

집단화를 이용한 음성외 표준 패턴신경에 관한 연구

* 김계국 ** 유영근 *** 김순환 * 윤재광
* 송준대 ** OPC *** 광운대

A study on creating Reference Pattern of speech by using the cluster

* K. K. KIM ** H. K. RYU *** S.H. KIM * J. K. YMOON
* Soong Jun Univ. ** OPC *** Kwang Woon Univ.

I. 서 론

요 약

본 특정 화자의 음성인식을 위해 150 숫자음에 대하여 10 개의 표준 패턴을 설정하는 데 목적을 두고 기술했다.
남성 화자 3인이 관숫자음(0-9)을 5 번씩 반복 발음한 150 음을 집단화하여 숫자음의 표준 패턴을 설정하였다.
특징 파라미터는 포르만트 주파수를 이용하였고 유클리드 거리 측정법을 유사도 비교에 사용하였다. 실험결과 85.3%의 인식률을 얻었다.

과학의 성장발전으로 기술혁명의 산업사회에서 통신매체의 역할은 지대해것으로 생각된다. 그러므로 한국어가 부강하기 위해서는 컴퓨터 시스템에 의한 정보의 전달은 본명이 연구되어야 할 것이다. 그러한 견지에서 21세기 화문이라고 일컬을 만큼 우리의 음성 인식에 대한 관심도는 참으로 높은 것으로 압고 있다.
본 연구의 목적은 한국어 음성을 인식하게 위하여 어떠한 집단화 알고리즘을 선택하고 집단화에 의하여 표준 패턴을 설정하고 보다 높은 인식율을 얻고자 하는 데 있다. 그림 1에서 음성인식의 기본 과정을 살펴보면, 우선 전처리된 음성데이터를 가지고 특징 파라미터를 추출하고 표준 패턴을 설정하여 유사성을 비교하고 인식하게 된다.

A B S T R A C T

This paper is of the pupose in the creation of ten reference patterns for 150 digits to recognize a isolated digits(0-9) which three speaker-independents pronounce.
The digit reference patterns are created from a statistical clustering analysis of speech date base consisted of 150 replications of each korean isolated digits.
They repeated five times by each of 3 male talkers.
Euclidean distance measure method is used for the comparison of pattern similarity.
The recognition accuracies for three male talkers is about 85.3 percent.

본 연구에서는 포르만트 주파수를 특징 파라미터로 이용하였으며 이것은 ARMA (Autoregressive Moving Average) 모델의 PSD (Power Spectrum Density) 추정 알고리즘을 이용하여 추출한 제1, 제2, 제3 포르만트로, 기존의 음성데이터를 이용하여 표준 패턴설정을 시도하였다.

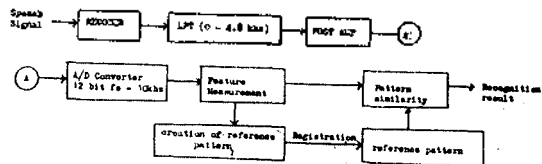


그림 1. 음성인식의 시스템 블록도
Fig. 1. Block diagram of speech recognition system

II. 본 론

1. 표준 패턴 설정 이론

본 연구에서는 한국어 단독 숫자들(0-9)을 3개의 집단화자가 각 숫자들에 대해서 5번씩 반복한 150음성을 집단화에 이용하였다. 표준 패턴은 각 집단을 대표할수 있는 단어로써 75개의 단어들중 어떠한 단어도 대체시도 최대거리가 최소가 될수 있도록 최적 기준을 만족할수 있는 단어가 결정되어야 한다. 그림2는 집단화의 일반적인 예를 보이고 있는데 이 그림에서 각음성 데이터의 특징파라미터집합을 X (token 이라칭함)로 나타내고 화자들간의 성대의 변화 특성 때문에 표준들이 서로 다른 영역에 분포되어있다.

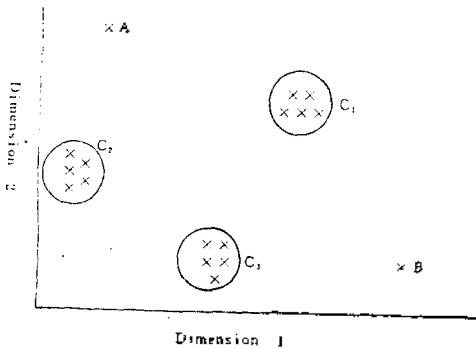


그림2. 집단화의 일반예

Fig. 2 Example Showing clustering of tokens

그림2에서는 3개의 집단이 이루어져있기 때문에 적어도 3개의 표준 패턴이 설정되어야 하며 3개의 집단과 일치할수 있는 두토큰 A와 B를 분리물이라고 하고 이 분리물은 3개의 집단중 어느것과도 한명거리를 이루지 못하고 있다. 음성데이터를 나타내는 J개의 토큰을 집합으로 표현하면 식(1)과 같다.

