



Universidade Federal de Pernambuco
Centro de Informática

A Fractional Gaussian Mixture Model for Speaker Verification

Final Term Paper

Eduardo Martins Barros de Albuquerque Tenório

March 3, 2015

Declaration

This paper is a presentation of my original research work, as partial fulfillment of the requirement for the degree in Computer Engineering. Wherever contributions of others are involved, every effort is made to indicate this clearly, with due reference to the literature, and acknowledgement of collaborative research and discussions.

The work was done under the guidance of Prof. Dr. Tsang Ing Ren, at Centro de Informática, Universidade Federal de Pernambuco, Brazil.

Eduardo Martins Barros de Albuquerque Tenório

In my capacity as supervisor of the candidate's paper, I certify that the above statements are true to the best of my knowledge.

Prof. Dr. Tsang Ing Ren

March 3, 2015

Acknowledgements

I am thankful to my parents, for the support and patience during the graduation,
To my adviser, Tsang Ing Ren, for the guidance,
To Cleice Souza, for the previous readings and help.

Live long and prosper

Vulcan salute

Abstract

TODO escrever o abstract após terminar tudo (após a conclusão)

Contents

1	Introduction	1
1.1	Speaker Recognition	1
1.2	Objectives	2
1.3	Document Structure	2
2	Speaker Recognition System	3
3	Feature Extraction	5
3.1	Mel-Frequency Cepstral Coefficient	6
3.1.1	The Mel Scale	6
3.1.2	Cepstrum	7
3.1.3	Extraction Process	7
4	Gaussian Mixture Model	9
5	Fractional Gaussian Mixture Model	11
6	Experiments	13
7	Conclusion	15
A	Codes	17

1. Introduction

The rise in popularity and naturality of computational systems in the everyday of modern life creates the need for easy and less invasive forms of authentication. While enter a hard to memorize password in a terminal still is the safest approach, voice biometrics presents itself as an alternative with continuing improvement. Also, speech is the most natural way humans have to communicate, being incredibly complex and with numerous specific details related to its producer [1]. Therefore, it is expected an increasing usage of vocal interfaces to perform actions such as authenticate in a system, command a machine, and identify who is talking and the content of the conversation.

Over the last decade, voice recognition technology has appeared in many commercial products (e.g. Google Now and Apple Siri) with relatively high popularity. Fingerprint biometrics is a reality for ATMs users, and retina scanners have been commercialized for some decades, then it is natural that the less invasive method of voice recognition be popularized for authentication processes in a near future.

Most of the commercial products based on voice technology are intended to perform **speech recognition** (determine *what* is being said) instead of **speaker recognition** (determine *who* is speaking). To achieve this goal, numerous voice processing techniques have become popular, e.g. Natural Language Processing (NLP), Hidden Markov Models (HMM) and Gaussian Mixture Models (GMM). Although all of these are interesting, the subject covered in this paper is a field of speaker recognition and only a small subset of techniques will be unraveled.

1.1 Speaker Recognition

As stated in [2], speaker recognition may be divided in two fields. The first is **speaker identification**, aimed to determine the identity of a speaker (using a speech signal) from a non-unitary set of speakers. This task is also named speaker identification in **closed set**. In the second field, **speaker verification**, the goal is to determine if a speaker is who he/she claims to be, not an imposter. As the set of imposters is unknown *a priori*, this is an **open set** problem. A special case is the **speaker identification in open set**, when an “imposter class” is added to the system in order to categorize all unmatched speakers found.

Restrictions on the type of text may be used. In **text-dependent** systems, the content of the speech is relevant to the evaluation, and the training and test utterances must contain the same text (but not the same intonation), e.g. a passphrase. In **text-independent** systems there is no restriction to the content of both sets. In this case the non-textual characteristics of the user’s voice is what have importance to the evaluator. These characteristics are presented in different sentences, usage of

different languages and even in gibberish.

This paper is focused in **text-independent speaker verification**, in other words, the determination of a user's claimed identity by analysis of his/her vocal characteristics with no predefined text to dictate. To achieve that, a GMM adapted from an UBM [3] is implemented. Also, an adaptation of the technique is proposed and evaluated using the theory of FCM presented in [4].

1.2 Objectives

The objectives of this study are:

- Analyse and evaluate the speaker verification system using adapted GMM proposed by [3];
- Propose and evaluate a new method derived from GMM, using the FCM theory proposed by [4];
- Conduct experiments and validation of the existent and the proposed methods.

1.3 Document Structure

Chapter 2 contains a brief historical context and some basic details about voice recognition, as well as the basic architecture for a speaker verification system. The feature extraction process is explained in chapter 3, from the reasons for its use to the chosen technique (MFCC). In chapter 4 is detailed the GMM and the UBM-GMM. Chapter 5 introduces FCM and the proposed FGMM. Experiments are described in chapter 6, as well as its results. Finally, chapter 7 concludes the study. Furthermore, this work contains an appendix with the most relevant pieces of the source code and some mathematics concepts used.

