

○ 임 제열, 안 수길
서울대학교 대학원 전자공학과

(A Study on the Recognition of Continuous Korean Digits using Phoneme Classification)

Jaeyeol RHEEM and Souguil ANN
Dept. of Electronics Eng, Seoul National Univ.

1. 서 론

오늘날 음성인식(speech recognition)은 급변하는 정보화 시대에 있어서 산업적으로 또는 공공적으로 광범위한 주목을 받고 있다. 외국의 경우에는 문장단위의 연속음 인식을 위해서 범용가격 연구 계획이 추진되어져와, 음성의 음향학적 측면외에도 언어학, 음성학, 문장학 등등의 언어 전반적 사항이 연구되어왔으며 인공지능(Artificial Intelligence)과의 결합도 시도되고 있다[1,2]. 그러나 국내의 연구의 경우 제한된 어휘수의 고립단어, 또는 연결단어 인식에 관한 연구는 활발하나, 연속음 인식을 위한 한국어전반에 관한 연구는 아직 초보적인 단계에 있다[3,4]. 따라서 한국어의 연속음 인식을 위한 기본적인 연구가 수행되어야 할 필요성이 있다.

본 논문은 한국어 연속음 인식의 필요성을 인식하고 제한된 대상인 한국어 연속 숫자음에 대한 인식 시스템의 구현을 시도 했다. 이는 한국어 전반에 관한 언어학, 음운학, 문장학 등등의 광범위한 지식없이 적용할 수 있는 소 대상이기 때문이다. 그 특징으로 강력한 에너지 곡선을 이용한 새로운 비음구간 점을 알고리즘과 음성 내의 음 구간 점을 알고리즘을 제시 했으며, optimal preemphasis factor의 이용과 기본주파수 검출시의 doubling과 halving 방지용 알고리즘을 사용했고, 포오반트 검출시 pole enhancement 방법을 적용했다. 그리고 화자 내부에서의 발생 방법의 변형에서 오는 변음을 음이교차 LPC 스펙트럼에서의 대역 에너지를 파라미터로 사용했다.

II. 연속음 인식 및 숫자음 구조

1. 연속음 인식

고립단어 인식을 위해서는 패턴 매칭(pattern matching) 방법이 주로 사용되고 있으나, 많은 단어를 인식할 때 그 계산량이 방대해지며, 조사 및 어미변화에 따른 대이러량 또한 커지고, 패턴의 표준화에 따른 문제점이 발생하게 된다. 따라서 음성 인식의 궁극적 목표인 자연스럽게 발생된 연속음성을 인식 또는 이해하기 위해서는 인식의 단위를 단어이하의 음소와 같은 미소 단위로 해야한다. 그리고 이러한 음소단위의 인식을 위해서는 개인자의 문체와 조음 결합이 가장 큰 어려움으로 지적되고 있다. 한편으로 무제한의 연속음음 인식을 하기 위해서는 단어 내부의 음소의 구별 외에도 문장으로서의 구조적 맥락을 이해하기위하여 음성학외에도 언어학, 음성학, 어휘론과 문법, 강세, 억양변화 등의 음조현상등 언어 전반 사항을 연구해야된다. 그림 2-1에 일반적인 연속음인식 시스템의 구조를 보았다. 신호 처리 부분은 음성의 음향학적 특징 파라미터를 추출하는 부분으로, 이때 사용되는 파라미터는 화자나, 발음시간, 발음속도 등의 영향을 적게 받는 안정된 것을 사용해야 한다. Bottom-Up Approach 부분은 추출된 파라미터를 이용하여 음성학적으로 분류하는 분류 인식에 해당된다. 음소단위인 경우 음소분류(phoneme classification)을 하게된다. 여기서서는 앞단에서 추출된 파라미터를 이용하여 특징 벡터를 추출한 후 분류 규칙(classification rule)에 의하여 유사성질을 지니는

음소단위의 유사군으로 분류된다. Top-Down Verification 부분은 앞단에서 얻은 음소 정보와 언어적인 문법구조를 이용하여 확인하는 부분이다. 이부분에서는 앞단의 분류된 음소 정보에 대하여 문법, 어휘, 발음 순환 등등에 관한 규칙을 이용하여, 수정 및 편집(editing)을 수행함으로써 최종인식을 하게된다.

