

Universidade Estadual Paulista  
Instituto de Biociências, Letras e Ciências Exatas  
Departamento de Ciência da Computação e  
Estatística

Luis Fernando Teixeira Silva

Um sistema para reconhecimento de comandos falados  
dependente do locutor

São José do Rio Preto - SP

2017

Luis Fernando Teixeira Silva

Um sistema para reconhecimento de comandos  
falados dependente do locutor

Monografia apresentada ao Programa de  
graduação em Ciência da Computação da  
UNESP para obtenção do título de Bacharel.

Orientador: Prof. Dr. Rodrigo Capobi-  
anco Guido

São José do Rio Preto - SP

2017

Ficha catalográfica elaborada pelo Serviço de Biblioteca do IBILCE/UNESP

Luis Fernando Teixeira Silva

titulo

titulo

titulo. / fulano de tal; orientador

Rodrigo Capobianco Guido. São José do Rio Preto, 2017.

xxx p.

Monografia (TCC

TCC

TCC, 2017.

1. Processamento de sinais. 2. Reconhecimento de locutor. 3. Acústica. 4. Energia. 5. Escala *Bark*. I. Capobianco Guido, Rodrigo, orient.

II. Título.





Dedico este trabalho a todos os meus familiares, em especial aos meus pais, Nilda, Luis Carlos e a minha irmã Ana Beatriz.

Dedico também esse trabalho para a minha namorada Cristiana Luiza.



# Agradecimentos

Primeiramente, gostaria de agradecer meus pais e minha madrinha, pois sem o apoio deles eu nunca teria conseguido ter acesso a um ensino de qualidade que o cursinho alternativo me proporcionou durante todo o ano de 2012. Foi graças a essas 3 pessoas que pude ingressar nessa linda universidade.

Gostaria de agradecer também a minha irmã que nos momentos mais difíceis da minha graduação me deu forças para continuar em frente e concluir minha formação de bacharel em ciência da computação. Agradeço também a todos os meus familiares que me apoiaram ao longo dessa jornada de 5 anos.

E também deixo um agradecimento especial a meus dois grandes amigos João Cesar Granville e Luiz Gustavo Caobianco que tornaram esses anos na universidade mais felizes. Agradeço também minha namorada por ter me auxiliado nesses dois últimos anos de universidade e por me dar forças a concluir o curso nessa etapa final.





*“No fim tudo dá certo, e se não deu certo é porque ainda não chegou ao fim.”*

**Fernando Sabino**



# Resumo

TAL, F. *titulo*. 2016. xxxp. TCC UNESP 2016.

Atualmente, ....

Palavras-chave: Processamento de sinais. Reconhecimento de locutor. Acústica. Escala *Bark*.



# Abstract

TAL, F. *titulo*. 2016. xxxp. TCC UNESP 2017.

Nowadays, ...

Keywords: Signal processing. Speaker recognition. Acoustics. Bark scale.



# Lista de Figuras

Figura 2.1 - Fisiologia da voz [19] . . . . .	32
Figura 2.2 - (a) Sinal na forma analógica. (b) Amostragem. (c) Quantização, extraído de [18] . . . . .	34
Figura 2.3 - [1] . . . . .	38





# Lista de Tabelas

Tabela 2.1 - Estrutura de um arquivo <i>WAVE</i> . . . . .	35
------------------------------------------------------------	----



# Lista de Abreviaturas

<b>WAVE</b>	<i>Waveform Audio File Format</i>
<b>PCM</b>	<i>Pulse-Code Modulation</i>
<b>IBM</b>	<i>International Business Machines</i>
<b>RIFF</b>	<i>Resource Interchange File Format</i>
<b>fmt</b>	<i>format</i>
<b>IEEE</b>	<i>Institute of Electrical and Eletronic Engineers</i>



# Sumário

<b>1</b>	<b>Introdução</b>	<b>25</b>
1.1	Introdução . . . . .	25
1.2	Objetivo . . . . .	25
1.3	justificativa . . . . .	26
1.4	Motivação . . . . .	26
1.5	Metodologia . . . . .	27
1.6	Exequibilidade . . . . .	28
1.7	Organização do trabalho . . . . .	29
<b>2</b>	<b>Revisão Bibliográfica</b>	<b>31</b>
2.1	A Voz Humana . . . . .	31
2.2	Sinais Digitais . . . . .	32
2.3	Arquivos Acústicos no Formato <i>WAVE</i> . . . . .	34
2.4	Reconhecimento de Padrões e Vetores de Características . . . . .	36
2.5	Energia . . . . .	37
2.6	Similaridade baseada em distancias . . . . .	38
2.7	Trabalhos Correlatos . . . . .	39
<b>3</b>	<b>Detalhamento do Trabalho Proposto</b>	<b>41</b>
3.1	Estrutura do Sistema . . . . .	41
3.2	Coleta e elaboração do banco de áudios . . . . .	41
<b>4</b>	<b>Testes e Resultados</b>	<b>43</b>
<b>5</b>	<b>Conclusões e Trabalhos Futuros</b>	<b>45</b>