2. Speaker Recognition System

3. Feature Extraction

As an acoustic wave propagated through space over time, the speech signal is not appropriate to be evaluated by the speaker verification system. In order to deliver decent outcomes, a good parametric representation must be provided to the system. This task is performed by the feature extraction process, which transforms a speech signal into a sequence of characterized measurements, i.e. features. As stated in [5], “the usual objectives in selecting a representation are to compress the speech data by eliminating information not pertinent to the phonetic analysis of the data, and to enhance those aspects of the signal that contribute significantly to the detection of phonetic differences”. According to [6] the ideal features should:

- occur naturally and frequently in normal speech;
- be easily measurable;
- vary highly among speakers and be very consistent for each speaker;
- not change over time nor be affected by the speaker’s health;
- be robust to reasonable background noise and to transmission characteristics;
- be difficult to be artificially produced;
- not be easily modifiable by the speaker.

Features may be categorized based on vocal tract or behavioral aspects, divided in (1) short-time spectral, (2) spectro-temporal, (3) prosodic and (4) high level [2]. Short-time spectral features are usually calculated using millisecond length windows and describe the voice spectral envelope, composed of supralaryngeal properties of the vocal tract, e.g. timbre. Prosodic and spectro-temporal occur over time, e.g. rhythm and intonation, and high level features occur during the conversation, e.g. accents.

The parametric representations evaluated in [5] may be divided into those based on the Fourier spectrum, Mel-Frequency Cepstrum Coefficients (MFCC) and Linear Frequency Cepstrum Coefficients (LFCC), and those based on the Linear Prediction Spectrum, Linear Prediction Coefficients (LPC), Reflection Coefficients (RC) and Linear Prediction Cepstrum Coefficients (LPCC). The better evaluated representation was the MFCC, with minimum and maximum accuracy of 90.2% and 99.4% respectively, leading to its choice as the parametric representation in this work.

3.1 Mel-Frequency Cepstral Coefficient

MFCC is a highly used parametric representation in the area of voice processing, due to its similarity with the mode the human ear operates. Despite the fact the ear is divided in three sections, i.e. outer, middle and inner ears, only the last is mimicked. The mechanical pressure waves produced by the triad hammer-anvil-stirrup are received by the cochlea (Fig. 3.1), a spiral-shaped cavity with a set of inner hair cells attached to a membrane (the basilar membrane) and filled with a liquid. This structure converts motion to neural activity through a non-uniform spectral analysis [7] and passes it to the pattern recognition in the brain.

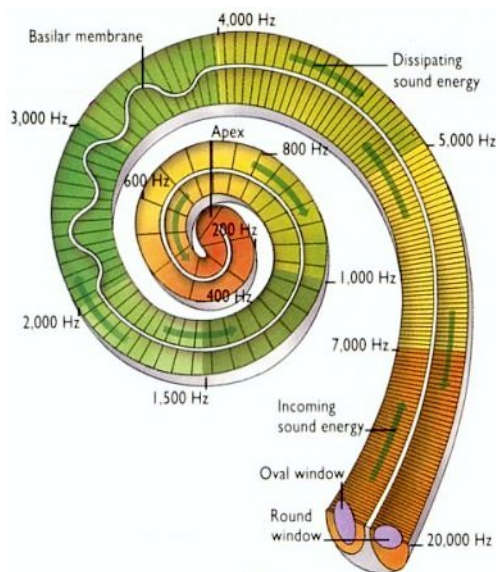


Figure 3.1: Cochlea divided by frequency regions.

A key factor in the perception of speech and other sounds is *loudness*, a quality related to the physical property of sound pressure level. Loudness is quantified by relating the actual sound pressure level of a pure tone (in dB relative to a standard reference level) to the perceived loudness of the same tone (in a unit called phons) over the range of human hearing (20 Hz–20 kHz) [7]. As shown in Fig. 3.2, a 100 Hz tone at 60 dB is equal in loudness to a 1000 Hz tone at 50 dB, both having the *loudness level* of 50 phons (by convention).

3.1.1 The Mel Scale

The mel scale is the result of an experiment conducted by Stevens, Volkman and Newman [9] intended to measure the perception of a pitch and construct a scale based on it. Each observer was asked to listen to two tones, one in the fixed frequencies 125, 200, 300, 400, 700, 1000, 2000, 5000, 8000 and 12000 Hz, and the other free to have its frequency varied by the observer for each fixed frequency of the first tone. An interval of 2 seconds separated both tones. The observers were instructed to say in which frequency the second tone was “half the loudness” of the first. A geometric mean was taken from the observers’ answers and a measure of 1000 mels was assigned to the frequency of 1000 Hz, 500 mels to the frequency sounding half as high (as determined by Fig. 1 in [9]) and so on.

