

유·무성음 및 묵음 식별에 관한 연구

김 명 환 유 영 근 김 순 협
 광운 대학 대학원 전자공학과 OPC 중앙연구소 광운 대학 전자계산기 공학과

A Study on the Voiced, Unvoiced and Silence Classification

김 명 환 유 영 근 김 순 협
 Kwang Woon University OPC R&I Kwang Woon University

(Abstract)

This paper reports on a Voiced-Unvoiced-Silence Classification of speech for Korean Speech Recognition.

In this paper, it is describe a method which uses a Pattern Recognition Technique for classifying a given speech segment into the three classes.

Best result is obtained with the combination using ZCR, Fl, Ep and classification error rate is less than 1%.

전체적인 식별시스템은 그림 1 과 같다.

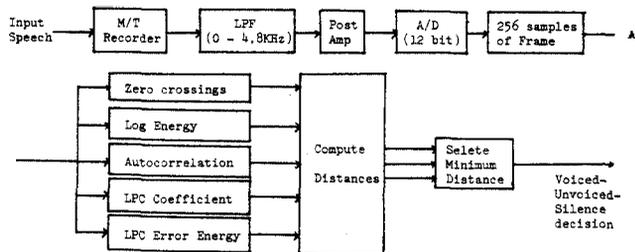


그림 1. 한국어 음성식별에 관한 시스템 블록도.

Fig. 1. A Block diagram of the Spoken Korean Speech Classification System.

1. 서론

주어진 음성마절을 유성음, 무성음, 묵음으로 식별하기 위하여 일반적으로 피치분석(Pitch analysis)으로 행하여져 왔으나 몇가지 큰 단점이 있다.

첫째, 유성음은 주기적인신호라는 단순한 특징에 가 본을 두었는데 성대진동의 갑작스러운 변화는 유성음 일지라도 비주기적인 음성마절을 만든다.

둘째, 음성합성분야에서는 피치분석이 가능하나,

음성분할(Speech Segmentation), 음성인식(Speech Recognition)과 같은 응용분야에서는 복잡성 뿐만아니라 유성음과 무성음 경계부분 에서는 매우 나쁜 수행을 초래한다.

따라서, 본 논문은 3가지음성신호를 식별하기 위하여 패턴인식방법을 사용하였다.

한편 본 연구에 사용된 분석하락이라는 음성신호의 영고 차율, 대수 에너지, 자동상관계수, 선형예측분석에서 얻은 첫번째 예측계수, 예측오차의 에너지이다.

그리고 유사도를 측정하기 위한 거리측정법은 Mahalanobis 거리측정법을 사용하였고 3가지 음성신호의 판정은 계산된 거리가 최소일때를 선택하여 결정하였으며 분석을 위한 음성데이터는 광운 대학 디지털 신호 처리연구실의 시스템에 의해 처리되었으며

2. 음성신호의 분석이론

가) 프레임 길이와 프레임 간격

프레임: 음성신호를 한번 처리하는 기본 단위

프레임길이: 프레임에 포함되어있는 샘플 수.

본 연구에는 프레임길이는 256 샘플

나) 영고 차율

$$Z_n = \sum_{m=-\infty}^{\infty} |\text{Sgn}[x(m)] - \text{Sgn}[x(m-1)]| \cdot W(n-m) \dots (1)$$

유성음은 스펙트럼의 하향때문에 에너지는 낮은

주파수에 집중되어 있다.

반면, 무성음은 에너지의 대부분이 높은주파수에 집중되어 있으며 높은주파수는 높은영고 차율을 낮은주파수는 영고 차율이 낮다.

그림 (2)는 각음성에 대한 영고 차율을 나타낸그림이다.

