

ارائه الگوریتمی مبتنی بر روش انرژی - انرژی برای تشخیص محدوده گفتار در محیط‌های نویزی

حمید دهقانی^۱، علی مهاجری^۲

۱- دانشگاه صنعتی مالک اشتر، hamid_deh@yahoo.com

۲- دانشجوی دانشگاه آزاد اسلامی واحد بوشهر، ali.mohajeri@gmail.com

چکیده

در یک سیستم بازشناسی گفتار کلمات گسسته یا پیوسته، بازدهی تشخیص محدوده گفتار، تاثیر مستقیمی روی کل سیستم بازشناسی دارد. معمولاً انرژی، و نرخ عبور از سطح صفر، به‌عنوان دو پارامتر در تشخیص نقاط ابتدا و انتهای گفتار استفاده می‌شوند. در محیط‌های نویزی پارامتر نرخ عبور از سطح صفر، غیرقابل اعتماد است و باعث خطا در کل سیستم بازشناسی می‌شود. در اکثر روش‌ها، نویز زمینه از روی فریم‌های اولیه گفتار تخمین زده می‌شود. در الگوریتم ارائه شده با استفاده از تغییراتی که در انرژی اصلاح شده تیگر می‌دهد، از آن در تشخیص نقاط اولیه ابتدا و انتهای گفتار استفاده می‌شود، سپس با استفاده از اصلاحات مؤثری که در الگوریتم انرژی - انرژی انجام شده، نقاط دقیقتر گفتار را می‌توان تشخیص داد. در این الگوریتم، حتی با وجود نویز در ابتدای گفتار، محدوده گفتار تشخیص داده می‌شود. از این الگوریتم، برای تشخیص محدوده گفتار های فارسی استفاده شده است. نتایج آزمایشات نشان‌دهنده بازدهی مناسب و قابل قبول روش ارائه شده در زمینه تشخیص گفتار از زمینه نویز و سکوت می‌باشد.

واژه‌های کلیدی

انرژی اصلاح شده تیگر، تشخیص نقاط ابتدا و انتهای گفتار، ویژگی انرژی - انرژی، محیط‌های نویزی، نویز نفس کشیدن

۱- مقدمه

انرژی در زبان فارسی هستند که بازشناسی آنها ممکن است همراه با خطا باشد. نرخ عبور از سطح صفر، در مورد حروف واکه دار، مقدار کمتری نسبت به سایشی‌های ضعیف بدون واکه و یا انفجاری‌ها دارد، که پس از تعیین نقاط اولیه گفتار، پارامتر مناسبی در تشخیص دقیقتر شروع و پایان گفتار است [۱]. در محیط‌های نویزی پارامتر نرخ عبور از سطح صفر، نمی‌تواند خیلی قابل اعتماد باشد. تیگر الگوریتم جدیدی را برای محاسبه کردن انرژی سیگنال پیشنهاد کرد که این الگوریتم تحت عنوان الگوریتم انرژی تیگر ارائه شد. با اصلاحاتی که در این الگوریتم انجام شد، جهت تشخیص محدوده

اولین گام در یک سیستم بازشناسی گفتار گسسته، آشکارسازی دقیق نقاط ابتدا و انتهای گفتار، از زمینه سکوت و نویز است. در سیستم‌های بازشناسی گفتار، تشخیص نادرست محدوده گفتار، اثرات منفی مانند افزایش خطای بازشناسی و افزایش حجم محاسبات در کل سیستم بازشناسی دارد. اکثر الگوریتم‌های تعیین محدوده گفتار، مبتنی بر اندازه‌گیری دو ویژگی نرخ عبور از سطح صفر و انرژی سیگنال صوتی است [۱،۳]. سایشی‌های ضعیف بدون صدا، نظیر /ف، س، ش، ه / انفجاری‌های ضعیف مانند /پ، ت، ک / و خیشومی‌های انتهایی مانند /م، ن / به‌عنوان بعضی واج‌های کم