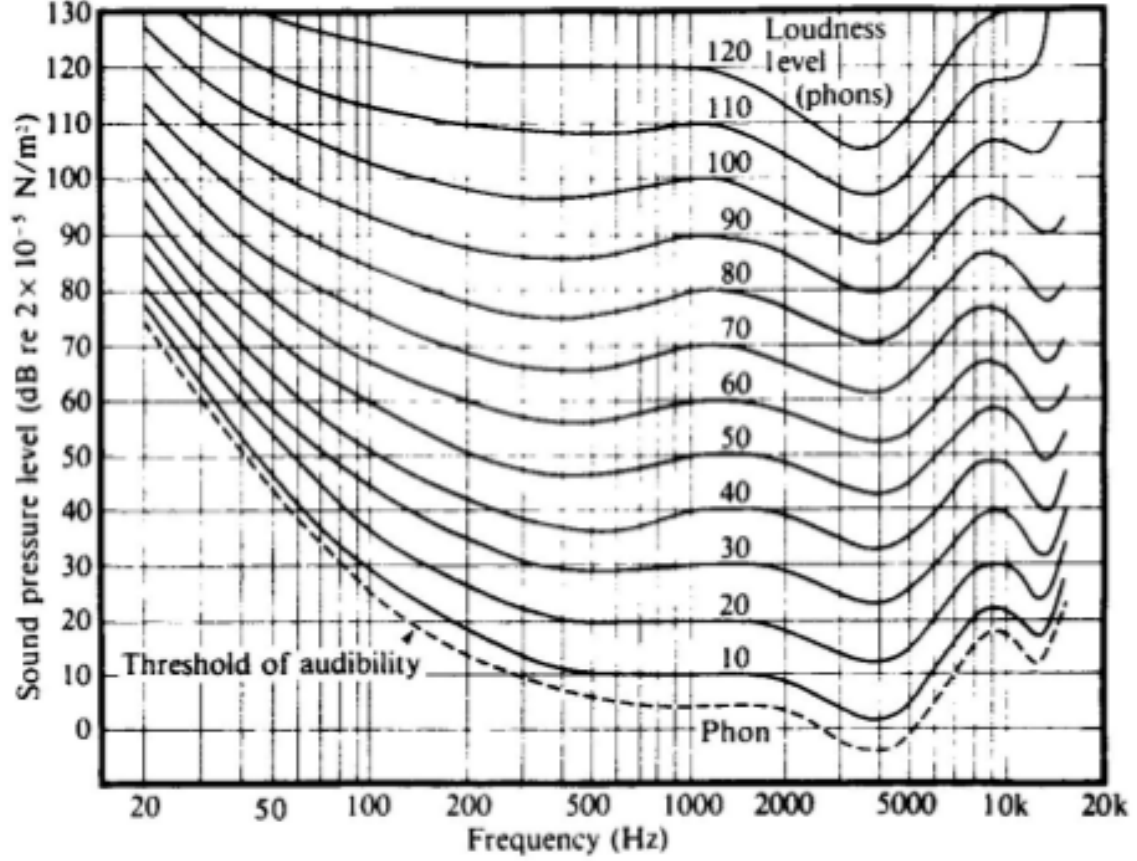


Figure 3.2: Loudness level for human hearing [8].

Decades after the creation of the mel scale, O'Shaughnessy [10] published an equation to convert frequencies in hertz to frequencies in mels.

$$f_{mel} = 2595 \log_{10} \left(1 + \frac{f}{700} \right) \quad (3.1)$$

Being logarithmic, the growth of a mel-frequency curve is slow with a linear growth of the frequency in hertz. Eq. 3.1 sometimes is used only for frequencies higher than 1000 Hz while the lower frequencies obey a linear function. In this work all conversions will use Eq. 3.1, as shown by Fig. 3.3.