다) 대수 에너지

대수 에너지는 샘플값 자승의 합에 대수를 취한것으로

$$LE_n = 10 * \log \sum_{m=-\infty}^{\infty} x^2(m) \cdot W(n-m) \dots (2)$$

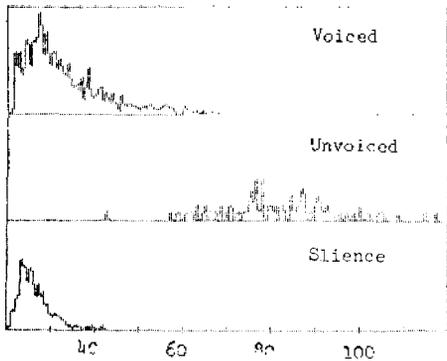


그림 2(a). 25.6ms 당 각음의 영고 차이의 통계분포

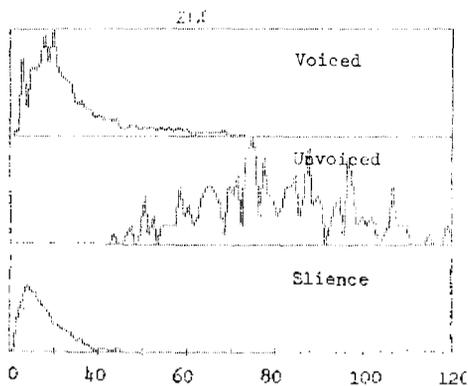


그림 2(b). 12.8ms 당 각음의 영고 차이의 통계분포

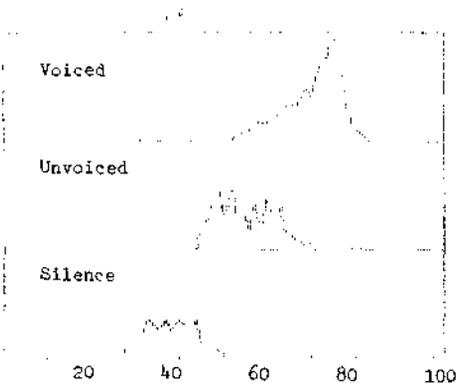


그림 3. 25.6ms 당 각음의 대수 에너지의 통계분포

다) 자동상관계수

$$C_1 = \frac{\sum_{m=1}^N x(m) \cdot x(m-1)}{\sqrt{\sum_{m=1}^N (x^2(m)) \sum_{m=0}^{N-1} (x^2(m))}} \dots (3)$$

자동상관계수는 인접한 샘플의 상관관계를 나타내며

정의에 의하여 1 - 1 사이에서 변화한다.

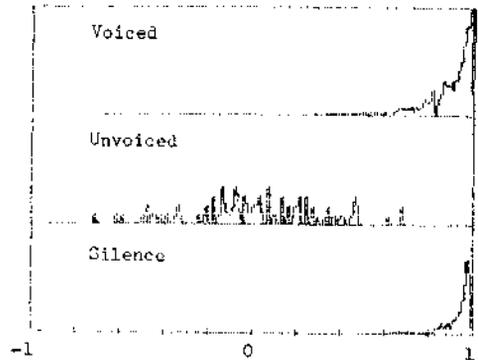


그림 4. 25.6ms 당 각음의 첫번째자동상관계수의 분포

라) 선형예측법

안음성행렬을 각기음성행렬들의 선형결합으로 근사할 수 있다는 데 있으며 선형시변 시스템을 특성화하는 매개변수를 추정하는데 확실하고 정밀한 방법이다.

$$s(n) = \sum_{m=1}^P a_m \cdot s(n-m) + G \sum_{l=0}^Q b_l \cdot u(n-l), b_0=1 \dots (4)$$

(1) 예측계수

평균 제곱오차 E_p 를 최소화 함으로써 얻을 수 있다.

$$E_p = \sum_{m=-\infty}^{\infty} e_n(m)^2 \dots (5)$$

$$= \sum_{m=-\infty}^{\infty} (s_n(m) - \sum_{k=1}^p a_k \cdot \hat{s}_n(m-k))^2 \dots (6)$$

여기서 E_p 가 최소가 되도록 하면

$$\frac{\partial E_p}{\partial a_i} = 0 \quad i=1,2,3,\dots,p \dots (7)$$

식(6),(7)을 풀면

$$\sum_{m=-\infty}^{\infty} s_n(m-1) \cdot s_n(m) = \sum_{k=1}^p a_k \sum_{m=-\infty}^{\infty} s_n(m-1) \cdot s_n(m-k) \dots (8)$$

(2) 예측오차의 에너지

$$E_p = \sum_{m=-\infty}^{\infty} s_n^2(m) - \sum_{k=1}^p a_k \sum_{m=-\infty}^{\infty} s_n(m) \cdot s_n(m-k) \dots (9)$$

3. 실험 알고리즘

확률밀도 함수 $P(X/W_1)$

$$P(X/W_1) = \frac{1}{(2\pi)^{n/2} |c_1|^{1/2}} \text{EXP} \left[-\frac{1}{2} (X - \hat{m}_1)^T c_1^{-1} (X - \hat{m}_1) \right] \dots (10)$$

