

GAUSSIAN MIXTURE MODEL BASED SPEAKER RECOGNITION USING AIR AND THROAT MICROPHONE

Aaliya Amreen.H

*Department of Computer Science and Engineering
B.S Abdur Rahman University
Vandalur 48
aaliyaamreen.amreen330@gmail.com*

Dr. Sharmila Sankar

*Department of Computer Science and Engineering
B.S Abdur Rahman University
Vandalur 48
hodcse@bsauniv.ac.in*

Abstract- Speaker recognition is a biometric recognition technique used to identify and verify a speaker from his/her speech data. Speaker recognition system uses mechanism to recognize the speaker by using the speaker's speech signal. Generally, speech information are recorded through the air microphone and these speech information are used as input for the speaker recognition system as they are prone to environmental background noise, the performance is enhanced by integrating an additional speech signal collected through a throat microphone along with speech signal collected from standard air microphone. The resulting signal is very similar to normal speech, and is not affected by environmental background noise. This paper is mainly focused on extraction of the Mel frequency Cepstral Coefficients (MFCC) feature from an air speech signal and throat speech signal to built Gaussian Mixture Model(GMM) based closed-set text independent speaker recognition systems.

Keywords: *Speaker Recognition, GMM, MFCC, Throat Microphone*

I. INTRODUCTION

Speaker recognition is a biometric procedure that uses an individual's speech for recognition purposes. The speaker recognition process specified by both the physical structure of an individual's vocal tract and the behavioural characteristics of the individual[1]. Speaker recognition technique uses the

speaker's voice to verify their identity and provides services such as voice dialling, database access services,[2] information services and voice mail. Speech is a complicated signal produced as a result which provides different levels of medium such as semantic, linguistic and acoustic.[4] Besides, there are speaker related differences which as a result specify a combination of anatomical differences that inherent in the vocal tract and the learned speaking habits of different individuals. In speaker recognition, all these signals are taken into account and used to discriminate between speakers. Speaker recognition can be classified into different categories such as Open Set vs. Close Set, Identification vs Verification and Text Dependent vs Text independent speaker recognition. Text dependent uses a constrained mode and Text independent uses an unconstrained mode[5]. In a system using text dependent speech, the individual utters either a fixed password or prompted phrase that is programmed into the system and this type of system can improve performance especially with cooperative users. A text independent system has no knowledge of the presenter's phrasing and is much more flexible in situations where the individual submits the sample which may be unaware of the collection or unwilling to cooperate, that presents a more difficult challenge[6]. The Figure 1 shows how the speaker recognition is classified based on the trained speakers in the system. An open set system can have any number of speakers[9] that are trained and registered, but in closed set system it can have only a fixed number of

users. Identification is the task of determining an unknown speaker's identity, but verification is the process of rejecting the identity claim of a speaker.

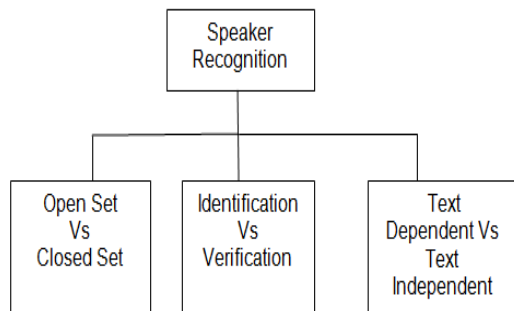


Figure1 Classification of Speakers

II. THROAT SPEECH AND AIR SPEECH

A throat microphone is a close contact microphone that absorbs vibrations directly from the wearer's throat in the form of single or dual sensors worn against the neck. These sensors are known as transducers that can pick up speech in extremely noisy or windy environments such as on a motorcycle, or in a war field[14]. An air microphone is a normal microphone which converts acoustic-to-electric transducer or sensor that formulates sound into an electrical signal. Microphones are used in many of the applications such as telephones, hearing aids, public address systems for concert halls and public events, two-way radios, megaphones, radio and television broadcasting, and in computers Throat microphones function well under noisy conditions,[10] but the functioning of normal microphone is not well because of high levels of background noise.