<b>Referências</b>	<b>45</b>
<b>Apêndice I - Gráficos das características extraídas</b>	<b>50</b>





# Capítulo 1

## Introdução

### 1.1 Introdução

*Speech recognition* ou reconhecimento de fala é uma área que ganhou destaque tanto no âmbito científico quanto no industrial. É comum encontrar nos *shoppings*, *smartphones*, *notebooks* sistemas que fazem uso da computação para poder realizar o reconhecimento de comandos, e essa realidade só tende a aumentar com os avanços da tecnologia em termos de *hardware* e *software*.

### 1.2 Objetivo

Este trabalho tem como objetivo implementar um algoritmo computacional desenvolvido em C/C++ para reconhecer comandos falados de modo *off-line* com locutor prédefinido, ou seja *speaker-dependent*. Esses comandos foram previamente gravados em arquivos no formato *WAVE* de 16 *bits* PCM.

## 1.3 justificativa

A área de *speech recognition* nos últimos anos tem ganhado um grande destaque no âmbito científico, o que pode ser facilmente verificado através de consultas nos mais diversos repositórios digitais. Além disso, a área de estudo do presente trabalho também tem ganhado destaque na indústria com a popularização dos *smartphones* com reconhecimento de voz, o que justifica ainda mais a elaboração desse projeto.

Ademais, convém ressaltar que o presente trabalho utiliza a computação junto com fórmulas matemáticas como parte de um ferramental para a resolução de um problema da área de *speech recognition*.

## 1.4 Motivação

Em 2008, é lançado o primeiro filme do Homem de Ferro, centrado no personagem Tony Stark. O filme além de despertar o interesse nas histórias em quadrinhos da Marvel, também prende a atenção dos cientistas da computação, uma vez que o personagem principal interage com uma inteligência artificial - Jarvis - que podia controlar a casa e a armadura do herói. Vale ressaltar que Jarvis além de fazer reconhecimento do locutor de modo *speaker-dependent*, também realiza síntese de voz. Já no ano de 2011, foi lançado a primeira versão do assistente pessoal do iOS, conhecido como Siri, que permite que o usuário execute determinadas funções do *smartphone* utilizando comandos falados. E mais recente, o criador do Facebook, decidiu no ano 2016, implementar seu próprio assistente pessoal, que funciona também de modo dependente do locutor, para auxiliar nas tarefas domésticas.

FALAR SOBRE ALGUM CONTEXTO CIENTIFICO TAMBEM.

E é a partir desse contexto que surgiu a inspiração para o desenvolvimento desse trabalho.

Inicialmente, o projeto tinha como proposta de conseguir auxiliar em determinadas funções de uma residência, além de executar comandos básico em um computador. Porém, o projeto necessitou de cortes em seu escopo para que se tornasse factível sua elaboração como forma de trabalho de conclusão de curso.

Assim, este projeto tem como objetivo iniciar o desenvolvimento de um futuro assistente pessoal, através da elaboração de um sistema para reconhecimento de comandos falados *speaker-dependent*.

## 1.5 Metodologia

Para a elaboração deste projeto foi determinado os seguintes 11 comandos:

- Bom dia, Logan;
- Bom noite, Logan;
- Oi, Logan;
- Como está o tempo hoje?;
- vai chover?;
- Abrir calculadora;
- Ver notícias;
- Pesquisar;
- Alarme;
- Calendário;
- Sair;

sendo que posteriormente foi realizada a gravação de 10 áudios para cada um dos 11 comandos referidos, totalizando 110 arquivos de áudio no formato MPEG-4. Tais arquivos foram convertidos para o formato *WAVE* de 16 *bits* PCM usando o programa *Audacity*. Vale ressaltar que todos os áudios foram gravados em um ambiente que proporciona-se um certo grau de isolamento sonoro, para assim se obter um som com menos ruído.

A partir dessa etapa inicial foi feita a extração dos dados brutos contidos nos arquivos *WAVE*. Para isso foi utilizada uma biblioteca fornecida pelo Prof.Dr.Rodrigo Capobianco Guido do Departamento de Ciência da Computação e Estatística (DCCE), IBILCE/Unesp. Tal biblioteca, escrita em C/C++, tem a função de separar o cabeçalho dos arquivos *WAVE*. A partir desse ponto, a biblioteca foi modificada para extrair os dados brutos e guardar as amplitudes dos sinais em arquivos de texto. Foi realizada a automação de todo o processo que exigia intervenção humana para a execução do algoritmo, como a passagem de áudios como parâmetro a cada nova execução, criação de arquivos para extração dos valores das amplitudes dos sinais, entre outros. Toda essa automatização foi implementada com a utilização de *scripts* escritos na linguagem *Shell script*.