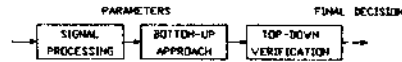


그림 2-1 연속음 인식 시스템의 구성

2. 한국어 숫자음의 음성학적 구조

본 논문에서 인식대상으로 한 숫자음은 '퐁', '일', ..., '구'의 10개의 단음절어음 연속해서 발음한 숫자이다. 그리고 숫자음 10개외에도 받도보 /에/발음도 포함시켰는데 이는 자연스러운 연속 숫자음을 얻고자 전화번호의 숫자들 인식 대상으로 삼아 전화번호의 '뿔' /에/로 발음을 했기때문이다. 10개의 숫자음을 발음기호와 함께 표 2-1에 나타냈다. 한국어 숫자음의 특징으로는 순수 모음으로 이루어진 숫자음 (이,오) 보다는 초성자음이 있는 숫자음(삼,사,칠,팔,구,풍)이 많아 상대적으로 초성 자음의 유성음화 현상이 적게 일어나며, 단음절(monosyllabic) 성격을 띄고 있고, 초성자음과 종성자음을 만 경우에는 모음을 구체적으로 알지 못해도 구별할 수 있고, 어느 한 음소운 강화적 침몰하지 않아도 추장할 수 있어 음소결합의 독립성이 작다는 점을 들 수 있다.

표 2-1 한국어 숫자음

숫자	분	특	발음 기호
일	관성모음 + 비음		/il/
이	관성모음		/i/
삼	무성음 + 관성모음 + 비음		/sam/
사	무성음 + 관성모음		/sa/
오	관성모음		/o/
육	관성모음 + 무성음		/yuk/
칠	무성음 + 관성모음 + 비음		/tʃil/
팔	무성음 + 관성모음 + 비음		/pʌl/
구	무성음 + 무성모음		/ku/
풍	무성음 + 무성모음 + 비음		/kʊŋ/

표 2-2 음소별 종류

음 소 명	종	류
모음(vowel)	/a/, /ɪ/, /o/, /u/, /ɔ/	(1, 1, 1, 1, 1)
류음(glide)	/j/	(2)
비음(nasal)	/m/, /n/	(2, 1)
피열음(plosive)	/p/, /t/	(1, 1)
마찰음(frictive)	/s/	(1)
마찰음(frictive)	/ʃ/	(1)
매체음(stop)	/k/	(60)

111. 음성 인식 시스템의 구성

1. 시스템의 구성

한국어 연속 숫자음을 인식하기 위한 전체 구성은 아래의 그림 3-1과 같다.

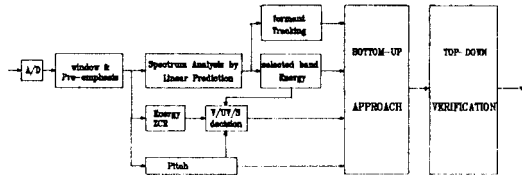


그림 3-1 시스템 구성

Microphone을 통해 입력된 음성 신호를 증폭기를 통과하여 증폭된 후, 4KHz의 차단 주파수를 갖는 저역 여파기를 통과하여 12-bit A/D 변환기에 의해서 10KHz의 표본화 주파수로 표본화된다. A/D 변환된 신호는 기본 신호 처리로 25.6ms의 Hamming window를 통과하여 optimal pre-emphasis 된다. 그리고 14차의 선형 예측 분석법이 매 5ms씩마다 수행되며, 그형에 주계수는 Melk(1/8)의 autocorrelation 방법을 사용했다. 각 신호 예측계수들은 256 point 40Hz resolution의 스펙트럼 분석되어 13상 대역의 에너지를 구하게 된다. 이때 고려된 에너지 대역은 다음과 같다.

ferg0	:	0	5000Hz
ferg1	:	300	5000Hz
ferg2	:	100	3000Hz
ferg3	:	640	2800Hz
ferg4	:	3700	5000Hz

기본 주파수의 추출 방법은 FFT기 방법을 수정하여 사용했으며, 포오펀트 추출은 McClellan[8,9]의 알고리즘을 수정하여 사용했다.

