



RELATIVE SPECTRAL ALGORITHM BASED VOICE RECOGNITION TECHNIQUES

Vijaya Babu Kuchipudi¹, Dr.HARSH PRATAP SINGH², Dr. Laxmaiah Mettu³

¹Research Scholar, Dept. of Computer Science and Engineering

Sri Satya Sai University of Technology and Medical Sciences,

Sehore Bhopal-Indore Road, Madhya Pradesh, India.

²Research Guide, Dept. of Computer Science and Engineering

Sri Satya Sai University of Technology and Medical Sciences,

Sehore Bhopal-Indore Road, Madhya Pradesh, India.

³Research Co-Guide, HOD. Dept. of Computer Science and Engineering

CMR Engineering College, Kandlakoya (V), Medchal, Hyderabad

3714

ABSTRACT

To find the trouble of secret phrase the executives and improve the convenience of authentication systems, biometric authentication has been broadly considered and has pulled in unique consideration in both scholarly world and industry. Numerous biometric authentication systems have been explored and grown, particularly for cell phones. The Voice is a signal of limitless data. Digital processing of speech signal is vital for rapid and exact programmed voice recognition technology. These days it is being utilized for medical care, communication military and individuals with handicaps in this manner the digital signal cycles, for example, Feature Extraction and Feature Matching are the most recent issues for investigation of voice signal. To remove important data from the speech signal, settle on choices on the cycle, and get results, the information should be controlled and examined. Fundamental technique utilized for removing the features of the voice signal is to discover the Mel frequency cepstral coefficients. Mel-frequency cepstral coefficients (MFCCs) are the coefficients that aggregately address the transient force range of a sound, in view of a linear cosine transform of a log power range on a nonlinear Mel size of frequency. This paper is



separated into two modules. Under the principal module feature of the speech signal are extricated as MFCC coefficients and in another module the non-linear grouping arrangement known as Dynamic Time Warping (DTW) presented by Sakoe Chiba has been utilized as features matching methods. Since clearly the voice signal will in general have distinctive transient rate, the arrangement is imperative to create the better presentation.

KEYWORDS: biometric authentication, password management, privacy protection, Mel frequency cepstral coefficient (MFCC), voice recognition.

1. INTRODUCTION

With the fast improvement of the Internet and cell phones, authentication systems have been broadly utilized in the Internet administration access and cell phone access for ensuring client gadgets, substance, and records. At the point when clients hold an ever increasing number of records, secret phrase the executives is getting genuinely troublesome practically speaking since it is ordinarily difficult to recollect different passwords for various framework gets to, particularly those with high security levels. To tackle this issue, biometrics were examined and applied in individual authentication because of their interesting qualities. Researchers have led broad and inside and out research on biometric authentication as of late. A few researchers zeroed in on explicit algorithms or structures utilized in biometric-based authentication. Biometric recognition techniques dependent on neural organizations by utilizing voice, iris, fingerprint, palm-print, face, and brought up expected approaches to improve these strategies. Multimodal biometric strategies are significantly more solid for developing a more secure authentication framework. They talked about such multimodal techniques as various sensors, numerous algorithms, different instances, different samples and half and half models.



Figure 1.1 Voice Recognition Systems

Voice or Speech is a characteristic method of correspondence for individuals; however sometimes it does not work for example impaired individual. In the course of the most recent multi decade, here is a need to empower human to speak with machines without playing out any content info. Here speech recognition is a technology can make competent crippled people to speak with machines with their inabilities. Speech recognition is a framework that utilized by the human to tune in, distinguish and comprehend what does the client need by talking. It is a transformation of speech to message in a framework. Speech recognition is the machine on the explanation or order of human speech to distinguish and comprehend and respond in like manner. It depends on the voice as the research object, it permits the machine to consequently distinguish and comprehend human communicated in language through speech signal processing and example recognition. The speech recognition technology is the cutting edge that permits the machine to transform the voice signal into the fitting content or order through the way toward recognizing and comprehension. Speech recognition is a cross-disciplinary and includes a wide reach. It has a cozy relationship with acoustics, phonetics, etymology, data hypothesis, and example recognition hypothesis and neurobiology disciplines. With the fast improvement of PC equipment and programming and data technology, speech recognition technology is slowly turning into a vital technology in the PC data processing technology. The objective in programmed speech recognition is to give a way to verbal human-to-machine correspondence. Albeit both speech coding and recognition include examination of the speech wave, the voice recognition issue is by a wide margin more troublesome. Utilizations of speech recognition advancements incorporate cycle mechanization, phone request, programmed banking, and secure voice access, to give some examples. Albeit the research in speech recognition is distribution driven. There is a huge need between research and business sending. Understanding speech requires the joining of various extraordinary and complex cycles, for example, signal processing (recognition of phonemes, syllables and words), syntactic parsing and semantic investigation. Language is a framework that empowers a speaker to utilize words that are typically gotten the hang of during adolescence. The qualities of speech sounds rely upon the specific human language or tongue. Essentially speech recognition is an example recognition issue. Speech Recognition Systems are largely delegated discrete or persistent systems that are speaker reliant, free or versatile. A speaker-subordinate framework necessitates that the client record an illustration of the word, sentence or expression before its being perceived by the framework for example the client prepares the framework. Some speaker-subordinate systems require just that the client record a subset of framework jargon to make the whole jargon unmistakable. A speaker-free framework doesn't need any chronicle before framework use. It is created to work for any speaker of a specific sort.