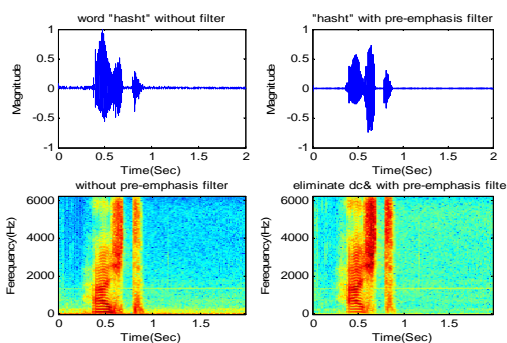
ضرب با یک جمع ساده و شیفت جایگزین شده است. تقسیم بر ۳۲ با شیفت دادن ۵ بیت به سمت راست ایجاد می‌شود [۶].

$$H(z) = 1 - az^{-1} \quad (1)$$

$$S_n = s_n - \frac{31}{32} s_{n-1}$$

$$S_n = s_n - (s_{n-1} - \frac{1}{32} s_{n-1}) \quad (2)$$

شکل (۱) شکل موج نطق هشت و طیف فرکانسی آن را قبل و بعد از حذف DC و استفاده از فیلتر پیش تأکید نشان می‌دهد، واضح است که نویز زمینه مخصوصاً نویز ۵۰-۶۰ هرتز حذف شده است [۲].



شکل ۱- شکل موج سمت چپ، نطق هشت و اسپکتروگرام آن را بدون فیلتر پیش تأکید نشان می‌دهد. - شکل موج سمت راست، نطق هشت و اسپکتروگرام آن با حذف DC و استفاده از فیلتر پیش تأکید نشان می‌دهد.

۳- توصیف الگوریتم

۳-۱- الگوریتم انرژی تیگر

تیگر در مدل‌سازی تولید گفتار، یک الگوریتم جدید برای محاسبه انرژی یک سیگنال پیشنهاد کرد که الگوریتم انرژی تیگر نامیده می‌شود [۲]. اگر نمونه‌های یک سیگنال به شکل $x_i = A \cos(\Omega i + \Phi)$ باشند، که در آن A دامنه سیگنال، Ω فرکانس گسسته زمان و Φ فاز اولیه سیگنال باشد. انرژی لحظه‌ای E_i برای نمونه‌های x_i در الگوریتم انرژی تیگر به صورت معادله (۳) تعریف شده است:

$$E_i = x_i^2 + x_{i+1}x_{i-1} \quad (3)$$

$$= A^2 \sin^2(\Omega) \quad (4)$$

$$E_i \cong A^2 \Omega^2$$

در معادله (۴) واضح است که عبارت انرژی، نه تنها به مربع دامنه سیگنال، بلکه به مربع مؤلفه‌های فرکانس نوسان نیز وابسته

گفتار مورد استفاده قرار گرفت. در الگوریتم انرژی تیگر، سیگنال گفتار نه تنها از دامنه، بلکه از فرکانس هم تأثیر می‌پذیرد و در قسمت‌هایی از گفتار که ضعیفتر هستند، ولی مؤلفه فرکانسی متفاوتی نسبت به نویز زمینه دارند، مؤثر عمل می‌کند [۲].

روش انرژی-انترویی، روی قسمت‌هایی از گفتار که تغییرات طیفی بالایی دارند، تأکید می‌کند. انترویی، با اندازه‌گیری اطلاعات ناخواسته، شیوه مؤثری در جایگزینی پارامتر نرخ عبور از صفر می‌باشد. در این مقاله با استفاده از روش مبتنی بر انرژی اصلاح شده تیگر، نقاط اولیه شروع و انتهای گفتار تعیین می‌شوند. سپس با تغییراتی در الگوریتم انرژی-انترویی و نرمالیزه کردن منحنی‌های اصلاح شده انرژی تیگر و انرژی-انترویی و تعیین آستانه‌های مناسب، می‌توان محدوده گفتار را بدون تخمین زدن از روی فریم‌های اولیه سیگنال به صورت دقیقتر تعیین نمود. در این مقاله این نوع اندازه‌گیری جدید انرژی در گفتار فارسی بررسی می‌شود. در قسمت دوم روی حذف نویز ناشی از بکارگیری مبدل A/D و تجهیزات ضبط صدا و غیره، همچنین استفاده از یک الگوریتم مؤثر و بهینه در پیاده‌سازی آن بحث می‌شود. در قسمت سوم روی توصیف الگوریتم، محاسبات الگوریتم اصلاح شده انرژی تیگر، الگوریتم انرژی-انترویی و همچنین پیشنهادهای انجام شده در بهبود الگوریتم و نحوه تعیین محدوده گفتار بحث شده است. بخش چهارم نتایج قابل قبول و ارزیابی پیاده‌سازی این الگوریتم‌ها را روی سیگنال گفتار فارسی مورد بررسی قرار داده است. در بخش پنجم نتیجه‌گیری‌های انجام شده، ارائه شده است.