Posteriormente a etapa de extração das amplitudes dos sinais digitalizados, foi realizada a extração das características dos áudios analisados. Essa etapa consiste na utilização do método A3 - descrito com maiores detalhes no capítulo 2 - para obter assim, vetores de características. Tal processo é essencial no presente trabalho, pois os valores obtidos no processo de extração são variáveis e demasiadamente grandes, sendo que o classificador exige valores menores e com tamanho fixo.

## 1.6 Exequibilidade

Para a elaboração desse projeto foram utilizadas ferramentas gratuitas como *Sublime*, Libre-Office, TeXstudio, *Audacity* e um computador pessoal. Além disso, foram também utilizados artigos científicos de acesso grátis como IEEE, e os acervo disponibilizado pela biblioteca física

localizada no Instituto de Biociências, Letras e Ciências Exatas da UNESP.

## 1.7 Organização do trabalho

A monografia está organizada a partir deste capítulo da seguinte forma:

- No Capítulo 2 apresenta-se uma série de trabalhos realizados na área *speaker-dependent* tanto a nível local quanto a internacional, exaltando a relevância da área no âmbito acadêmico. Além disso, é apresentado também neste capítulo toda a fundamentação teórica do trabalho, definindo e exemplificando os principais conceitos utilizados na elaboração do deste projeto.
- No Capítulo 3 apresenta-se uma breve descrição do estado do atual do trabalho e também um cronograma para finalização.



## Capítulo 2

# Revisão Bibliográfica

### 2.1 A Voz Humana

A voz é uma característica única dos seres humanos, que foi desenvolvida a partir da necessidade do homem em socializar, e assim poder se comunicar entre seus semelhantes. Em outras palavras, essa nossa característica foi adquirida através de um processo evolutivo, assim sendo não somos dotados de um aparelho específico para fala, apenas adaptamos um conjunto de órgãos da parte superior do aparelho digestivo e do respiratório para produção de fonemas.

De modo geral, podemos dividir o processo de síntese de voz em três partes:

1. **Pulmões** - São eles os responsáveis por gerar o fluxo de ar, que é ocasionado através da contração do órgão pelo diafragma. Tal fluxo, funciona como um "combustível para a voz.

#### 2. **Aparelho digestivo**

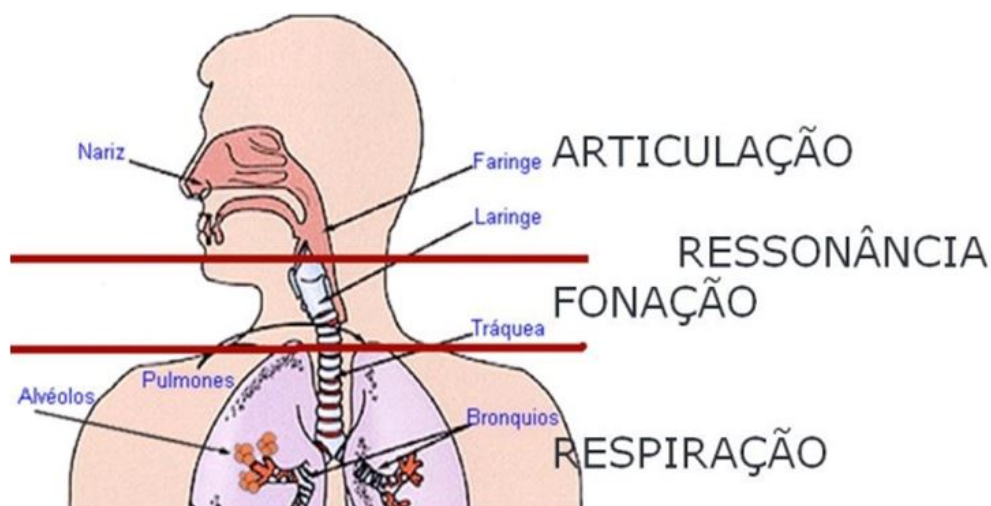
- 2.1. **Laringe** - Local onde são encontrada as pregas vocais, que são dois músculos que tem a responsabilidade de modelar a voz através de sua contração e/ou relaxamento. Quando o fluxo de ar que sai dos pulmões passa faringe, as cordas vocais vibram com maior ou menor intensidade, logo pode-se perceber que a voz humana é formada pela interação de duas principais forças: A força com que o fluxo de ar sai dos

pulmões e força muscular da laringe.

