

포르만드 주파수를 이용한 한국어 음성의
자동 인식에 관한 연구

83306

○ 김 순 협 *
광운대학 전자계산기공학과 교수
박 규 태 **
연세대학교 전자계산기공학과 교수

A Study on the Automatic Recognition of
Korean speech by Formant Frequency.

(Abstract)

In Speech signal processing, ARMA spectral estimation method is used. It has been demonstrated that the ARMA model provides better spectral estimation than the more specialized AR model and MA model.

Dynamic program is used to achieve time alignment. Speech sound similarity is defined to be proportional to the distance separating two sound in a Vector space defined by ARMA model.

As a result, the recognition rate of 97.3% for three speaker is obtained.

<본문>

(1) 본 연구의 목적은 자동회귀 이동 평균 모형 (Autoregressive Moving Average model ; 이하 'ARMA' 모형이라 함) 방식을 이용하여 한국어의 음성을 분석하고, 단독 숫자 음을 대상으로 한 인식 방법과 인식 시스템을 제시하는데 있다.

(2) 시간 영역에서의 음성 신호 분석 방법은 영고차율과 대수 에너지를 이용하였다. 영고차율은 분석 구간 폴라입체의 좌푯값이 영점(zero) 초과 고차하는 회수를 말하며 일반식은 $Z = \sum_{m=1}^N [1 - sgn(x_{m+1}) \cdot sgn(x_m)] / N$ (1)

이고. 여기서 $sgn[x(m)] = \begin{cases} 1, & x(m) > 0 \\ -1, & x(m) < 0 \end{cases}$

으로 표시한다.

대수 에너지는 이산계(discrete system)에서 $D = \sum_{m=0}^{N-1} x^2(m)$ 으로 표현되며 무정음으로 부터 유정음을 구분하는데 유용하다.

(3) DP 정합법의 기본 원리는 시험 패턴 (Test pattern)과 표준 패턴 (Reference pattern)을 비교하여 두 패턴 사이의 유사성을 결정하는 것이다. 유사도는 패턴 사이의 시간적 정규화와 거리 계산 등으로 이루어진다. 한편, 근본적인 음성 패턴의 시간적 구조에는 연속성, 단조증가, 그리고 음성학 과정이 변화 속도에 따른 조건이 따른다. 이러한 조건들은 warping 함수에 대한

① 끝점 제약 또는 초기 조건 ② 부분적 경로의 연속성 제약 ③ 전체 경로의 제약으로 만족될 수 있고, DTW 알고리즘의 최적 warping 경로를 위한 거리 측정 (distance measurement)은 다음과 같은 식의 일반적 형태를 갖는다.

$$D(i(k), j(k)) = \frac{\sum_{k=1}^L d[i(k), j(k)] \cdot w(k)}{N(w)} \quad \dots (3)$$

(4) 음성 신호의 PSD 추정 알고리즘에서 음성 생성 과정에 대한 정확한 모델을 ARMA 모형으로 조사하였다.

$$y_k = -\sum_{i=1}^p a_i y_{k-i} + \sum_{j=0}^q b_j u_{k-j} \quad \dots \dots \dots (4)$$

$$\text{이고 } \text{이} \rightarrow y(m) = \sum_{i=0}^p b_i u(m-i) - \sum_{j=1}^q a_j y(m-j) \dots (5)$$

로 변형 가능하다. 여기서 출력 임펄스
의 $y(1), y(2), \dots, y(m)$ 은 시계열의 자동
상관함수는 $r_y(m) = E\{y(k) \cdot y^*(m-k)\} \dots (6)$
로 되며 이식을 푸리에 변환하면 PSD를
추정할 수 있다.

$$\text{즉, } P_y(w) = \sum_{m=-\infty}^{\infty} r_y(m) e^{-jwm} \dots \dots (7) \text{이다.}$$

이것을 유리화 스펙트럼 모형으로 표시하면

$$P_y(w) = \left| \frac{b_0 + b_1 e^{jw} + \dots + b_p e^{jpw}}{1 + a_1 e^{-jw} + \dots + a_q e^{-jqw}} \right|^2 \dots \dots (8)$$

이 된다. 여기서 $b_0 = 1$ 이며 이 식으로부터
 a_k 와 c_k 의 계수를 결정하여 음성 신호의
PSD를 추정했다.

(5) 본 연구에서 사용된 절차 방법은
음성 신호를 $4.8 [kHz]$ 의 저역 통과 여과
기(Low Pass Filter)에 통과 시킨 다음
 12 bit A/D 변환기 의 입력 경계 전압이
되도록 증폭기를 구성하여 컴퓨터에 연결하
여 인식시켰다.

이상의 방법으로 성인 화자 3인에 의해
발성된 숫자들에 대해 97.3%의 인식률
을 얻었다.

앞으로는 화자 수와 어휘 수에 상관없는
음성 인식을 위해서 음성의 음소 분석과 음소
인식 설계가 반드시 선행되어야 할 것이다.