식(10)에서 정규밀도함수의 지수 항을 단조증가함수인 자연대수로 취하고 최적결정함수 $\hat{d}_1(X) = P(X/W_1)P(W_1)$ 에 대입하면 다음 식과 같다.

$$\hat{d}_1(X) = \ln P(W_1) - \frac{1}{2} \ln |C_1| - \frac{1}{2} \ln 2\pi - \frac{1}{2} \left\{ (X - \hat{m}_1)^t C_1^{-1} (X - \hat{m}_1) \right\} \dots \dots \dots (11)$$

또 위식의 처음 3항들은 측정치 벡터와 무관하므로 다음 식과 같이 간단하게 쓸 수 있다.

$$\hat{d}_1(X) = (X - \hat{m}_1)^t C_1^{-1} (X - \hat{m}_1) \dots \dots \dots (12)$$

위식은 측정치 벡터 X 와 평균 벡터 \hat{m}_1 사이의 거리인데 통계적인 특징이 분명히 그려지기 때문에 유사도를 측정하는데 유용한 식이며 이를 Mahalanobis 거리 측정법이라고 한다.

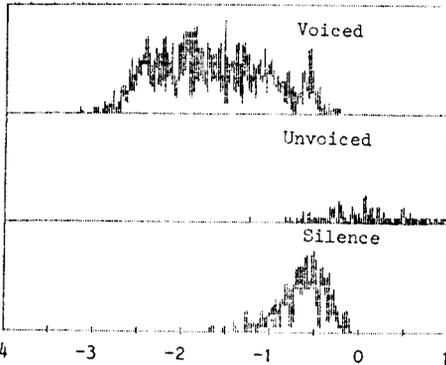


그림 5. 첫번째 예측 계수에 관한 통계분포

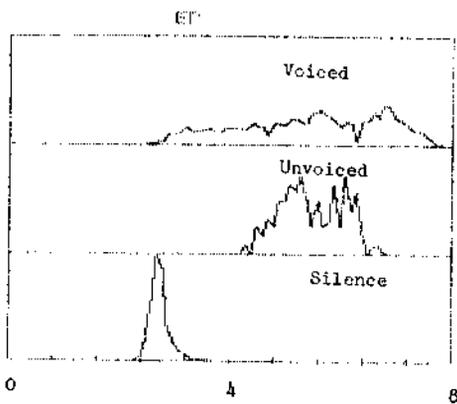


그림 6. 예측오차의 에너지에 관한 통계분포

4. 실험 및 결과

성인남성 3명이 발성한 한국어 숫자음 (/영/-/구/)를 7번 반복한 음성중 임의로 선정된 154개숫자음을 대상으로 실험하였으 며 전체적인 식별시스템은 그림 7과 같다.

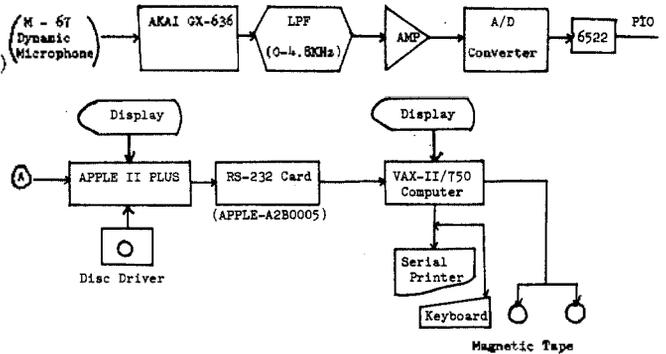


그림 7. 한국어 음성분석 및 식별에 관한 하드웨어 시스템.

Fig. 7. A Hardware System for Korean Speech analysis and Classification.

(1) 포준집합 설정

본 연구는 학자독립방법으로 음성을 식별하였기 때문에 포준집합설정에는 2명의 남성화자가 발성한 숫자음 파형을 순수식별하여 3가지음성신호부류로 구분하고 각부류에 대한 평균치,포준편차,공분산 행렬을 구하였다. 다음 그림은 포준집합에 사용한 음성파형이다.

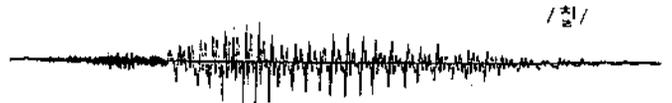


그림 8. 포준집합에 사용한 숫자음 /칠/의 음성파형

표 1. 포준집합에서의 평균치,포준편차,공분산행렬

Table 1. Means, Standard Deviations, and Covariance Matrices in Reference set.

