

## Speaker Recognition Using Mel-Frequency Cepstrum Coefficients and Artificial Neural Network for Security Systems

Ahmet Küçüker<sup>1</sup>, Can Yüzkollar<sup>2</sup>, Ahmet Sansli<sup>2</sup>, Fatih Sen<sup>2</sup>

<sup>1</sup>Sakarya University, Electronic Engineering Department, Esentepe Campus,

<sup>2</sup>Sakarya University, Computer Engineering Department, Esentepe Campus,  
Sakarya, Turkey

**Abstract:** In this study; the sample sounds of the speakers are transferred to computer for identifying the speaker. Study is based on three steps; In first step analog sound that taken is converted to data to be processed and determine own characteristic of each sound. in second step converted sounds called data are used to train neural network. at last step trained neural network is test by the test sounds and obtained istatistical informations.

**Keywords:** Neurel Network, Speaker Identification, Speech Recognition, Mel Cepstrum Coefficients

## Güvenlik Sistemleri İçin Mel Frekans Kepstrum Katsayıları ve Yapay Sinir Ağları Kullanılarak Konuşmacı Tanıma

**Özet:** Bu çalışmada konuşmacı ses örnekleri bilgisayara alınarak bu konuşmacıların kimliklerinin belirlenmesi sağlanmıştır. Çalışma üç aşamadan oluşmaktadır. İlk aşama konuşmacı sesinin dijital ortamda veriye dönüştürülmesi işlenebilir hale getirilerek özellik çıkartılması ikinci aşama elimizdeki konuşmacı ses örnekleriyle yapay sinir ağının eğitilmesi ve son aşama da eğitilmiş sinir ağına test verileri göndererek istatistik bilgileri elde edilmesidir.

**Anahtar Kelimeler:** Sinir Ağları, Konuşmacı kimliklendirme, Konuşmacı tanıma, Mel Kepstrum Katsayıları

Reference to this paper should be made as follows (bu makaleye aşağıdaki şekilde atıfta bulunulmalı):

A. Kucuker et al, 'Speaker Recognition Using MFCC and Artificial Neural Network For Security Systems' , Elec Lett Sci Eng , vol. 2 (2) , (2006), 8-15

### 1 Giriş

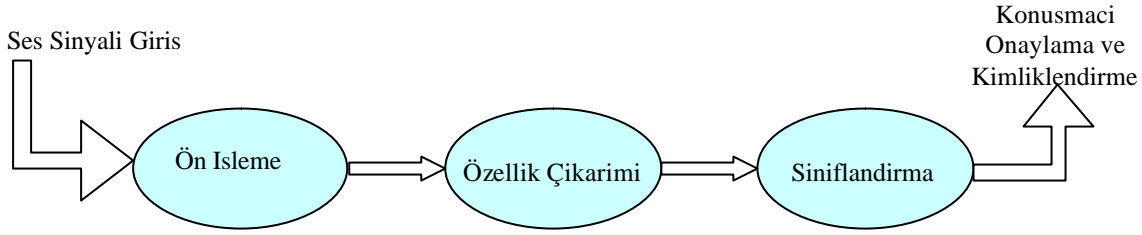
Otomatik olarak tanımlama veya doğrulama için kullanılan, insana ait, ölçülebilir, essiz fizyolojik ölçüler vardır. Bu ölçülerin kaybolması, unutulması ve bir baskısı tarafından kullanılması, fiziksel saldırılar dışında pratikte imkânsız kabul edildiği gibi taklit edilmeleri de çok zordur. Gelistirilen teknolojileri kullanarak her ortamda kişinin el, parmak izi, ses, göz, imza, retina, kulak şekli, DNA, klavyeye bası deseni, koku ve yürüyüş kistaslarının ölçülerek tanınması mümkündür. [1] Dolayısıyla konuşmacı tanımlama bu alanda kullanılabilir nitelik kazanmaktadır. Böyle bir sisteme ihtiyaç duyulmasının sebebi insanların parola hatırlamalarının gerekliliğini ortadan kaldırmak ya da parolanın çalınması ihtimaline karşı ek bir tedbir almak da denilebilir. Bu yüzden günden güne ASR(Automatic Speaker Recognition) otomatik konuşmacı tanıma sistemleri geliştirilmektedir. . Konuşmacı tanıma sistemlerinin güvenlik dışında da birçok uygulama alanı vardır. Konuşma tanıma sistemleri sayesinde telefonda kredi kartı numarası tuşlamak yerine numarayı sesli söyleyebiliyor veya arabada sesli komutlarla bilgiye ulaşabiliyoruz. [2]

