

FREQUENCY ANALYSIS OF SPEECH SIGNALS FOR DEVANAGARI SCRIPT USING FFT

UMESH KUMAR GUPTA¹ & R. K. PRASAD²

¹M-Tech Student, Department of Electronics Engineering, B.V.D.U. College of Engineering, Pune, India

²Department of Electronics Engineering, B.V.D.U. College of Engineering, Pune, India

ABSTRACT

This paper aims to discuss the implementation of an isolated word Automatic Speech Recognition system (ASR) for an Indian regional language Devnagari script (HINDI). Devnagari vowels are playing the vital role in pronunciation of any word. Each vowel is classified as starting, middle and end according to the duration of occurrences in the word. The Devnagari script having 12-vowels and 34-consonants are used in some Indian language like Hindi. Sound samples from multiple speakers were utilized to extract different features. Initial processing of data, i.e., normalizing and time-slicing was done using a combination of Simulink and MATLAB. Afterwards, the same tools were used for calculation of Fourier descriptions and correlations. The correlation allowed comparison of the same words.

So the frequency has been calculated in statistical manner and generates a table between amplitude and frequencies. Mean and standard deviation such a system can be potentially utilized in implementation of a voice-driven help setup at call centres of commercial organizations operating in India and other foreign region. The implementation, experiments and result discussions are also existence. The paper also describes the role of each HTK tool, used in various phases of system development, by presenting a detailed architecture of an ASR system developed using HTK library modules and tools

KEYWORDS: Correlation, Feature Extraction, Fourier Descriptors and Spoken Hindi Words

INTRODUCTION

Fundamental frequency estimation has been a popular topic in many fields of research. Such as speech synthesis, speech processing, speaker identification etc. The Devnagari vowels and numerals cannot pronounce two ways but it can be pronounced only one way e.g. Devnagari 12-vowels are classified with the phonetic transcription structure of phonemes according to organ used in produce the sound.

Devnagari is based on phonetics principles which are considered as Place of articulation (POA) vowels. These Devnagari vowels having Frequency analysis of speech signals are estimated in noisy environment (original signals) for analysis and synthesis. The original speech signals are unbalanced to adjustment of an interval with help of some feature extraction techniques or use Sound Forge 9.0 software. The initial objective is to estimating the pitch of Devnagari vowels and numerals with noisy environments speech signals. When one looks at a person, car or house, one's brain tries to match the incoming plot with hundreds (or thousands) of plot that are already stored in memory.

In the speech recognition research literature, no work has been reported on Devnagari speech processing and numerals. So we consider our work to be the first such attempt in this direction. The process involves extraction of some distinct characteristics of individual words by utilizing Fourier transforms and their correlations. The system is speaker-independent and is moderately tolerant to background noise.

DEVNAGARI VOWELS

The 12-Devnagari vowels are categorised as per IPA (International Phonetics Association) as shown in Table-2. These are used for the speech analysis and synthesis purpose. It describes in different categories such as follows:

Short Vowels

The short vowel is a single vowel (V) in a short word or syllable, that vowel usually makes a short sound. These short vowels usually appear at the beginning of the word or between two consonants.

E.g. the short vowels represent character in Marathi and in Hindi.

Long Vowels

The long vowels a short word or syllable ends with a vowel-consonant (VC). The 'a' at the end of the word is silent. Long vowels when the word or syllable has a single vowel and the vowel appears at the end of the word or syllable, the vowels usually represent makes the long sound in Hindi.

Conjunct Vowels

The conjunct vowels are combination of short and long vowels. These phonemes are produced in Hindi e.g. as shown in Table-2.

Nasal Vowel

A nasal vowel is produced with a low tune so that air pressure through nose as well as mouth. The term "nasal" is slightly air pressure which does not come exclusively out of the nose in nasal vowels.

Visarg Vowel

The Visarg symbol is used rarely in Devnagari. The visarg is pronounced as the voiceless sound after the vowels. E.g.in Hindi.

