



Speech Recognition Techniques: A Review

Er. Shweta Doda

Department of Computer Science and Engineering
JCD Vidyapeeth, Sirsa, India

Er. Rajni Mehta

Assistance Professor in Dept. Computer Science
JCD Vidyapeeth, Sirsa, India

Abstract: This paper presents detailed review of Speech Recognition. Speech recognition is essential for a computer to reach the goal of natural human-computer communication. Speech recognition has gained a lot of interest in the researchers from various fields. A wide variety of approaches have been proposed to recognize isolated words. The objective of this review paper is to summarize some well-known methods used in various stages of speech recognition system such as Dynamic Time Warp and Hidden Markov Model approaches for isolated speech recognition.

Keywords: Speech Capturing, Feature Extraction, Vector Quantization, DTW, HMM

I. INTRODUCTION

Speech is the most natural and efficient way to exchange information for human beings. To make a real “intelligent computer”, it is important that the machine can hear, understand, and act upon spoken information, and also speak to complete the information exchange. Speech Recognition (is also known as Automatic Speech Recognition (ASR) or computer speech recognition) is the process of converting a speech signal to a sequence of words, by means of an algorithm implemented as a computer.

Structure of Basic Speech Recognition System:

Speech Capturing

The input speech can be captured with the help of microphone. The sound card in the computer changes the analog signal into digital signal. The sound card can record the sound and can also play it. Using Window’s MCI (Media Control Interface) commands, the sampling frequency and sample size can be controlled.

Preprocessing

Once sound-capturing process is complete, the speech is available in continuous samples. Now our next step is pre-process these samples to make available for feature extraction and recognition. It involves following steps.

- 1 Background Noise and Silence Removing
- 2 Preemphasis filter
- 3 Blocking into Frames
- 4 Windowing

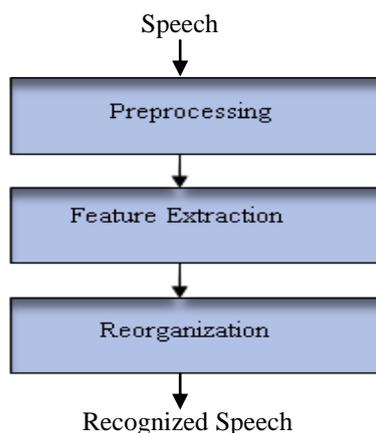


Fig.1 General Structure of Speech Recognition System

Feature Extraction

Feature extraction is the process of parameterization of the speech i.e. representation of speech utterance in terms of feature vectors, which can be used for the recognition purpose. These feature vectors should not change with the speaker i.e. the features should be same for the same utterance by different speakers. These features can be extracted by using several methods, for example digital filter, Fourier Transformation or Linear Predictive Coding. Linear Predictive Coding is the most powerful speech feature extraction (e.g., autocorrelation, cepstral coefficient etc.) technique. It

provides extremely accurate estimates of speech parameters, and is relatively efficient for computation. Feature extraction can be subdivided into three basic operations: spectral analysis, parametric transformation, and statistical modeling.

Recognition

This phase is divided into two parts: one is training and other is testing. The training phase of a recognition system is similar to the learning process of a baby. A child should experience a phenomenon many times and with a wide variability before being able to recognize it. The current speech recognition technology does not allow real-time implementation of models comparable to human complexity. This means that the variability of speech must be limited to achieve proper results. In the testing phase, an unknown utterance is scored over the reference patterns. The word corresponding to the reference pattern closest to the unknown pattern is the word recognized.

II. MODELS OF SPEECH RECOGNITION SYSTEM

1. Dynamic Time Wrap (DTW)

DTW is a method that allows a computer to find an optimal match between two given sequences. It is template based approach. In order to understand DTW, two concepts need to be dealt with:

Features- the information in each signal has to be represented in some manner.

Distances- some form of metric has to be used in order to obtain a match path.

There are two types of distances:

Local Distance: The computational difference between a feature of one signal and a feature of the other is called Local distance.

Global Distance: The overall computational difference between an entire signal and another signal of possibly different length is called Global distance.

DTW based on two concepts:

Symmetrical DTW: Speech is a time-dependent process. Several utterances of the same word are likely to have different durations, and utterances of the same word with the same duration will differ in the middle, due to different parts of the words being spoken at different rates.

- Matching paths cannot go backwards in time;
- Every frame in the input must be used in a matching path;
- Combine local distance scores by adding to give global distance.

It is known as Dynamic Programming (DP). When applied to template-based speech recognition, it is often referred to as Dynamic Time Warp (DTW). DP is guaranteed to find the lowest distance path through matrix, while minimizing the amount of computation.

Asymmetrical DTW:

Each frame of the input pattern is used only once and only once. This means that dispense with template-length normalization and it is not required to add the local distance in twice for diagonal path transitions. This approach is referred to as asymmetric dynamic programming.