۳-۲ فیلتر پیش تأکید

ابتدا جبران‌سازی DC را که ممکن است از طریق میکروفن یا هر مبدل دیگری ایجاد شده باشد، با گرفتن میانگین سیگنال گفتار و کم کردن آن از کل سیگنال حذف می‌شود. جهت حذف نویز زمینه تداخلی که فرکانس آن زیر ۱۰۰ هرتز می‌باشد و همچنین پهن کردن طیف فرکانس سیگنال گفتار، از یک فیلتر دیجیتال FIR مرتبه اول استفاده شده می‌شود [۲]. یک روش بهینه و مؤثر در حجم محاسبات آن به صورت معادلات می‌باشد. محاسبات فیلتر پیش تأکید مرتبه اول با انجام تغییرات اندک، ساده‌تر می‌شود، ضریب a معمولاً عددی بین ۰.۹ و ۱ انتخاب می‌شود، این ضریب در شبیه‌سازی‌های انجام شده در این مقاله ۰.۹۷ (معادل با ۳۱/۳۲) در نظر گرفته شده است. معادله (۱) یک فیلتر دیجیتال مرتبه اول را نشان می‌دهد که در آن s_n نمونه n ام در یک دنباله می‌باشد. این معادله را می‌توان با معادله (۲) جایگزین نمود که در آن عملیات

تأثیر قابل ملاحظه‌ای در بالا بردن دامنه قسمت‌های ضعیف گفتار، در منحنی انرژی اصلاح شده تیگر و منحنی EEF، دارد و قادر است نقاط دقیقتری را در محیط‌های نویزی بدست آورد. عملکرد آن در شکل‌های بخش شبیه‌سازی مشاهده می‌شود. در این مقاله با یک فیلتر میان‌گذر FIR با مرتبه ۵۰ و فرکانس قطع بین ۲۵۰ هرتز و ۳۷۵۰ هرتز که بسیار شبیه باند تلفن است و تقریباً تمامی اطلاعات گفتار را دارد، سیگنال گفتار فیلتر شده است [۴].

۳-۲-۳- محاسبات انرژی اصلاح شده تیگر

این الگوریتم به صورت خلاصه در ادامه بیان شده است، ابتدا سیگنال گفتار فیلتر شده با فریم‌های ۲۰ میلی ثانیه و ۵۰ در صد روی هم‌افتادگی فریم‌بندی می‌شود. که W به عنوان طول ویندو و S_i به عنوان نمونه گفتار i ام نام‌گذاری می‌شود. انرژی اصلاح شده تیگر از لحاظ طیفی قوی‌تر از معادله (۳) است. تبدیل فوریه سریع FFT در هر فریم محاسبه شده است:

$$X(w) = \sum_{i=-\infty}^{\infty} S_i e^{-jwi} \quad (5)$$

که در آن $X(w)$ تبدیل دنباله S در حوزه فرکانس می‌باشد. سپس دامنه‌های این طیف، با مجذور مؤلفه‌های فرکانس وزن‌دهی می‌شود. در این قسمت بجای محاسبه طیف توان [۴]، مطابق معادله (۶) از طیف نمونه‌های سیگنال گفتار استفاده شده است، علت استفاده از آن به این دلیل است که باعث تقویت دامنه قسمت‌های ضعیف گفتار می‌شود. این رویکرد یکی از نوآوری‌های این مقاله محسوب می‌شود:

$$f_i = w_i^2 X(w_i) \quad (6)$$

در انتها انرژی اصلاح شده تیگر T_i به صورت ریشه جمع مؤلفه‌های فرکانس محاسبه می‌شود (معادله (۷)).