1.1. Basic Principle of Voice or Speech Recognition

The speech recognition framework is an example recognition framework, including feature extraction, design matching and the reference model library. The obscure voice through the mouthpiece is transformed into an electrical signal on the contribution of the distinguishing proof framework, the first after the pre-treatment. The framework sets up a voice model as per the human voice attributes dissects the information voice signal and concentrates the necessary features on this premise; it builds up the necessary layout of the speech recognition. The PC is utilized in the recognition interaction as indicated by the model of the speech recognition to think about the voice layout put away in the PC and the attributes of the info voice signal. Search and matching systems to recognize the ideal scope of the info voice coordinates the layout.

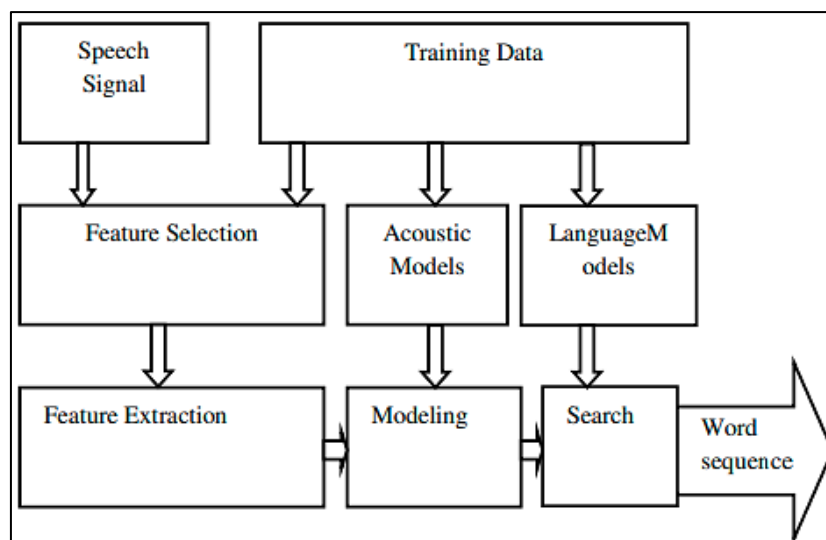


Figure 1.2 Principle of Speech Recognition

According to the definition of this template through the lookup table can be given the recognition results of the computer.

2. LITERATURE REVIEW

M.A.Anusuya et al (2019), This paper presents a short overview on Automatic Speech Recognition and talks about the significant topics and advances made in the previous 60 years of research, to give a mechanical viewpoint and an enthusiasm for the essential advancement that has been cultivated in this significant zone of speech correspondence. Following quite a while of research and improvement, the exactness of programmed speech recognition stays one of the significant research difficulties. The plan of Speech Recognition framework requires cautious considerations to the accompanying issues: Definition of different sorts of speech classes, speech portrayal, feature

extraction methods, speech classifiers, and information base execution assessment. The issues that are existing in ASR and the different methods to take care of these issues developed by different research labourers have been introduced in a sequential request. Thus creators trust that this work will be a commitment nearby speech recognition. The target of this survey paper is to sum up, think about a portion of the notable strategies utilized in different phases of speech recognition framework, and recognize research point and applications, which are at the bleeding edge of this energizing and testing field.

