

음성 신호를 사용한 감정인식의 특징 파라미터 비교

Comparison of feature parameters for emotion recognition using speech signal

| | |
|----------|---|
| 저자 | 김원구 |
| 저널명 | 電子工學會論文誌. Journal of the Institute of Electronics Engineers of Korea. SP, 신호처리 |
| 발행기관 | 대한전자공학회 |
| NDSL URL | http://www.ndsl.kr/ndsl/search/detail/article/articleSearchResultDetail.do?cn=JAKO200311921927764 |
| IP/ID | 1.212.198.211 |
| 이용시간 | 2018/04/23 17:46:51 |

저작권 안내

- ① NDSL에서 제공하는 모든 저작물의 저작권은 원저작자에게 있으며, KISTI는 복제/배포/전송권을 확보하고 있습니다.
- ② NDSL에서 제공하는 콘텐츠를 상업적 및 기타 영리목적으로 복제/배포/전송할 경우 사전에 KISTI의 허락을 받아야 합니다.
- ③ NDSL에서 제공하는 콘텐츠를 보도, 비평, 교육, 연구 등을 위하여 정당한 범위 안에서 공정한 관행에 합치되게 인용할 수 있습니다.
- ④ NDSL에서 제공하는 콘텐츠를 무단 복제, 전송, 배포 기타 저작권법에 위반되는 방법으로 이용할 경우 저작권법 제136조에 따라 5년 이하의 징역 또는 5천만 원 이하의 벌금에 처해질 수 있습니다.

論文2003-40SP-5-8

음성 신호를 사용한 감정인식의 특징 파라미터 비교 (Comparison of feature parameters for emotion recognition using speech signal)

金 元 九 *

(Weon-Goo Kim)

요 약

본 논문에서는 음성신호를 사용하여 인간의 감정을 인식하기 위한 특징 파라미터 비교에 관하여 연구하였다. 이를 위하여 여러 가지 감정 상태에 따라 분류된 한국어 음성 데이터 베이스를 이용하여 얻어진 음성 신호의 피치와 에너지의 평균, 표준편차와 최대 값 등 통계적인 정보 나타내는 파라미터와 음소의 특성을 나타내는 MFCC 파라미터가 사용되었다. 파라미터들의 성능을 평가하기 위하여 문장 및 화자 독립 감정 인식 시스템을 구현하여 인식 실험을 수행하였다. 성능 평가를 위한 실험에서는 운율적 특징으로 피치와 에너지와 각각의 미분 값을 사용하였고, 음소의 특성을 나타내는 특징으로 MFCC와 그 미분 값을 사용하였다. 벡터 양자화 방법을 사용한 화자 및 문장 독립 인식 시스템을 사용한 실험 결과에서 MFCC와 델타 MFCC를 사용한 경우가 피치와 에너지를 사용한 방법보다 우수한 성능을 나타내었다.

Abstract

In this paper, comparison of feature parameters for emotion recognition using speech signal is studied. For this purpose, a corpus of emotional speech data recorded and classified according to the emotion using the subjective evaluation were used to make statical feature vectors such as average, standard deviation and maximum value of pitch and energy and phonetic feature such as MFCC parameters. In order to evaluate the performance of feature parameters speaker and context independent emotion recognition system was constructed to make experiment. In the experiments, pitch, energy parameters and their derivatives were used as a prosodic information and MFCC parameters and its derivative were used as phonetic information. Experimental results using vector quantization based emotion recognition system showed that recognition system using MFCC parameter and its derivative showed better performance than that using the pitch and energy parameters.

Keyword : 감정 인식, 음성 파라미터, MFCC

* 正會員, 群山大學校 電子情報工學部

(School of Electronic and Information Eng. Kunsan National University)

※ 본 연구는 정보통신부 정보통신연구진흥원에서 지원하고 있는 정보통신기초연구지원사업의 연구결과입니다.

接受日字:2002年5月21日, 수정완료일:2003年9月17日

I. 서 론

인간의 감정 상태를 인지하는 시스템의 개발은 보다 고차원적인 휴먼-컴퓨터 인터페이스 기능을 갖는 제품을 가능하게 한다. 인간의 감정 정보를 얻는 방법은 얼굴 표정, 음성, 몸 동작, 심장 박동수, 체온, 혈압 등을

통하여 다양하게 얻을 수 있고, 응용 분야에 따라 달라질 수밖에 없다. 이 중에서 음성을 이용한 시스템의 경우 센서가 신체 부위에 직접 단지 않거나, 전화와 같이 반드시 음성을 이용하여야 하는 시스템에 사용할 때 이용될 수 있다.