$$\Omega = \{x_1, x_2, \dots, x_J\} \quad (1)$$

이들 토큰들중 임의의 x_i 와 x_j 토큰 간의 차는 거리 d_{ij} 로 측정하고 식(2)와 같이 표현된다.

$$d_{ij} = d(x_i, x_j) = \frac{1}{K} \sum_{k=1}^K d(k, \omega_k) \quad (2)$$

여기서 K는 모든 x_i 의 프래임수를 나타내며 $d(k, \omega_k)$ 는 Itakura가 제안한것으로 x_i 의 k 번째 프래임과 x_j 의 $\omega(k)$ 번째 프래임들간의 거리를 의미한다.

또 $\omega(k)$ 는 x_i 와 x_j 의 DTW 정향으로 얻어진 워핑함수라고 한다.

식(1)의 Ω 는 M개의 집단으로 이루어진 어떠한 반복방음한 특징같은 음의 단어를 나타내고있으며 집단의 수 M은 실험과정에서 적절히 고러되어야 하고 일반적으로 식(3)과 같이 나타내고 있다.

$$\Omega = \bigcup_{i=1}^M \omega_i \quad (3)$$

여기서 ω_i 는 i 번째 집단을 의미하며 ω_i 에 속해있는 토큰의 개수를 m_i 로 결정하고 ω_i 를 대표할수 있는 중심점 (cluster center point) 을 $x_p^{(i)}$ 로 정의하고 $x_p^{(i)}$ 는 반드시 i 번째 집단의 토큰들중 하나가 결정되어야 한다.

$$x_p^{(i)} \in \omega_i \quad (4)$$

표준 패턴설정에 이용되는 일반적인 집단화 알고리즘에는 Chainmap, SNM (Shared nearest neighbour), K-means iteration, ISODATA (Iterative Self Organizing DATA Analysis Technique A) 등이 있다.

본 연구의 사용된 알고리즘은 다음과 같다.

1. 임의의 최단 표준 패턴 알고리즘을 사용하여 최단 표준 패턴을 얻는다.

2. 임의의 15개의 표준 패턴을 임의의 길이로 임의의 길이로 임의의 표준 패턴을 생성한다.

3. 임의의 길이로 임의의 표준 패턴을 생성한다.

4. 다시 임의의 길이로 임의의 표준 패턴을 생성한다.

5. 임의의 길이로 임의의 표준 패턴을 생성한다.

6. 임의의 길이로 임의의 표준 패턴을 생성한다.

7. 임의의 길이로 임의의 표준 패턴을 생성한다.

8. 임의의 길이로 임의의 표준 패턴을 생성한다.

9. 임의의 길이로 임의의 표준 패턴을 생성한다.

10. 임의의 길이로 임의의 표준 패턴을 생성한다.

11. 임의의 길이로 임의의 표준 패턴을 생성한다.

12. 임의의 길이로 임의의 표준 패턴을 생성한다.

13. 임의의 길이로 임의의 표준 패턴을 생성한다.

14. 임의의 길이로 임의의 표준 패턴을 생성한다.

15. 임의의 길이로 임의의 표준 패턴을 생성한다.

16. 임의의 길이로 임의의 표준 패턴을 생성한다.

17. 임의의 길이로 임의의 표준 패턴을 생성한다.

18. 임의의 길이로 임의의 표준 패턴을 생성한다.

19. 임의의 길이로 임의의 표준 패턴을 생성한다.

20. 임의의 길이로 임의의 표준 패턴을 생성한다.