Manoj Kumar Sharma et al (2012), Speech is a characteristic mode to cooperate with others. With speech, we can communicate our words to other people. Speech recognition is a way or technology where the assertions or orders of human speech to comprehend and respond as needs be. Speech recognition permits machining framework to transform the approaching speech signals into orders through the way toward distinguishing and comprehension. It additionally makes the regular voice correspondence work. Primary Goal of speech recognition is to accomplish better language correspondence among man and machine. So it is an extraordinary technology of human machine interface. The paper depicts the speech recognition technology advancement is all fundamental standards, strategies and arrangement of this technology. Precision of various techniques for speech technology is furnished to figure out strategies with their exhibition viewpoint.

3. RESEARCH METHODOLOGY

With the quickly developing ubiquity and usefulness of voice-driven IoT gadgets, the capability of voice-based assaults turns into a non-unimportant security hazard. As talked about, an assault may prompt extreme misfortunes, e.g., a robber could go into a house by deceiving a voice-based keen lock or an assailant could make unapproved buys and MasterCard charges utilizing a voice-based framework. Such assaults can be exceptionally basic and regularly troublesome or even difficult to distinguish by people and voice assaults can be covered up by different sounds or installed into sound and video accounts. Further, it is likewise extremely simple to scale up such assaults, e.g., a concealed vindictive sound example in a YouTube video could all the while target a great many gadgets. Although the usage of existing assault strategies might be altogether different, their objectives are the equivalent: producing a signal that drives a voice controlled framework to execute a particular malevolent order that the client can't distinguish or perceive. In the accompanying areas, we initially present delegate cutting edge assault approaches as per the kind of usage. We at that point further examine the positives and negatives of each approach and how they identify with one another. The aggressor execution examined in this segment is assessed and revealed by the first



distribution. Because of fast changes of cloud-based systems, the assailant execution is likewise prone to change over the long run.

3.1. Attack Classification Based On Implementation

1) Basic Voice Replay Attack: It is generally realized that voice controlled systems are helpless against voice replay assaults, i.e., an aggressor can replay a formerly recorded voice to cause a framework to play out a particular activity as shown beforehand with the mainstream Amazon Alexa technology . An inadequacy of the fundamental voice replay assault is that it is not difficult to distinguish and in this manner has a restricted useful effect. In any case, as shown later in this part, voice replay assaults are the premise of other further developed and risky assaults.

2) Operating System Level Attack: Contrasted with essential voice replay assaults, a working framework (OS) level assault misuses weaknesses of the OS to make the assault self-set off and more impalpable. Agent instances of this are the assault, GVS-Attack, and the methodology introduced. In the creators, propose a malware that gathers a client's voice and afterward plays out a self-replay assault as a foundation administration. In, the creators further confirm that the underlying mouthpiece and speaker can be utilized at the same time and that the utilization of the speaker does not need client authorization on Android gadgets. They exploit this and propose a zero-consent malware, which persistently dissects the climate and behaviours the assault once it finds that no client is close by. The assault utilizes the gadget's implicit speaker to replay a recorded or engineered speech, which is then acknowledged as a real order. This self-set off assault is subsequently more hazardous and down to earth. While the client can in any case distinguish this assault, the creators call attention to that if the malware has high authorizations; it is even feasible for it to import a sound document to the amplifier without playing it, which can make the assault totally quiet. In, the creators dissect the authorization weakness to the voice assault in detail and propose a way to deal with sidestep the consent the executives of the Android framework. The creators likewise locate that some vindictive activities require a numerous progression order and further propose an intelligent assault that can execute further developed orders.

3) Hardware Level Attack: An equipment level assault replays an engineered non-speech simple signal rather than human voice. The simple signal is painstakingly planned by the qualities of the equipment (e.g., the simple digital converter). The signal is indistinct, yet can be changed over into a genuine digital speech signal by the equipment. Agent approaches are the Dolphin assault and the IEMI assault. In, the creators use the non-linearity of a Micro Electro Mechanical Systems (MEMS) amplifier over ultrasounds and effectively produce indiscernible ultrasound signals that can be



acknowledged as authentic objective orders. Producing such ultrasound signals requires an extraordinary gadget that incorporates a regulator, a speaker, and a ultrasonic transducer. The longest assault distance is 175cm. In, the creators exploit the way that a wired amplifier able earphone can be utilized as a mouthpiece and a FM radio wire all the while and exhibit that it is conceivable to trigger voice orders distantly by discharging a painstakingly planned indistinct AM-balanced signal. This assault is just powerful when the wired earphone is connected to the gadget. A limit of equipment level assaults is that creating the assault signal requires uncommon gadgets, and furthermore a few preconditions should be met. While the engineered signal is unintelligible to the client, the client may in any case see the signal generator.