In advanced throat microphone it is capable of picking up the whispered speech to perform well in all kind of environments, but considering the normal microphone it does not be able to proceed any in any kind of noisy environment.. Throat microphones have advantages and disadvantages, based on their appearance and usage. The biggest advantage is that throat microphones are not affected virtually and insensitive to noise. For example a motorcycle rider while he communicates

the disturbances is generated inside and outside the helmet. The placement of a throat mic and its shape makes it ideal for not only motorcycling, but many other applications[8]. Next the disadvantages include cost, as this type of microphone tends to be more expensive to develop and produce. Figure 2 and 3 depicts the images of throat microphone and normal microphone.



Figure 2 Throat Microphone



Figure 3 Standard microphone

A .FEATURES FOR SPEAKER RECOGNITION

The speech signal is a form that can be represented by a sequence of feature vectors in order to apply mathematical tools without the loss of generality. Most of these features are also used for speaker dependent and speaker independent recognition systems that rely on real life systems.

B.EXTRACTION METHODS

Various methods available for feature extraction are

- a) Mel-Frequency Cepstral Coefficients (MFCC),
- b) Real Cepstral Coefficients (RCC),
- c) Linear Prediction Coding (LPC),
- d) Linear Predictive Cepstral Coefficients (LPCC)
- e) Perceptual Linear Predictive Cepstral Coefficients (PLPC).

a) Mel-Frequency Cepstral Coefficients (MFCC)

MFCC is one of the most popular technique and commonly used in most of the applications of speech signal for feature extraction.[5] It is based on the human peripheral auditory method. According to human perception the frequency contents of sounds does not follow a linear scale. It is mainly used in speaker or speech recognition systems.

b) Real Cepstral Coefficients (RCC)

RCCs are signals that are transformed from the time domain to the frequency domain by applying a Fast Fourier Transform (FFT) to each frame. The results of this logarithm specify inverse Fast Fourier transform (IFFT)[7] and then it is applied to get the real Cepstrum of the signal, and the equation for the result is given below

$$\text{Real Cepstrum} = \text{IFFT}(\log(\text{FFT}(s(n))))$$

c) Linear Prediction Coding (LPC)

This technique analysis the speech signal by estimating the formants. LPC is a form that removes the effects of formants and calculates the intensity and frequency of the remaining buzz from the speech signal[10]. The procedure that removes the formants is called Inverse filtering and the remaining signal is specified as the Residue for the current speech signal. In this LPC method, each sample of the speech signal is formulated as a linear combinational transformation for the previous samples and hence it is called a linear predictor or linear predictive coding.

d) Linear Predictive Cepstral Coefficients (LPCC)

It is also a technique with widely used extracted features from speech signal. In this process it specify LPC parameters that can be effectively used to energy and frequency spectrum. The base of explaining acoustic signals spectrum, modeling [12] and pattern recognition is given by the set of increasing logarithm which restrains the fast change of frequency spectrum, and hence it is more centralized for short-time character. One of the most common short term spectral measurements are used for LPC that are derived from Linear Predictive cepstral coefficients (LPCC) with regression coefficients[11]. LPCC is the process which shows the differences of the biological structure of human vocal tract that is computed through iteration from the LPC Parameters to the LPC Cepstrum.

e) Perceptual Linear Predictive Cepstral Coefficients (PLPC).

This technique is based on the magnitude spectrum of the speech analysis window. Other techniques such as MFCC and LPC are cepstral techniques while this PLPCC is a temporal technique[11]. There are certain steps which are followed to calculate the coefficients of the PLPCC they are described below in steps

Step-1 compute the power spectrum of a windowed speech.

Step-2 For frequency sampling of 8kHz it performs grouping of the results to 23 critical bands using bark scaling.

Step-3 It helps in simulating the power law hearing and carries out loudness, equalization and cube root compression.

Step-4 This technique performs inverse Fast Fourier Transform (IFFT).