### 2 Konuşmacı Tanıma Sistemi

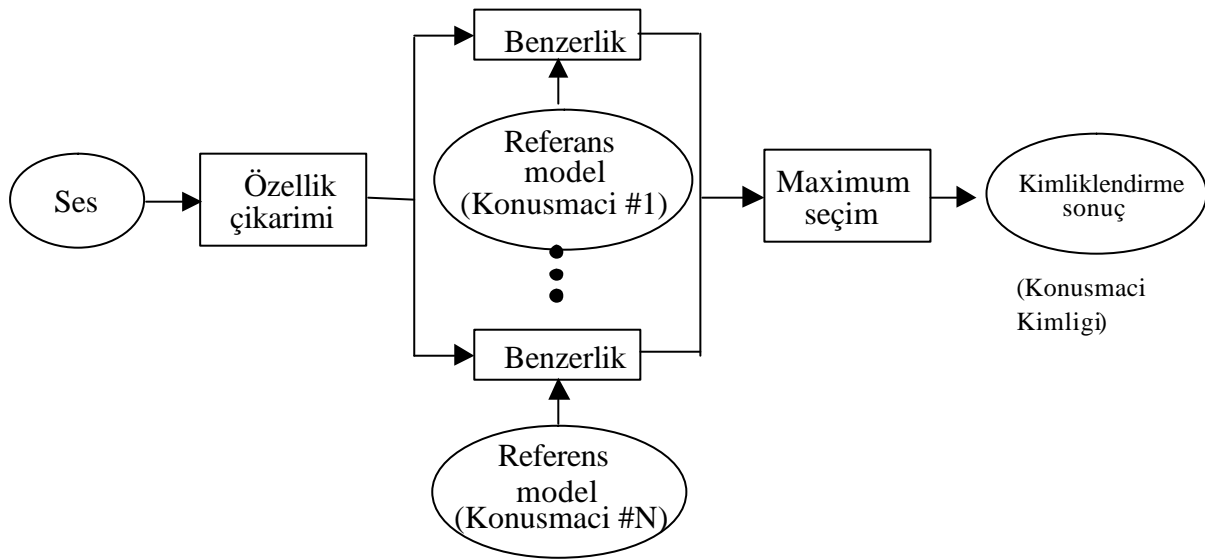
İnsanın ses özelliklerini barındıran ses dalgaları mikrofon vasıtasıyla elektrige çevrilerek bilgisayara veri olarak alınmaktadır. Birçok araştırmacı istatistiksel özellikler kullanarak konuşmacı tanıma üzerine çalışma yapmıştır [3]. Konuşmacı tanıma sistemini fonksiyon yönünden iki başlık altında sınıflandırabiliriz. İlki konuşmacıyı onaylama ikincisi konuşmacıyı

kimliklendirir. [4]. Uygulamamızda Mel Frekans Kepstrum Katsayilari yöntemiyle konusmacinin seslerinden gerekli özellikler çıkarılmış, çıkarılan bu özelliklere Yapay zeka tekniği uygulanarak konusmacinin taninması sağlanmıştır.