2. 음성음/무성음/무음구간 결정

무음(silence) 구간과 음성 구간의 구별은 대역 에너지 ferg1을 이용한다. 이외권이 300에서 5000Hz의 대역을 이용하는 것은 voiced bar, 또는 voiced stop은 무음 구간에 포함시키기 위해서이다. ferg1의 함수에서 음성 신호 시작 전후의 각 100ms씩 동안의 배경 잡음에 대한 평균 값을 얻은 후, 이 값에 대해 3dB 이상의 frame에 대해서는 음성구간으로 간주하며 그 미만에 대해서는 무음 구간으로 간주한다. 이때 무음 구간 중에서 pitch가 0이 아닌 경우를 특히 voiced silence라고 명명한다.

$$\text{threshold} = \text{background noiselevel} + 3\text{dB} \quad (3-1)$$

음성 구간에 대해서는 pitch의 존재에 따라서 간단히 음성음 구간과 무성음 구간을 구분한다.

IV. 사전 분할 및 음소 분류 (Bottom-Up Approach)

Bottom-Up approach는 그림 4-1과 같이 사전 분할(primary segmentation) 부분과 음소 분류(phoneme classification) 부분으로 구성되어있다.

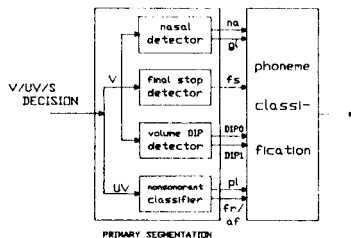


그림 4-1 Bottom-Up Approach

1. 비음 구간의 검출 및 분류

일반적으로 연속음 인식에서는 3절에서 정의된 volume DIP을 이용하여 비음 및 기타 유성 지음 또는 낮은 에너지의 무음 구간을 검출한다[11]. 그러나, 한국어 숫자음의 경우 조성자음의 영향으로 연속된 다음 단어와의 조음현상이 잘 일어나지 않을므로써 DIP으로 검출하지 못하는 구간이 자주 발생하게 된다. 이것을 방지하고자 시간영역에서의 비음구간 검출기를 선행 시켰다.

그림 4-2에서와같이 기 frame 마다의 에너지 차를 간략화 시켜서 얻은 패턴으로 비음 구간을 검출할 수 있다. smoothing은 3 point Hanning window(1/4, 1/2, 1/4)를 사용했으며, 에너지 차는 식(4-1)로 정의되고, median smoothing[12]을 이용했다. 그림 4-3에는 비음 구간의 패턴과 확인 알고리즘을 제시했으며 그림 4-4에는 그 결과를 보였다.

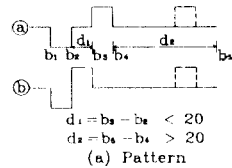
$$\text{ergd}(i) = \text{out}(0, i) - \text{out}(0, i-1) \quad (4-1)$$

비음 패턴으로 추출된 부분은 비음과 유음 그리고 기타 유성음 구간으로 확인 알고리즘을 통해 1차 확인된후 음소분류를 하게된다. 1차 확인서의 파라미터는 다음과 같다.

