



مجموعه مقالات - ۵
مخابرات (الف)

دانشگاه صنعتی شریف
دانشکده مهندسی برق

ICEE - 97

پنجمین کنفرانس مهندسی برق ایران

۱۷-۱۹ اردیبهشت ۱۳۷۶



جداسازی دو سیگنال صحبت درهم

مسعود بابایی زاده و محمود تیانی

دانشکده برق دانشگاه صنعتی شریف

چکیده: در این مقاله روش جدیدی برای جداکردن دو سیگنال صحبت مخلوط ارائه شده است. بر خلاف روشهای قبلی موجود، که دو صحبت را در حوزه فرکانس از هم جدا می کردند، کلیه محاسبات این روش در حوزه زمان انجام می گردد. این روش بر مبنای قضیه جالب و ساده ای قرار دارد، که آن هم برای نخستین بار در این مقاله بیان می شود. برای آنالیز و سنتز صحبت نیز، از مدل سینوسی صحبت استفاده شده است.

کلمات کلیدی: سیگنال صحبت، جداسازی، مدل سینوسی صحبت، حداقل مربعات.

۱- مقدمه

مسئله جداکردن دو سیگنال صحبت مخلوط، با وجود بیش از بیست سال تحقیق، هنوز به جواب مطلوبی منتهی نشده است. شکل موج پیچیده سیگنال صحبت، و طبیعت غیرایستای آن، باعث عدم موفقیت کلیه تلاشهای انجام شده برای حل این مسئله گردیده است. نخستین بار شیلدز و اپنهایم^۱ به این مسئله پرداختند و استفاده از فیلتر شانهای را برای حل این مسئله بررسی نمودند [۱]. ایده انتخاب دامنه هارمونیکها برای حل این مسئله، نخستین بار توسط پارسنز^۲ مطرح گردید [۲] و بعد

¹ V.C. Shields and A.V. Oppenheim

² T.W. Parsons

توسط دیگران کاملتر شد [۳ و ۴ و ۵]. کواتیری و دانیسویچ^۳ استفاده از مدل سینوسی صحبت را در حل این مسأله آزمودند [۶] و برخی کار آنها را دنبال کردند [۷ و ۸]. ویژگی همه این روشها آنست که عمل جداسازی، در حوزه فرکانس انجام می شود و محدود بودن طول فریمها (ناشی از طبیعت غیرایستای صحبت)، باعث کاهش قدرت تفکیک در حوزه فرکانس و همپوشانی پیکها و بالاخره دشواری عمل جداسازی می گردد.

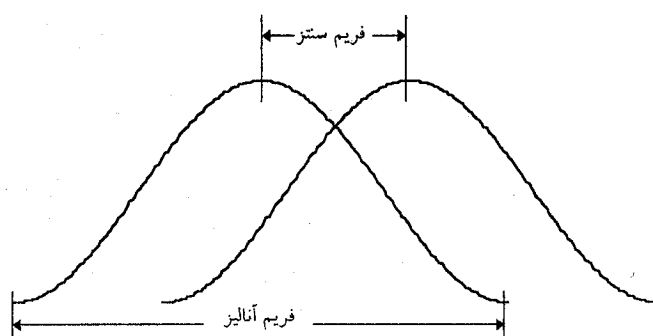
روش معرفی شده در این مقاله، عمل جداسازی را در حوزه زمان انجام می دهد و در نتیجه مشکل قدرت تفکیک کم در حوزه فرکانس، تا حد زیادی کاهش می یابد. همچنین، قضیه ای مطرح خواهیم کرد که با استفاده از آن و پذیرفتن درصدی خطا، می توان حجم محاسبات را به مقدار زیادی کاهش داد. نتایج بدست آمده نیز حاکی از موفقیت این روش در عمل هستند.