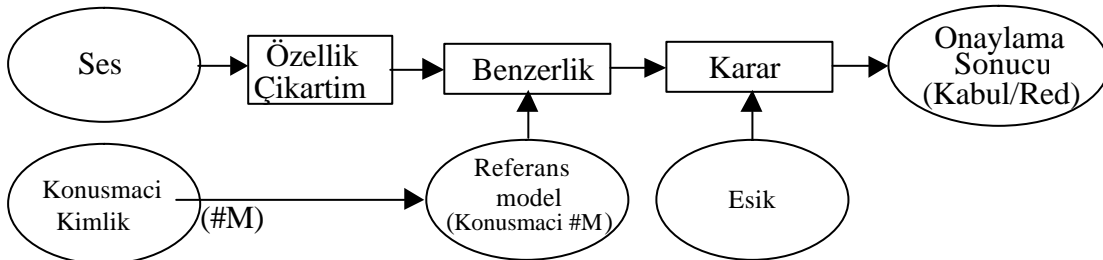
Konusmaci onaylama, konusmacinin kimlik talebini onaylama ya da onaylamama işlemidir. Konusmaci onaylamanın iki çeşit hata kriteri vardır. Yanlış Kabul(False Accept) ve Yanlış Red (False Reject) Bu iki kriterin kesisi mi ise “ Esit Hata Oranı (Equal Error Rate) “ olarak adlandırılır.Uygun bir tanıma sistemi için Esit Hata Oranı ; Yanlış Kabul ve Yanlış Red arasında uzlastirici görev yapmalıdır.[5] Konusmaci kimliklendirme sunulan insan seslerinin hangi konusmaciya ait olduğunun belirlenmesi işlemidir. Temel bir konusmaci tanıma sistemi şekilde gösterildiği gibidir. (Sekil 1)



Sekil-1 Temel Konusmaci Tanıma Sistemi Blok diyagramı



(a) Konusmaci Kimliklendirme



(b) Konusmaci Onaylama

Sekil 2 Konusmaci Onaylama ve Kimliklendirme

Konusmaci Tanima sistemini fonksiyon yönünden iki bölüme ayrıldığı gibi method olarak da iki bölüme ayrılmaktadır. Bunlardan ilki *Konusma Metninden Bagimsiz* konusmaci tanima digeri *Konusma Metnine Bagimli* Konusmaci tanima olarak adlandırılır [4]. Çalışmamızda ikinci metod kullanılmıştır.

### 3 Özellik Çıkarım

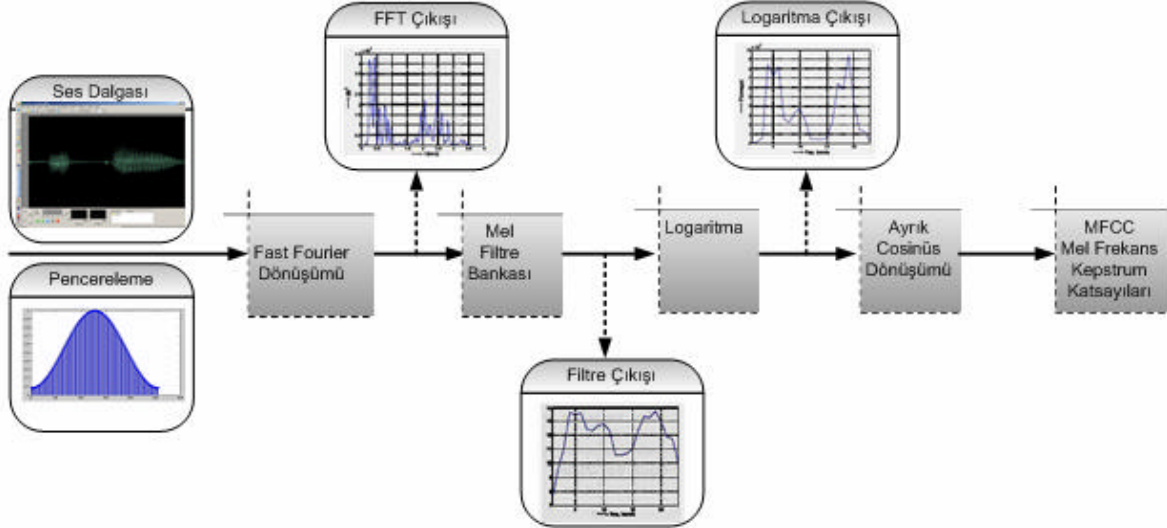
Bu çalışmada özellik çıkarım yöntemlerinden olan mel-frekans Cepstrum katsayıları yönteminden faydalanılmıştır. [6 , 7] *Cepstrum* kavramı ilk olarak 1963'de *Bogert, Healy* ve *Tukey* tarafından kullanılmıştır. *Cepstrum*, homomorfik sinyal işleme teknikleri içinde yer alır. [6]. Homomorfik sistemler doğrusal olmayan sistemlerin bir sınıfı olarak kabul edilirler. Doğrusal sistemler homomorfik sistemlerin özel bir durumudur. Sesli ifade bağlamında kullanılan homomorfik sistem,  $S(f)=V(f).G(f)$  şeklinde ifade edilebilir. Burada  $s(n)$  sesli ifadeyi,  $v(n)$  girtlagi, yani, sesin izlediği yolu,  $g(n)$  ise asıl ses sinyalini, yani ses telleri tarafından üretilen ve değişime uğramamış ses sinyalini temsil eder. Bu şekilde ifade edildiğinde girtlagin etkisini asıl sesteki homomorfik sinyal işleme teknikleriyle ayırmak mümkün olmaktadır. *Cepstrum* değerlerinin hesaplanması için (Es. 1.1) eşitliği kullanılır. [7]

