# Learning Assistant in Educational Field Using Automatic Speech Recognition

**Prajakta Kotwal, Prof. M. R. Dixit**

*Abstract-* -Automatic Speech recognition is the translation of spoken words into text. It takes speech data as input and divides it into small time domain frames. Speech signal processing considering speech signals stationary for a small time interval. From point of view speech signals are divided into small units Morphims or Phonims. Any speech data can be sorted as word uttered followed by voice and silence intervals. Voice activity detection can be are employed to detect voiced and unvoiced part of speech. Speech processing consists of speech recognition, speech synthesis, speaker recognition, understanding of speech with reference to context, speech coding, speech enhancement, speech transmission, speech to text conversion & text to speech conversion etc. In general speech to text conversion system will convert input speech data to output text data. If the input speech data is inappropriate with some errors then there is a possibility to get incorrect output data.

The proposed system contains options for correction of inappropriate input data so that the output text and speech data produce and pronounce is correct. The proposed system will be employed as learning assistance in educational field for students to learn correct pronunciation of words. The proposed system will also help tourists for conversation in local language.

Keywords: - HMM (Hidden Markov Model), GMM(Gaussian Mixture Model), ANN(Artificial Neural Network), VAD(Voice Activity Detection)

## I. INTRODUCTION

The task of speech recognition is to recognize input speech. Currently used Speech Recognition software available in market takes speech as input though it is correct in pronunciation or not and proceeds on it to produce output which may not give us 100% accurate result .Obviously its depends on input. So it is essential to develop a system which will check incoming data is correction pronunciation or not. Here proposed system will do this job. Speech recognition systems have a wide range of applications from the relatively simple isolated-word recognition systems. Like any pattern recognition problem,

Manuscript received November 15, 2013

PrajaktaKotwal, Research Candidate, Kolhapur Institute of Technology, Kolhapur, Maharashtra,India.
(e-mail:- kotwalprajakta77@gmail.com)
Prof.M.R.Dixit,Dept.Of Electronics, Kolhapur Institute of Technology,Kolhapur,Maharashtra, India.
(e-mail:- mrdixit@rediffmail.com)

the fundamental problem in speech recognition is the speech pattern variability. In general the sources of speech variability are as follows: Duration variability, Accent, Speaker variability, Noise etc. So proposed system based on speaker independent system and recording of input is taken in noise-free environment only.

## II. THEORETICAL ANALYSIS

### A. MFCC Calculation:

Generally speaking, a conventional automatic speech recognition (ASR) system can be organized in two blocks: the feature extraction and the modeling stage.The feature extraction is usually a non-invertible (lossy) transformation. Making an analogy with filter banks, such transformation does not lead to perfect reconstruction, i.e., given only the features it is not possible to reconstruct the original speech used to generate those features. Computational complexity and robustness are two primary reasons to allow loosing information. Increasing the accuracy of the parametric representation by increasing the number of parameters leads to an increase of complexity and eventually does not lead to a better result due to robustness issues. The greater the number of parameters in a model, the greater should be the training sequence.

Speech is usually segmented in frames of 20 to 30 ms, and the window analysis is shifted by 10 ms. Each frame is converted to 12 MFCCs plus a normalized energy parameter. The first and second derivatives (D's and DD's) of MFCCs and energy are estimated, resulting in 39 numbers representing each frame. Assuming a sample rate of 8 kHz, for each 10 ms the feature extraction module delivers 39 numbers to the modeling stage..

### B. Pitch Calculation:

The human perception of the frequency contents of sound for speech signals does not follow linear scales. Thus for each tone with an actual frequency measured in Hz. A subjective pitch is measured on a scale .A person pitch originates in the vocal cord and the rate at which the vocal folds vibrates is the frequency of the pitch .e.g. when the vocal folds oscillates at 300 times per seconds, they are said to be producing a pitch of 300Hz.

## III. PROPOSED METHODOLOGY

### A. OBJECTIVE OF THE PROJECT :-

a) Applying Voice Activity Detection (VAD) and removal of silence segments.
b) Estimation of Mel – Frequency Cepstrum Coefficient (MFCC) .
c) To develop Hidden Markov Model.
d) To develop Gaussian Mixture Model.
e) To develop artificial neural network.
f) Correct the inappropriate pronounce word and replace it by appropriate pronunciation.

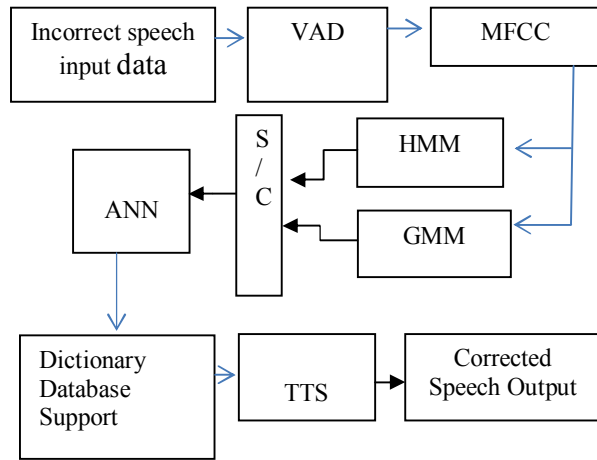### B. METHODOLOGIES OF IMPLEMENTATION :-



**Fig1 :** Methodology of Implementation

Fig 1.- Shows design of a system which convert incorrect speech input data into correct speech output data. The proposed system includes workflow as – recording of speech input data using sound recorder of PC and recording tools. Speech input data can be recorded over appropriate frequency range which can be applied to voice activity detection (VAD) also known as speech activity detection. This technique can be used in speech processing in which presence or absence of human speech is detected. Output speech signal from VAD is use for calculating Mel Frequency.