3.1.2 Cepstrum

3.1.3 Extraction Process

Pre-emphasis

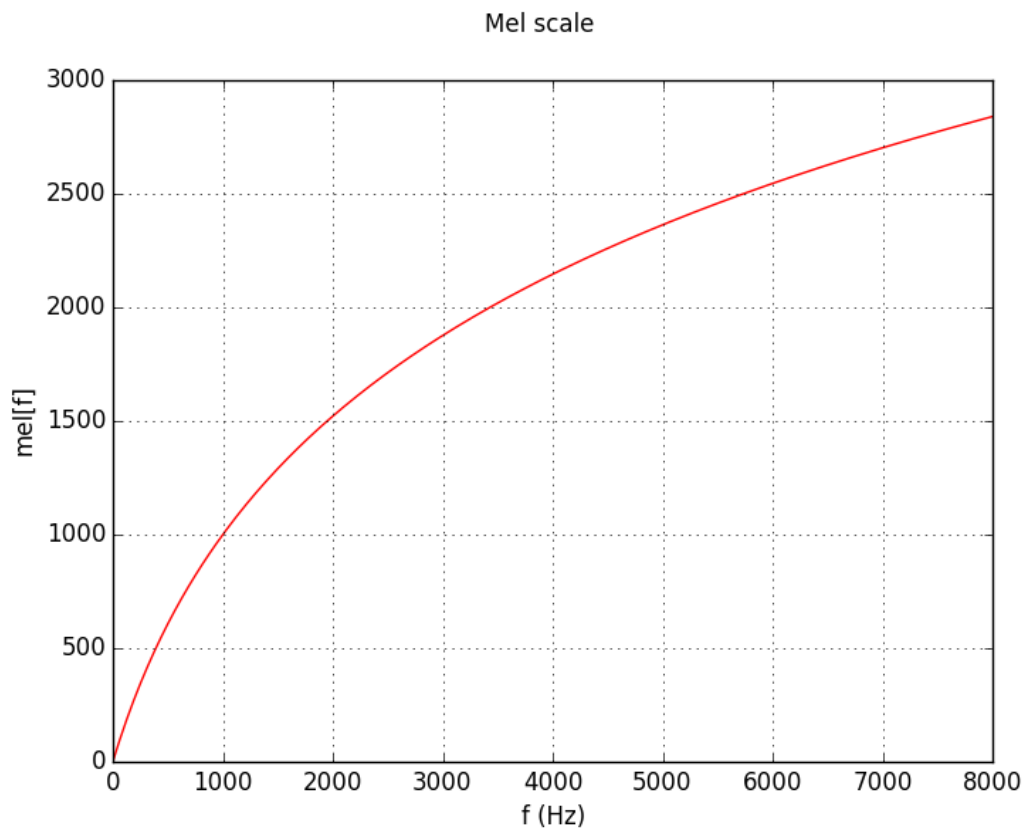


Figure 3.3: The logarithm curve of the mel-frequency.

4. Gaussian Mixture Model

5. Fractional Gaussian Mixture Model

6. Experiments

7. Conclusion

TODO escrever a conclusão após terminar tudo (antes do abstract)

A. Codes

Bibliography

- [1] Frédéric Bimbot, Jean-François Bonastre, Corinne Fredouille, Guillaume Gravier, Ivan Magrin-Chagnolleau, Sylvain Meignier, Teva Merlin, Javier Ortega-García, Dijana Petrovska-Delacrétaz and Douglas A. Reynolds. “A Tutorial on Text-Independent Speaker Verification”. In: *EURASIP Journal on Applied Signal Processing* 4 (2004), pp. 430–451.
- [2] Hector N. B. Pinheiro. *Sistemas de Reconhecimento de Locutor Independente de Texto*. Final Term Paper. Universidade Federal de Pernambuco, 2013.
- [3] Douglas A. Reynolds, Thomas F. Quatieri and Robert B. Dunn. “Speaker Verification Using Adapted Gaussian Mixture Models”. In: *Digital Signal Processing* 10.1 (2000), pp. 19–41.
- [4] Chaobang Gao, Jiliu Zhou and Qiang Pu. “Theory of fractional covariance matrix and its applications in PCA and 2D-PCA”. In: *Expert Systems with Applications* 40.13 (2013), 5395–5401.
- [5] Steven B. Davis and Paul Mermelstein. “Comparison of Parametric Representations for Monosyllabic Word Recognition in Continuously Spoken Sentences”. In: *IEEE Transactions on Acoustics, Speech, and Signal Processing* ASSP-28.4 (1980), pp. 357–366.
- [6] Jared J. Wolf. “Efficient acoustic parameters for speaker recognition”. In: *Journal of the Acoustical Society of America* 51 (1972), pp. 2044–2056.
- [7] Lawrence R. Rabiner and Ronald W. Schafer. “Introduction to Digital Speech Processing”. In: *Foundations and Trends in Signal Processing* 1.1-2 (2007), pp. 1–194.
- [8] Harvey Fletcher and Wilden A. Munson. “Loudness, Its Definition, Measurement and Calculation”. In: *Bell Telephone Laboratories* 12.4 (1933), pp. 82–108.
- [9] Stanley S. Stevens, John Volkman and Edwin B. Newman. “A Scale for the Measurement of the Psychological Magnitude Pitch”. In: *The Journal of Acoustical Society of America* 8.3 (1937), pp. 185–190.
- [10] Douglas O’Shaughnessy. *Speech Communications: Human and Machine*. Addison-Wesley, 1987.