$$c(n) = \frac{1}{N_s} \sum_{k=0}^{N_s-1} \log_{10} |S(k)| e^{(2\pi / N_s)kn}, \quad 0 \leq n \leq N_s - 1 \quad (1.1)$$

Burada  $c(n)$ ,  $n$ . *Cepstrum* olarak adlandırılır.  $S(k)$ , *Cepstrum* değerinin ait olduğu frekans aralığı için alınan *fourier* dönüşümünü belirtir.  $N_s$  o andaki çerçevenin boyunu gösterir. Burada dikkat edileceği gibi  $c(0)$  doğrudan o andaki DFT (**Discrete Fourier Transform**) spektrum değerini gösterir. Sesli ifadenin gürültüden ayırt edilmesi için önce spektrumun logaritması alınır ve ters *fourier* dönüşümü yapılır. Bu şekilde belirlenen *Cepstrum* değerleri *fourier* dönüşümünden türetilmiş *cepstral* katsayılar (*fourier transform derivated cepstral coefficients*) olarak adlandırılır. Burada hesaplanan *Cepstrum* değeri sesli ifade tanımında kullanılan önemli bilgileri elde etmede etkilidir. *Fourier* dönüşümünde kullanılan frekans *mel* skalasında örneklenirse elde edilen *Cepstrum* değerleri *Mel Cepstrum* değerleri olarak adlandırılır . [7]

Mel frekans cepstral özellikleri için her pencereye sirasiyla su adımlar uygulanır[5]. (Sekil 3)

- Ayrik Fourier Dönüşümü alınır
- Ayrik Fourier Dönüşümü katsayıları Mel filtre bankasının genlik frekans cevabına göre ağırlıklandırılır.
- Logaritmik Enerjileri Hesaplanır
- Ayrik Kosinüs Dönüşümü bulunur



Sekil 3 Özellik Çıkarım Adımları Seması

*Mel Cepstrum* değerlerini hesaplamak için Mel filtre bankası yöntemi (1.2) formülüyle uygulanır. (Davis and Mermelstein). Burada  $X_k$  mel skalasında k. bant geçişli filtrenin uygulanmasıyla elde edilen spektrum değeridir. Mel Filtre bankasının amacı duyma mekanizmasının kritik bant filtrelerini simülasyonunu yapmaktır. Filtreler Mel Skalasında düzenli bir şekilde yerleşen üçgensel filtrelerdir. [7 , 10 , 11 ]. Üçgensel filtre çıkışları  $Y(i)$ ;  $i = 1; \dots; M$  logaritma kullanarak sıkıştırılmış ve ayrık kosinüs dönüşümü uygulanmıştır. [8].

$$melc(n) = \sum_{i=0}^M \log Y(i) \cos \left[ \frac{pn}{M} \left( i - \frac{1}{2} \right) \right] \quad (1.2)$$

Bu fonksiyonun sonucunda elde edilen katsayılar ile sesimizin artık yapay sinir ağına girebilecek şekilde yani işlenebilecek şekilde getirildiğini söyleyebiliriz.