2. Hidden Markov Model (HMM)

It is a mathematical approach to recognize speech. It is a doubly embedded stochastic process with an underlying stochastic process that is not directly observable (it is hidden) but can be observed only through another set of stochastic processes that produce the sequence of observations. Stochastic modeling entails the use of probabilistic models to deal with uncertain or incomplete information. In speech recognition, uncertainty and incompleteness arise from many sources; for example, confusable sounds, speaker variability, contextual effects, and homophones words. Thus, stochastic models are particularly suitable approach to speech recognition. Hidden Markov Model is a collection of states connected by transitions. Each transition carries two sets of probabilities: Transition Probability: Which provides the probability for taking this transition, Output Probability: Which defines the conditional probability of emitting each output symbol from a finite alphabet given that a transition is taken. Problems of HMM Evaluation, Decoding Problem, The Learning Problem.

III. LITERATURE REVIEW OF SPEECH RECOGNITION

The review process was adopted by surveying the research in last few years for extraction of information about some issues. Various research articles were reviewed to cover the review of speech recognition technique.

F. Itakura[February 1975] In this found that a new measure of distance for all pole model of speech has been derived on the basis of the likelihood ratio criteria and is applied to automatic recognition of isolated words. An algorithm to find the best match between the input pattern and a reference pattern is derived. In which the dynamic programming technique is used in conjunction with a sequential decision scheme. The system is implemented on a DDP-516 computer to recognize 200 isolated words. The validity of the scheme has been confirmed experimentally. Further work is in progress to test the system for a greater number of talkers and for telephone connection switched over greater distances.

J. K. Baker[February 1975] This paper describe that termination that the hidden articulatory Markov model as an alternative or companion to standard phone-based HMM models for speech recognition. Found that either in noisy

conditions, or when used in tandem with a traditional HMM, a hidden articulatory model can yield improved WER results. Also shown that the HMM is able to reasonably estimate articulator motion from speech. There are a number of avenues to improve this work. In the future, the plan to add more articulatory knowledge, with rules for phoneme modification that arise as a result of physical limitations and shortcuts in speech production, as was done in (Erler 1996) (for example, vowel nasalization). Such rules may help speech recognition systems in the presence of strong co articulation, such as in conversational speech.

D. Raj Reddy [April 1976] In this paper stated that the focus has been to review research progress, to indicate the areas of difficulty, why they are difficult, and how they are being solved. The past few years have seen several conceptual and scientific advances in the field. For the first time use the available extensive analysis of connected speech. Know connected speech recognition is not impossible. The role and use of knowledge are better understood. Almost all systems use knowledge to generate hypotheses and/or verify them. Error and ambiguity can be handled within the framework of search. Stochastic representations and dynamic programming provide a simple and effective solution to the matching problem.

F. Jelinek [April 1976] This paper presented that a new approach for visual speech recognition based on a data driven lip model and HMMs. Experiments have demonstrated high recognition performance using very low dimensional shape information only. The recognition task described is relatively simple because it only consists of four word classes and only deals with isolated words. Nevertheless, recognition tests were speaker independent and have demonstrated high recognition accuracy and generalization ability of the system. More extensive tests with more speakers and sub word classes are necessary to estimate the discrimination ability of shape features for all phonemes. The results are not as good as with 89.58% correct and which was about equivalent to the performance of untrained humans performing the same task. The ability to locate and track lips accurately opens several other potential applications, as example model based image coding, facial animation, facial expression recognition and audio-visual person identification.

Yoseph Linde, Andres Buzo and Robert M. Gray [January 1980] This paper stated that an efficient and intuitive algorithm is presented for the design of vector quantizers based either on a known probabilistic model or on a long training sequence of data. The basic properties of the algorithm are discussed and demonstrated by examples. Quite general distortion measures and long block lengths are allowed, as exemplified by the design of parameter vector quantizers of ten-dimensional vectors arising in Linear Predictive Coded (LPC) speech compression with a complicated distortion measure arising from LPC analysis that does not depend only on the error vector. The hidden-articulator Markov model (HAMM) and have implemented it using HMMs.

Bing-Hwang Juang, David Y. Wong, and H. Augustine, Jr. Gray [April 1982] found that Analytical as well as experimental comparisons between vector and scalar quantization have been presented in detail. It was shown that vector quantization performs a multidimensional clustering process which effectively eliminates unnecessary model spectra. Detailed comparisons between vector and scalar quantization results show that the spectral distortion fluctuates less from frame to frame in vector quantization as compared to scalar quantization.

Lawrence R. Rabiner, and B.H. Juang [January 1986] to search out that to present the theory of hidden Markov models from the simplest concepts (discrete Markov chains) to the most sophisticated models (variable duration, continuous density models). The purpose to focus on physical explanation of the basic mathematics; hence to avoid long, drawn out proofs and/or derivations of the key results, and concentrated primarily on trying to interpret the meaning of the math, and how it could be implemented to illustrate some applications of the theory of HMMs to simple problems in speech recognition, and pointed out how the techniques could be (and have been) applied to more advanced speech recognition.