3.2. RASTA (Relative Spectral Algorithm)

RASTA or Relative Spectral Algorithm as it is known is a strategy that is created as the underlying stage for voice recognition. This technique works by applying a band pass channel to the energy in every frequency sub-band to streamline momentary commotion varieties and to eliminate any steady balance. In voice signals, fixed commotions are regularly identified. Fixed clamours are commotions that are available for the full time of a specific signal and does not have lessening feature. Their property does not change over the long run. The expectation that should be made is that the commotion shifts gradually regarding speech. This makes the RASTA an ideal instrument to be remembered for the underlying phases of voice signal sifting to eliminate fixed clamours. The fixed commotions that are recognized are clamours in the frequency scope of 1Hz - 100Hz.

3720

3.3. Formant Estimation

Formant is one of the significant parts of speech. The frequencies at which the resounding pinnacles happen are known as the formant frequencies or just formants. The formant of the signal can be gotten by dissecting the vocal plot frequency reaction. Figure 1.3 shows the vocal plot frequency reaction. The x-hub addresses the frequency scale and the y-pivot addresses the size of the signal. As it very well may be seen, the formants of the signals are named F1, F2, F3 and F4.



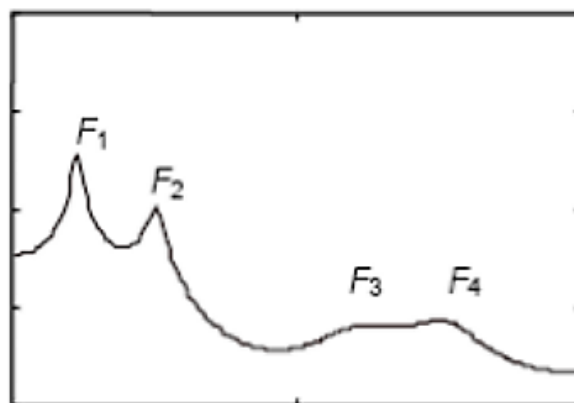


Figure 1.3 Formant Estimation

To acquire the formant of the voice signals, the LPC (Linear Predictive Coding) strategy is utilized. The LPC (Linear Predictive Coding) strategy is gotten from the word linear expectation. Linear forecast as the term infers is a sort of numerical activity. This numerical capacity that is utilized in discrete time signal gauges the future qualities dependent on a linear capacity of past samples.

3.4. RASTA-LPC and DWT Implementation

To execute the framework, a specific philosophy is actualized by breaking down the voice signal to its estimate and detail. From the guess and detail coefficients that are separated, the approach is actualized to complete the recognition interaction. The proposed strategy for the recognition stage is the factual count. Four distinct sorts of measurable estimations are completed on the coefficients. The factual figuring's that are completed are mean, standard deviation, change and mean of total deviation. The wavelet that is utilized for the framework is the wavelet as that this wavelet has an exceptionally close connection with the voice signal. This is resolved through various preliminary and blunders. The coefficients that are extricated from the wavelet deterioration measure is the second level coefficients as the level two coefficients contain the majority of the related information of the voice signal. The information at more elevated levels contains almost no measure of information considering it unusable for the recognition stage. Henceforth for introductory framework execution, the level two coefficients are utilized.

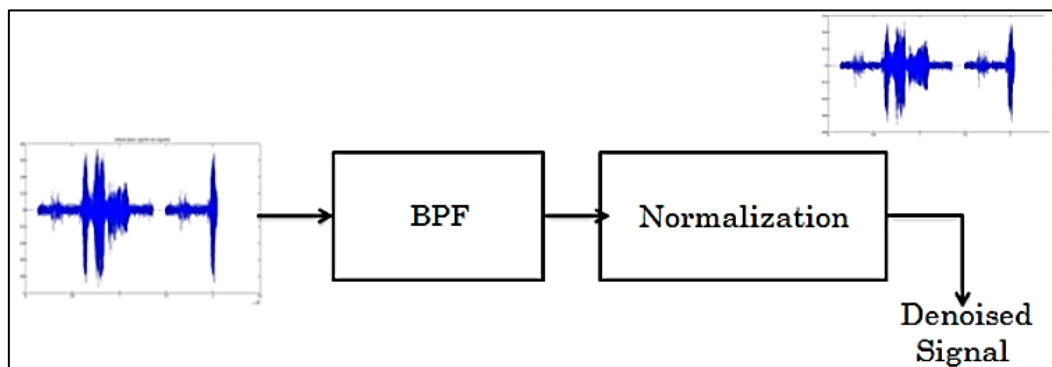


Figure 1.4 Block Diagram of Rasta Process

The coefficients are further edge to eliminate the low connection esteems, and utilizing this coefficients factual calculation is completed. The measurable calculation of the coefficients is utilized in examination of voice signal along with the formant assessment and the wavelet energy. All the separated data acts like a 'fingerprint' for the voice signals. The level of check is determined by contrasting the current qualities signal qualities against the enlisted voice signal qualities.