- 2.2. **Articuladores** - Eles são localizados na cavidade bucal e são responsáveis por fazer a articulação do fluxo de ar vindo dos pulmões. Sendo que os principais articuladores são: lábios, língua, dentes e mandíbula.

Vale ressaltar que a partir da característica vibratória das cordas vocais, é possível determinar a frequência da voz humana, o que nos torna capaz de fazer a distinção de gêneros e até mesmo reconhecimento de voz e comandos.

**Figura 2.1** – Fisiologia da voz [19]



## 2.2 Sinais Digitais

Os dados podem ser representados de duas maneiras: digital ou analógica. Sendo a primeira forma uma onda eletromagnética que pode assumir infinitos valores ao longo do tempo e um excelente exemplo desse tipo de sinal, é a voz humana. Já a segunda forma tem um número finito de valores discretos no tempo. Em outras palavras, o sinal digital assume valores em determinados instantes de tempo e sua representação pode ser baseada na discretização de um sinal analógico.



Para se realizar uma análise em um sinal analógico em um computador, precisamos necessariamente fazer sua conversão para a forma digital, uma vez que os computadores possuem uma capacidade finita de processamento, o que o permite tratar apenas valores discreto, não sendo possível assim processar os infinitos pontos de um sinal analógico. Esse processo de conversão é denominado de digitalização, que é realizado através das 3 etapas. Na figura 2.1 é apresentado o processo de amostragem e quantização que serão comentados a seguir.

1. **Amostragem:** Consiste no processo de obter amostras de um sinal analógico em instantes de tempo com espaços iguais. Vale ressaltar que para se obter sucesso nesse processo de amostragem, é necessário seguir o Teorema de Nyquist. Tal teorema define que para realizar a amostragem de um sinal, é necessário no mínimo ter uma taxa de amostra duas vezes maior que a máxima frequência do sinal analógico. Caso não se obtenha tal taxa de amostragem ocorrerá um efeito chamado de *aliasing*, que consiste na medição errônea de uma frequência mais alta como sendo mais baixa. Um bom exemplo da aplicação de tal teorema é na digitalização da voz humana, que possui uma frequência máxima de 4 mil Hertz. Logo, para realizar amostragem desse sinal, seria necessário no mínimo 8 mil amostras por segundo.
2. **Quantização:** Esse processo consiste atribuição de valores discretos para as amplitudes do sinal analógico, que pertencem a um intervalo contínuo de valores. Cada amplitude é alocada ao nível de valor discreto mais próximo, o que diminui assim os erros absolutos. Vale ressaltar que o número de níveis é definido pelo número de *bits* que serão utilizados na etapa posterior de codificação, que é definido como sendo  $2^n$ , onde  $n$  é o número de bits que será usado. De modo geral, no sinal de áudio é utilizada a quantização de 16 *bits*, em outras palavras, as amplitudes de cada amostra podem assumir  $2^{16} = 65536$  valores.

Para realizar a quantização, pode-se utilizar duas técnicas diferentes para diminuir a quantidade de erros do processo:

- 2.1. **Quantização uniforme:** Essa técnica é aplicado em sinais que possuam uma pequena diferença entre sua amplitude máxima e mínima. O processo consiste em

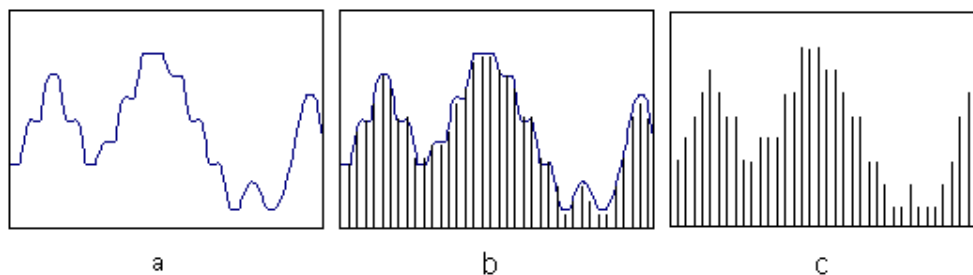
dividir de forma igualmente espaçados os intervalos de amplitude, ou seja de forma uniforme.

- 2.2. **Quantização não uniforme:** É utilizada para sinais que possuam um range dinâmico alto, ou seja existe uma grande diferença entre a amplitude máxima e mínima. Esse processo consiste em encontrar valores adequados para região do sinal, o que exige diferentes quantizadores. Uma solução alternativa para esse tipo de abordagem, seria comprimir o sinal, e logo após utilizar a técnica de quantização uniforme.

Vale ressaltar que cada técnica descrita anteriormente dependerá do sinal a ser quantificado.