Mel Frequency ceptrum coefficient is based on human peripheral auditory system. The human perception of the frequency contents for sound of speech signal does not follow a linear scale. Thus for each tone with an actual frequency measured in Hz, a subjective pitch is measured on a scale called melscale. Isolated speech recognition can be performed by hidden markov model and Gaussian mixture model ,from them both will compare and the best output that model will be used to give input data for artificial neural network, which will classify speech segments as voiced, unvoiced, nasal / frative / plosive etc.The available development tools like Colea, Dragon, Natural Reader will be applied for speech to text conversion. The output speech will be represent the appropriate or corrected utterance of the input speech data.

Output signal based on MFCC coefficient.
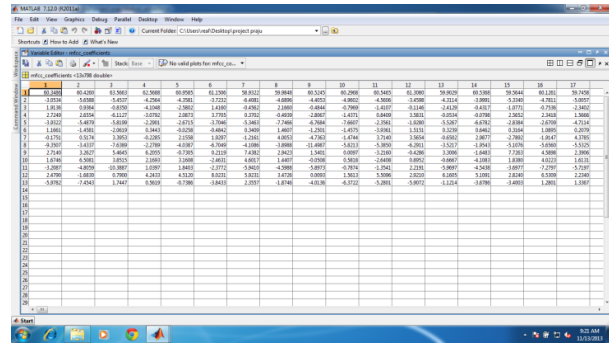
**BASIC RESULS:**
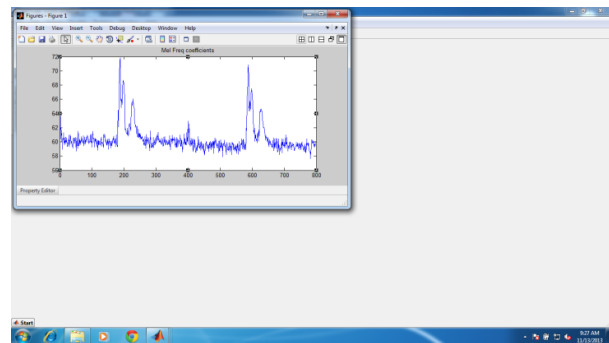


Fig 2 :MFCC Extraction



Fig 3: Signal Plotting

Pitch calculation=253 for first correct input voice.

### IV. CONCLUSION / FUTURE SCOPE

Feature extraction is done by MFCC (Mel Frequency Cepstral Coefficient) which represents audio based on perception of human ear.This result is useful in speech recognition application. The same techniques can be used in different applications. There are lots of techniques likes GMM, HMM, DTW, ANN etc .which can be used as per requirement of application. Dictionary can be modified when it needs to change.

### REFERENCES

[1] Nelson Morgan,"Deep and Wide: Multiple Layers in Automatic Speech Recognition",IEEE Transactions on audio,speech and language processing,vol.20, no.1, January 2012.
[2]ArchanaShende,Subhash Mishra, Shiv Kumar, "Comparison of different parameters used in GMM based automatic speaker recognition", International Journal of Soft Computing and Engineering (IJSCE)

Volume-1, Issue-3, July 2011

[3]Mohamad Adan AL - ALaoui, Lina AL-Kanj, JimmyAzar,and Elias Yaacoub ,"Speech recognition using Artificial Neural Networks and Hidden Markov Model", IEEE multisciplinary engineering education magazine.vol.3, no.3.september 2008.

[4]Chulhee Lee, Donghoon Hyun, Euisun Choi, Jinwook Go, and ChungyongLee, "Optimizing Feature Extraction for Speech Recognition",IEEE transactions on speech and audio processing, vol. 11, no. 1, january 2003.

[5]Harry Printz and Isabel Trancoso ,"Editorail", IEEE transactions on speech and audio processing, vol. 10, no. 8, november 2002.

[6] Alexandros Potamianos, Member, IEEE, and Petros Maragos, "Time-Frequency Distributions for Automatic Speech Recognition", IEEE transactions on speech and audio processing, vol. 9, no. 3, march 2001.

[7] VibhaTiwari,"MFCC and its applications in speaker recognition", International Journal on Emerging Technologies 1(1): 19-22(2010).

[8]D.B.Paul,"Speech Recognition using Hidden Markov Model",The Lincoln Laboratory Journal vol.3, no.1,1990.

[9]Lawrence R.Rabiner, "A-Tutorial on Hidden Markov Models and selected applications in speech recognition",vol.77,no.2, Feb 1999.

[10]Yang Liu,"Enriching Speech Recognition with Automatic Detection of Sentence Boundaries and Disfluencies"IEEE Transactions on audio, speech and language processing, vol. 14, Sept 2006