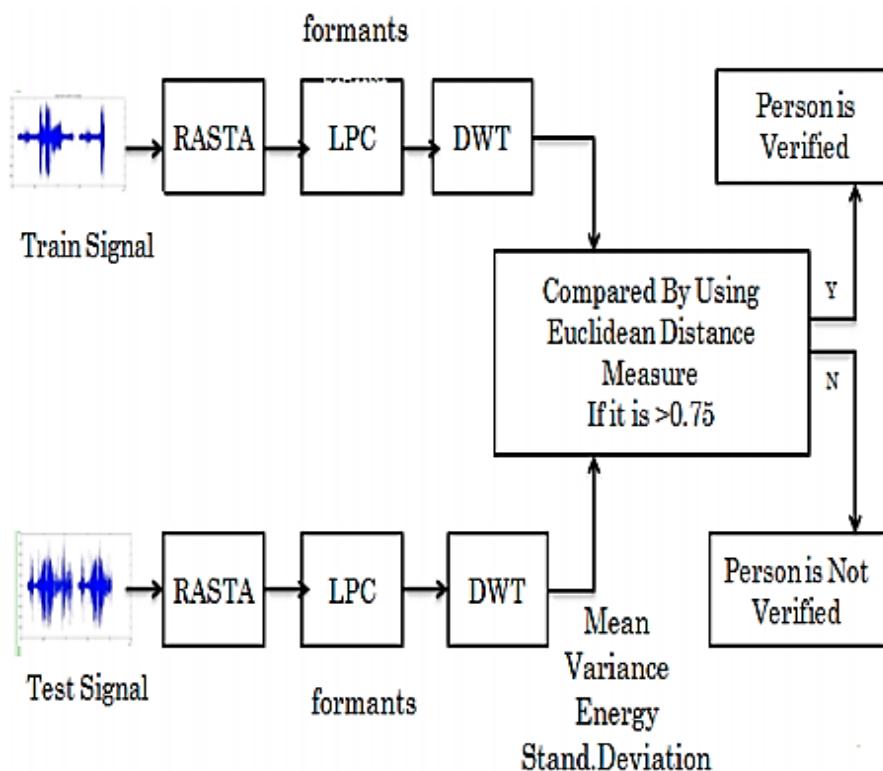


Figure 1.5 Block Diagram of Proposed Text Dependent Speaker Identification System

Verification % = (Test value / Registered value) x100. Between the tested and registered value, whichever value is higher is taken as the denominator and the lower value is taken as the



numerator. Figure 9 shows the complete flowchart, which includes all the important system components that are used in the voice verification program.

4. SIMULATION RESULTS

In this section, experimental results have been shown for various voice test signals with LPC and proposed algorithms. All the experiments have been done in MATLAB 2011a version with 4GB RAM and i3 processor for speed specifications.

$$\text{Mean} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\text{Std. deviation} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \mu)^2}$$

$$\text{Variance} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \mu)^2$$

$$\text{Energy} = T \sum_{i=1}^n x^2(i)$$

Figure 1.6 has shown the performance comparison of proposed and LPC in terms of recognition accuracy with statistical parameters. Finally, LPC achieved 66.66% accuracy where the proposed algorithm achieved almost 90% accuracy.