$$T_i = \sqrt{\sum_{k=1}^K f_k} \quad (7)$$

۳-۳-۳- محاسبات پیشنهادی انرژی-انترپوی

در ادامه محاسبات ویژگی‌های انرژی-انترپوی خلاصه می‌شود و به بعضی از جزئیات پیاده‌سازی‌های انجام شده، اشاره می‌شود. ابتدا برای هر فریم i ، انرژی با جمع مجذور تمام نمونه‌های فریم محاسبه می‌شود:

$$E_i = \sum_{k=1}^K (S_{ik}^2) \quad (8)$$

که در آن S_{ik} نمونه k ام از فریم i ام می‌باشد. برای هر فریم، سیگنال از حوزه زمان به حوزه فرکانس نگاشت داده می‌شود که در

است. یک اصلاح‌سازی روی الگوریتم تیگر انجام شده، که در معادله (۳) بجای محاسبه کردن انرژی لحظه‌ای برای هر نمونه سیگنال، محاسبات برای هر فریم انجام داده می‌شود، که تحت عنوان اندازه‌گیری انرژی تیگر روی فریم بکار برده می‌شود [۲]. با انجام محاسبات زیر:

۱- محاسبه کردن طیف توان برای هر فریم

۲- وزن‌دهی هر نمونه از طیف توان هر فریم با مجذور مؤلفه‌های

فرکانس

۳- گرفتن ریشه مربع طیف توان داده شده انرژی برای هر

فریم [۲].

الگوریتم انرژی تیگر در تشخیص نقاط ابتدا و انتهای گفتار بکار برده می‌شود. مخصوصاً وقتی که یک سایشی ضعیف نظیر / ف س ش ... یا انفجاری مانند / ب پ ت ... وجود داشته باشد، که دامنه ضعیف اما فرکانس بالاتری دارند.

۳-۲-۳- الگوریتم ویژگی انرژی-انترپوی (EEF)

اگرچه انرژی اصلاح شده تیگر در حضور نویزهایی از قبیل نویز سخنان نامفهوم و نویز زمینه موزیک خوب عمل می‌کند، ولی در برابر نویزهای غیرایستاد مانند صداهای مکانیکی، از قبیل نویز ناشی از صدای باز و بسته شدن درب و نویز ناشی از صدای لرزش موتور، دچار خطا می‌شود [۴]. الگوریتم ویژگی‌های انرژی-انترپوی که به اختصار EEF نامیده می‌شود، قابل اعتمادتر از روش‌های مبتنی بر انرژی خالص است و مقاومت بیشتری در برابر انواع نویزهای با شدت بالا دارد. در این روش در ابتدا، انرژی و انترپوی هر فریم محاسبه می‌شود، سپس با حذف مقدار میانگین از هر فریم، مقادیر تنظیم شده به صورت نقطه به نقطه در هم ضرب می‌شود که با عمل ضرب، روی نواحی گفتاری تأکید و نواحی غیرگفتاری تضعیف می‌شود و برای تشخیص محدوده گفتار پارامتر مناسبی می‌باشد [۴].

۳-۳-۳- توصیف الگوریتم در تشخیص محدوده گفتار

در این قسمت با استفاده از فیلتر پیش تأکید و فیلتر میان‌گذر و انجام تغییراتی در دو الگوریتم ذکر شده بالا، از آن به صورت بهینه، برای سیگنال گفتار فارسی استفاده شده است و در بخش‌های زیر به صورت خلاصه به محاسبات آن اشاره می‌شود.