3. **Codificação:** Etapa que consiste na modificação de um sinal para uma aplicação específica, como por exemplo o armazenamento ou transmissão de dados.

**Figura 2.2** – (a) Sinal na forma analógica. (b) Amostragem. (c) Quantização, extraído de [18]



## 2.3 Arquivos Acústicos no Formato *WAVE*

*Waveform audio file format* é a abreviação de *WAVE* ou simplesmente *WAV*, que é um tipo de formato de arquivo de áudio que foi desenvolvido pela *Microsoft* em conjunto com a *IBM*. O formato *WAVE* é amplamente utilizado em uma variedade de trabalhos, sejam eles científicos ou profissionais, visto que o formato permite uma fiel representação dos dados digitalizados, uma vez que os dados digitalizados podem ser armazenados sem sofrer obrigatoriamente um

processo de compressão, o que evita perdas. Porém, devido a essa característica o *WAV* ocupa muito mais espaço que os demais formatos de arquivos de áudios.

Na Tabela 2.1 pode-se ver a estrutura de um arquivo *WAVE*. Basicamente o arquivo é dividido em 2 grandes blocos, sendo o primeiro bloco um cabeçalho *RIFF* e o segundo bloco é dividido em dois sub-blocos, sendo um com informações referentes ao formato *WAVE* e o outro com os dados do áudios.

**Tabela 2.1** – Estrutura de um arquivo *WAVE*

Classe	Posição ( <i>bytes</i> )	Tamanho ( <i>bytes</i> )	Descrição
Cabeçalho	0	4	Apresenta o identificador do cabeçalho - "RIFF".
Cabeçalho	4	4	Tamanho do arquivo sem o identificado do cabeçalho.
Cabeçalho	8	4	Mostra o identificador <i>WAVE</i> .
Formato	12	4	Mostra o identificador do segundo bloco - "fmt".
Formato	16	4	Tamanho do bloco sem o identificador.
Formato	20	2	Mostra se o arquivo é do tipo PCM ou se tem alguma compressão.
Formato	22	2	Mostra a quantidade de canais.
Formato	26	4	Apresenta o valor da taxa de amostragem.
Formato	30	4	Apresenta a taxa de <i>bytes</i> .
Formato	32	2	Demonstra a quantidade de <i>bytes</i> para uma amostra.
Formato	34	2	Demonstra a quantidade de <i>bits</i> para cada amostra.
Dados	36	4	Apresenta o identificador do terceiro bloco - "data".
Dados	40	4	Mostra o tamanho do bloco sem o identificador.
Dados	44	4	Demonstra os dados reais da música.

Vale ressaltar que os valores mais comuns para cada amostra de um arquivo *WAVE* pode ser 8 *bits* ou 16 *bits*. Um áudio de 8 *bits* significa que o valor da amplitude do sinal, de cada amostra, pode ser representado por 256 valores, sendo 127 positivos e 128 negativos. Já para um arquivo 16 *bits* a amplitude do sinal pode ser representado por 65536 valores, com 32757 positivos e 32768 negativos. Para *WAVE* de 16 *bits* é utilizada a codificação de complemento de 2 para representar o valor da amplitude do sinal. Assim, o valor do *bit* mais significativo representa se o sinal é negativo ou positivo.

Nesse trabalho foi utilizado o formato *WAV* de 16 *bits* PCM (*Pulse-code Modulation*) que não utiliza compressão, para se obter assim uma melhor qualidade na elaboração deste projeto final. Foi fornecida uma biblioteca escrita em C/C++ pelo orientador para isolar o primeiro bloco referente ao cabeçalho *RIFF* e o sub-bloco de formato *WAV* dos dados brutos, que contém

as amplitudes dos sinais de voz digitalizados.

## 2.4 Reconhecimento de Padrões e Vetores de Características

Reconhecimento de Padrões também conhecido como *Pattern Recognition* ou simplesmente como PR, é uma área de estudo que tem como objetivo classificar objetos em determinadas classes ou categorias com base na extração de características relevantes de tais objetos. De modo geral, as etapas do PR pode ser dividida da seguinte maneira:

- **Extração das características:** Principal etapa de todo o processo de *Pattern Recognition*, pois é nesse momento que é realizada a redução de dimensionalidade do objeto. Caso esse processo seja mal elaborado, ocorrerá perdas de características significativas, que podem tornar a classificação dispendiosa. Portanto, é necessário ter um conhecimento específico sobre o problema, para poder realizar a redução de dimensionalidade sem que ocorra perdas de informações relevantes do objeto e também que diminua o esforço computacional. Nessa etapa é realizada também a extração das características do objeto, logo é interessante realizar uma seleção de modo que tais características sejam relevantes na classificação.
- **Classificação do objeto:** Etapa que determina os procedimentos para realizar a classificação do objeto em uma determinada classe. Diferentemente da fase anterior, o classificador pode ser definido independente do problema, uma vez que os métodos usados são os mesmos.