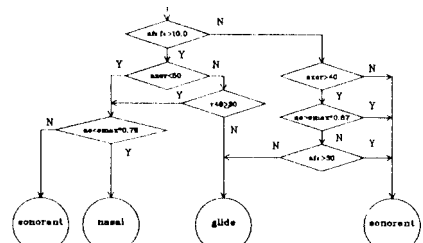
af214	:	ferg4와 ferg2의 평균 비음
azet	:	평균 양피치
fac	:	ferg3와 ferg4의 평균 비음
ac	:	평균 에너지
emax	:	입력 신호의 최대 에너지
af4	:	ferg4의 평균 에너지



그림 4-2 비음구간 검출도



(a) Pattern



(b) Check Algorithm

그림 4-3 비음구간 패턴 및 확인 알고리즘

2. 종성 내파음구간의 검출 및 분류

종성 내파음(final impulsive stop) 구간의 검출도 비음 검출과 같이 간략화된 에너지차 패턴을 이용하여 구할 수 있다. 특히 종성 내파음 뒤에는 반드시 무음 구간이 오며, 에너지 차가 현격하다[13]는 사실을 이용하여 1차 확인후 음소 분류할 계획이 된다.

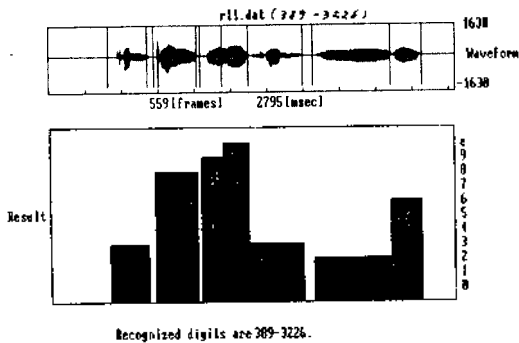


그림 5-1 최종 인식 결과

표 5-1 음소 분류 결과

분류된 음소	
음소	a i u o ju /r/ fs k ŋ ɳ l p k m ts/ e
a	25 3 3
i	9 3 7
u	2 1 9 1
o	10 3 2
ju	6 1 1
fs-k	6 4 1 1
ŋ	8 16
ɳ	3 13
l	14 18
p	7
k	
m	
ts	

표 5-2 최종 인식 결과

	데이터	인식 결과	비교
1	309-3226	309-3226	
2	446-0909	446-0909	
3	878-2670	878-2670	
4	434-3505	434-3509	
5	879-5290	879-5290	
6	482-6556	482-2555-6	→ 2
7	386-2824	38 2824	5 →
8	878-2637	878-2037	
9	647-4173	64 4173	7 →
10	355-0109	355-0109	
11	446-0909	446-0909	

2. 실험 결과 및 검토

음소 분류 결과 부성음구간은 모두 옳게 분류되었으며, 비음 검출기는 모든 비음을 검출하는 외에 류음 및 유성을 부분 검출했으나 뒤에 오는 음소 분류시 모두 옳바로 판정되었다. 종성 내파음의 경우에도 모두 바르게 판정되었다. 그러나 모음의 판정에는 음소 분류시 많은 오�판정을 했다. 특히 한 모음구간이 여러 (1)P로 나뉘어 길 경우에는 한 모음 내에서도 여러 모음으로 판명되는 경우가 다소 있었으며 종성 자음으로 비음이나 류음이 오는 경우에는 종성 모음을 구별하지 못하는 경우가 있었다. 그래서 /a/ 모음의 경우 삼파 램의 비음 및 류음에 의하여 미확인 모음(/r/)으로 분류된 경우가 10% 발생했으며, 특히 램의 류음 /m/의 영향으로 류음(/l/)으로 분류된 경우가 10% 발생했다. 각 음소의 분류결과는 표 5-1에 있다.

최종인식 결과 11개의 전화 번호에있는 숫자음 77개에 대해서는 2개의 오인식이 발생해 97.4%의 인식률을 얻었으며 /에/ 모음을 고려한 전체 88개의 음절에 대해서는 4개의 오인식이 발생해 95.5%의 인식률을 얻었다. 표 5-2에있는 Top-Down approach후의 최종 인식결과를 보면 주된 오인식이 모음 /에/에 있음을 알 수 있다. 이것은 /에/ 모음의 앞 또는 뒤에 종성 모음이 올 경우에, /에/ 모음이 종성 모음과 유사해지거나 또는 종성 모음이 /에/ 모음과 유사하게 되어 /에/모음과 종성 모음의 중간적인 성질을 띄게 되기 때문이다. 한편 77개의 숫자음 중 2개의 오인식의 경우 '7'가 /치레/로 발음된 경우와 '5'가 /에/로 판정된 경우이므로 순수 숫자음만을 가지고 실험할 경우 더 높은 인식률이 예상된다.

