



Universidade Federal de Pernambuco
Centro de Informática

Gaussian Mixture Models for Text-Independent Speaker Recognition

Final Term Paper

Eduardo Martins Barros de Albuquerque Tenório

March 3, 2015

Declaration

This paper is a presentation of my 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 family, 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 suggestions,
To Sérgio Vieira and James Lyons, for clarify many of my questions.

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 Systems | 5 |
| 2.1 | Basic Concepts | 5 |
| 2.1.1 | Utterance | 5 |
| 2.1.2 | Features | 5 |
| 2.2 | Speaker Identification | 6 |
| 2.2.1 | Training | 6 |
| 2.2.2 | Test | 6 |
| 2.3 | Speaker Verification | 7 |
| 2.3.1 | Likelihood Ratio Test | 7 |
| 2.3.2 | Training | 7 |
| 2.3.3 | Test | 8 |
| 3 | Feature Extraction | 9 |
| 3.1 | Mel-Frequency Cepstral Coefficient | 10 |
| 3.1.1 | The Mel Scale | 10 |
| 3.1.2 | Cepstrum | 11 |
| 3.1.3 | Extraction Process | 11 |
| 4 | Gaussian Mixture Models | 13 |
| 5 | Experiments | 15 |
| 6 | Conclusion | 17 |
| A | Results from Experiments | 19 |
| B | Codes | 21 |

1. Introduction

The increasing popularity and the intensive usage of computational systems in the everyday of modern life creates the need for easier and less invasive forms of user recognition. While entering a hard-to-memorize password in a terminal and identifying a person by placing a human finger to listen to telephone calls are the status quo for respectively authentication and identification, voice biometrics presents itself as a continuing improvement alternative. Passwords can be forgotten and people are biased and unable to be massively trained, but the unique characteristics of a person's voice combined with an automatic speaker recognizer (ASR) outperform any "manual" attempt.

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 that an increasing usage of vocal interfaces to perform actions such as computer login, voice search (e.g., Apple Siri, Google Now and Samsung S Voice) and identification of speakers in a conversation and its content. At present, fingerprint biometrics is adopted in several solutions (e.g., ATMs [2]), authentication through facial recognition comes as built-in software for average computers and iris scan was adopted for a short time by United Kingdom and permanently by United Arab Emirates border controls [3, 4]. These examples indicate a near future where biometrics are common, with people talking to the computer and receiving concise answers, and cash withdrawals allowed via a combination of speaker verification, corrected captcha dictated and other techniques.

Current commercial products based on voice technology (e.g., Dragon NaturallySpeaking, KIVOX and VeriSpeak) are usually intended to perform either **speech recognition** (*what* is being said) or **speaker recognition** (*who* is speaking). Voice search applications are designed to determine the content of a speech, usually with no concern about who the speaker is or if there is more than one, while computer login and telephone fraud prevention supplement a memorized personal identification code with speaker verification [5], not interested in the message spoken. Few applications perform both processes, such as automatic speaker labeling of recorded meetings, that transcribes what each person is saying. To achieve this goal, numerous voice processing techniques have become known in industry and academy, e.g., Natural Language Processing (NLP), Hidden Markov Models (HMMs) and Gaussian Mixture Models (GMMs). Although all of these are interesting state-of-the-art techniques, the subject covered by this paper is the area of speaker recognition and only a small subset of these techniques will be unraveled.

1.1 Speaker Recognition

As stated in [6], speaker recognition may be divided in two subareas. The first is **speaker identification**, aimed to determine the identity of a speaker from a non-unitary set of known speakers. This task is also named speaker identification in **closed set**. In the

second, **speaker verification**, the goal is to determine if a speaker is who he or she claims to be, not an imposter. As the set of imposters is unknown, this is an **open set** problem. An intermediate task is **open set identification**, when an “unmatched class” is added in order to categorize all unknown speakers.