۲- مدل سینوسی صحبت

در این مدل، سیگنال صحبت در هر لحظه، بصورت جمع تعدادی سینوسی با دامنه، فاز و فرکانسهای مختلف مدل می شود [۹]، یعنی:

$$x(t) = \sum_{k=1}^M a_k \cos(\omega_k t + \phi_k) \quad (1)$$

یکی از مزایای این مدل، مستقل بودن طول فریمهای سنتز و آنالیز از یکدیگر است (شکل ۱). هنگام تجزیه صحبت، پارامترهای مدل در طول یک فریم آنالیز محاسبه شده و به نقطه مرکزی فریم نسبت داده می شود. سپس در هنگام سنتز، با توجه به معلوم بودن پارامترها در دو انتهای فریم سنتز، مقادیر آنها در نقاط میانی درونیابی می گردد. نشان داده شده است [۹] که برای درونیابی مقادیر دامنه، می توان از درونیابی خطی استفاده کرد، در حالیکه برای محاسبه فاز و فرکانس در نقاط میانی، باید از یک چندجمله ای درجه سه بهره برد.



شکل (۱): فریمهای سنتز و آنالیز در مدل سینوسی صحبت

توجه شود که برای نواحی واکدار صحبت، فرکانسهای این مدل، هارمونیکهای فرکانس اصلی (گام) صحبت هستند.

³ T.F. Quatieri and R.G. Danisewics

مدل فوق را می‌توان به سادگی برای دو صحبت نیز تعمیم داد [۶]. در اینحالت داریم:

$$x[n] = \sum_{k=1}^{M_a} a_k \cos(\omega_{a,k}n + \varphi_{a,k}) + \sum_{k=1}^{M_b} b_k \cos(\omega_{b,k}n + \varphi_{b,k}) \quad (2)$$

استخراج فرکانسها از روی صحبت مخلوط، عملی دشوار است و در این مقاله (همچون بیشتر مقاله‌های قبلی)، بمنظور شبیه‌سازی روش ارائه شده، آنها را از قبل معلوم فرض می‌نماییم.

۳- جداسازی دو صحبت مخلوط

برای جداکردن دو صحبت از یکدیگر، قضیه زیر را بکار می‌گیریم [۱۰]:

قضیه: فرض کنید سیگنال $x[n]$ جمع تعدادی سیگنال سینوسی با فرکانسهای متفاوت باشد، یعنی:

$$x[n] = \sum_{k=1}^{M_a} a_k \cos(\omega_{a,k}n + \varphi_{a,k}) + \sum_{k=1}^{M_b} b_k \cos(\omega_{b,k}n + \varphi_{b,k})$$

علاوه بر آن فرض کنید فرکانس تعدادی از سینوسی‌های تشکیل دهنده سیگنال فوق، مثلاً $\omega_{a,k}$ را می‌دانیم و می‌خواهیم دامنه و فازهای آنها را بدست آوریم. در این صورت اگر سیگنال زیر را با دامنه و فازهای مجهول تشکیل دهیم:

$$y[n] = \sum_{k=1}^{M_a} a'_k \cos(\omega_{a,k}n + \varphi'_{a,k})$$

و مقادیر این دامنه و فازهای مجهول را بگونه‌ای تعیین کنیم که میانگین $|x[n] - y[n]|^2$ حداقل شود، آنگاه همان مقادیر دامنه و فاز واقعی را بدست خواهیم آورد، یعنی:

$$a'_k = a_k, \quad \varphi'_{a,k} = \varphi_{a,k}, \quad k = 1, \dots, M_a$$

با توجه به قضیه بالا، برای یافتن یکی از صحبتها، کافی است یک ترکیب خطی از فرکانسهای تشکیل دهنده آن صحبت را، با دامنه و فازهای مجهول، تشکیل داده و مقادیر این دامنه و فازها را بگونه‌ای تعیین کنیم که میانگین مربعات خطای سیگنال حاصل با سیگنال مخلوط حداقل شود. در اینصورت طبق قضیه بالا، سیگنال بدست آمده، همان صحبت مورد نظر خواهد بود.

برای محاسبه این دامنه و فازها، از فرم تریبیتی^۴ سیگنال سینوسی استفاده کرده و ابتدا سعی می‌کنیم خطا را صفر نماییم.