이러한 감정 관련 연구는 일부이긴 하지만 이미 제품에 적용되어 상용화되는 단계에 이르고 있다. 이러한 제품으로 1999년 일본 소니사가 개발하여 시판한 오락용 애완견 로봇 AIBO가 있다. AIBO는 6가지 감정을 포함하는 감정 모델을 적용하여 주인과의 관계에 의해서 감정상태가 변화하고 그에 따라 반응하도록 만들어졌으며 현재는 보다 다양한 기능을 갖는 로봇으로 개량되어 판매되고 있다. IBM에서는 Blue Eyes 프로젝트를 통하여 차세대 감정 인식 제품을 개발하고 있으며 표정을 이용한 감정 인식시스템과 생체 신호를 이용한 감정 인식 마우스 등을 상용화하는 단계에 이르렀다. 또한 미국 MIT대 미디어랩의 Rosalind W. Picard는 감정적 착용 컴퓨터를 개발하여 감정과 관련된 연구 결과의 실용화 가능성을 밝게 하고 있다^[1].

인간 감정과 음성의 상관관계에 대한 연구는 화자의 감정 상태를 반영하는 음성 신호의 요소로서 시간에 따른 피치 값의 변화를 제시하고 다양한 감정 상태에 따라 피치 값이 변화하는 차이점을 보여 줌으로써, 현재까지의 감정인식 연구에 많은 역할을 하였다^[2]. 또한 감정이 포함된 음성을 합성하는 연구에서는 화자의 감정을 반영하는 요소로서 발음 속도, 피치 평균, 피치 변화 범위, 발음 세기, 음질, 피치의 변화, 발음법 등이 있으며, 기쁨, 화남, 슬픔, 두려움, 혐오감 등의 주요 감정들을 표현할 때 이들 요소들이 감정 합성에 미치는 차이점들을 정리하였다^[3]. 이후 감정 반영에 주요한 역할을 하는 음성 요소를 변화시켜 다양한 감정이 포함된 음성을 합성하는 연구에서 합성된 감정 음성에 대한 주관적 감정 평가 결과가 약 50% 정도의 인식률을 나타내는 연구 결과를 발표하였다^[4].

감정 인식 기술은 음성 중에 포함된 감정에 작용하는 파라미터에 관한 연구와 이러한 파라미터를 사용하여 감정을 구분하는 패턴 인식 기술로 구분된다. 감정 인식에 주로 사용되는 파라미터로는 피치와 에너지가 대표적이며 이러한 파라미터의 평균, 표준편차, 최대값, 최소값, 변화율, 지속시간 등을 사용하여 감정 인식을 수행한다. 또한 패턴 인식 기술로는 MLB(Maximum-Likelihood Bayes), KR(Kernel Regression), KNN(K-

Nearest Neighbor) 분류기 등 기본적인 패턴 인식 기법들이 사용되었고 음성과 표정을 동시에 이용한 감정 인식에 관한 연구도 수행되었다^[5-10], 한편 국내에서는 외국의 연구 동향과 다르게 아직 감정인식에 대한 연구가 초기 단계에 있다. 국내 일부 대학에서 감정 인식에 대한 연구가 진행되어 피치와 에너지를 사용한 감정 인식 실험이 이루어졌다^[11].

본 논문에서는 음성신호를 사용한 인간의 감정인식을 위한 특징 파라미터의 비교에 관하여 연구하였다. 이를 위하여 여러 가지 감정 상태에 따라 분류된 한국어 음성 데이터 베이스를 이용하여 감정 인식 및 음성 신호 처리에 널리 사용되고 있는 음성 신호의 피치와 에너지의 평균, 표준편차와 최대 값 등 통계적인 정보 나타내는 파라미터와 음소의 특성을 나타내는 MFCC 파라미터의 성능을 비교 평가하였다.