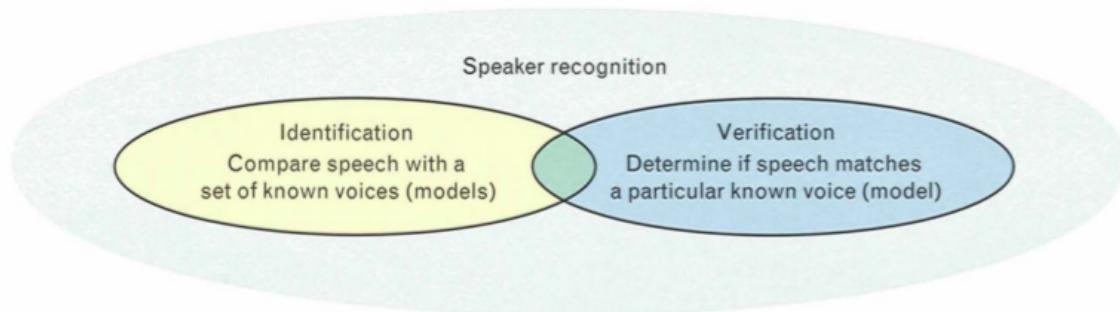


Figure 1.1: Speaker identification and speaker verification are different, but not entirely [5].

The text used may be constrained, such as by type (e.g., digits and letters) and/or by number of words used (e.g., one word or sentences). In **text-dependent** systems the content of the speech is relevant to the evaluation, and the testing texts must belong to the training set (not necessarily be the entire set) [7]. A change in the training text demands an complete new training section. **Text-independent** systems have no restrictions to the message in both sets, with the non-textual characteristics of the user’s voice (e.g., pitch and accent) being the important aspects to the evaluator. These characteristics are presented in different sentences, usage of foreign languages and even gibberish. Between the extremes in constraints falls the **vocabulary-dependent system**, which constrains the speech to come from a limited vocabulary (e.g., digits) from which test words or phrases are selected (e.g., “two” or “one-two-three”) [5].

The focus of this paper is in **text-independent speaker recognition** and to achieve that, Gaussian Mixture Models are used.

1.2 Objectives

The objectives of this study is to implement an ASR that executes the listed actions:

- From a group of enrolled speakers identify who produced a given speech signal, for all speakers in the group;
- Determine if a speaker is the claimed enrolled speaker or an imposter, given the speech signal produced. This experiment is performed for a group of enrolled speakers and for a group of imposters.
- Analyze the performance for different number of mixtures and features.

1.3 Document Structure

Chapter 2 contains basic information about speaker recognition, as well as the basic architecture of speaker identification and verification systems. The feature extraction process

is explained in Chapter 3, from the reasons for its use to the chosen technique (Mel-Frequency Cepstral Coefficient, MFCC). In Chapter 4 the GMM and the Universal Background Model (UBM) are detailed. Experiments are described in Chapter 5, as well as its results. Finally, Chapter 6 concludes the study. Furthermore, this work contains an appendix with all results from the experiments performed (Section A) and the most relevant pieces of the implemented code (Section B).

2. Speaker Recognition Systems

The process of speaker recognition lies on the field of pattern classification, with the speaker's utterance (a speech signal) as input for a classifier. This decision may be, given a speech signal Y produced by a speaker S and a set $\mathcal{S} = \{S_1, \dots, S_S\}$ of enrolled users, “identify S as S_i if $i = \arg \max_j P(S_j|Y)$ ”. This is a case of speaker identification and the output is a S_i from \mathcal{S} . Another type of decision is “accept S as S_i if $P(S_i|Y) \geq \alpha$ ”, where S_i is the claimed identity of S . This is a speaker verification decision. Both are covered in this chapter.

2.1 Basic Concepts

Before explain the architecture of both types of recognizers, the elucidation of some basic concepts is necessary.

2.1.1 Utterance

An utterance is a piece of speech produced by a speaker. It may be a word, a statement or any vocal sound. The terms *utterance* and *speech signal* sometimes are used interchangeably, but from herenow speech signal will be defined as an utterance recorded, digitalized and ready to be processed. An example is shown in Fig. 2.1.