Step-5 In this step one is to perform LP analysis by Levinson-Durbin algorithm. And the next step is to convert LP coefficients into cepstral coefficients. And hence the relationship for frequency in Bark and frequency in Hz is given as

$$f(\text{bark}) = 6 * \arcsin(h(f(\text{Hz})/600))$$

III ARCHITECTURAL DIAGRAM

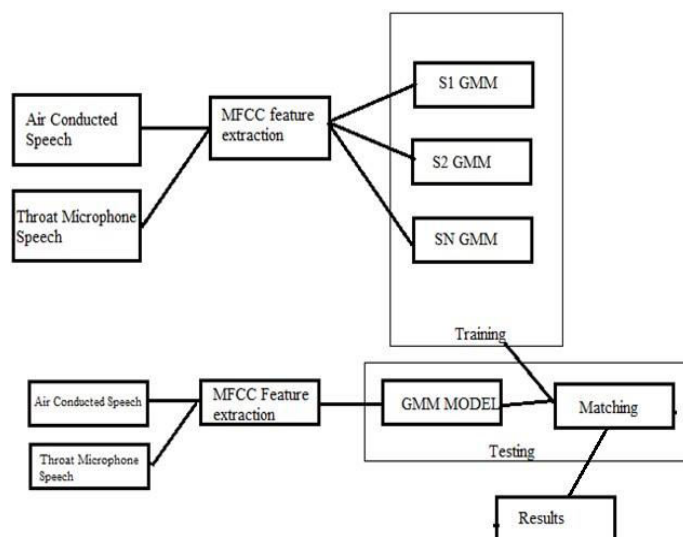


Figure 4 Combination of Air phone and Throat Phone

As such of all other pattern recognition systems this speaker recognition systems also has two phases such as training and testing. Training is the process which specifies familiarity to the system with the voice characteristics of the speakers who are all registered[5]. Testing is the process with actual recognition task. The above block diagram of Figure 4 specifies that in training phase the speech signals from two kinds of phones are depicted and features of those signals are extracted with the help of MFCC and then all the speech signals are modelled using GMM[11]. In testing phase again the speech signals and features for those signals are gathered and extracted accordingly and then the next step of modelling happens then in this step a new formation such as matching the models, and at last results are produced by specifying he is the recognized speaker.

IV GAUSSIAN MIXTURE MODEL

Definition of GMM specifies that it is the density function with probability parameters that are represented as a weighted sum of Gaussian component densities,[4] It is also a form of parametric model for probability distribution with continuous measurements in biometric system. And the estimation is done for training data using the iterative Expectation-Maximization (EM) algorithm or Maximum Posteriori (MAP)

from a well-trained prior model. It is also a form of non-parametric methods for speaker identification .In this feature vectors are provided in d-dimensional feature space for clustering as they are related to Gaussian distribution,[2] which specifies that each corresponding cluster can be seen as Gaussian probability distribution with features belonging to the clusters of probability values. Below Figure 5 depicts the GMM model with its corresponding feature space. The usage of Gaussian mixture density for speaker identification is motivated by two facts. They are:

- Individual Gaussian classes are represented as the set of acoustic classes. These acoustic classes formulate vocal tract information.
- Gaussian mixture densities provide free approximation to distribute all feature vectors in multi-dimensional feature space.

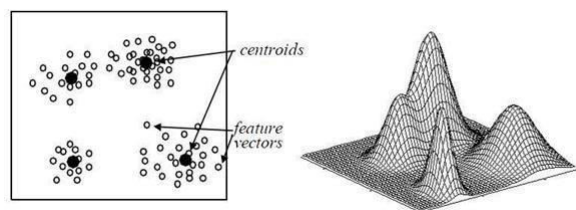


Figure 5 GMM model showing a feature space and corresponding Gaussian model.

V BUILDING SPEAKER RECOGNITION SYSTEM

a) RECORDING SPEECH

A standard air conduction microphone is used to record the speech. The recorded speech is stored as ‘.wav’ format. Various speeches from different speaker are recorded to train the speaker recognition system. The recorded speech is not clean and contains background noise due to microphone recording. The collected speech requires pre processing to make it suitable for the feature extraction to which sample features are generated based on two techniques of pre processing to make the recorded speech to work properly without occurring any error rate to do this speech data has to be normalized .

b) PRE-PROCESSING

It is the process of removing the unwanted or channel error disturbances to make the sample to process without error rate, During pre-processing it helps in lowering of high frequency energies that are not useful for creating GMM model.