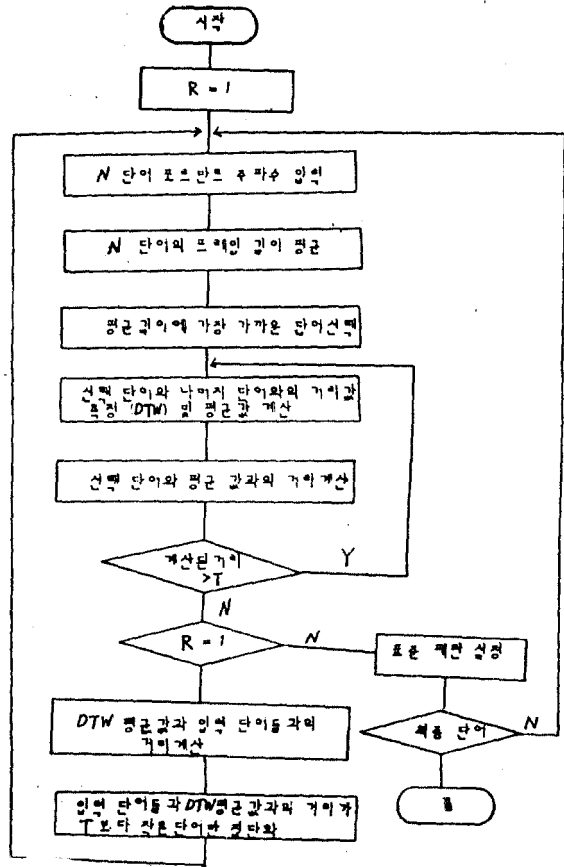


그림 3. 임의의 최단 표준 패턴 알고리즘에 사용된 알고리즘의 흐름도
 Fig. 3. Flow chart of algorithm used to reference pattern are then into clusters

본 연구의 사용된 알고리즘은 다음과 같다.

1. 임의의 최단 표준 패턴 알고리즘을 사용하여 최단 표준 패턴을 얻는다.

2. 임의의 15개의 표준 패턴을 임의의 길이로 임의의 길이로 임의의 표준 패턴을 생성한다.

3. 임의의 길이로 임의의 표준 패턴을 생성한다.

4. 다시 임의의 길이로 임의의 표준 패턴을 생성한다.

5. 임의의 길이로 임의의 표준 패턴을 생성한다.

6. 임의의 길이로 임의의 표준 패턴을 생성한다.

7. 임의의 길이로 임의의 표준 패턴을 생성한다.

8. 임의의 길이로 임의의 표준 패턴을 생성한다.

9. 임의의 길이로 임의의 표준 패턴을 생성한다.

10. 임의의 길이로 임의의 표준 패턴을 생성한다.

11. 임의의 길이로 임의의 표준 패턴을 생성한다.

12. 임의의 길이로 임의의 표준 패턴을 생성한다.

13. 임의의 길이로 임의의 표준 패턴을 생성한다.

14. 임의의 길이로 임의의 표준 패턴을 생성한다.

15. 임의의 길이로 임의의 표준 패턴을 생성한다.

16. 임의의 길이로 임의의 표준 패턴을 생성한다.

17. 임의의 길이로 임의의 표준 패턴을 생성한다.

18. 임의의 길이로 임의의 표준 패턴을 생성한다.

19. 임의의 길이로 임의의 표준 패턴을 생성한다.

20. 임의의 길이로 임의의 표준 패턴을 생성한다.

2. DTW 알고리즘 원리

인간의 음성은 위치, 강도, 에너지 측정, 그리고
신경 예측 계수, 포르만트 주파수 분산등의 특징
추출로 특징 벡터를 표현할수 있다.

보존 패턴과 시험 패턴 의 특징 벡터는

$$R(m) = [R_1, R_2, \dots, R_j, \dots, R_N] \quad (8)$$

$$T(n) = [T_1, T_2, \dots, T_k, \dots, T_N] \quad (9)$$

본 연구에서는 신경 예측법에 의한 MPA영역도 추출
안 제1, 제2, 제3 포르만트 주파수를 보존 패턴과 시험
패턴의 특징 벡터로 이용하였다.

1) 거리 측정법

시험패턴과 보존 패턴간의 두 특징 벡터의 차는
두 벡터간의 거리 측정으로 구할수 있다.

$$d(i, j) = d(i(R), j(R)) = \|T_i - R_j\| \quad (10)$$

최적 왜임 경로를 구하기 위하여, 이용되는 전체거리
는 다음과 같다.

$$D(i(R), j(R)) = \sum_{k=1}^K d(i(R), j(R)) \cdot W(R) \quad (11)$$

$W(R)$ 는 k 번째 소구간 경로의 왜임 함수를 뜻하며
 $N(W)$ 는 왜임 함수 W 의 함수인 정규화 요소이며 최
적 경로는 전체거리를 최소로 하는 경우로서 식(12)
와 같이 구할수 있다.

$$D_T = \min_{(i(R), j(R))} [D(i(R), j(R))] \quad (12)$$

본 연구에서는 f_1, f_2, f_3 의 포르만트들 특징 벡터로
이용하고 식(15)의 Euclidean 거리 측정법을 도입
하였다.

$$d = \left[\sum_{p=1}^3 \left(\frac{f_{ip}^{(T)} - f_{jp}^{(R)}}{\bar{f}_K} \right)^2 \right]^{\frac{1}{2}} \quad (13)$$

$$\text{단, } \bar{f}_K = \begin{cases} 500 & \text{for } K=1 \\ 1500 & \text{for } K=2 \\ 2500 & \text{for } K=3 \end{cases}$$

$f_{ip}^{(T)}$ 는 시험 패턴의 특징 벡터이며 $f_{jp}^{(R)}$ 는 보존
패턴의 특징 벡터를 말하며, P 는 J 개의 포르만트를
말한다.