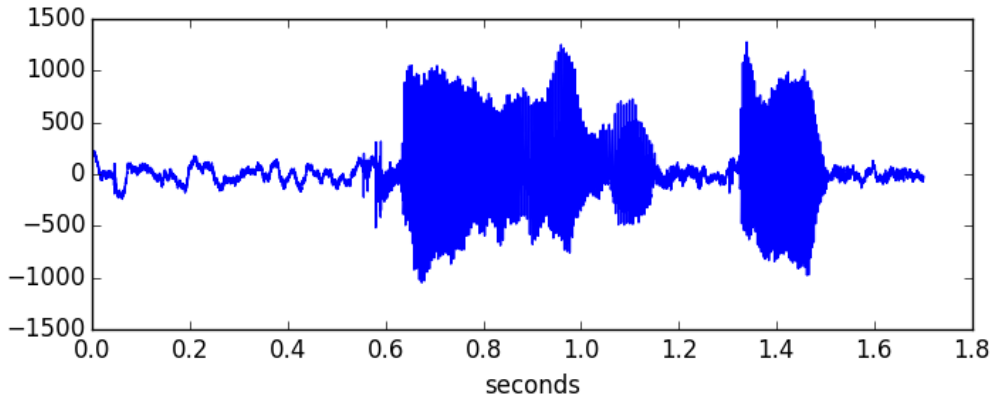


Figure 2.1: Speech signal for utterance “karen livescu”, from the MIT dataset [8].

2.1.2 Features

The raw speech signal is unfit for usage by a recognition system. For a correct processing, the unique features from the speaker's vocal tract are extracted, what reduces the number of variables the system needs to deal with (leading to a simpler implementation) and

performs a better evaluation (prevents the curse of dimensionality). Due to the stationary properties of the speech signal when analyzed in a short period of time, it is divided in overlapping frames of small and predefined length, to avoid “loss of significancy” in the features [9, 10]. This extraction is executed by the MFCC algorithm, explained in details in Chapter 3.

2.2 Speaker Identification

In speaker identification, the objective of the system is to assign an identity from a set of enrolled speakers to the so far unknown speaker, using a speech signal produced by him or her. This system contains a set \mathcal{S} and is fed with features $\mathbf{X} = \{x_1, \dots, x_T\}$ extracted from the speech \mathbf{Y} produced by the speaker \mathcal{S} , the one to be identified.

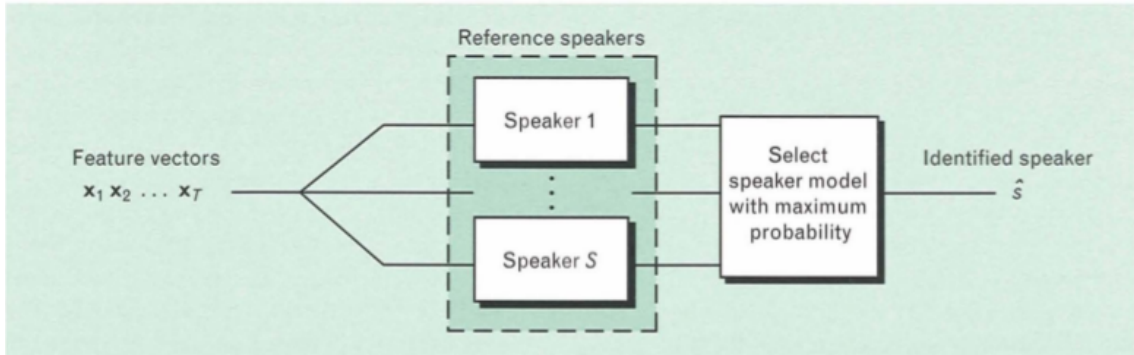


Figure 2.2: Speaker-recognition system for identification [5].

The system is the closed set speaker identification described in Section 1.1, and it classifies any speaker, be it an enrolled one or not. Unenrolled speakers are incorrectly identified as the speaker which model returns the highest probability. The equation to determine the identity of \mathcal{S} is

$$\mathcal{S}_i \text{ if } i = \arg \max_j P(\mathcal{S}_j | \mathbf{X}), \quad (2.1)$$

and as all speakers have the same prior probability, Eq. 2.1 is reduced to

$$\mathcal{S}_i \text{ if } i = \arg \max_j p(\mathbf{X} | \mathcal{S}_j). \quad (2.2)$$

2.2.1 Training

The use of \mathcal{S}_j to represent the speaker is inaccurate, because it maintains the problem in a high level of abstraction. To present a more concrete solution, it is necessary to use a model λ_j for each \mathcal{S}_j . Each λ_j is trained independently until a stop condition is fulfilled.