1) Frame Blocking

In this process continuous speech signal are taken which is divided into frames of some N samples, with adjacent frames being separated by some M samples with the value M less than that of N. [11] Considering the first frame with N samples and second frame with M samples overlapping takes place by N - M samples. This step continues until all the speech is accumulated for using one or more frames. For example values of M and N are said to be specified as N = 256 and M = 128 respectively. The N's value is taken as 256 it specifies that speech signal is assumed to be periodic. As of frame length 256 which is helpful in proceeding fast implementation of Discrete Fourier Transform and Fast Fourier Transform.

2) Windowing

It is a process of taking a small subset of a larger data set, for processing and analysis. A naive approach, the rectangular window, involves simply truncating the data set before and after the window, while not modifying the contents of the window at all.[10] The next step is to window each individual frame to minimize the signal discontinuities at the beginning and end of each frame. The below Figure 6 shows how Hamming Window is related in terms of co-efficient and gain.

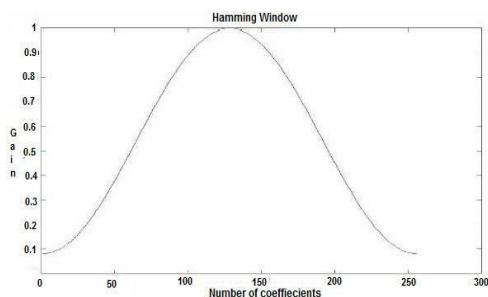


Figure 6 Hamming Window

c) MEL-FREQUENCY CEPSTRAL COEFFICIENT(MFCC)

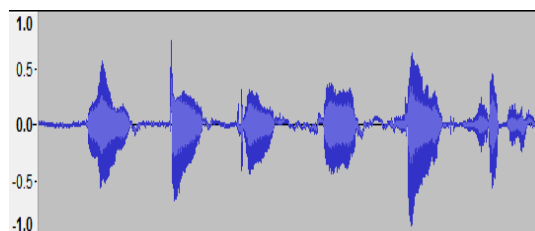
MFCC Calculation is the second feature extraction method to be used, for which it is based on the known variation of the human ear's critical bandwidths with frequency. Filters are spaced linearly at low frequencies and logarithmically at high frequencies that has been used to capture the phonetically important characteristics of speech.[6] This is expressed in mel-frequency scale, which is linear frequency spacing below 1000 Hz and a logarithmic spacing above 1000 Hz.

d) DATA COLLECTION PROCESS

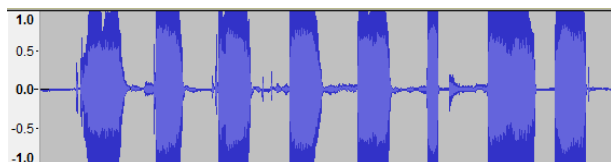
The speech samples data are recorded for processing in the laboratory environment. The speech sample data from each volunteer speaker is collected one by one using both the TM and NM. Before the actual recording to check speaker adaptability, a 10 sec sample recording is performed to check whether the microphones are placed properly.[2] Volunteer speakers are asked to speak in their natural voice and each recording is checked using audacity to see whether there is any background noise in the signal. If the waveform for the speech signal is good then it is saved. Otherwise, the volunteer speaker is asked to re-record the same content.

e) SPEECH DATA SAMPLES

Below Figure 7 shows the waveform of speech data a) shows the five minutes speech recorded data using normal microphone b) shows the five minutes recorded data using throat microphone.



a) Sample air conducted speech



b) Sample throat conducted speech

Figure 7 Speech Sample data collection with two Microphones

VIEVALUATION

A.EVALUATION

There are two steps in which speech signal can be evaluated for the creation of models using GMM.

Step-1 Record and play

A command-line sound file recorder which supports several file formats and ALSA soundcard driver with multiple soundcards and multiple devices . It will record a 10-second WAV file with DAT quality on your available soundcard (hw: 0, 0). DAT quality is defined as stereo digital audio record with a 48 kHz sampling rate and 16-bit resolution and play the recorded sample.