	Zero Crossing	Log Energy	First autocorrelation	First LPC	LPC Log Error
1) Voiced (/1-3)					
Mean	25.7442	2.6285	0.897	-1.653	5.14539
Standard Deviation	17.9159	0.4014	0.12053	0.5084	1.26081
Covariance Matrix	1	0.01802	-0.9219	0.5210	0.01583
	0.01802	1	-0.3305	0.9219	0.87403
	-0.9219	-0.3305	1	-0.6270	-0.4409
	0.5210	0.9219	-0.6270	1	0.3842
	0.01583	0.87403	-0.4409	0.3842	1
2) Unvoiced (/4-7)					
Mean	12.1379	1.0714	0.05174	-0.0246	1.48372
Standard Deviation	10.1141	0.11732	0.36664	0.41399	0.56652
Covariance Matrix	1	0.5681	-0.8301	0.5200	0.6103
	0.5681	1	-0.3109	-0.091	0.9601
	-0.8301	-0.3109	1	-0.7353	-0.3731
	0.5200	-0.091	-0.7353	1	-0.0561
	0.6103	0.9601	-0.3731	-0.0561	1
3) Silence (/8-9)					
Mean	1.62905	0.3268	0.9801	-0.5819	2.7581
Standard Deviation	0.60751	1.39429	0.6215	0.25317	0.1916
Covariance Matrix	1	-0.210	-0.560	0.03120	0.6511
	-0.210	1	0.6403	-0.6016	0.3603
	-0.560	0.6403	1	-0.4056	-0.0784
	0.03120	-0.6016	-0.4056	1	-0.40902
	0.6511	0.3603	-0.0784	-0.40902	1

(2) 결과 및 고찰

(가) 1개의 파라메타로 식별하였을 때

Measured parameter	V = 0	V = 3	V = 5	Total	Error Rate
CRP	13	490	66	570	14.35
LE	143	14	71	228	6%
P1	4	677	51	732	18.55
P1	65	471	64	540	16.15
Sp	455	28	0	543	14.15

표 2. 각각의 파라메타로 식별하였을 때의 전체식별오차
Table 2. Total Classification error using each parameter.

(나) 2개의 파라메타로 식별하였을 때

Measured parameter	V = 0	V = 3	V = 5	Total	Error Rate
CRP, LE	54	46	5	105	2.85
CRP, P1	20	555	1	576	15.05
CRP, Sp	15	272	16	293	8.65
LE, P1	31	13	2	46	1%
LE, Sp	296	17	6	319	3.05
P1, Sp	60	42	6	108	3.85
CRP, P1	32	544	16	612	16.25
CRP, Sp	2	245	5	252	3.25
P1, Sp	560	47	3	610	16.15

표 3 2개의 파라메타를 조합하여 식별 하였을 때
Table 3 Total Classification error using two parameter.

(다) 3개의 파라메타로 식별하였을 때

Measured parameter	V = 0	V = 3	V = 5	Total	Error Rate
CRP, LE, P1	37	50	5	86	2.55
CRP, LE, Sp	61	45	6	112	3%
CRP, P1, Sp	2	54	2	58	1.85
LE, P1, Sp	33	17	10	60	1.8%
CRP, LE, Sp	5	71	2	78	2.15
CRP, P1, Sp	1	31	7	39	0.95%
LE, P1, P1	28	34	4	66	1.6%
LE, P1, Sp	1	47	2	50	1.5%
LE, P1, Sp	130	32	7	169	15.55
LE, P1, Sp	7	53	3	63	1.8%

표 4. 3개의 파라메타를 조합하여 식별하였을 때 전체오차
Table 4. Total Classification error using three parameters.

(라) 4.5개의 파라메타로 식별하였을 때

Measured parameter	V = 0	V = 3	V = 5	Total	Error Rate
CRP, LE, P1, Sp	4	14	7	25	2.6%
CRP, LE, P1, Sp	3	31	2	36	2.1%
CRP, P1, Sp	5	56	6	67	1.55
CRP, LE, P1, Sp	1	34	4	39	2.1%
CRP, LE, P1, Sp	4	17	3	24	2.4%
CRP, LE, P1, Sp	3	24	1	28	2.6%

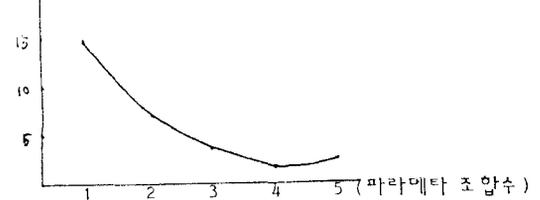
표 5. 4.5개의 파라메타로 식별하였을 때 전체오차
Table 5. Total Classification error using four or five parameters.