2.2.2 Test

Once all λ_j 's are trained, the system is able to be used for classification. Eq. 2.2 is then redefined as

$$\mathcal{S}_i \text{ if } i = \arg \max_j p(\mathbf{X} | \lambda_j), \quad (2.3)$$

and the unknown speaker receives the identity of the enrolled speaker for which \mathbf{X} maximizes the **likelihood** of λ_j (see Fig. 2.2).

2.3 Speaker Verification

In speaker verification, \mathcal{S} claims to be a particular \mathcal{S}_i from \mathcal{S} . The strength of this claim resides on how similar the features \mathbf{X} are to the features from \mathcal{S}_i “memorized” by the system. Then a simple equation

$$p(\mathbf{Y}|\mathcal{S}_i) \begin{cases} \geq \alpha, & \text{accept } \mathcal{S}, \\ < \alpha, & \text{reject } \mathcal{S}, \end{cases} \quad (2.4)$$

should be enough (again considering all speakers equally probable). However a subset of enrolled speakers may have vocal similarities, leading to a misclassification of one enrolled speaker as another (a false positive). To reduce the error rate, the system must decide not only if a speech signal came from the claimed speaker, but also if it came from a set composed of all other enrolled speakers.

2.3.1 Likelihood Ratio Test

Given the vector of features \mathbf{X} , and assuming it was produced by only one speaker, the detection task can be restated as a basic test between two hypotheses [11]:

H_0 : \mathbf{X} is from the claimed speaker \mathcal{S}_i ;

H_1 : \mathbf{X} is not from the claimed speaker \mathcal{S}_i .

The optimum test to decide which hypothesis is valid is the **likelihood ratio test** between both likelihoods $p(\mathbf{X}|H_0)$ and $p(\mathbf{X}|H_1)$

$$\frac{p(\mathbf{X}|H_0)}{p(\mathbf{X}|H_1)} \begin{cases} \geq \theta, & \text{accept } H_0, \\ < \theta, & \text{reject } H_0, \end{cases} \quad (2.5)$$

where the decision threshold for accepting or rejecting H_0 is θ . Applying the logarithm, the behavior of the likelihood ratio is maintained and Eq. 2.5 is replaced by the log-likelihood ratio

$$\Lambda(\mathbf{X}) = \log p(\mathbf{X}|H_0) - \log p(\mathbf{X}|H_1). \quad (2.6)$$

2.3.2 Training

Once the features are extracted from the speech signal, they are used to train the models λ_{hyp} and $\lambda_{\overline{hyp}}$ for H_0 and H_1 , respectively. A high-level demonstration of the training of λ_{hyp} (mathematical representation of \mathcal{S}_i) is shown in Fig. 2.3.

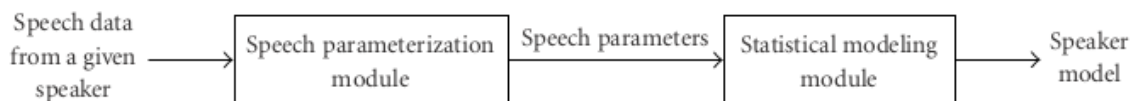


Figure 2.3: The statistical model of \mathcal{S} is created from the speech signal \mathbf{Y} [1].

Due to λ_{hyp} be a model of \mathcal{S}_i , the features used for training (i.e., estimate $p(\mathbf{X}|\lambda_{hyp})$) are extracted from speech signals produced by \mathcal{S}_i . The model λ_{hyp} , however, is not well-defined. It should be composed of the features extracted from speech signals from all other speakers except \mathcal{S}_i , but creating a single λ_{hyp} for each speaker is complicated and with no expressive gain. Instead, what is normally done is use all speakers to generate a background model λ_{bkg} [12], in which the weight of each \mathcal{S}_i is minimized.