Table 1: Range of Human Speech

Gender	Fundamental Frequency (F0) Min Hz	Fundamental Frequency (F0) Max Hz
Male	80	200
Female	150	350

Table 2: Devnagari Vowels Classified into Five Types

Type of Devnagari Vowels	1	2	3	4
SHORT	अ	इ	उ	-
LONG	आ	ई	ऊ	-
CONJUN-CT	अ+इ=ए	अ+ई=ऐ	अ+उ=औ	अ+उ=औ
NASAL	अं	-	-	-
VISARG	अः	-	-	-

Table 3: Hindi Character Set

Vowels	अ आ इ ई उ ऊ ऋ ए ऐ ओ औ व अः
	a ā i ī u ū r e ai o au aā ah
Gutturals (कवर्ग)	क ख ग घ ङ
	ka kha ga gha ŋa
Palatals (चवर्ग)	च छ ज झ ञ
	ca cha ja jha ña
Cerebrals (टवर्ग)	ट ठ ड ढ ण
	ṭa ṭha ḍa ḍha ṇa
Dentals (तवर्ग)	त थ द ध न
	ta tha da dha na
Labials (पवर्ग)	प फ ब भ म
	pa pha ba bha ma
Semi-Vowels	य र ल व
	ya ra la va
Sibilants	श ष स
	sa pa sa
Aspirate	ह
	Ha

SPEECH PRODUCTION

The theoretical section pretends to give an essential background about the speech analysis involved in recognition tasks, in order to understand the basic principles in which the procedures and implementations carried out during this Master Thesis are based on the theoretical section is divided into three sections.

In the first one, the speech signal and its characteristics are described; the second one is an introduction to frontend analysis for automatic speech recognition, where the important feature vectors of speech signal are explained; and the third is an approach of distance measures based on spectral measures for speech processing.

Essential Features of the Human Vocal Tract

Figure 3 portrays a medium section of the speech system in which we view the anatomy midway through the upper torso as we look on from the right side. The gross components of the system are the lungs, trachea (windpipe), larynx (organ of speech production), pharyngeal cavity (throat), oral or buccal cavity (mouth), and nasal cavity (nose). In technical discussions, the pharyngeal and oral cavities are usually grouped into one unit referred to as the vocal tract, and the nasal cavity is often called the nasal tract. Accordingly, the vocal tract begins at the output of the larynx (vocal cords, or glottis) and terminates at the input to the lips.

The nasal tract begins at the velum and ends at the nostrils. When the velum (a trapdoor-like mechanism at the back of the oral cavity) is lowered, the nasal tract is acoustically coupled to the vocal tract to produce the nasal sounds of speech. Air enters the lungs via the normal breathing mechanism. As air is expelled from the lungs through the trachea, the tensed vocal cords within the larynx are caused to vibrate by the airflow.

The airflow is chopped into quasi-periodic pulses, which are then modulated in frequency in passing through the throat, the oral cavity, and possibly nasal cavity. Depending on the positions of the various articulators (i.e., jaw, tongue, velum, lips, mouth), different sounds are produced.

Figure –3 Schematic view of human speech production mechanism a simplified representation of the complete physiological mechanism for creating speech is shown in Figure 3. The lungs and the associated muscles act as the source of air for exciting the vocal mechanism. The muscle force pushes air out of the lungs (shown schematically as a piston pushing up within a cylinder) and through the trachea.