Step-2 Creation of MFCC features

The speech recognition cannot process directly on speech waveforms it needs tools that has to be represented in a more compact and efficient way. This step is called “pre processing”:

- a) The signal is segmented in successive frames , overlapping with each other.
- b) Each frame is multiplied by a Hamming windowing function
- c) A MFCC feature is extracted from each windowed frame.

The conversion of waveform with a series of feature vectors is done with the HCopy of HTK tool:

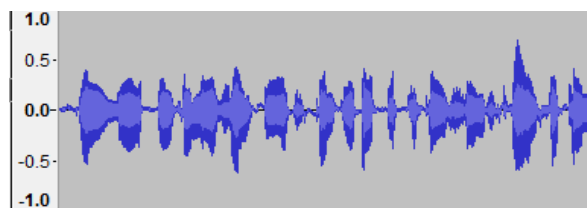
To do this process a configuration file has been set to which the parameters of the acoustical coefficient extraction can be done using analysis.conf.targetlist.txt specifies the name and location of each waveform to process, along with the name and location of the target

coefficient files. And below explains the MFCC configuration file. With such a configuration file, an MFCC feature extraction is performed .For each signal frame, the following coefficients are extracted:

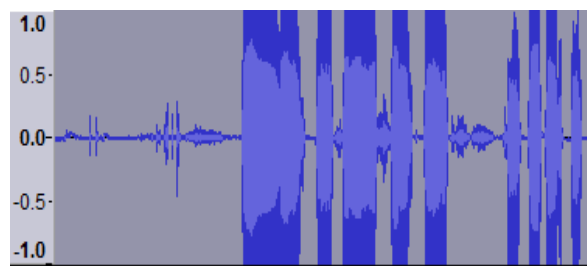
- The 12 first MFCC coefficients $[c_1, \dots, c_{12}]$ (since $NUMCEPS = 12$)
 - The “null” MFCC coefficient c_0 , which is proportional to the total energy in the frame (suffix “_0” in TARGETKIND)
 - 13 “Delta coefficients”, estimating the first order derivative of $[c_0, c_1, \dots, c_{12}]$ (suffix “_D” in TARGETKIND)
 - 13 “Acceleration coefficients”, estimating the second order derivative of $[c_0, c_1, \dots, c_{12}]$ (suffix “_A” in TARGETKIND)
- Altogether, a 39 coefficient feature vector is extracted from each signal frame.