É nessa etapa que o classificador aprende a distinguir a qual classe o objeto pertencerá, de modo que esse aprendizado pode ser classificado em duas principais categorias:

- **Aprendizagem supervisionada:** Seu funcionamento é baseado em exemplos de entradas e saídas previamente conhecidos, o que permite assim, aprender uma re-

gra genérica para classificar as entradas posteriores. Nesse trabalho é utilizada a distância Euclidiana para tomar a decisão de classificação dos comandos falados.

- **Aprendizagem não supervisionada:** Ao contrário da aprendizagem supervisionada, nessa técnica não é utilizado nenhum conhecimento prévio para se classificar os objetos.

Vale ressaltar que o conceito de característica definido aqui é entendido como qualquer medida útil que possa ser extraída no processo de identificação de um padrão. Nesse trabalho a característica escolhida para fazer a classificação dos sinais foi sua energia. Ao final do processo de extração todos os áudios foram convertidos em vetores de características com tamanhos fixos.

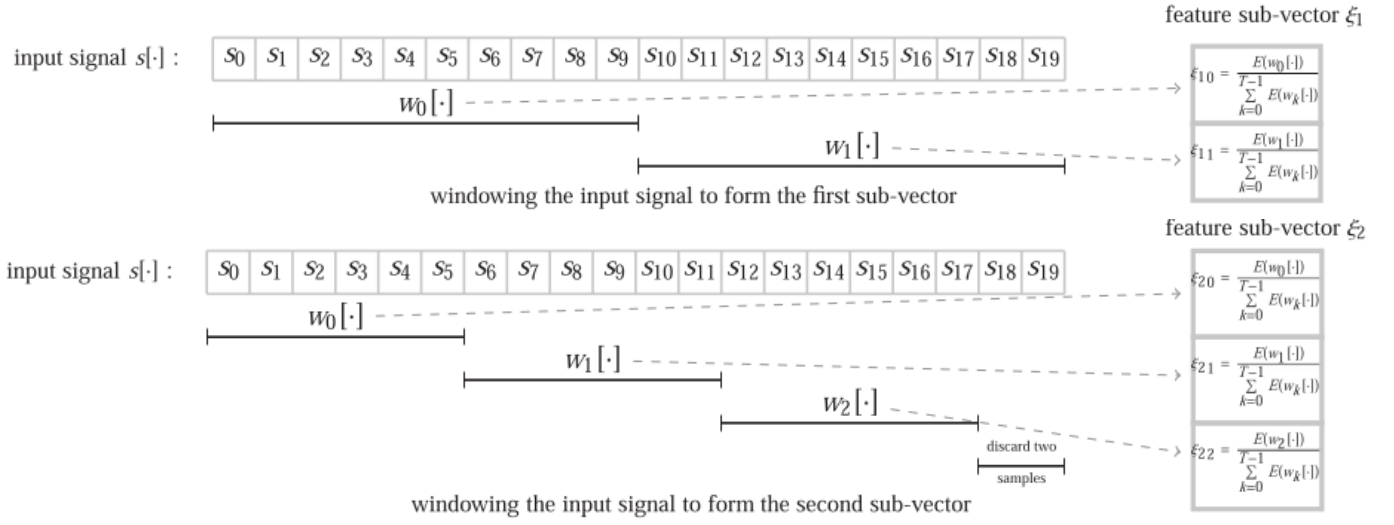
## 2.5 Energia

A definição de energia está relacionada ao conceito de conseguir realizar trabalho. Neste projeto, será considerada energia a capacidade das estruturas vocálicas e dos pulmões de produzir um sinal acústico. Na equação 2.1 definida a energia total  $E(s[.])$ , de um dado sinal de áudio digitalizado  $s[.]$ , de tamanho  $M$ .

$$E(s[.]) = \sum_{i=0}^{M-1} (S i)^2 \quad (2.1)$$

Para realizar a captura das características dos áudios foi utilizado o método A3, que foi referenciado pelo próprio orientador em (1), conforme é mostrado na figura 2.3. Para realizar a extração das características, tal método se baseia no em determinar tamanhos ou áreas proporcionais para atingir níveis predefinidos da energia do sinal que se encontra em análise. A3 é ideal para avaliar os níveis de energia de um sinal de voz digitalizado que foi gerado por um agente.