파라미터들의 성능을 평가하기 위하여 벡터 양자화 방법을 사용한 문장 및 화자독립 감정 인식 시스템을 구현하여 인식 실험을 수행하였다. 성능 평가를 위한 실험에서는 운율적 특징으로 피치, 에너지와 각각의 미분 값을 결합하여 사용하였고, 음소의 특성을 나타내는 특징으로 MFCC와 그 미분 값을 결합하여 사용하였다.

II. 감정 인식 알고리즘

1. 감정 인식을 위한 특징 파라미터

감정 인식에 사용되는 음성의 특징 파라미터로는 운율적 특징으로 피치와 에너지에 관한 파라미터가 주로 사용된다. 음성 특징 파라미터는 음성신호의 단구간에 대해 구한 피치와 에너지 값으로부터 피치 평균 (pitch mean), 피치 표준편차 (pitch standard deviation), 피치 최대 값 (pitch maximum), 에너지 평균 (energy mean), 에너지 표준편차 (energy standard deviation) 등의 통계적 정보를 감정 인식을 위한 특징으로 사용하였다^[12].

MFCC(Mel Frequency Cepstral Coefficient) 파라미터는 음소의 특성을 나타내는 특징으로 음성 인식에 널리 사용되고 있다. 이러한 파라미터는 같은 음소라도 포함된 감정에 따라 음소의 형태가 다르다는 점에서 감정인식에도 사용될 수 있다.

2. 패턴 인식 알고리즘

2.1 KNN(K-Nearest Neighbor Classifier)

KNN 분류기는 기준 패턴의 분포 함수를 사용하는 대신에 미리 구하여 놓은 각각의 기준 패턴과의 거리를 계산하여 가장 가까운 기준패턴의 클래스를 입력 패턴의 클래스로 결정하는 방법이다^[26]. 여기서 입력 패턴과 기준 패턴간의 거리는 특정한 거리 측정 방법을 사용하여 구하며 최소 거리는 계산된 거리 측정의 결과가 가장 작은 것을 의미한다. 기준 패턴 생성 방법은 적은 수의 패턴으로 클래스를 잘 표현할 수 있어야 한다. 일반적으로 기준 패턴 생성 방법으로는 k-means 알고리즘과 LBG 알고리즘이 많이 사용된다^[13, 14]. 거리 측정 방법은 가장 기본적인 유클리디안 거리측정(euclidean distance) 이외에도 음성 인식에서 사용되고 있는 많은 방법들이 사용될 수 있다^[12]. 유클리디안 거리 측정 방법은 다음과 같다.

$$d(a, b) = \sqrt{\sum_{i=1}^P (b_i - a_i)^2} \quad (1)$$

여기서 a 와 b 는 두 벡터이며 P 는 벡터의 차수이다. 유클리디안 거리측정은 가장 널리 사용되는 거리 측정 함수로서 상이성이 큰 특징을 강조하는 성질이 있다. 사전에 클래스마다 기준이 되는 기준 패턴을 생성한 후 KNN 분류기는 전체 기준 패턴 중에서 미지의 입력 패턴 x 로부터 가장 가까운 거리에 있는 K 개의 패턴을 x 의 K-NN이라 하며, K-NN규칙은 패턴 x 의 K-NN의 각 요소가 어느 클래스에 가장 많이 속하는가를 조사하여, 그 클래스를 x 의 클래스로 결정한다. K 가 2 이상인 경우에는 K-NN 규칙은 $K=1$ 인 척보다 많은 정보를 참조하기 때문에 인식결과가 보다 양호하다고 알려져 있지만, 패턴의 분포에 따라서 그 반대의 경우도 있다.

2.2 벡터 양자화를 이용한 인식기

벡터 양자화(VQ : Vector Quantization)를 이용한 인식 방법은 인식 대상마다 집단화(clustering)를 통하여 코드북을 만든 후 인식 시에 양자화 오차를 계산하여 가장 적은 오차를 갖는 코드북을 입력 대상으로 인식하는 방법으로 주로 음성인식 초기단계에 사용되었고 문장독립 화자 인식에도 사용되어 왔다.

벡터 양자화를 이용한 인식 시스템의 블록도는 <그림 1>과 같다. 학습 과정에서는 각 감정마다 학습 데이

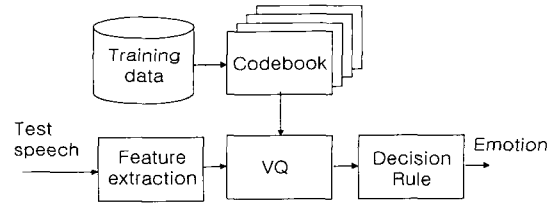


그림 1. 벡터 양자화를 이용한 감정인식 시스템 블록도
Fig. 1. Block diagram of the emotion recognition system using vector quantization.