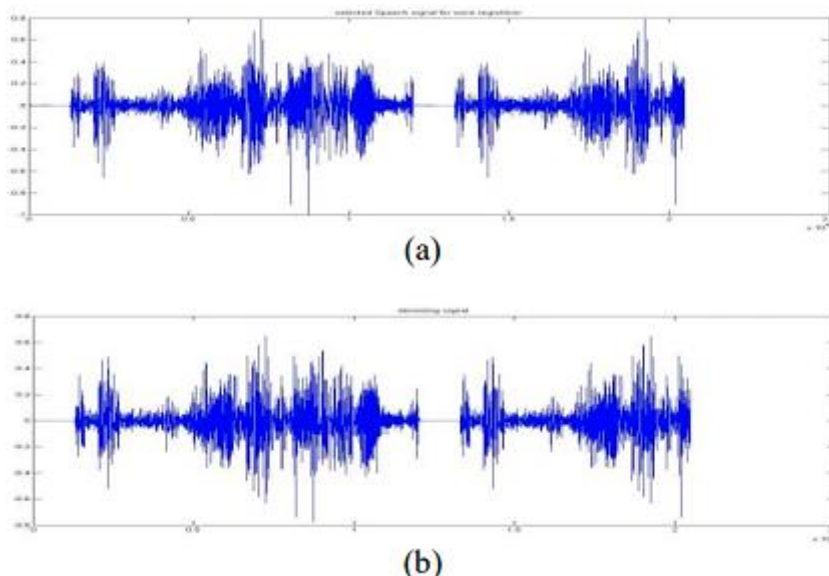


Figure 1.6 (a) Original voice signal (b) De-noised signal for training

Text dependant Speaker Recognition system used to check the character of an individual dependent on their own speech signal utilizing the factual calculation, formant assessment and wavelet energy.



By utilizing the fifty preloaded voice signals from six people, the check tests have been conveyed and an exactness pace of roughly 90 % has been accomplished by proposed algorithm where the LPC has accomplished just 66.66%. By noticing the reproduction results on different speech, signals with various speakers we can infer that the proposed algorithm precision has been improved when contrasted with LPC.

CONCLUSION

In this paper, given a survey of Speech recognition. The territory of Speech recognition is consistently changing and improving. Speech recognition technology is skilled to make conceivable to speak with handicapped people. It makes control of digital system. In future, tremendous prospects to improve the territory of speech recognition technology. By improving of speech, recognition can offer better types of assistance for handicap people. Speech recognition can give a safe climate to our system by voice authentication. Various strategies and their exactness additionally organized that shows the utilization of HMM and ANN model is a lot more extensive utilized techniques for persistent speech recognition measure. Later on, the rightness of speech recognition and the nature of speech will be more improve that is makes correspondence so natural and dependable for everyone including impair people. Future systems should be more effective and fit contrast with conventional systems. **Future scope:** The universe of Speech recognition is quickly changing and developing. Early uses of technology have made fluctuating levels of progress. The guarantee for what's to come is essentially better for pretty much every speech recognition technology zone, with more heartiness to speakers, foundation clamour and so on This will eventually prompt solid, vigorous voice interfaces to each media transmission administration that is offered, accordingly making them generally accessible.

3724

REFERENCES

- [1] Preeti Saini, Parneet Kaur, "Automatic Speech Recognition: a review," International Journal of Engineering Trends and Technology- Volume issue - 2013
- [2] Santosh K.Gaikwad, Bharti W.Gawali, "A review on speech recognition technique," International Journal of computer applications (0975 – 8887) Volume 10– no.3, November 2010
- [3] Manoj Kumar Sharma, Omendri kumari, "Speech Recognition: A Review" 2012
- [4] Santosh K.Gaikwad, Bharti W.Gwali and Pravin Yannawar, "A Review on Speech Recognition Technique", International Journal of Computer Applications (0975 – 8887) Volume 10– No.3, November 2010



- [5] Wiqas Ghai, Navdeep Singh, "Literature review on automatic speech recognition," international journal of computer applications (0975 – 8887) volume 41– no.8, march 2012
- [6] Lin-shan Lee and Yi-cheng Pan, "Voice-based Information Retrieval How far are we from the text based information retrieval", IEEE ASRU 2009.
- [7] H. S. Bhatt, S. Bharadwaj, R. Singh, and M. Vatsa, "Recognizing surgically altered face images using multiobjective evolutionary algorithm," IEEE Trans. Inf. Forensics Security, vol. 8, no. 1, pp. 89–100, Jan. 2013.
- [8] M.A.Anusuya, S.k.katti "speech recognition by machine: A review" international journal of computer science and Information security 2009
- [9] G. Riccardi and d. Hakkani-tür, "active and unsupervised learning for Automatic speech recognition," in proc. Euro speech, 2003.
- [10] Ranu Dixit, Navdeep Kaur, "Speech Recognition Using Stochastic Approach: A Review," International Journal of Innovative Research in Science, Engineering and Technology Vol. 2, Issue 2, February 2013