Kai-Fu Lee, Hsio-Wuen Hon, and Raj Reddy [January 1990] concluded that SPHINX-a hidden Markov model based system for large-vocabulary speaker-independent continuous speech recognition. On the one hand, HMM's perform better with detailed models. On the other hand, HMM's need considerable training. This need is accentuated in large-vocabulary speaker-independence, and discrete HMM's. However, given a fixed amount of training, model specificity and model trainability pose two incompatible goals. More specificity usually reduces trainability, and increased trainability usually results in over generality.

Joseph W. Picone [September 1993] termination that several popular signal analysis techniques in a common framework that emphasized the importance of accurate spectral analysis and statistical normalization. When viewed in this common framework, the differences amongst these competing approaches seem small when compared to the enormous challenges, still face in the speech recognition problem. All approaches share some important basic attributes: time-derivative information, perceptually motivated transformations, and parameter normalization.

L. R. Bahl, P. F. Brown, P. V. de Souza, R. L. Mercer and M.A. Picheny [October 1993] stated that a new technique for constructing Markov models for the acoustic representation of words is described. Word models are constructed from models of sub-word units called fenones. Fenones represent very short speech events, and are obtained automatically through the use of a vector quantizer. The fenonic base form for a word-i.e., the sequence of fenones used to represent the word-is derived automatically from one or more utterances of that word. Since the word models are all composed from a small inventory of sub-word models, training for large-vocabulary speech recognition systems can be accomplished with a small training script. A method for combining phonetic and fenonic models is presented.

Sahar E. Bou-Ghazale and John H. L. Hansen [May 1998] found that novel modeling approach for speech parameter variations under stress using HMM's. The variations in overall pitch contour, voiced duration, and overall spectral contour were modeled for angry, loud, and Lombard effect speech. The models were trained with the variations, referred to as perturbations, in speech parameters from neutral to each stressed condition as opposed to training with actual speech

parameters. A pitch perturbation model was developed for each stressed condition using a three-state single-mixture Gaussian HMM.

IV. CONCLUSION

In this research we study about the speech recognition, speech capturing, preprocessing and how the speech recognized, for this we study two techniques for speech recognition i.e. Hidden Markov model and dynamic time wrap technique. We try to find the working of both the techniques for speech recognition and on the basis of our study we find that how both the techniques are works and how they recognize the speech and extract the words. We conclude that Hidden Markov model is less expensive in compare to Dynamic time wrap techniques but the performance of Hidden Markov Model is somewhat poorer than the Dynamic Time Wrap based recognizer appears to be primarily because of the insufficiency of the Hidden Markov Model training data. The accuracy of dynamic time wrap technique is more accurate than the hidden Markov model.

REFERENCES

- [1] F. Itakura, *Minimum prediction residual principle applied to speech recognition*, IEEE Transactions on Acoustics, Speech, and Signal Processing, Vol. ASSP-23, pp. 67-72, February 1975.
- [2] J. K. Baker, *The DRAGON system - An overview*, IEEE Transactions on Acoustics, Speech, and Signal Processing, Vol. ASSP-23, February 1975, pp. 24-29.
- [3] D. Raj Reddy, *Speech Recognition by Machine: A Review*, Proceedings of the IEEE, Vol. 64, No. 4, April 1976, pp. 501-531.
- [4] F. Jelinek, *Continuous Speech Recognition by Statistical Methods*, Proceedings of IEEE, Vol. 64, April 1976, pp. 532-556.
- [5] Yoseph Linde, Andres Buzo and Robert M. Gray, *An Algorithm for Vector Quantizer Design*, IEEE Transaction on Communications, Vol. COM-28, No. 1, January 1980, pp. 84-95.
- [6] Bing-Hwang Juang, David Y. Wong, and H. Augustine, Jr. Gray, *Distortion Performance of Vector Quantization for LPC Voice Coding*, IEEE Transaction on Acoustics, Speech, and Signal Processing, Vol. ASSP-30, No. 2, April 1982, pp. 294-303.
- [7] Lawrence R. Rabiner, and B.H. Juang, *An Introduction to Hidden Markov Models*, IEEE ASSP Magazine, Vol. 3, No. 1, January 1986, pp. 4-16.
- [8] Kai-Fu Lee, Hsio-Wuen Hon, and Raj Reddy, *An Overview of the SPHINX Speech Recognition System*, IEEE Transaction on Acoustics, Speech, and Signal Processing, Vol. ASSP-38, No. 1, January 1990, pp. 35-45
- [9] Joseph W. Picone, *Signal Modeling Techniques in Speech Recognition*, Proceedings of the IEEE, Vol. 81, No. 9, September 1993, pp. 1214-1245.
- [10] L. R. Bahl, P. F. Brown, P. V. de Souza, R. L. Mercer and M.A. Picheny, *A Method for the Construction of Acoustic Markov Models for Words*, IEEE Transaction on Speech and Audio Processing, Vol. 1, No. 4, October 1993, pp. 443-452.
- [11] Sahar E. Bou-Ghazale and John H. L. Hansen, "HMM-Based Stressed Speech Modeling with Application to Improved Syndisertation and Recognition of Isolated Speech Under Stress", IEEE Transaction on Speech and Audio Processing, Vol. 6, No. 3, May 1998, pp. 201-216.