터를 집단화하여 코드북을 만들고 인식 단계에서는 입력 음성을 각각의 코드북으로 양자화 한 후 양자화 오차를 계산하여 그 오차가 가장 적은 코드북의 감정을 입력 음성의 감정으로 결정한다. 양자화 이러한 방법은 입력 문장의 시간적인 변화에는 상관없이 동작하므로 이러한 특징을 이용하여 문장독립 감정 인식 시스템에 응용할 수 있다. 즉 감정의 구분된 학습데이터를 사용하여 감정별 코드북을 만들어 인식에 사용하는 것이다.

III. 실험 및 결과

1. 감정 인식 시스템 구성

감정 인식 시스템 구현하기 위해서는 DB 구축 과정, 특징 추출 과정, 학습 및 인식 과정으로 구성된다. 특징 추출 과정에서 음성으로부터 감정 인식을 위하여 필요한 정보를 얻어내고, 이러한 정보를 이용하여 학습 과정에서 기준패턴을 생성하고, 인식 과정에서 결정법칙을 이용하여 최소 거리나 최대 확률을 갖는 기준패턴으로 인식을 한다. 기본적인 인식 시스템은 <그림 2>와 같다.

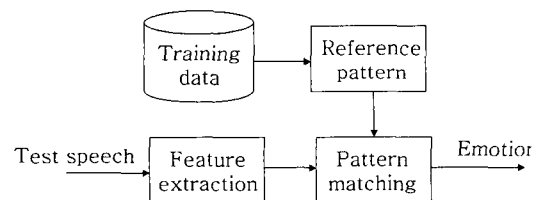


그림 2. 감정 인식 시스템 블록도
Fig. 2. Block diagram of emotion recognition system.

2. 특징 추출

본 연구에서는 화자독립-문장독립형 인식 시스템을 구현하기 위해서는 한국어 감정 음성 DB를 직접 구축하여 실험하였다. 구축한 DB의 데이터를 이용한 특징

추출 과정은 다음과 같다. 전처리를 통하여 16KHz로 샘플링하고, 고주파 성분을 보강한다. 이렇게 샘플링된 신호를 20 msec씩 프레임별로 나누어 분석하여 특징벡터를 구한다. 본 연구에서는 음성의 특징벡터를 음소군 특징벡터와 감정 특징벡터로 구분하였는데, 음소군 특징벡터는 발성기관의 해부학적인 차이나 발성기관의 조음 방법 차이에서 나타나는 음소특징을 추출한 MFCC, 델타 MFCC와 같은 특징벡터이고, 감정 특징벡터는 감정의 표현에 기여하는 피치, 델타 피치, 델타 델타 피치, 에너지, 델타 에너지, 델타 델타 에너지 등으로 구성된 특징벡터이다. 지금까지 감정과 음성과의 상관관계에 대한 연구에 따르면, 운율적 요소 즉 감정 특징벡터가 감정을 표현하는데 많은 영향을 끼친다고 알려져 있다. 부가적인 특징 추출 방법으로 벡터 양자화 기법을 사용하였는데, 이 방법에서는 N개의 클러스터링 방법으로 다차원 코드 벡터로 이루어진 코드북을 구성한 후, 특징벡터들을 코드북과 비교하여 가장 근접한 코드 벡터의 인덱스를 특징벡터로 사용하였다.

학습 및 인식 과정은 기준 패턴이나 기준 확률모델을 구하기 위해 사용되는 패턴 인식 기법에 따라 분류기가 달라진다. kNN 분류기는 제안된 방법과 비교하기 위하여 구성하였고, 음성의 감정 특징만을 이용하였다. 또한 피치, 델타 피치, 델타 델타 피치, 에너지, 델타 에너지, 델타 델타 및 MFCC, 델타 MFCC를 파라미터로 사용하여 벡터 양자화를 사용한 감정 인식 실험도 수행하였다.