VI. 결론

본 논문은 한국어 연속음중 제한된 대상인 숫자음에 대한 음소단위의 인식을 시도 했다. 인식기의 구성은 일반적인 연속음인식 시스템의 구성을 따랐으며 optimal pre-emphasis factor를 파라미터의 하나로 사용했고, SIFT 알고리즘과 McCandless의 포먼트 추적 알고리즘을 수정해서 적용했다. 그리고 간략화된 에너지차의 패턴을 이용하는, 한국어 숫자음에 맞는 비음구간의 검출 방법과 종성 내파음의 검출 방법을 제안했다. 그리고 /에/로 한 화자가 발음한 전화번호를 대상으로 인식 실험한 결과 숫자음만의 경우 97.4%, /에/를 포함한 경우 95.5%의 인식률을 얻었다. 본 논문은 한국어 연속음 인식 시스템의 구성을 위한 앞으로의 시스템구현에 발판이 될 수 있겠으며, 특히 새로이 제안된 비음 및 종성 내파음의 검출 방법은 앞으로도 유용하게 사용될 수 있으리라 본다. 그리고 앞으로의 과제에는 한국어의 모든 음소를 분류하기 위한 분류 파라미터의 개발과 그 시스템의 구현, 그리고 무제한의 화자 인식을 위한 정규화의 방법 연구등이 될 것이다.

참고 문헌

- [1] Wayne A. Lea, Trends in Speech Recognition, Prentice Hall, 1980.
- [2] Jean-Paul Haton, Automatic Speech Analysis and Recognition D. Reidel Publishing Company, 1981.
- [3] 이 철희, "한국어 연속음 중 단모음 인식에 관한 연구," 서울대학교 대학원 공학석사학위 논문, 1986.
- [4] 이용주, 김경태, 차광현, "한국어 단모음의 성별, 연령별 특징 변화 및 인식," 음성통신 및 처리기술 Workshop 논문집, 1987.
- [5] D.Ter Haar, Mechanisms of Speech Recognition, International series in natural Philosophy, Vol.85, Pergamon Press.
- [6] J.D.Markel and A.H.Gray, Jr., Linear Prediction of Speech, Springer Verlag Berlin Heidelberg, New York, 1976.
- [7] Markel, J.D., "The SIFT Algorithm for Fundamental Frequency Estimation," IEEE Trans. AU-20, 367-377, 1972.
- [8] Stephanie S.McCandless, "An Algorithm for Automatic Formant Extraction Using Linear Prediction Spectra," IEEE Trans. ASSP-22, No.2,135-141, April, 1974.
- [9] Stephanie Seneff, "Modifications to Formant Tracking Algorithm of April 1974," IEEE Trans. ASSP-22, pp.192-193, April, 1976.
- [10] 배명진, 임재열, 안수길, "음성 발생 모델로부터 G-peak를 이용한 음성 에너지 추출에 관한 연구," KIEE, Vol. 24, pp. 12-17, 1987.
- [11] Clifford J. Weinstein, Stephanie S.McCandless, LEE F.Mondshein, and Victor W.Zue, "A System for Acoustic Phonetic Analysis of Continuous Speech," IEEE Trans. ASSP-23, No.1, pp.54-67, Feb., 1975.
- [12] H.R.Rabiner, M.R.Sambur, and C.E.Schmidt, "Applications of a Nonlinear Smoothing Algorithm to Speech Processing," IEEE Trans. ASSP-23, No.6, pp.552-557, Dec., 1975.
- [13] 배명진, 배명진, 안수길, "음성 신호에서의 종성내파음 구간 검출에 관한 연구," 음성통신 및 처리기술 Workshop 논문집, 1987.