۳-۳-۱- فیلتر پیش تأکید

با حذف DC و استفاده از یک فیلتر پیش تأکید FIR مرتبه اول، ضمن حذف بعضی نویزها، طیف سیگنال نرم تر می‌شود که

دارند [۵]. با جستجو از $tb1$ طبق تعریف زیر نقطه شروع مشخص می‌شود:

$$tb = \operatorname{argmin}_i \{ EEF1(i) \geq Tresh_B, tb1 \leq i \leq N \} \quad (14)$$

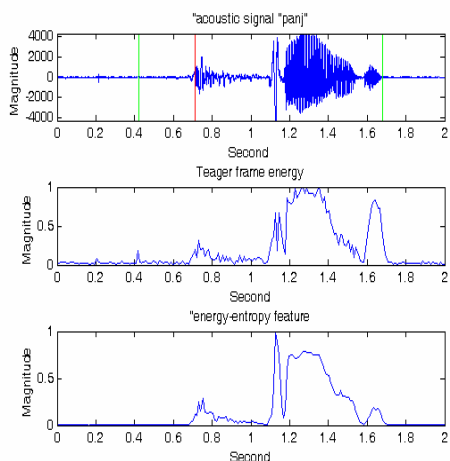
با جستجو از $te1$ طبق تعریف زیر نقطه انتهای گفتار مشخص می‌شود:

$$te = \operatorname{argmax}_i \{ EEF1(i) \geq Tresh_E, tel \leq i \leq N \} \quad (15)$$

۴- آزمایش‌های انجام شده

لازم به ذکر است که پیاده‌سازی‌های انجام شده در محیط مطلب انجام گرفته است و داده‌های مورد استفاده توسط تیم مجری جمع‌آوری گردیده است.

الگوریتم توصیف شده بالا با یک منبع اطلاعاتی گفتار تشکیل شده از اعداد صفر تا نه فارسی که توسط ۱۰ گوینده، ۵ گوینده مرد و ۵ گوینده زن، که هر کدام ۲ بار کلمات صفر تا ده را تکرار کرده‌اند تست شده است و در جدول ۱ آورده شده است. میزان دقت الگوریتم با گوش‌دادن متوالی به اجرای گفتار و استفاده از نرم‌افزارها تعیین شده است. فرکانس نمونه‌برداری ضبط شده ۲۲۵۰۰ هرتز و طول فریم ۲۰ میلی ثانیه که شامل ۴۵۰ نمونه در هر فریم همراه با ۵۰ درصد روی هم افتادگی است در شکل (۲) تشخیص نادرست محدوده گفتار، در حضور نویز ناشی از نفس کشیدن و صدای لب‌ها و کلیک‌های دهان واضح است.



شکل ۲- تشخیص نادرست نقطه شروع نطق "پنج" در حضور نویز نفس کشیدن و نویز کلیک دهان، به دلیل عدم استفاده از فیلتر پیش تأکید (محدوده گفتار بین ۰.۶۷۲ تا ۱.۰۶۸ ثانیه است).

معادله (۵) محاسبات آن انجام شده است [۴]. سپس بجای تخمین زدن تابع چگالی احتمال، با نرمالیزه کردن طیف نمونه‌ها، بهتر است تابع چگالی احتمال با مجذور طیف نمونه‌های سیگنال تخمین زده شود، که تاثیر آن در بالابردن دامنه قسمت‌های ضعیف گفتار در منحنی انرژی-انترپیی است:

$$P_i = \frac{x(w_i)^2}{\sum_{k=1}^K x(w_k)^2} \quad (9)$$

$x(w_i)$ طیف انرژی مؤلفه فرکانس w_i ، چگالی احتمال، k تعداد کل نقاط FFT در هر فریم است. انترپیی هر فریم H_i برای هر فریم i به صورت زیر تعریف می‌شود:

$$H_i = \sum_{k=1}^K P_k \log(P_k) \quad (10)$$

در اینجا بنابر الگوریتم ویژگی انرژی [۴] که معمولاً از میانگین ۱۰ فریم اول جهت کاهش نویز زمینه استفاده می‌شود، صرف نظر می‌شود. با استفاده از نمودار انرژی تیگر محدودده گفتار، بدون تخمین نویز زمینه محاسبه می‌شود و در انتها ویژگی EE به صورت زیر تعریف می‌شود [۵].