3. 데이터 베이스

데이터 베이스를 구성하기 위해서는 사용 용도를 고려한 감정 선정, 문장 선정, 녹음 대상 선정, 녹음 환경, DB 규모 등의 결정 작업이 필요하다. 본 연구에서는 인간의 주요 감정인 기쁨, 슬픔, 화남의 3가지 감정과 이들의 기준이 되는 평상 감정을 포함한 4가지 감정을 인식 대상 감정으로 결정하였다. 또한 데이터베이스로 구성할 문장은 화자독립 및 문장독립형 감정 인식 시스템 개발에 사용될 수 있도록 하기 위하여 다음과 같은 사항을 고려하였다. 첫 번째, 하나의 문장을 3가지 감정(기쁨, 슬픔, 화남)상태로 표현하기에 용이한 문장을 선정하였다. 즉 문장의 내용이 중성적인 것을 선정하여 여러 가지 감정을 모두 가능한 것을 선정하였다. 두 번째, 자연스런 대화체 문장을 선정하였다. 문장에 감정이 포함되기 위해서는 단어보다는 문장 단위로 표

현하는 것이 보다 자연스럽고 감정을 표현하기에 용이하다. 세 번째, 전체적으로 모든 음소를 골고루 포함하도록 문장을 구성하였다. 또한 제시된 문장이외에 감정에 수반될 수 있는 보조적인 효과음은 배제하였다.

음성의 녹음은 평소 감정 표현을 훈련하는 아마추어 연극단원 남/녀 각 15명을 대상으로 하였고, 모든 참여자에 대해서 표준어 사용여부 및 감정 표현능력을 심사하여 선별하였다.

녹음작업은 조용한 사무실 환경에서 이루어졌고, DAT를 이용하여 녹음되었다. 이후 16kHz 16비트로 양자화 하여 PC에 저장 후 각 음성을 수작업으로 앞 뒤 약 50ms정도의 마진을 두고 분리하여 저장하였다. 각 화자는 45개의 문장을 4가지 감정으로 녹음하였고 녹음동안에 감정 표현이 미흡하다고 판단된 경우에는 다시 녹음을 하였다.

본 연구를 위하여 사용된 데이터의 규모는 5400(30명×4감정×45문장×1회)문장이다. 향후 실험에서는 제작된 DB 중 감정이 적절히 반영되었다고 판단되는 문장을 선별하는 주관적 평가를 거쳐 선택하였다. 구축된 DB가 화자의 감정을 어느 정도로 정확히 반영하는지를 판단하기 위해서 평소 음성 신호 처리 실험에 숙련된 연구원들을 대상으로 주관적 평가를 실시하였다. 주관적 평가는 5400문장을 문장 당 10명이 청취한 후 평가를 하였다. 주관적 평가를 통한 데이터 베이스의 감정 평가는 다음 <표 1>와 같다.

표 1. 주관적 평가 결과

Table 1. The result of subjective evaluation.

| 점수 | 문장수 | 백분율(%) | 누적백분율(%) |
|----|------|--------|----------|
| 10 | 2063 | 38.20 | 38.20 |
| 9 | 1031 | 19.09 | 57.30 |
| 8 | 592 | 10.96 | 68.26 |
| 7 | 415 | 7.69 | 75.94 |
| 6 | 296 | 5.48 | 81.43 |

표에서 알 수 있듯이 화자가 의도적으로 특정한 감정을 담아 발음한 음성이라 실제로 청취자들은 다르게 느낄 수 있는 것이다. 본 데이터 베이스를 만드는데 참가한 화자들은 다년간 연구 경험을 통하여 감정조절이 잘 되는 화자임에도 불구하고 10점(10명이 청취하여 10명 모두 정답인 경우)을 받은 문장은 전체의 38.2%였다.

4. 실험 결과

주관적 평가에서 100%의 정답률을 가진 데이터만을 선별하여 실험하였고, 20명의 화자(남성 10명, 여성 10명)는 학습 데이터, 10명의 화자(남성 5명, 여성 5명)를 인식 데이터로 사용하였다. 또한 총 45문장 중에 35문장은 학습에 나머지 10문장은 인식에 사용하여 화자 및 문장독립 감정인식 실험을 수행하였다.

