE Transaction on Speech and audio processing, VOL. 12, NO. 5, September 2004.
- [5] M. L. Seltzer, R. M. Stern, "Subband parameter optimization of microphone arrays for speech recognition in reverberant environments", Proceedings of the International Conference on Acoustics, Speech, and Signal Processing, Hong Kong, July 2003.
- [6] M. Berouti, R. Schwartz, J. Makhoul, "Enhancement of speech corrupted by acoustic noise", Proceedings of the IEEE International Conference on Acoust, Speech and Signal Processing, pp. 208-211, Apr. 1979.
- [7] M. Bijankhan, J. Sheikhzadegan, "FARSDAT-The Farsi Spoken Language Database", proceeding of the 5th International conference of speech science and technology, vol. 2, pp. 826-831, 1994.

زیرنویس‌ها

-
- ¹ Musical
² Seltzer
³ Filter-and-Sum
⁴ Beamforming
⁵ Berouti