$$EEF_i = \sqrt{1 + |E_i \times H_i|} \quad (11)$$

۳-۳-۴- مکان یابی مناطق کم انرژی

انرژی تیگر تحت عنوان T و انرژی-انترپیی EEF را در فاصله بین صفر و یک نرمالیزه و به ترتیب $T1$ و $EEF1$ می‌نامیم. اولین سطح تصمیم‌گیری مربوط به $T1$ است که با تعریف سطح آستانه $Tresh_B1$ نقطه ابتدای گفتار و آستانه $Tresh_E1$ نقطه اولیه انتهای گفتار تعیین می‌شوند:

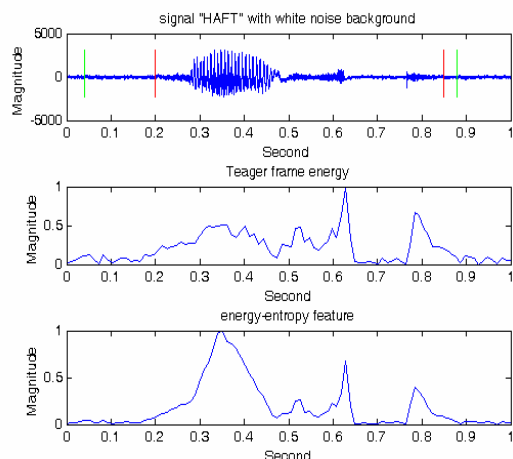
$$tb1 = \operatorname{argmin}_i \{ T1(i) \geq Tresh_B1, 1 \leq i \leq N \} \quad (12)$$

$$te1 = \operatorname{argmin}_i \{ T1(i) \geq Tresh_E1, N \leq i \leq 1 \} \quad (13)$$

مقداری که برای آستانه‌ها بکار برده می‌شود مقدار ثابتی است که در اینجا ۰.۹، برای آستانه‌های $Tresh_B1$ و $Tresh_E1$ بکار رفته است. بنابراین انرژی فریم تیگر، به عنوان تخمینی از ابتدا و انتهای گفتار است و از ویژگی انرژی-انترپیی طبق توصیف زیر جهت تعیین نقاط واقعی استفاده شده است.

۳-۳-۵- تشخیص نهایی محدوده گفتار

برای منحنی نرمالیزه شده $EEF1$ دو آستانه یکی برای محدوده شروع گفتار $Tresh_B$ و دیگری برای محدوده انتهای گفتار $Tresh_E$ در نظر گرفته می‌شود، این دو آستانه مقدار خیلی کمتری



شکل ۵- تشخیص صحیح محدوده گفتار در نطق "هفت" در حضور نویز سفید.

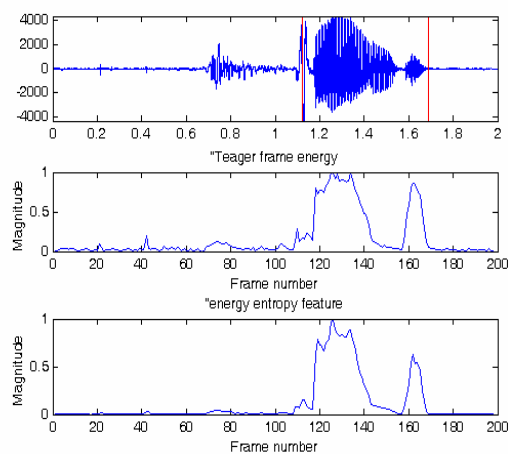
جدول ۱- بازدهی الگوریتم، در تشخیص ابتدا و انتهای اعداد صفر تا نه.

