

**ANALIZA CEPSTRALĂ ÎN SISTEME DE RECUNOAȘTERE A VORBIRII****Vladimir DUBINEANSCHI***Universitatea Agrară de Stat din Moldova*

Automatic Speech Recognition (ARS) has progressed considerably over the past several decades, but still has not achieved the potential imagined at its very beginning. Almost all of the existing applications of ASR systems are PC based. This publication describes the possibility of using spectral analyses in the speech recognition.

**Material și metodă**

Recunoașterea vorbirii este la ora actuală un domeniu în creștere datorită succeselor din ultimii ani în prelucrarea numerică a semnalelor, în dezvoltarea de aplicații cu procesoare specializate și în rapida dezvoltare a tehnicii de calcul. Sistemele actuale de recunoaștere a vorbirii se situează deocamdată în limite restrânse ale parametrilor caracteristici și dedicate unor aplicații specifice. Din punctul de vedere al dimensiunii vocabularului și al modului de vorbire, sistemele de recunoaștere cu performanțe acceptabile se împart în trei categorii principale.

- sisteme cu vocabular mic (10-100 cuvinte);
- sisteme cu vocabular mediu și mare și vorbire izolată (10 000-20 000 cuvinte);
- sisteme cu vocabular mediu și vorbire conectată sau continuă, restrictivă la un domeniu de aplicabilitate (1 000-5 000 cuvinte).

Cele mai multe sisteme realizate aparțin claselor sistemelor mici și mijlocii cu recunoașterea vorbirii izolate. Sistemele de recunoaștere a vorbirii continue, în marea lor majoritate, există doar în formă experimentală, în condiții de laborator. Chiar și sistemele utilizate în practică, cele pentru vorbirea izolată sau conectată, nu sunt destul de robuste la zgomotul mediului în care funcționează și la variabilitatea vorbirii. Toate sistemele dau performanțe mai bune, dacă numărul de utilizatori este mai redus și dacă cei care folosesc sistemul sunt cei a căror voce s-a folosit pentru învățarea sistemului. Performanțele se degradează semnificativ, dacă vorbitorii se schimbă sau dacă sistemul este folosit cu alte cuvinte decât cu cele pentru care a fost antrenat. Caracteristicile principale ale unui sistem de recunoaștere automată a vorbirii, fără a aminti parametrii și metodele specifice prin care s-a implementat, sunt următoarele:

- dimensiunea vocabularului, adică numărul de cuvinte capabil să le recunoască
- monolocator sau multilocutor (aici se poate preciza și sexul vorbitorilor)
- vorbirea izolată sau continuă
- condiții de zgomot și robustețea sistemului
- domeniul de aplicabilitate
- timpul de operare, care poate fi în timp real, cu întârziere sau off-line
- procentajul de recunoaștere
- costul.

**Rezultate și discuții**

Semnalul vocal reprezintă convoluția în timp a semnalului excitație produs de vibrația corzilor vocale și răspunsul în timp al filtrului reprezentat de tractul vocal. Astfel, supus analizei vom avea un semnal convolutat. Separarea celor două semnale, precum și analiza parametrilor în domeniul temporal este imposibilă. Prin urmare, problema care se pune este trecerea semnalului în alt domeniu în care componentele se combină liniar, caz în care separarea ar fi simplă. Un astfel de domeniu este domeniul frecvența. Trecerea din domeniul temporal în domeniul frecvența se face cu ajutorul transformatei Fourier. În urma unei transformări liniare a unui semnal combinat liniar se obține o combinație tot liniară. În consecință, dacă un semnal este obținut prin combinația liniară a două semnale de frecvențe diferite, în domeniul frecvența combinând liniar componentele, vom obține semnalul ca o combinație liniară a componentelor spectrale. Pentru a demonstra că cele două componente sunt ușor de separat în domeniul frecvența, vom lua un exemplu în care un semnal de joasă frecvență  $x(n)$  este perturbat de zgomotul de înaltă frecvență  $w(n)$ .  $y(n) = x(n) + w(n)$ . Aplicând transformata Fourier în domeniul frecvența, vom avea câte o componentă separată la frecvența celor două

semnale. Prin eliminarea frecvențelor nedorite și transformarea inversă înapoi în domeniul timp se obține semnalul curățat de zgomot. În cazul semnalului vocal componentele nu sunt combinate liniar. Astfel, pentru a separa liniar componentele trebuie trecut și semnalul vocal într-un domeniu, unde componentele se combină liniar. Fie  $s(n)$  semnalul vocal obținut prin convoluția excitației corzilor vocale  $e(n)$  și a răspunsului la impuls al funcției de transfer a tractului vocal  $\Phi(n)$ .

$$s(n) = e(n) \otimes \Phi(n) \tag{1}$$

Aplicând asupra acestei egalități transformata Fourier, convoluția devine înmulțire și, prin urmare, vom avea:

$$F\{s(n)\} = F\{e(n) \otimes \Phi(n)\} = F\{e(n)\} \cdot F\{\Phi(n)\} = E(\omega) \cdot \Phi(\omega) \tag{2}$$