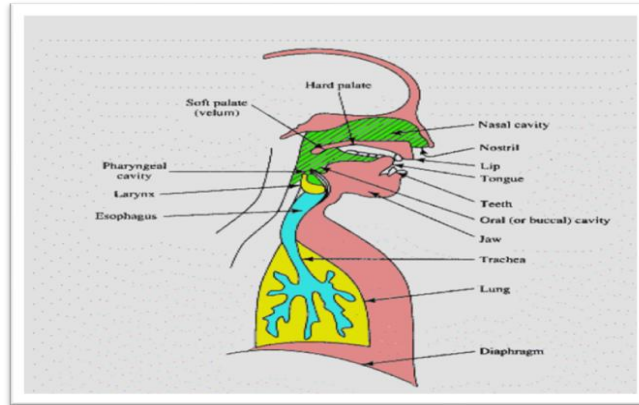


Figure 1: Schematic View of Human Speech Production Mechanism

When the vocal cords are tensed, the airflow causes them to vibrate, producing so-called voiced speech sounds. When the vocal cords are relaxed, in order to produce a sound, the air flow either must pass through a constriction in the vocal tract and thereby become turbulent, producing so-called unvoiced sounds, or it can build up pressure behind a point of the total closure within the vocal tract, and when the closure is opened, the pressure is suddenly and abruptly release, causing a brief transient sound.

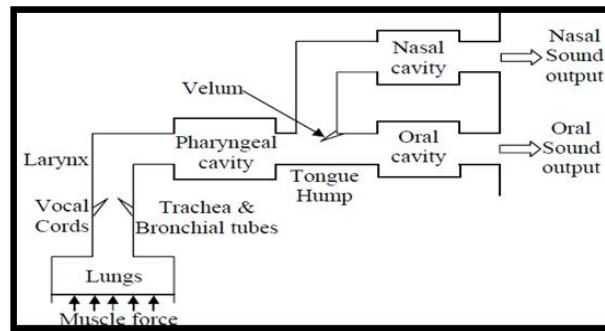


Figure 2: Mechanism for Creating Speech

SPEECH MODELLING USING AVERAGE ENERGY IN THE ZEROCROSSING INTERVAL

The speech production model suggests that the energy of the voiced speech is concentrated about 8 kHz, where as in the case of unvoiced speech, most of the energy is found at higher Frequencies. Since high frequency implies high zerocrossing rate and low frequency implies low zerocrossing rate, there is strong correlation between zerocrossing rate and energy distribution with frequency. This motivates us to model the speech signal using average energy in zerocrossing interval of the signal. Consider the speech segment shown in Figure 2. The ZC_i^k shows the i th zerocrossing and ZC_{i+1}^k shows the $i+1$ th zerocrossing of k th observation window. The time interval between these two points is called i th zerocrossing interval T_i^k in the k th observation window.

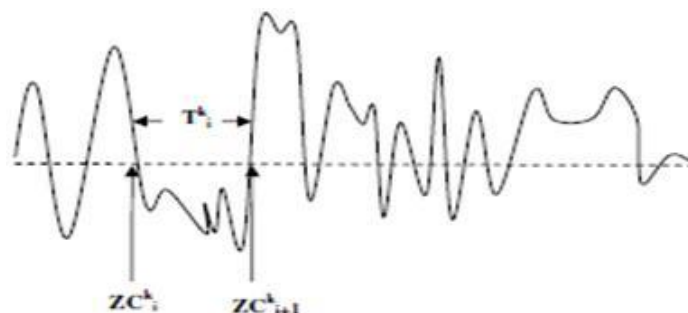


Figure 3: Speech Segment in Kth Observation

The average energy in the i th zerocrossing interval can be obtained by the expression:-

$$E_i^k = \frac{1}{T_t^k} \int_{z_c^k}^{z_c^{k+1}} X^2(t) dt$$

E_i^k Is the average energy of the signal in T_t^k th zerocrossing interval of k th observation window and $X(t)$ is the instantaneous signal amplitude. The aim of the present study is to find a robust coefficient for speech recognition application using the average energy in the zerocrossing interval (AEZI). An XY plot is generated by plotting index number of zero crossing intervals along X axis and Average Energy in the Zerocrossing Interval (AEZI) along Y axis. Figure 4 represents the average energy in the zerocrossing interval vs index number of the zerocrossing interval for the Hindi script.