بازدهی	تعداد خطا	تعداد فایلها	گوینده	روش
%۸۶	۱۶	۱۰۰	مرد	انرژی- ZCR
			زن	
%۹۲	۸	۱۰۰	مرد	انرژی-انترپوی
			زن	
%۹۸	۲	۱۰۰	مرد	انرژی-انترپوی
			زن	
%۹۷	۳	۱۰۰	زن	پیشنهادی

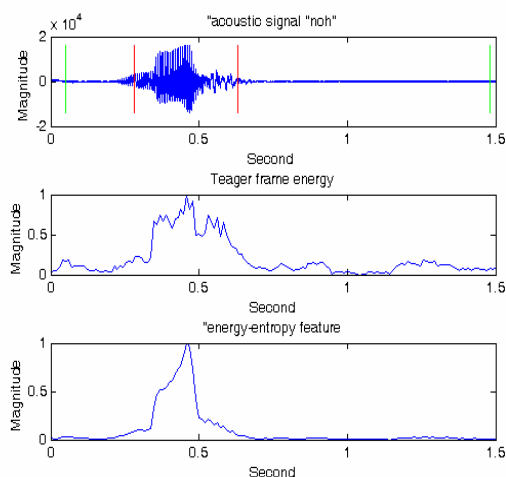
۵- نتیجه گیری

واج های خیلی ضعیف ابتدا و انتهای گفتار در محیط های نویزی به سهولت قابل تشخیص نیست. با استفاده از تغییراتی در انرژی اصلاح شده تیگر، و همچنین اصلاحاتی در روش انرژی-انترپوی، واج های ضعیفی که در ابتدا و یا انتهای گفتار قرار می گیرند، بهتر آشکار می شوند. نتایج آزمایشات نشان می دهد که الگوریتم پیشنهادی در این مقاله، با استفاده از فیلتر پیش تأکید قبل از فیلتر میان گذر و تغییراتی در الگوریتم تیگر و الگوریتم اصلاح شده انرژی-انترپوی و استفاده از آستانه های ثابت، بازدهی قابل قبولی در تشخیص نقاط ابتدا و انتهای گفتار دارد، مخصوصاً وقتی که واج های ضعیف در شروع یا انتهای گفتار قرار گیرد. این الگوریتم بدون تخمین نویز، در محیط هایی با نویز بالا، بازدهی مناسبی از خود نشان داده است.

در منحنی نرمالیزه شده انرژی-انترپوی شکل (۳)، حذف نویز ناشی از نفس کشیدن را که در فاصله ۰.۷ تا ۱.۱ ثانیه است به وضوح دیده می شود. شکل (۴) یک نمونه تشخیص صحیح نقطه شروع و پایان نطق "نه" را با وجود نویز در ابتدای گفتار را نشان می دهد. شکل (۵) تشخیص صحیح کلمه هفت را در حضور نویزهای غیرایستاد نظیر نویز سفید نشان می دهد.



شکل ۳- تشخیص صحیح نقاط ابتدا و انتهای نطق "پنج" بعد از استفاده فیلتر پیش تأکید و حذف نویز نفس کشیدن و نویز دهان و لب ها و متعادل شدن شکل موج در منحنی انرژی-انترپوی.



شکل ۴- تشخیص نقاط ابتدا و انتهای عدد "نه" با وجود نویز در ابتدای گفتار.

۶- مراجع

- [1] L.R. Rabiner, M.R. Sambur; “**An Algorithm for Determining the Endpoints of Isolated Utterances**”, Bell Sjisl. Tech. J., Vol. 54, pp.297-315, 1975.
- [2] [G.S. Ying, C.D. Mitchell, L.H. Jamieson; “**Endpoint Detection of Isolated Utterances Based on A Modified Teager Energy Measurement**”. In Proc. IEEE ICASSP-92, pp.732-735, 1992.
- [3] H. Qiang, Z. Youwei; “**On Prefiltering and Endpoint Detection of Speech Signal**”, Proceedings of ICSP 1998 , pp.749-752, 1998.
- [4] L.S.Huang, C.H.Yang; “**A Novel Approach to robust speech endpoint detection in car environments**”, ICASSP-2000, Vol. 3, PP.1751-1754, 2000
- [5] L. Gu and S.A. Zahorian; “**A New Robust Algorithm for Isolated Word Endpoint Detection**”, IV-4161 ICASSP, 2002.
- [6] W.Han, C.F.Chan, C.S.Choy, K.P.Pun; “**An Efficient MFCC Extraction Method in Speech Recognition**”, ISCAS 2006. Volume , Issue , 4 pp, 2006.