4.1. KNN 분류기를 사용한 성능평가

KNN 분류기는 기존의 감정 인식 알고리즘으로 제안된 알고리즘과 비교하기 위하여 실험되었다. KNN의 경우 특징 파라미터로 피치 평균, 피치 표준편차, 피치 최대값, 에너지 평균, 에너지 표준편차를 사용하였다. KNN을 이용한 실험에서 LBG 군집화 알고리즘을 사용하여 감정별로 기준패턴을 생성하고 기준 패턴과의 거리측정을 위해 유클리디언 거리를 사용하였다. 코드북의 크기를 8, 16, 32, 64로 바꾸어 실험한 결과 인식률은 약 37.58 ~ 46.44%의 인식률을 보였으며 그 중 32일 때의 결과는 <표 2>와 같다. 클러스터 크기의 변화에 따른 인식률 편차는 인식률 대비 5% 미만으로 크기를 최적화함에 따른 인식률 향상은 크게 기대되지 않았다.

표 2. KNN 분류기를 이용한 감정 인식 시스템의 인식 성능(%)

Table 2. The recognition performance of emotion recognition system using KNN classifier(%).

| 감정 | 평상 | 기쁨 | 슬픔 | 화남 |
|----|------|------|------|------|
| 평상 | 32.1 | 21.4 | 25.0 | 21.4 |
| 기쁨 | 18.2 | 72.7 | 9.1 | 0.0 |
| 슬픔 | 20.0 | 20.0 | 40.0 | 20.0 |
| 화남 | 18.2 | 36.4 | 4.5 | 40.9 |
| 평균 | 46.4 | | | |

4.2. 벡터 양자화를 이용한 인식기의 성능평가

피치, 델타 피치, 델타 델타 피치, 에너지, 델타 에너지, 델타 델타 및 MFCC, 델타 MFCC를 파라미터로 하여 각 감정별로 집단화(clustering)를 통한 코드북을 만든 후 입력을 테스트 입력을 양자화하여 최소의 거리를 갖는 코드북을 입력의 감정으로 인식하는 인식 시스템을 구성하여 성능을 평가하였다. <표 3>은 각종 파라미터에 따른 인식 성능과 그때 사용된 코드북의 크기

를 나타낸다. 여기서 사용된 파라미터의 기호는 다음과 같다.

| | |
|-----|----------------|
| P | : 피치 |
| DP | : 델타 피치 |
| DDP | : 델타 델타 피치 |
| E | : 에너지 |
| DE | : 델타 에너지 |
| DDE | : 델타 델타 에너지 |
| M | : 멜 캡스트럼 |
| DM | : 델타 멜 캡스트럼 |
| DDM | : 델타 델타 멜 캡스트럼 |

표 3. 벡터 양자화를 이용한 감정 인식 시스템의 성능 평가

Table 3. Performance evaluation of emotion recognition system using vector quantization.

| 파라미터 | 코드북 크기 | 인식률(%) |
|----------|--------|--------|
| P | 32 | 42.24 |
| DP | 64 | 42.24 |
| DDP | 64 | 38.39 |
| E | 256 | 41.38 |
| DE | 128 | 33.62 |
| DDE | 256 | 23.28 |
| M | 64 | 67.24 |
| DM | 256 | 56.03 |
| DDM | 256 | 56.03 |
| P+DP | 256 | 42.24 |
| P+DP+DDP | 64 | 45.69 |
| DP+DDP | 64 | 40.52 |
| E+DE | 128 | 46.55 |
| E+DE+DDE | 128 | 51.72 |
| DE+DDE | 4 | 41.38 |
| M+DM | 256 | 73.28 |
| M+DM+DDM | 127 | 71.55 |

표에서 알 수 있듯이 가장 우수한 성능을 나타낸 것은 MFCC와 델타 MFCC를 결합한 MFCC+DMFCC로 73.28%의 인식 성능을 나타내었다. 피치와 에너지는 화자중속 또는 문장중속 형태의 시스템에서는 비교적 우수한 성능을 나타내지만 문장독립 및 화자독립 감정 인식 시스템에서는 40~50%정도의 낮은 인식 성능을 나타내고 있다. 이러한 것은 시스템의 형태가 문장독립 및 화자독립 감정 인식 시스템이기 때문으로 코드북에 다양한 화자와 다양한 문장이 포함되어 있기 때문이다.