2.3.3 Test

As seen in Eq. 2.5, the decision process is based on a function *Score*. Replacing each H_j for its corresponding model, the likelihood of a λ_j given \mathbf{X} can be written as

$$p(\mathbf{X}|\lambda_j) = \prod_{t=1}^T p(\mathbf{x}_t|\lambda_j). \quad (2.7)$$

Using the logarithm function, Eq. 2.7 becomes

$$\log p(\mathbf{X}|\lambda_j) = \frac{1}{T} \sum_{t=1}^T \log p(\mathbf{x}_t|\lambda_j), \quad (2.8)$$

where the term $\frac{1}{T}$ is used to normalize the log-likelihood to the duration of the speech signal. That said, the likelihood ratio given by Eq. 2.6 becomes

$$\Lambda(\mathbf{X}) = \log p(\mathbf{X}|\lambda_{hyp}) - \log p(\mathbf{X}|\lambda_{bkg}), \quad (2.9)$$

and the speaker is accepted if $\Lambda(\mathbf{X}) \geq \log \theta$, for an arbitrary value of θ .

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 [9], “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 [13] 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 [14]. 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 [9] 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 [10] and passes it to the pattern recognition in the brain.

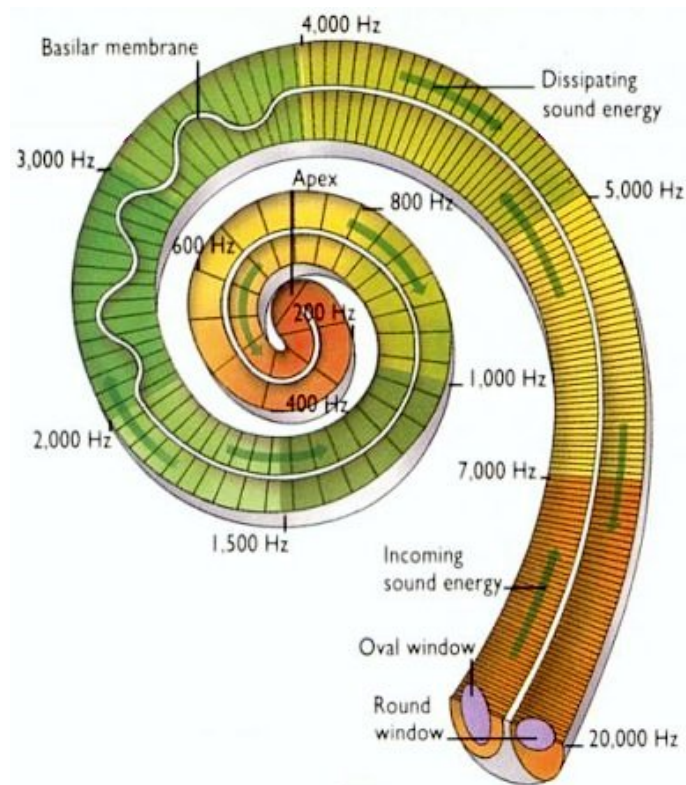


Figure 3.1: Cochlea divided by frequency regions [15].

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) [10]. 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, Volkmann and Newman [17] 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