در اینصورت به دستگاه معادلات زیر خواهیم رسید:

⁴ Quadrature Form

$$\mathbf{H} \cdot \mathbf{c} = \mathbf{s}, \quad (3)$$

که در آن:

$$\mathbf{H} = \begin{bmatrix} 1 & 1 & \dots & 1 & \dots & 0 & \dots & 0 & \dots \\ 1 & \cos(\omega_{a,1} \cdot 1) & \dots & \cos(\omega_{a,2} \cdot 1) & \dots & \sin(\omega_{a,1} \cdot 1) & \dots & \sin(\omega_{a,2} \cdot 1) & \dots \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots & \ddots & \vdots & \ddots \\ 1 & \cos(\omega_{a,1} \cdot (N-1)) & \dots & \cos(\omega_{a,2} \cdot (N-1)) & \dots & \sin(\omega_{a,1} \cdot (N-1)) & \dots & \sin(\omega_{a,2} \cdot (N-1)) & \dots \end{bmatrix} \quad (4)$$

و:

$$\mathbf{c} = \begin{bmatrix} \alpha_0 \\ \alpha_1 \\ \alpha_2 \\ \vdots \\ \beta_1 \\ \beta_2 \\ \vdots \end{bmatrix}, \quad \mathbf{s} = \begin{bmatrix} x[0] \\ x[1] \\ \vdots \\ x[N-1] \end{bmatrix} \quad (5)$$

دستگاه بالا را باید به مفهوم میانگین مربعات حل کنیم، یعنی \mathbf{c} را طوری بدست آوریم که $\|\mathbf{H} \cdot \mathbf{c} - \mathbf{s}\|^2$ حداقل شود. می توان ثابت کرد [۱۱] که در این صورت \mathbf{c} از رابطه زیر بدست می آید:

$$\mathbf{c} = (\mathbf{H}^* \mathbf{H})^{-1} \mathbf{H}^* \cdot \mathbf{s} \quad (6)$$

پس از محاسبه \mathbf{c} ، داریم:

$$\begin{aligned} \alpha_k &= \sqrt{\alpha_k^2 + \beta_k^2} \\ \varphi_{a,k} &= \tan^{-1}(\beta_k / \alpha_k) \end{aligned} \quad (7)$$

و در نتیجه صحبت اول بازسازی می شود.

۴- درونیابی چندفریمی

گاهی اوقات برخی از فرکانسهای صحبت دوم، بیش از اندازه به فرکانسهای صحبت اول نزدیک می شوند و این پدیده باعث ایجاد خطا در محاسبه دامنه و فاز متناظر می گردد. برای رفع این مشکل، فاصله فرکانسهای دو نفر را همواره تحت مراقبت قرار می دهیم و در صورت وقوع پدیده فوق، فرکانس مربوطه را به نشانه غیرقابل اعتماد بودن دامنه علامت گذاری می کنیم. در نهایت هنگام بازسازی صحبت، وقتی به چنین بلوکهایی رسیدیم، بجای آنها یک فریم سست را در نظر

بگیریم، فریمهای قبلی و بعدی را نیز تا جاییکه به دامنه قابل اعتمادی برسیم با هم به عنوان یک فریم سنتز طولانی تر در نظر گرفته و درونیایی پارامترها را در این فریم طولانی انجام می‌دهیم. به این عمل درونیایی چندفریمی می‌گوییم. استفاده از این روش باعث افزایش کیفیت صداها می‌گردد.

۵- نتایج

با استفاده از روش مطرح شده، سیستمی برای جداکردن دو صحبت پیاده‌سازی گردید. برای این منظور از صحبتها بطور جداگانه و با فرکانس ۱۰ KHz نمونه‌برداری کردیم. سپس فرکانسهای مدل سینوسی آنها محاسبه گردیده و در فایل‌های جداگانه‌ای ذخیره شدند. همچنین گام صحبتها را نیز بطور جداگانه محاسبه نمودیم.