Vale ressaltar que A3 defini um nível crítico de energia que varia de 0 a 100%. Em um sinal de áudio, A3 extrai um vetor de características dividindo o áudio em partes proporcionais ao valor definido pelo nível crítico - parâmetro  $C$  - de modo que, cria uma janela de tamanho  $C\%$

**Figura 2.3** – [1]

e extrai a característica dessa faixa de valores do sinal digitalizado, conforme é mostrado na equação 2.2, onde  $\epsilon$  é uma característica extraída,  $E(W_0[\cdot])$  a energia obtida do início do sinal até o ponto final da janela definido por  $C$ ,  $E(W_k[\cdot])$  é a energia total do sinal. Tal processo é repetido até que se atinja o montante total de energia do áudio, de modo que, defini-se assim um vetor de características.

$$\epsilon = (E(W_0[\cdot]) / (\sum_{k=0}^{T-1} E(W_k[\cdot]))) \quad (2.2)$$

## 2.6 Similaridade baseada em distancias

Quando precisamos agrupar objetos, normalmente utilizamos como quesito a similaridade ou dissimilaridade entre tais objetos. A primeira mede o quanto dois ou mais objetos são parecidos entre si, enquanto que a segunda analisa o quanto dois ou mais objetos são diferentes.

Em específico para avaliar similaridades em determinados objetos, podemos usufruir de uma série de métodos, sendo um dos mais conhecidos a distância Euclidiana, que pode ser calculada com base na equação 2.3, onde  $P = \{P_1, P_2, \dots, P_n\}$  e  $Q = \{Q_1, Q_2, \dots, Q_n\}$  são dois vetores de

características e a distância  $d$  será o nível de similaridade entre os dois vetores. A similaridade entre dois objetos será máxima quando  $P = Q$ , caso seja igual à zero, será mínima.

$$d(P, Q) = \sqrt{(P_1 - Q_1)^2 + (P_2 - Q_2)^2 + \dots + (P_n - Q_n)^2} = \sqrt{\sum_{i=1}^n (P_i - Q_i)^2} \quad (2.3)$$

Outro método para medir similaridade é através do uso da distância absoluta. Tal técnica utiliza o módulo da diferença entre dois números para realizar a medição, conforme é mostrado na equação 2.4. Caso a diferença seja igual a zero, temos que os dois objetos comparados são idênticos, caso contrário, serão diferentes.

$$d = |X - Y| \quad (2.4)$$

O classificador que é elaborado nesse presente trabalho, foi implementado com o método de distância Euclidiana, visto que tal métrica oferece além de simplificação, ótimos resultados para reconhecimento de padrões em sinais acústicos (2).

## 2.7 Trabalhos Correlatos

A área de *speech recognition* tem sido bastante estudada em diversos trabalhos. Nas monografias de conclusão de curso [2, 4, 5, 6, 12, 13, 14] é discutido temas sobre o reconhecimento de vocabulário restrito. Já nos artigos [15, 16] os esforços são voltados para o reconhecimento de um amplo vocabulário, tal como idiomas.

Em [2], o autor utiliza conceitos como energia, limiar *hard*, níveis críticos de energia, para elaborar um sistema de reconhecimento de vocabulário restrito implementado na linguagem C/C++, com um classificador baseado em distância Euclidiana.

Em [4], o autor utiliza duas abordagens diferentes para realizar o reconhecimento de voz, sendo a primeira baseada em *pattern-matching* e a outra em *knowledge-based*. Ao final do

trabalho pode-se concluir e verificar que a segunda abordagem apresentou melhores resultados, pois utilizava inteligência artificial.

Particularmente em [5], o autor propõem um sistema para reconhecimento de voz independente do locutor. Para isso, o autor utiliza conceitos como *Zero Crossing Rate* e energia para realizar o pré-processamento, extração de características e um classificador também baseado em distância Euclidiana.

Em [7], o autor desenvolve um sistema de reconhecimento de voz baseado em automação residencial via wireless. Para a elaboração desse projeto, o autor utiliza o dispositivo HM2007L IC, para assim realizar o reconhecimento de comandos simples, como "acender", "apagar", "acender a luz", etc.

Em [8], o objetivo do autor é melhorar a relação homem-máquina, e com isso, desenvolve um método para reconhecimento de voz que uso de características *voice-likeness*.

Em [9] o autor tem como objetivo desenvolver um sistema para *smartphone* multifuncional que além de realizar o reconhecimento de voz, também a faça a conversão do diálogo reconhecido para texto.

Em [13], o autor utiliza conceitos como MFCC, HMM (*Hidden Markov Models*), para reconhecer um restrito vocabulário que contem dígitos de 0 a 9.

Já em [16], é realizada uma comparação de desempenho de GMMs - Modelos de Misturas Gaussianas - com DNNs - Redes Neurais Profundas - no reconhecimento de um amplo vocabulário Chines.