DATA ACQUISITION AND PROCESSING

One of the obvious methods of speech data acquisition is to have a person speak into an audio device such as microphone or telephone. This act of speaking produces a sound pressure wave that forms an acoustic signal. The microphone or telephone receives the acoustic signal and converts it into an analog signal that can be understood by an electronic system. Finally, in order to store the analog signal on a computer, it must be converted to a digital signal.

The data in this paper is acquired by speaking Hindi Word and numeral into a microphone connected to Windows-7 based PC. The data is saved into '.wav' format files by the using of MATLAB. The sound files are processed after passing through a (Simulink) filter, and are saved for further analysis such as FFT. We recorded the data form speakers who spoke the same word set, i.e. Devnagari Script & numerals.

In general, the digitized speech waveform has a high dynamic range, and can suffer from additive noise. So first, a Simulink model was used to extract and analyze the acquired data; see Figure 1.

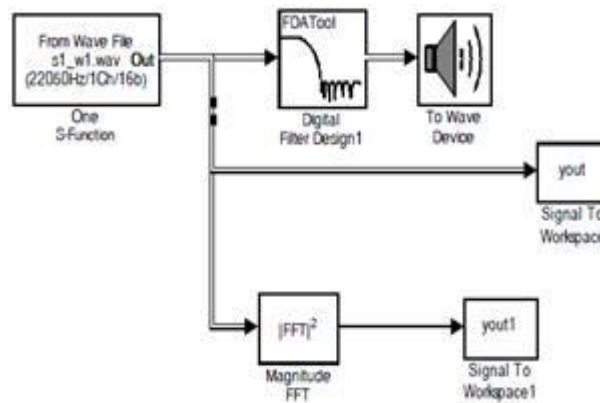


Figure 4: Simulink Model for Analyzing Hindi Data and Numerals

The Simulink model, as shown in Figure 4, was developed for performing analysis such as standard deviation, mean, autocorrelation, magnitude of FFT, data matrix correlation. We also tried a few other statistical techniques.

We would also like to mention that we had started our experiments by using Simulink, but soon found this GUI-based tool to be somewhat limited because we did not find it easy to create multiple models containing variations among them. This iterative and variable-nature of models eventually led us to MATLAB’s (text-based) .m files. We created these files semi-automatically by using a Hindi-language script; the script was developed specifically for this purpose.

Three main data pre-processing steps were required before the data could be used for analysis.

Pre-Emphasis

By pre-emphasis, we imply the application of a *normalization* technique, which is performed by dividing the speech data vector by its highest magnitude.

Data Length Adjustment

FFT execution time depends on exact number of the samples (N) in the data sequence $[x_k]$, and that the execution time is minimal and proportional to $N \log_2(N)$, where N is a power of two. Therefore, it is often useful to choose the data length equal to a power of two.

Endpoint Detection

The goal of endpoint detection is to isolate the word to be detected from the background noise. It is necessary to trim the word utterance to its tightest limits, in order to avoid errors in the modeling of subsequent utterances of the same word. As we can see from the upper part of Figure 5, a threshold has been applied at both ends of the waveform. The front threshold is normalized to a value that all the spoken numbers trim to a maximum value. These values were obtained after observing the behaviour of the waveform and noise in a particular environment. We can see the difference in frequency characteristics of the words.

Fourier Transform

The MATLAB algorithm for the two dimensional FFT routine is as follows:

```
fft2(x) =fft (fft (x));
```

Thus the two dimensional FFT is computed by first computing the FFT of x , that is, the FFT of each column of x , and then computing the FFT of each row of the result. Note that as the application of *fft2* command produced even symmetric data, we only show the lower half of the frequency spectrum in our graphs.

Correlation

Calculations for correlation coefficients of different speakers were performed. As expected, the cross-correlation of the same speaker for the same word did come out to be 1. The correlation matrix of a spoken number was generated in a three-dimensional form for generating different simulations and graphs.