VARIOUS SAMPLES FOR PROCESSING WITH BOTH MICROPHONES

SAMPLE-I

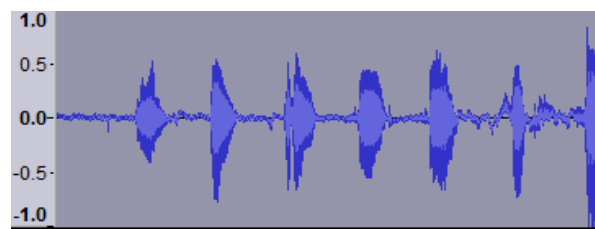


a) Normal Microphone TIMIT Sentence-1

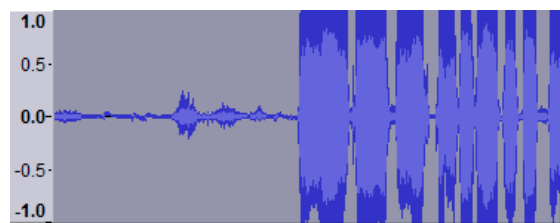


b) Throat Microphone TIMIT Sentence-1

SAMPLE-II

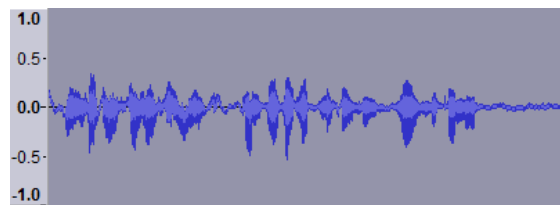


a) Normal Microphone TIMIT Sentence-2

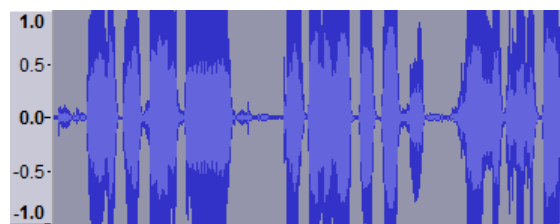


b)Throat Microphone TIMIT Sentence-2

SAMPLE-III



a)Normal Microphone TIMIT Sentence-3



b)Throat Microphone TIMIT Sentence-3

VII CONCLUSION

In this paper, the air conducted speech signals are recorded from different speakers using a standard microphone and also with the another mode of speech recording using a throat microphone. The recorded speech signals are preprocessed in which the speeches that are recorded from air microphone has much background noise while the throat microphone is free from background noise and then both speeches are made suitable for the feature extraction process. From these preprocessed speech signals, mfcc features were successfully extracted for generating the GMM speaker model.

REFERENCES

- [1] Jia-Ching Wang, Yu-Hao Chin, Wen-Chi Hsieh, Chang-Hong Lin, Ying-Ren Chen, Siahaan.E, "Speaker Identification With Whispered Speech for the Access Control System", Automation Science and Engineering, IEEE Transactions , vol 12, no 4, pp. 1191-1199, 2015.
- [2] Kawthar Yasmine Zergat and Abderrahmane Amrouche, "New Scheme based on GMM-PCA-SVM Modeling for Automatic Speaker Recognition", International Journal of Speech Technology, vol 17, no 4, pp. 373-381, 2014.
- [3] Maxim Sidorov, Alexander Schmitt, Sergey Zablotskiy and Wolfgang Minker, "Survey of Automatic Speaker Identification Methods", Proceedings of the Ninth International Conference on Intelligent Environments, pp. 236-239, 2013.
- [4] Cherifa S. and Messaoud R, "New Technique to use the GMM in Speaker Recognition System (SRS)", International Conference on Computer Applications Technology, pp. 1-5, 2013.
- [5] Seiichi Nakagawa, Longbiao Wang and Shinji Ohtsuka, "Speaker Identification and Verification by Combining MFCC and Phase Information" IEEE Transactions on Audio, Speech and Language Processing, vol. 20, no. 4, 2012.
- [6] Rishiraj Mukherjee, Tanmoy Islam and Ravi Sankar, "Text Dependent Speaker Recognition using Shifted MFCC", Proceedings of IEEE Southeast Conference, pp. 1-4, 2012.
- [7] Homayoon Beigi, "Fundamentals of Speaker Recognition" Springer Publications, pp. 75-84, 2011.
- [8] Douglas A. Reynolds, Thomas F. Quatieri, and Robert B. Dunn, "Speaker Verification Using Adapted Gaussian Mixture Models", Digital Signal Processing, vol. 10, no. 1-3, July 2010.

- [9] Tomi Kinnunen and Haizhou Li, "An Overview of Text-Independent Speaker Recognition: From Features to Supervectors", Elsevier Journal of Speech Communication, vol. 52, no 1, pp. 12-40, 2010.
- [10] Marcos Faundez-Zanuy and Enric Monte-Moreno, "State-of-the-Art in Speaker Recognition" IEEE Aerospace and Electronic Systems Magazine, vol. 20, no 5, pp. 7-12, May 2005.
- [11] Joseph P. Campbell Jr., "Speaker Recognition: A Tutorial", Proceedings of the IEEE, vol. 85, no 9, pp. 1437-1462, September 1997.
- [12] Robin King, "New Challenges in Automatic Speech Recognition and Speech Understanding", IEEE TENCON, Conference on Speech and Image Technologies for Computing and Telecommunication, pp. 287-294, 1997.
- [13] Weng, Z., Li, L., & Guo, D, "Speaker recognition using weighted dynamic MFCC based on GMM", Proceedings - 2010 International Conference on Anti-Counterfeiting, Security and Identification, pp.285–288, 2010.
- [14] M. Arun Marx, G.Vinoth, A. Shahina, A. Nayeemulla Khan, "Throat Microphone Speech Corpus for Speaker Recognition",MES Journal of Technology and Management,pp.16-20,2008.