Özellik çıkarımında farklı bir yaklaşım olarak Yapay Sinir Ağları da kullanılmaktadır [12, 13]

#### 4- Yapay Sinir Ağları

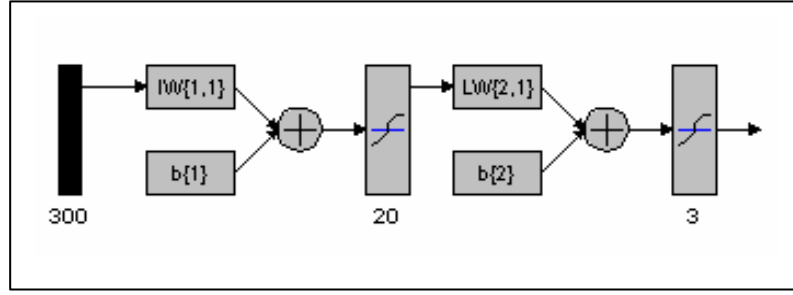
Yapay sinir ağları (YSA) günümüzde bilgi sınıflama ve bilgi yorumlamanın içinde bulunduğu değişik problemlerin çözümünde kullanılmaktadır [9]. Özelliklerini elde ettiğimiz ses sinyallerinin tanınması da bu kapsama girmektedir. Çalışmamızda genelleme özelliğini kullanarak bozulmuş seslerin de tanınması üzerinde durulmuştur.

Yapay sinir ağları şöyle sınıflandırılmaktadır:

- Tek katmanlı ileri beslemeli ağlar : Katmanlı modellerdeki en basit ağ tipi olup bir çıktı katmanı ve buna bağlı bir girdi katmanı bulunmaktadır
- Çok katmanlı ileri beslemeli ağlar : Çalışmamızda bu yöntemi seçmemizin sebebi tek katmanlı ağlara göre daha karmaşık problemlere çözüm getirebilmesidir. Dezavantajı ise Tek katmanlı ileri beslemeli ağlara göre eğitimi daha uzun sürmektedir [10]

Tek katmanlı ağlardaki girdi ve çıktı katmanından başka, bir yada daha fazla sayıda gizli katman içeren ağlara çok katmanlı ağ denir. Dis dünya tarafından doğrudan müdahale

edilmediği için gizli katman adı verilen katmanda bulunan birimlere de gizli birimler adı verilir ( Şekil 4).



Şekil 4 Çok katmanlı ileri beslemeli ağ modeli

Geri yayımlı öğrenme metodu olarak isimlendirilen eğitim algoritmasının temeli ilk olarak Werbos'un Harvard Üniversitesi'nde verdiği doktora tezinde ileri sürülmüş [11], daha sonra Parker tarafından MIT'nin bir teknik raporunda ele alınmış [12] ve Rumelhart ve arkadaşları tarafından da popüler hale getirilip uygulanabilir sekile dönüştürülmüştür.

İleri besleme aşamasında eğitim için kullanılan girdi, sisteme beslenir ve bunun sonuçları her bir katmanı geçerek çıktı katmanına kadar gelir ve girdiye karşılık bir çıktı elde edilir. Bu aşama sırasında ağ üzerindeki ağırlık değerleri sabit tutulur.

İkinci aşama olan geri besleme aşamasında, elde edilen çıktı ile hedef çıktı arasındaki farktan hata sinyali elde edilir ve bu sinyal ağ yapısında geriye doğru yayılır. Bu geri yayılım aşamasında, oluşan hatayı minimuma indirecek şekilde ağırlık değerleri güncellenir.[13] Bu durumda sistemin eğitilmesi uzun süre almakla beraber, eğitilmiş bir sistemden bilgi alınması çok hızlı olmaktadır

## 5 Uygulama

Uygulamada konuşmaciyi tespit etmek için kapi kelimesi kullanılmıştır. Kapi kelimesi çok basit bir kelime gibi görünmesine rağmen YSA nin tanımı açısından büyük zorluk içermektedir. Fonetik olarak söylenisi rahat olduğu için herkes birbirine yakın şekilde telaffuz etmektedir. Bu sebeple bu kelimeyi tanıyabilen bir sistem diğer kelimeleri daha rahat öğrenecektir. Örneğin Sıcak kelimesi gibi bir kelime kullanılmış olsaydı S harfi ve K harfinin sonda kullanılması kisten kişiye daha çok farklılık gösterdiğinden tanıma işlemi çok daha kolaylaşmış olacaktı. [14]