RELATED WORK

Short Time Fourier Transform

Short-time analysis and synthesis enables us to represent the spectra of signals with spectral profiles that change over time (which is the case for most “interesting” 1- dimensional signals such as speech and music). We can think of STFT as multiplying the signal $x[n]$ by a short-time window that is centred on the time frame n . The segment of the signal contained in the window is analyzed using the DFT, which implies the evaluation of the Time-Frequency representation at a set of discrete frequencies where

$$X[n, k] = \sum_{m=-\infty}^{\infty} x[m]w[n-m]e^{-j\omega_k m} = \sum_{m=-\infty}^{\infty} x[m]w[n-m]e^{-j2\pi nk/N}$$

ANALYSIS & RESULTS

We observed that Fourier descriptor feature was independent for the spoken Devnagari Script and numerals with the combination of the Fourier transform and correlation technique commands used in MATLAB, a high accuracy recognition system can be realized. Recorded data was used in Simulink model for introductory analysis.

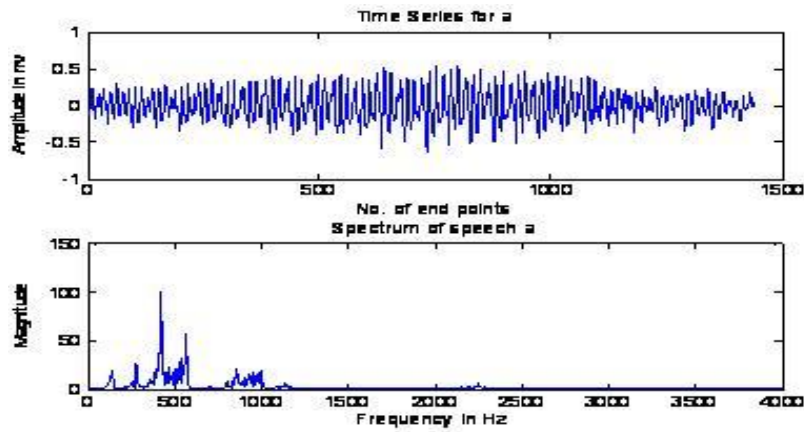


Figure 5: The Fft Waveform of the Word अ in Devnagari Script

X = 1500, It's having 1500 numbers of data points. It's denoted by X. and having a 5 peaks values for each & every word same for अ in Devnagari script.