MFCC의 경우에는 오히려 피치나 에너지의 영향보다는 각 감정상태에서 발음한 음성의 스펙트럼 차이를 표현하기 때문에 인식 성능이 더 우수한 것으로 판단된다. 다음 <그림 3, 4, 5>는 각각의 파라미터와 코드북 크기에 따른 인식 성능을 보여준다. <그림 3>는 피치, 델타 피치, 델타 델타 피치 및 이들을 결합한 것을 사용한 결과로 전반적으로 인식 성능이 낮은 편이다. 이러한 것은 피치 값과 그 변화 특성이 화자의 감정 상태뿐만 아니라 다른 화자나 문장의 내용에 따라서도 크게 영향받는다는 것을 나타낸다. <그림 4>은 에너지, 델타 에너지, 델타 델타 에너지 및 이들을 결합한 것을 사용한 결과로 피치 파라미터 보다는 우수하지만 전반적으로 낮은 인식 성능을 보인다. <그림 5>는 MFCC, 델타 MFCC, 델타 델타 MFCC 및 이들을 결합한 것을 사용한 결과로 이들중 가장 우수한 결과를 나타낸다. 이러한 결과는 감정 상태에 따른 음성 신호의 차이를 나타내는 MFCC 파라미터가 문장 및 화자독립 시스템인 경우에 보다 적합하다는 것을 나타낸다.

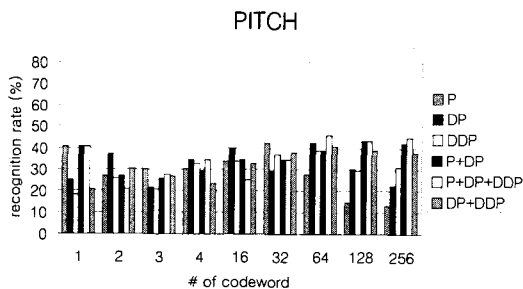


그림 3. 피치 파라미터와 코드북 크기에 따른 인식 성능
Fig. 3. Recognition performance according to the pitch parameters and codebook size.

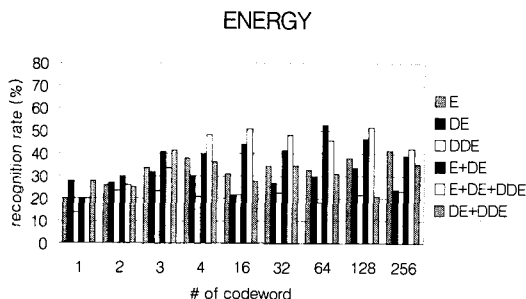


그림 4. 에너지 파라미터와 코드북 크기에 따른 인식 성능
Fig. 4. recognition performance according to the energy parameters and codebook size.

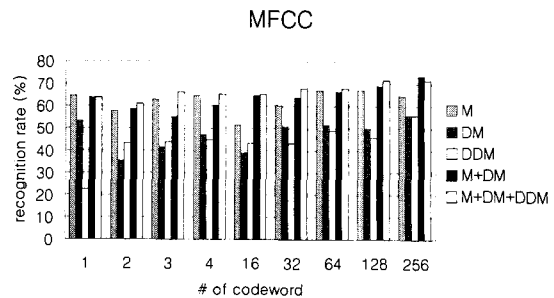


그림 5. MFCC 파라미터와 코드북 크기에 따른 인식 성능
Fig. 5. recognition performance according to the MFCC parameters and codebook size.

IV. 결 론

본 논문에서는 음성신호를 사용한 인간의 감정인식을 위한 특징 파라미터의 비교에 관하여 연구하였다. 이를 위하여 여러 가지 감정 상태에 따라 분류된 한국어 음성 데이터 베이스를 이용하여 감정 인식 및 음성 신호 처리에 널리 사용되고 있는 음성 신호의 피치와 에너지의 평균, 표준편차와 최대 값 등 통계적인 정보 나타내는 파라미터와 음소의 특성을 나타내는 MFCC 파라미터의 성능을 비교 평가하였다.

이러한 목적을 위하여 본 연구에서는 다양한 입력문장에 담긴 화자의 감정을 인식할 수 있는 문장독립(text-independent) 및 화자 독립 감정 인식 알고리즘으로 벡터 양자화 방법을 사용한 시스템을 구성하여 사용하였다. 인식을 위한 대상 감정은 평상, 기쁨, 슬픔, 화남의 4가지를 사용하였고 주관적 평가를 통해서 정답률이 100%인 문장만을 선별하여 실험에 사용하였다.