Seçtiğimiz bu kelime için her bir konuşmacıdan 10 örnek alınarak toplam 30 ses kaydı kullanılmıştır. Alınan veriler Matlab ortamında *voicebox toolbox* kullanılarak MFCC fonksiyonu ile katsayılar elde edilmiştir [15]. Elde edilen 25 satır 12 sütunlu bu öznelik vektörlerinin yapay sinir ağına girebilmesi için gereken matris dönüşümleri yapılmış ve 300\*30 luk eğitim seti (aas1 olarak adlandırıldı) hazırlanmıştır. Daha sonra her bir konuşmacıdan 5'er test sesi kaydedilerek gerekli matrisler elde edilmiştir.

Sinir ağı aşağıdaki Matlab kodlarıyla oluşturulmuştur. [18, 20]

```
net = newff(minmax(aas1),[20 3],{'tansig' 'tansig'},'trainscg');
```

Egitim kodlari [18, 20]

```
net.trainParam.epochs = 1250;
net.trainParam.goal = 1e-003;
net.trainParam.min_grad = 1e-07;
net.trainParam.time = Inf;
net.trainParam.show = 1;
net = train(net,aas1,cikiss);
```

Elimizdeki test seslerinin kime ait oldugunu belirlemek için asagidaki kod kullanilmistir.

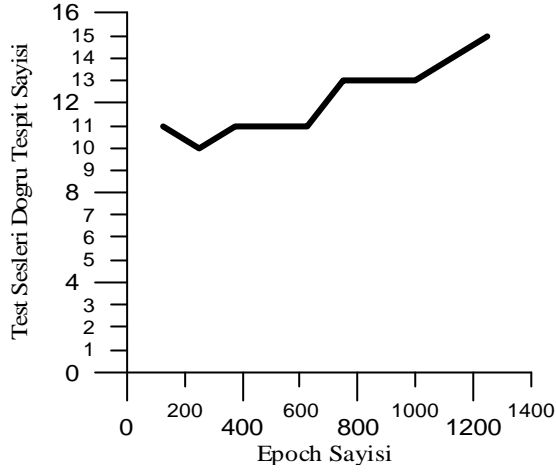
```
sonuç=sim(net,ses1)
```

## 6 Sonuç

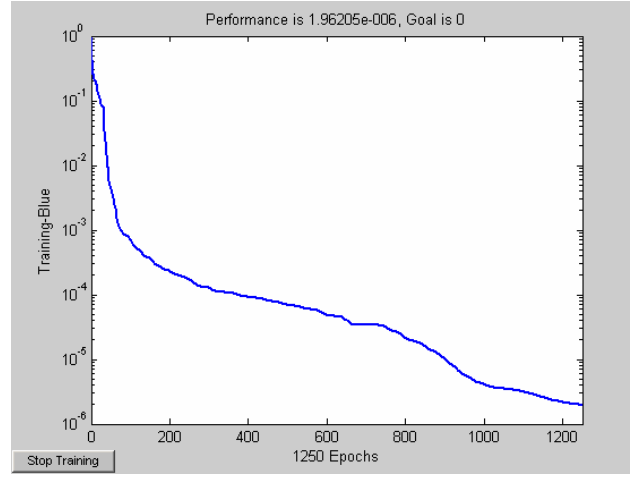
Egitimli sinir agina gönderilen test seslerine verilen cevaplar kaydedilerek ilgili dogruluk oranlarinin epochs sayısına göre dagilimi asagida gösterilmistir. (Tablo 1)

**Tablo 1. Dogruluk oranlarinin epochs sayısına göre dagilimi**

Epoch Sayisi	125	250	375	500	625	750	875	1000	1125	1250
Test Sesleri										
Ak1	+	+	+	+	+	+	+	+	+	+
Ak2	+	+	+	+	+	+	+	+	+	+
Ak3	+	+	+	+	+	+	+	+	+	+
Ak4	+	+	+	+	+	+	+	+	+	+
Ak5	+	+	+	+	+	+	+	+	+	+
As1	+	+	-	-	-	-	-	-	+	+
As2	+	+	+	+	+	+	+	+	+	+
As3	+	+	+	+	+	+	+	+	+	+
As4	+	+	+	+	+	+	+	+	+	+
As5	-	-	-	-	-	-	-	-	-	+
C1	-	-	-	-	-	+	+	+	+	+
C2	-	-	+	+	+	+	+	+	+	+
C3	-	+	-	-	-	+	+	+	+	+
C4	+	+	+	+	+	+	+	+	+	+
C5	+	+	+	+	+	+	+	+	+	+