در مرحله اول با استفاده از فرکانسهای مدل سینوسی هرکدام از صحبتها، یکبار بدون استفاده از روش درونیایی چندفریمی و بار دیگر با استفاده از آن، مخلوط آنها را جداکردیم. سپس چندین نفر، به عنوان داور، به این صحبتها گوش کرده و دو نمره، یکی برای میزان جداسازی صحبتها، و دیگری برای کیفیت صحبتهای سنتز شده ثبت کردند. همگی افراد به میزان جداسازی نمره کامل (۴) دادند. میانگین نمره داده شده به کیفیت صحبتها، در حالت استفاده از درونیایی چندفریمی، ۲/۹ و در حالت استفاده نکردن از آن، ۲/۱ بود. این نکته، کارایی درونیایی چند فریمی را در بهبود کیفیت، نشان می‌دهد.

در مرحله بعد بجای معلوم بودن همه فرکانسهای مدل سینوسی، تنها فرکانس اصلی (گام) آنها معلوم فرض شد و بقیه فرکانسهای مدل، هارمونیکهای آن فرض گردید. نتیجه حاصل از بررسی این صحبتها نشان داد که میزان جداسازی از ۴ به ۲/۳ افت پیدا کرده است.

در مرحله سوم، بجای استفاده از قضیه مطرح شده، ترکیب خطی از کلیه فرکانسها بکار برده شد (حجم محاسبات در اینحالت حدود ۸ برابر بیشتر است). هیچکدام از داوران تفاوتی بین کیفیت صحبتهای بازسازی شده با صحبتهای مرحله اول، تشخیص ندادند و این، کارایی قضیه مطرح شده را در کاهش حجم محاسبات نشان می‌دهد.

۶- خاتمه

در این مقاله پس از معرفی مدل سینوسی، روشی برای جداسازی دو صحبت درهم معرفی گردید. ویژگی این روش، انجام همه محاسبات در حوزه زمان است. همچنین قضیه‌ای برای افزایش سرعت محاسبه بیان شد. در نهایت یک مطالعه subjective برای بررسی کارایی روش انجام گردید که این مطالعه، کارایی روش معرفی شده را بخوبی نشان داد.

مراجع

- [1] V.C. Shields, "Separation of Added Speech Signals by Digital Comb Filtering," M.S. Thesis, MIT, 1970.
- [2] P.W. Parsons, "Separation of Simultaneous Vocalic Utterances of Two Talkers," Ph.D. Thesis, Department of Elec. Eng., Polytech. Inst. of New York, 1975.
- [3] B.A. Hanson and D.Y. Wong, "The Harmonic Magnitude Supression (HMS) Technique for Intelligibility Enhance," 1984, 18A.5.1-18A.5.4.
- [4] D.G. Childers and C.K. Lee, "Co-Channel Speech Separation," ICASSP 1987, 6.4.1-6.4.4.
- [5] J.A. Naylor and C.K. Lee, "Co-Channel Speech Separation," ICASSP 1987, 6.12.1-6.12.4.
- [6] T.F. Quatieri and R.G. Danisewics, "An Approach to Co-Channel Talker Interference Supression Using a Sinusoidal Model for Speech," IEEE Trans. on ASSP, Vol. 38, pp.56-69, January 1990.
- [7] F.M. Silva and L.B. Almeida, "Speech Separation by Means of Stationary Least Squares Harmonic Estimation," ICASSP 1990, pp. 809-812.
- [8] J.A. Naylor and J. Porter, "An Effective Speech Separation System which Requires No Priori Information," ICASSP 1991, pp. 937-940.
- [9] R.J. McAulay and T.F. Quatieri, "Speech Analysis/Synthesis Based on A Sinusoidal Representation," IEEE trans. on ASSP, vol. 34, pp. 744-754, Aug. 1986.
- [10] بابایی زاده، مسعود، "جهدسازی دو سیگنال صحبت مخلوط"، رساله کارشناسی ارشد، دیماه ۱۳۷۵، دانشکده برق دانشگاه صنعتی شریف.
- [11] Kay, S.M., *Modern Spectral Estimation*, Prentice Hall, 1988.