전반적으로 많은 파라메타를 조합하여 식별하였을 때 3% 보다도 적은 식별오차를 얻었다. 오차의 대부분은 끝점부분에서 유성음 구간을 독음으로 식별하였고, 유성음과 무성음의 경계부분에서도 오차가 발생하였다.

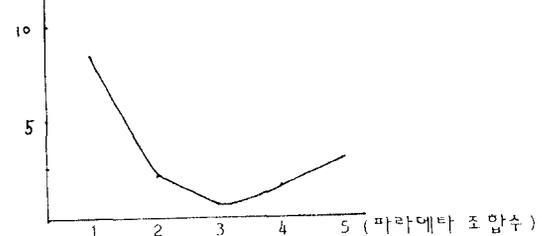
각 파라메타조합수와 평균오차율 및 최소오차율을 살펴 보면

파라메타 조합수	평균오차율	최소오차율
1	14.6 %	8 %
2	7.18 %	2 %
3	3.3 %	0.95 %
4	2.2 %	1.5 %
5	2.6 %	2.6 %

표 6. 파라메타 조합수와 평균오차율 및 최소오차율 (평균오차율)



(최소오차율)



5. 결론

본 논문은 한국어 음성인식을 위하여 음성구간을 유성음, 무성음, 묵음으로 식별하기 위하여 패턴인식방법 (Pattern Recognition Technique)을 사용하였다. 첫째, 특징집합은 음성분석에서 일반적으로 사용되는 음성신호의 영고 차를, 대수 에너지, 첫번째 예측 계수, 첫번째 작동 상관 계수, 예측 오차의 에너지를 사용하였으며, 둘째, 포론 집합은 3가지 음성신호 부류에 대하여 손수 식별한 음성 데이터들로 부터 얻은 평균 벡터와 분산을 통계적으로 처리 하였으며 셋째, 유사도 거리측정법은 각 파라메타들의 상관관계 및 통계적인 특징이 분명히 그려지는 Mahalanobis 거리측정법을 사용하였으며 마지막으로, 3가지 음성신호 부류 판정은 각 음성에 대한 연습 집합과 포론 집합과의 계산된 거리가 최소일 때를 선택하여 결정 하였다. 그리고 각 파라메타들을 여러가지로 조합하여 식별한 결과 영고 차를, 선형 예측 분석에서 얻은 첫번째 예측 계수, 예측 오차의 에너지를 파라메타로 수행하였을 때 다른 조합 방식보다도 가장 좋은 결과를 얻었으며 그때의 식별 오차는 1%보다도 적었다.

참고 문헌

- (1) L.Rabiner, M.Sambur, and C.Schmidt, "Applications of a nonlinear smoothing algorithm to speech processing," IEEE Trans. Acoust., Speech, Signal processing, vol. ASSP-23, PP. 552-557, Dec. 1975.
- (2) B.S. Atal and L.R. Rabiner, "A pattern recognition approach to voiced-unvoiced-silence classification with applications to speech recognition," IEEE Trans. Acoust., Speech, Signal processing, vol. ASSP-24, pp. 201-221, June 1976.
- (3) J.J. Dubnowski, R.W. Schafer, and L.R. Rabiner, "Real-time digital pitch detector," IEEE Trans. Acoust., Speech, Signal processing, vol. ASSP-24, pp 2-8, Feb. 1976
- (4) L.J. Siegel and K. Steiglitz, "A pattern classification algorithm for the voiced/unvoiced decision," in 1976 ICASSP Proc., 1976, pp. 326-329.
- (5) L.F. Lamel, L.R. Rabiner, A.E. Rosenberg, and J.G. Wilpon, "An improved ending point detector for isolated word recognition," IEEE Trans. Acoust., Speech, Signal processing, vol. ASSP-29, pp. 777-785, August, 1981.
- (6) H.F. Silverman and N.R. Dixon, "A comparison of several speech-spectra classification methods," IEEE Trans. Acoust., Speech, Signal processing, vol. ASSP-24, August, 1976.
- (7) L.J. Siegel and A.C. Bessey, "Voiced/Unvoiced/Mixed excitation classification of speech," IEEE Trans. Acoust., Speech, Signal processing, vol. ASSP-30, No. 3 Jun. 1982.