2) 실험 결과

본 연구에서는 단어도 두개 이상의 모음을 하나의
집단으로 처리하기로 하고 실험 하였다.

특 같은 반복음 15 단어의 대하여 단 한개의 보존 패
턴만 실험 해주고 인식에 사용 하였으며
다른 모음들과 혼동 하지 못하이 어느 집단에도 포함
되지 못한 단일 모음은 본리물로 처리하였다.

표 7은 본 연구에 사용된 150 숫자음을 집단화한
결과를 제시하고있다. 일반적으로 본리물이 아닌
모든 모음들은 하나의 집단으로 묶을 수 있었으며
세일론 집단에는 포함 15개중 15모음이 모두 속해
있음을 알 수 있었다. 앞으로는 집단화 임계치에
대한 연구가 계속되어 좀 더 바람직한 집단화가 이루어
져야 함것으로 생각된다. 실험에 이용한 150단어
중 22개의 단어가 잘못 인식되어 85.3%의 인식 결과
를 얻었다.

표 7. 집단화 결과

숫자음	집단수	본리물수	오인수
일	1	1	14
일	1		1
이	1		15
삼	1	2	1
사	1		15
오	1		1
육	1		15
일	1		15
팔	1		15
구	1	1	14

표 2, 인식 결과

		Recognized word										
		영	입	이	삼	사	오	육	십	판	구	
Spoken word	영	10				2				1	2	
	입		15									
	이	1	14									
	삼				11					1	3	
	사					15						
	오			1			12	1				1
	육	1	1	1				12				
	십								13			
	판										15	
	구						2	1	1			11

III. 결 론

본 연구에 제안된 알고리즘을 가지고 보른 패턴
 심정에 이용한 결과는 85.3%의 인식을 얻었다.
 본 논문을 쓰면서 이를 했던 것은 한국어 음성
 인식을 위해서 보른 패턴 심정에 관한 분석 자료가
 없었다는 것이 큰 어려움이었다. 인식 결과 22개 음성
 이 잘못 인식 되었으며, 이것은 150 단어를 대하여
 10개의 보른 패턴을 심정비로 재확인하고 집단화
 한데 이유가 있다. 150단어에 대해서 20개 정도의
 보른 패턴을 허용하고, 집단을 추적하는 데 있어
 패턴을 심정하였으면 더욱 높은 인식효과를 얻을 수
 있을 것으로 생각된다.
 특히 앞으로는 집단화에 활용되지 못한 본 특성의
 처리 방안이 연구되어야 할 것이며 본 연구에서는 단
 독 숫자음의 보른 패턴 심정에 그쳤으나 앞으로 어휘가
 지 단어를 인식하기 위해서 한국어 특성에 맞는
 통계적인 패턴 심정 알고리즘이 개발되어야 할 것
 이다. 그러나 여기에 제시된 알고리즘도 잘 이용만
 된다면 화자수와 어휘에 관계없이 음성 인식에 기여
 할 수 있을 것으로 본다.

참 고 문 헌

1. Warren J.G.Kootz. "A banch and bound clustering Algorithm."
 IEEE Transaction computers 1975.9.
2. Helmuth spath, Ellis Horwood Limited.
 "Cluster Analysis Algorithms for Data reduction and classification of objects".
3. 김순환, "한국어 음성의 분석과 자동인식에
 관한 연구" 연세대학교원 박사학위논문
 1982. 12.
4. Hiroaki Sakoe. " Dynamic Programming
 Algorithm optimisation for spoken word
 recognition".
 IEEE Trans. Vol. Assp-26 No.1, 1978.2.
5. L.R.Rabiner. " Speaker-Independent
 Recognition of Isolated Words using
 Clustering Techniques".
 IEEE Vol.ASSP-27,NO.4, 1979.8
6. R.s.Jarvis. "Clustering using a simila-
 rity Measure Based on Shared Nearest
 Neighbors". IEEE Vol c-22.No.11.
 1973.11.
7. L.R.Rabiner. " On creating Reference
 templates for speaker independent
 recognition of isolated word".
 IEEE, Vol.ASSP-26, No.1, 1978.2.