Sekil 5 – Tablo 1 Değerlerinin Grafiği



Sekil 6 – Eğitim Performans Grafiği

Eğitim sesleri test sesi olarak kullanıldığında tanıma işlemi 125 epochs tan itibaren tüm sesler için doğru sonuç vermiştir.

## 7 Tartışma ve Öneriler

Bu çalışmada üç konuşmacının ses örnekleri bilgisayara alınarak Mel Frekans Kepstrum Katsayıları yöntemiyle konuşmacının seslerinden gerekli özellikler çıkarılmış, çıkarılan bu özelliklere Yapay zeka tekniği uygulanarak konuşmacının tanınması sağlanmıştır. Tanıma sonuçlarına bakıldığında başarı oranının % 100 olması bu sistemin kullanılabilir bir sistem olduğunu ortaya koymuştur. Fakat ses taklidi düşünüldüğünde bu sistemin parmak izi tanıma vb. sistemlerle birlikte kullanımı sistemin kullanılabilirlik oranını artıracaktır.

## Referanslar

- [1] Simon Liu, Mark Silverman. "A Practical Guide to Biometric Security Technology," *IT Professional*, vol. 03, no. 1, pp. 27-32, January/February, 2001.
- [2] Hemphill, C. T., Agarwal, R., Muthusamy, Y. K., and Gong, Y. "Voice-Driven Information Access in the Automobile", *IEEE Vehicular Technology Society News*, August, 8-11, 2000
- [3] FURUI S., "Comparison of Speaker Recognition Methods Using Statistical Features and Dynamic Features", *IEEE Transaction on ASSP*, Vol. 29, No. 3, pp. 342-350, June 1981.
- [4] İNAL M., FATİHOĞLU Y. S. " Self Organizing Map And Associative Memory Model Hybrid Classifier For Speaker Recognition" *Kocaeli 2002*
- [5] FARRELL, K. R., MAMMONE, R. J., and ASSALEH, K. T., Jan. 1994. *Speaker Recognition Using Neural Networks and Conventional Classifiers*, *IEEE Trans. On Speech and Audio Processing*, Vol.2, No.1, part II.
- [6] A.V.Oppenheim and R.W.Schafer, *Discrete-Time Signal Processing*, Prentice Hall, Englewood Cliffs, NJ, 1989
- [7] MENGÜSOĞLU, Erhan " Rule Based Design And Implementation Of A Speech Recognition System For Turkish Language" 1999
- [8] HUANG, X., Acero, A., and Hon, H.-W. *Spoken Language Processing: a Guide to Theory, Algorithm, and System Development*. Prentice-Hall, New Jersey, 2001.

[9] ELMAS, Ç., 2003, “*Yapay Sinir Aglari (Kuram,Mimari,Uygulama)*”, Seçkin Yayinlari.

[10] AYDIN Ö. “*Yapay Sinir Aglarini Kullanarak Bir Ses Tanima Sistemi Gelistirilmesi*” Yüksek Lisans Tezi, Edirne 2005

[11] WERBOS, P. J., 1974, “*Beyond Regression: New Tools for Prediction and analysis in the behavioral sciences*”

[12] PARKER, D. B., 1985, “*Learning-logic: Casting the Cortex of the Human Brain in Silicon*”, Technical Report TR-47, Center for Computational Research in Economics and Management Science, MIT, Cambridge, MA.

[13] TEMURTAS Dr.F. “*Yapay Sinir Aglari Ders notlari*” SAÜ 2005

[14] KUCUKER Dr.P. “*Fonetik Ders Notlari*” SAÜ 2004

[15]. MATLAB® Documentation (2002) *Neural Network Toolbox Help, Version 7.0, Release*