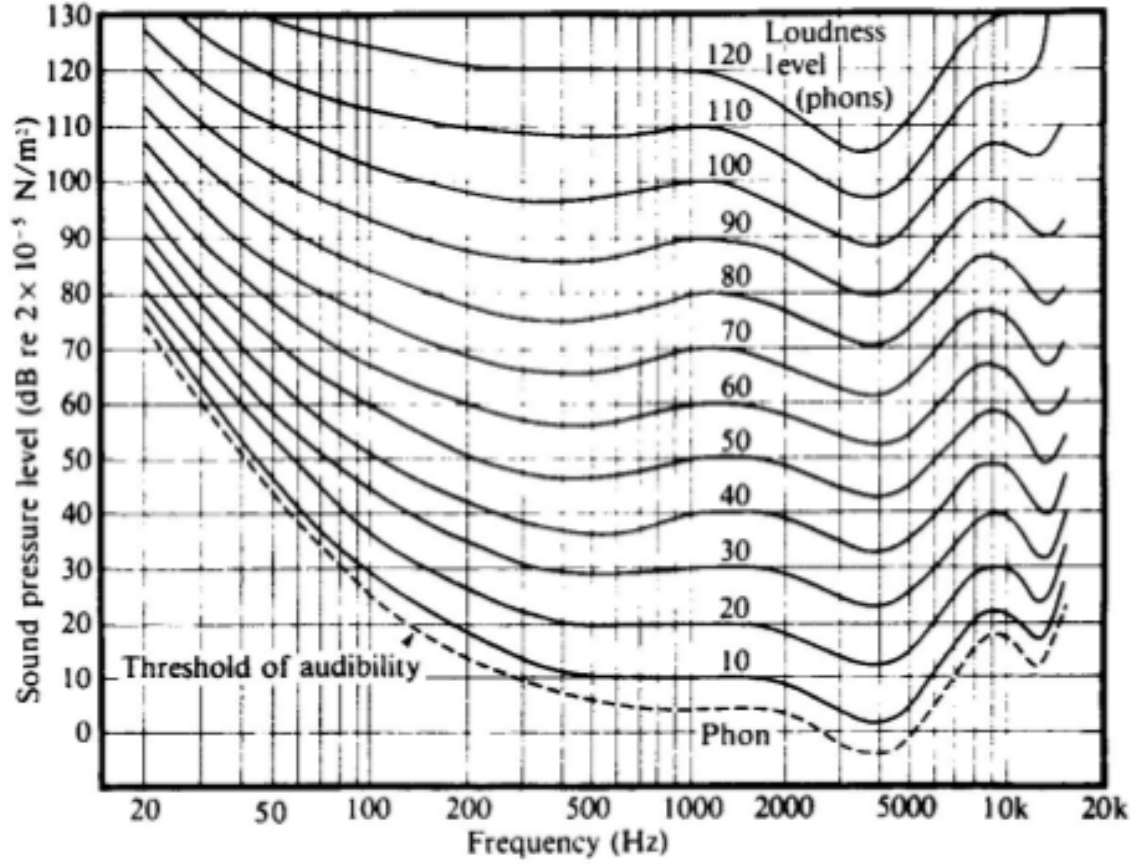


Figure 3.2: Loudness level for human hearing [16].

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 [17]) and so on.

Decades after the creation of the mel scale, O’Shaughnessy [18] 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