Pentru o combinare liniară avem nevoie de adunare. Trecerea din multiplicativ în aditiv se face prin logaritmare și, prin urmare, vom avea:

$$C(\omega) = \log[F\{s(n)\}] = \log[E(\omega) \cdot \Phi(\omega)] = \log[E(\omega)] + \log[\Phi(\omega)] \tag{3}$$

Am ajuns într-un domeniu în care componentele sunt combinate liniar. De acum încolo vom considera  $C(\omega)$  ca semnal. Pentru a putea trece într-un domeniu în care componentele  $C(\omega)$  sunt ușor de separat, putem aplica iarăși transformata Fourier, trecând în domeniu „frecvența”:

$$c(n) = \frac{1}{2\pi} \cdot \int_{-\pi}^{\pi} C(\omega) \cdot e^{-jn\omega} \cdot d\omega \tag{4}$$

Funcția  $C(\omega)$  fiind o funcție reală pară, relația (4) este echivalentă cu:

$$c(n) = \frac{1}{2\pi} \cdot \int_{-\pi}^{\pi} C(\omega) \cdot e^{jn\omega} \cdot d\omega, \tag{5}$$

deoarece

$$\int_{-\pi}^{\pi} C(\omega) \cdot e^{-jn\omega} \cdot d\omega = \int_{-\pi}^{\pi} C(\omega) \cdot (\cos(n\omega) - j \cdot \sin(n\omega)) \cdot d\omega = \int_{-\pi}^{\pi} C(\omega) \cdot \cos(n\omega) \cdot d\omega$$

$$\int_{-\pi}^{\pi} C(\omega) \cdot e^{jn\omega} \cdot d\omega = \int_{-\pi}^{\pi} C(\omega) \cdot (\cos(n\omega) + j \cdot \sin(n\omega)) \cdot d\omega = \int_{-\pi}^{\pi} C(\omega) \cdot \cos(n\omega) \cdot d\omega$$

unde

$$\int_{-\pi}^{\pi} C(\omega) \cdot \sin(n\omega) \cdot d\omega = 0$$

Aplicând amândouă relațiile, se obține ca rezultat semnalul de start. În cazul în care din domeniul spectral se alege aplicarea transformării Fourier directe, domeniul în care se ajunge nu are semnificație fizică. Dacă se alege transformarea Fourier inversă, se ajunge într-un domeniu temporal. Acest domeniu temporal nu mai este identic cu domeniul temporal de unde am pornit, de aceea noțiunile utilizate pentru caracterizare se modifică (noțiuni de natură spectrală într-un domeniu temporal). Astfel, noul domeniu va avea denumirea de *domeniu cepstral*. Denumirile din noul domeniu se determină prin corespondența cu domeniul spectral, însă inversând literele din prima silabă. Deci, se definesc următoarele denumiri: *spectrum*  $\Leftrightarrow$  *cepstrum*, *frecvență*  $\Leftrightarrow$  *cvefrență*, *armonica*  $\Leftrightarrow$  *ramonica*, *filtrare*  $\Leftrightarrow$  *liftrare*. În acest domeniu cepstral, componentele combinate ale semnalului original pot fi separate prin metode simple, cum ar fi liftrarea.

Schema-bloc a întregului proces este prezentată în figura următoare:

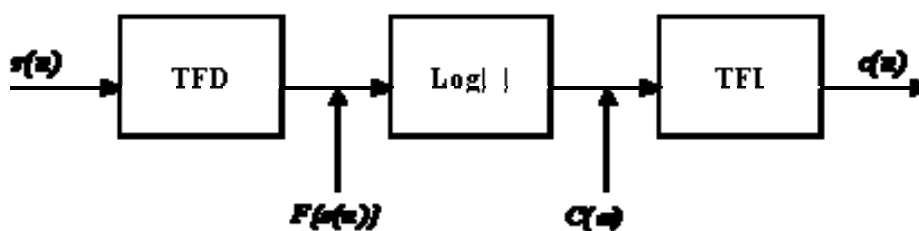


Fig.1. Schema-bloc de determinare a cepstrului.

Spectrul semnalului vocal este compus prin multiplicare dintr-o *anvelopă* de variație lentă și *pulsuri*, o componentă cu variație rapidă. Spectrul logaritmic are forma asemănătoare, însă el reprezintă combinația liniară a celor două componente. În cepstrum (spectrul spectrului) se pot vedea separate cele două componente, una de joasă cvefrență, rezultată din componenta lent variabilă a anvelopei, și componenta de cvefrență înaltă, rezultată din componenta rapid variabilă a spectrului. Distanța în cvefrență dintre două impulsuri, sau distanța de la origine la primul impuls, se poate interpreta ca și perioada acestor impulsuri, deci se poate obține frecvența fundamentală de rezonanță a corzilor vocale. Această metodă este folosită pentru determinarea frecvenței fundamentale și pentru determinarea caracterului vocalic sau nevocalic al cadrului de semnal analizat.