성능 평가를 위한 실험에서는 운율적 특징으로 피치, 에너지와 각각의 미분 값을 결합하여 사용하였고, 음소의 특성을 나타내는 특징으로 MFCC와 그 미분 값을 결합하여 사용하였다. 알고리즘의 특성상 KNN의 경우 입력 신호에 대한 피치 평균, 피치 표준편차, 피치 최대 값, 에너지 평균, 에너지 표준 편차를 이용하여 KNN 방법은 46.44%의 인식률을 나타내었다. 인식 실험에서 문장 종속이나 화자 종속인 경우에는 우수한 성능을 나타내는 것으로 알려진 피치 및 에너지 파라미터는 문장 독립 및 화자 독립인 경우에는 성능이 많이 저하되는 것을 알 수 있다. MFCC와 델타 MFCC를 결합한 파라미터로 하여 크기가 256개인 코드북을 사용한 경우 약 73.3%의 인식 성능을 나타내었다. 이러한 것은

기준에 사용하던 KNN 방법보다 우수한 성능을 나타내고 구조도 간단한 장점이 있다.

참 고 문 헌

- [1] Rosalind W. Picard, Affective Computing, The MIT Press 1997.
- [2] C. E. Williams and K. N. Stevens, "Emotions and speech: Some acoustical correlates", Journal Acoustical Society of America, Vol. 52, No. 4, pp. 1238-1250, 1972.
- [3] Lain R. Murray and John L. Arnott, "Toward the simulation of emotion in synthetic speech: A review of the literature on human vocal emotion", Published in J. Acoust. Soc. Am., pp. 1097-1108, Feb. 1993.
- [4] Janet E. Cahn, "The generation of affect in synthesized speech", Journal of the American Voice I/O Society, Vol. 8, pp. 1-19, July 1990.
- [5] Frank Dellaert, Thomas Polzin, Alex Waibel, "Recognizing emotion in speech", Proceedings of the ICSLP 96, Philadelphia, USA, Oct. 1996.
- [6] Thomas S. Huang, Lawrence S. Chen and Hai Tao, "Bimodal emotion recognition by man and machine", ATR Workshop on Virtual Communication Environments - Bridges over Art/Kansei and VR Technologies, Kyoto, Japan, April 1998.
- [7] K. R. Scherer, D. R. Ladd, and K. E. A. Silverman, "Vocal cues to speaker affect: Testing two models", Journal Acoustical Society of America, Vol. 76, No. 5, pp. 1346-1355, Nov. 1984.
- [8] Michael Lewis and Jeannette M. Haviland, Handbook of Emotions, The Guilford Press, 1993.
- [9] D. Roy and A. Pentland, "Automatic spoken affect analysis and classification", in Proceedings of the Second International Conference on Automatic Face and Gesture Recognition, pp. 363-367, Killington, VT, Oct. 1996.
- [10] Jun Sato, and Shigeo Morishima, "Emotion Modeling in Speech Production using Emotion Space", Proceedings of the IEEE International Workshop 1996, pp. 472-477, IEEE, Piscataway, NJ, USA, 1996.
- [11] 강봉석, 음성 신호를 이용한 감정 인식, 석사학위논문, 연세대학교, 1999년 12월
- [12] L. R. Rabiner and B. H. Juang, Fundamentals of speech recognition, Prentice-Hall Inc., 1993.
- [13] R.O. Duda, and P.E. Hart, Pattern classification and scene analysis, John Wiley & Sons Inc., 1973.
- [14] Earl Gose, Richard Johnsonbaugh, and Steve Jost, Pattern Recognition and Image Analysis, Prentice Hall Inc., 1996.

저 자 소 개



金元九(正會員)

1983년 3월~1987년 2월 : 연세대학교 전자공학과 학사. 1987년 9월~1989년 8월 : 연세대학교 전자공학과 석사. 1989년 9월~1994년 2월 : 연세대학교 전자공학과 박사. 1994년 9월~현재 : 군산대학교 전

자정보공학부 부교수. 1998년 9월~1999년 9월 : Bell Lab, Lucent Technologies(USA) 객원연구원. <주관심 분야 : 음성 및 디지털 신호처리, 음성 인식, 감정 인식, 음성 통신 등임>