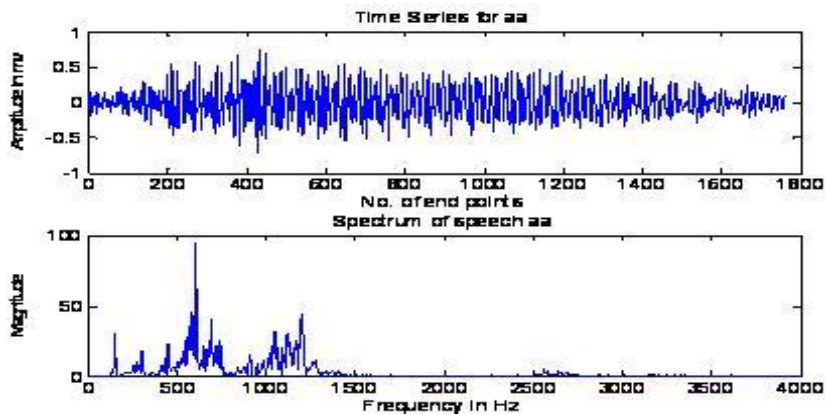


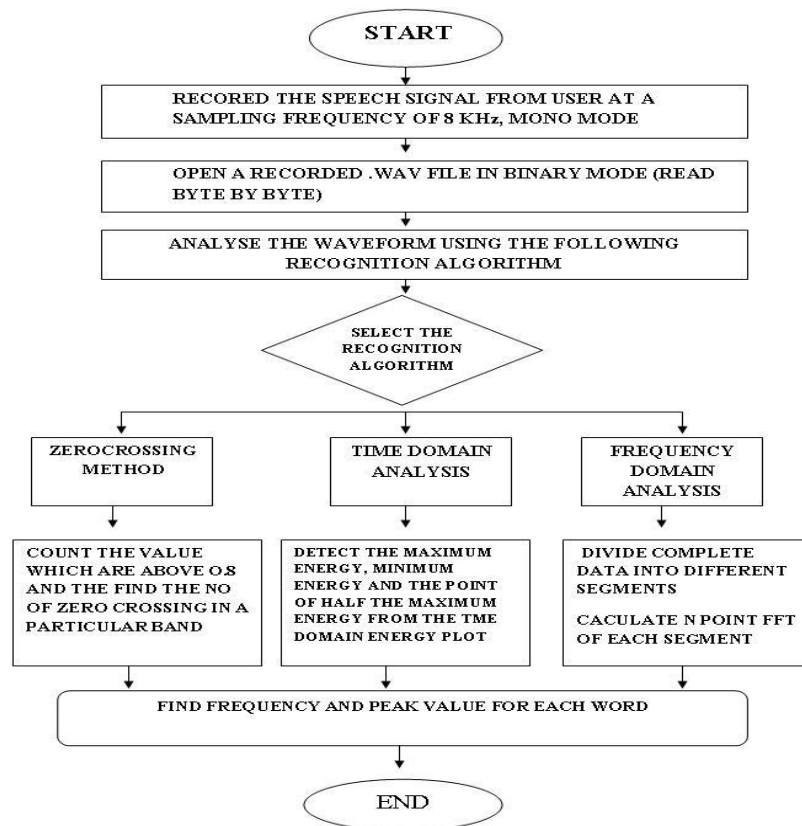
Figure 6: The Fft Waveform of the Word आ in Devnagari Script

X = 1800, It's having 1800 numbers of data points. It's denoted by X. and having a 5 peaks values for each & every word same for आ in Devnagari script.

Table 4: Peaks and its Corresponding Frequencies

Sr.No	Speech Word	Peak	Frequency in (Hz)
1	FOR A	P1	100.9565 F1
		P2	71.1906 F2
		P3	57.5883 F3
		P4	46.6103 F4
		P5	37.3028 F5
2	FOR AA	P1	95.4134 F1
		P2	77.7759 F2
		P3	70.3393 F3
		P4	46.3746 F4
		P5	44.5413 F5

FLOW CHART FOR FREQUENCY ANALYSIS OF SPEECH SIGNALS



CONCLUSIONS AND FUTURE WORK

In conclusion, an efficient, abstract and fast ASR system for regional languages like Hindi is need of the hour. The work implemented in the paper is a step towards the development of such type of systems. The work may further be extended to large vocabulary size and to continuous speech recognition. As shown in results, the system is sensitive to changing spoken methods and changing scenarios, so the accuracy of the system is a challenging area to work upon. Hence, various Speech enhancements and noise reduction techniques may be applied for making system more efficient, accurate and fast.

REFERENCES

1. S K Hasnain, Perez Akhter, "Digital Signal Processing, Theory and Worked Examples", January 2007.
2. Samuel D Stearns, Ruth A David, "Signal Processing Algorithms in MATLAB," Prentice Hall, 1996.
3. S K Hasnain, Nighat Jamil, "Implementation of Digital Signal Processing real time Concepts Using Code Composer Studio 3.1, TI DSK TMS 320C6713 and DSP Simulink Blocksets," IC-4 conference, Indian Navy Engineering College, Goa, Nov. 2007.
4. M. Habibullah Pagarkar, Lakshmi Gopalakrishnan, et.al. "Language Independent Speech Compression using Devnagari Phonetics", 2002.
5. D. O'Shaughnessy, "Interacting with Computers by Voice-Automatic Speech Recognitions and Synthesis", (Invited Paper), Proceedings of the IEEE, Vol. 91, No. 9, 2003, pp. 1272-1305.