$$c_s(n) = c_e(n) + c_r(n) \quad (6)$$

Bineînțeles, această metodă poate fi aplicată pentru determinarea frecvenței fundamentale numai în cazul în care cele două componente combinatorice sunt destul de distanțate între ele. În cazul semnalelor de origine vocalică (vocale), metoda se justifică, însă în cazul sunetelor nevocalice (consoane), metoda este inadecvată, pentru că impulsurile, dacă există, se confundă cu zgomotul cepstral de înaltă cvefrență. Pe baza semnalului cepstral se pot obține, independent, spectrul de variație lentă sau cel de variație rapidă, prin operația de liftrare și transformare în domeniul frecvența. De exemplu, dacă se dorește anvelopa spectrală, din relația (6) se elimină  $c_r(n)$  și se aplică transformata Fourier directă.

### Concluzii

Odată determinați coeficienții cepstrali pentru cadrele dintr-un semnal vocal, se pune problema comparării acestora cu alt set de coeficienți ai cuvintelor etaloane din vocabularul analizat. Compararea se face calculând distanța între cele două cepstre. Pentru aceasta folosim relația lui Parseval, conform căreia energia semnalului periodic, calculată din eşantioane, este egală cu energia calculată din componentele spectrale înmulțite cu un factor. Considerând semnalul „temporal” cel spectral, se poate scrie egalitatea:

$$\sum_{n=1}^{2\pi} c_n^2 = \frac{1}{2\pi} \cdot \int_{-\pi}^{\pi} (\log(|H(\omega)|))^2 d\omega \quad (7)$$

$$s(n) = \sum_{p=1}^L a_p \cdot s(n-p) + U \cdot u(n) \quad (8)$$

În formula

demonstrată anterior se consideră factorul de amplificare  $G=1$ , iar suma coeficienților cepstrali făcându-se de la  $n=1$ , dat fiind faptul că ceilalți coeficienți sunt nuli. Distanța dintre cele două spectre se poate calcula ca fiind energia conținută în spectrul obținut prin scădere componentă cu componentă, care este identică cu o distanță euclidiană:

$$E_{\Delta} = \frac{1}{2\pi} \cdot \int_{-\pi}^{\pi} |H_1(\omega) - H_2(\omega)|^2 d\omega \quad (9)$$

Prin logaritmare, egalitatea (9) devine:

$$E_{\Delta} = \frac{1}{2\pi} \cdot \int_{-\pi}^{\pi} \log(|H_1(\omega) - H_2(\omega)|)^2 \cdot d\omega \quad (10)$$

Din (9) se observa că energia  $E_{\Delta}$  este egală cu energia semnalului cepstral diferențial, care, fiind de natură temporală, se poate calcula și ca energia semnalului diferență:

$$E_{\Delta} = \sum_{n=1}^{2\pi} (c_{1n} - c_{2n})^2 \quad (11)$$

Se poate scrie relația:

$$\sum_{n=1}^{2\pi} (c_{1n} - c_{2n})^2 = \frac{1}{2\pi} \cdot \int_{-\pi}^{\pi} \log(|H_1(\omega) - H_2(\omega)|)^2 \cdot d\omega \quad (12)$$

Prin urmare, distanța euclidiană între două semnale cepstrale este o distanță care poate exprima o diferență spectrală. În practică distanța se calculează între coeficienți până la ordinul  $P$ . Generalizând, distanța cepstrală se definește:

$$d_{cep} = \sqrt{\sum_{p=1}^P (c_{1p} - c_{2p})^2} \quad (13)$$

A fost studiat efectul ponderării în calculul distanței cepstrale. Ponderarea îmbunătățește performanțele calculului distanței, ale recunoașterii, deoarece coeficienții cepstrali de ordin diferit au importanță diferită din punctul de vedere al tractului vocal. Analiza cepstrală este o analiză specială, prin care efectul compus al excitației și al funcției de transfer a tractului vocal se separă cu o eficiență mai bună, ca în cazul analizei spectrale sau liniar predictive.

#### Bibliografie:

1. Обжелян Н.К., Трунин–Донской В.Н. Речевое общение в системах „Человек – ЭВМ”. - Кишинёв: Штиинца, 1985.
2. Винцюк Т.К. Анализ, распознавание и интерпретация речевых сигналов. - Киев: Наукова Думка, 1987.
3. Rabiner L.R., Juang B.H. Fundamentals of speech recognition. - Prentice Hall, 1993.
4. Speech Analysis FAQ <http://svr-www.eng.cam.ac.uk/~ajr/SA95/SpeechAnalysis.html>
5. Microsoft Speech API Help. <http://research.microsoft.com>.
6. Методы автоматического распознавания речи: В 2-х книгах. Пер. с англ. / Под ред. У.Ли. - Москва: Мир, 1983.

*Prezentat la 31.01.2011*