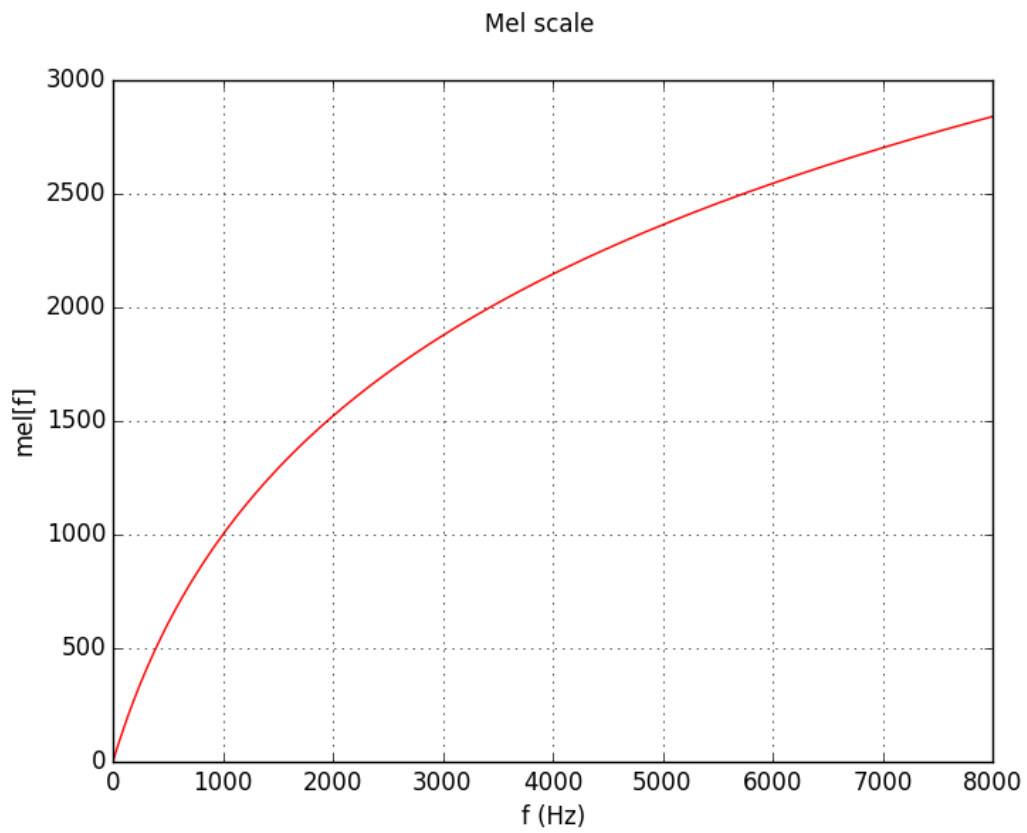


Figure 3.3: The logarithm curve of the mel-frequency.

4. Gaussian Mixture Models

5. Experiments

6. Conclusion

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

A. Results from Experiments

B. Codes

References

- [1] Frédéric Bimbot et al. “A Tutorial on text-independent speaker verification”. In: *EURASIP Journal on Applied Signal Processing* 4 (2004), pp. 430–451.
- [2] P.T. Wang and S.M. Wu. “Personal fingerprint authentication method of bank card and credit card”. Pat. US Patent App. 09/849,279. 2002. URL: <https://www.google.com/patents/US20020163421>.
- [3] M. Angela Sasse. “Red-Eye Blink, Bendy Shuffle, and the Yuck Factor: A User Experience of Biometric Airport Systems”. In: *Security & Privacy, IEEE* 5.3 (2007), pp. 78–81.
- [4] Ahmad N. Al-Raisi and Ali M. Al-Khoury. “Iris recognition and the challenge of homeland and border control security in UAE”. In: *Telematics and Informatics* 25.2 (2008), pp. 117–132.
- [5] Douglas A. Reynolds. “Automatic Speaker Recognition Using Gaussian Mixture Speaker Models”. In: *The Lincoln Laboratory Journal* 8.2 (1995), pp. 173–192.
- [6] Douglas A. Reynolds and William M. Campbell. “Springer Handbook of Speech Processing”. In: ed. by Jacob Benesty, M. Mohan Sondhi, and Yiteng Huang. Berlin: Springer, 2008. Chap. Text-Independent Speaker Recognition, pp. 763–780.
- [7] Martial Hébert. “Springer Handbook of Speech Processing”. In: ed. by Jacob Benesty, M. Mohan Sondhi, and Yiteng Huang. Berlin: Springer, 2008. Chap. Text-Dependent Speaker Recognition, pp. 743–762.
- [8] Ram H. Woo, Alex Park, and Timothy J. Hazen. “The MIT Mobile Device Speaker Verification Corpus: Data Collection and Preliminary Experiments”. In: *Odyssey 2006: The Speaker and Language Recognition Workshop, San Juan, Puerto Rico, June 28-30, 2006*. IEEE, 2006, pp. 1–6.
- [9] 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.
- [10] 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.
- [11] Douglas A. Reynolds. “Speaker identification and verification using Gaussian mixture speaker models”. In: *Speech Communication* 17.1 (1995), pp. 91–108.
- [12] D. A. Reynolds. “Comparison of background normalization methods for text-independent speaker verification”. In: *Proceedings of the European Conference on Speech Communication and Technology*. 1997, 963–966.

- [13] Jared J. Wolf. “Efficient acoustic parameters for speaker recognition”. In: *Journal of the Acoustical Society of America* 51 (1972), pp. 2044–2056.
- [14] Hector N. B. Pinheiro. *Sistemas de Reconhecimento de Locutor Independente de Texto*. Trabalho de Graduação. Universidade Federal de Pernambuco, 2013.
- [15] Ethan. *Don’t you hear that?* May 10, 2010. URL: <http://scienceblogs.com/startswithabang/2010/05/10/dont-you-hear-that/>.
- [16] Harvey Fletcher and Wilden A. Munson. “Loudness, Its Definition, Measurement and Calculation”. In: *Bell Telephone Laboratories* 12.4 (1933), pp. 82–108.
- [17] 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.
- [18] Douglas O’Shaughnessy. *Speech Communications: Human and Machine*. Addison-Wesley, 1987.