Em [17], com o objetivo de melhorar a eficácia do reconhecimento de um vocabulário Russo amplo, foi utilizado técnicas *knowledge-based* e estatísticas para modelagem de fonemas.



## Capítulo 3

# Detalhamento do Trabalho Proposto

### 3.1 Estrutura do Sistema

### 3.2 Coleta e elaboração do banco de áudios

Inicialmente, foram definidos os 11 comandos que serão reconhecidos pelo sistema, conforme já mencionado. Os comandos foram gravados 10 vezes em diferentes dias e horários, para obter se assim, uma melhor veracidade e fidelidade à voz do locutor, pois o mesmo pode sofrer alterações significativas com base na variação do seu humor ou estado físico. Como também já mencionado, a gravação dos arquivos de áudio foi realizada em um ambiente fechado, para diminuir assim, a probabilidade de ruídos nos sinais. Todos os arquivos foram gravados no formato MPEG-4, que é o formato disponível no *software* de gravação de áudio do *Windows 10*, e posteriormente foram convertidos para o formato *WAVE* de 16 *bits* PCM, com o auxílio do editor de áudio *Audacity*. Ao total, foram gravados e convertidos 110 arquivos de áudio.



## **Capítulo 4**

### **Testes e Resultados**

bla bla bla...



## **Capítulo 5**

### **Conclusões e Trabalhos Futuros**

Neste trabalho, ...



## Referências

- 1 GUIDO, R. C. A tutorial on signal energy and its applications. *Neurocomputing*, v. 179, p.264-282, 2016.
- 2 DORDAN, M, K. Verificação de Locutores dependente do discurso baseada na Evolução do Esforço Vocálico. IBILCE, UNESP, São José do Rio Preto (Trabalho de Conclusão de Graduação em Ciência da Computação), 2015.
- 3 SILVA, V. F. Identificação e Classificação de Gêneros Musicais com Abordagens Múltiplas de Reconhecimento de Padrões. IBILCE, UNESP, São José do Rio Preto (Trabalho de Conclusão de Graduação em Ciência da Computação), 2014.
- 4 CAOBIANCO, A. P. Comparação de Abordagens *PatternMatching* e *Knowledge-Based* para Reconhecimento de Locutor Dependente de Texto. IBILCE, UNESP, São José do Rio Preto (Trabalho de Conclusão de Graduação em Ciência da Computação), 2011.
- 5 MASCHIO, J. V. D. Projeto e Implementação Acústico-Computacional de Palavras Isoladas. IBILCE, UNESP, São José do Rio Preto (Trabalho de Conclusão de Graduação em Ciência da Computação), 2017.
- 6 LOUREIRO, W. d. F. Reconhecimento de Comandos Falados Português Brasileiro com Parâmetros de Frequência e Tempo. IBILCE, UNESP, São José do Rio Preto (Trabalho de Conclusão de Graduação em Ciência da Computação), 2015.
- 7 PAUL, A.; PANJA, M.; BAGCHI, M.; DAS, N.; MAZUMDER, R. M.; GHOSH, S. Voice Recognition Based Wireless Room Automation System. 2016 International Conference on Intelligent Control Power and Instrumentation, 2016.
- 8 SAKANO, T.; KIGAWA, T.; SUGIMOTO, M.; KUSUNOKI, F.; INAGAKI, F.; MIZOGUCHI, H. Laughing Voice Recognition Using Periodic Waveforms and Voice-likeness Features. Proceedings of the 2016 IEEE International Conference on Robotics and Biomimetics Qingdao, China, 2016.

- 9 CHERN, A.; TSAO, Y.; CHANG, R.; HSIU-WEN, C.; YING-HUI, L. A Smartphone-Based Multi-Functional Hearing Assistive System to Facilitate Speech Recognition in the Classroom. Disponível em: <<http://ieeexplore.ieee.org/document/7938619/>>. Acesso em: 05 jun. 2017.
- 10 RABINER, L. R.; JUANG, B. H. Fundamentals of Speech Recognition. Prentice Hall, 1993.
- 11 DUDA, R. O.; HART, P. E.; STORK, D. G. Pattern Classification. 2. ed. Wiley-Interscience, 2000.
- 12 THEODORIDIS, S.; KOUTROUMBAS, K. Pattern Recognition. 4. ed, Academic Press, 2008.
- 13 SILVA, A. G. d. Reconhecimento de Voz para Palavras Isoladas. Monografia (Graduação em Engenharia da Computação) - UFPE, Recife, 2009.
- 14 LOUZADA, J. Reconhecimento Automático de Fala por Computador. Monografia(Graduação em Ciência da Computação) - PUC, Goiás, 2010.
- 15 Piccoli, E. E. M. Reconhecimento de Padrão em Áudio no Formato Wave. IBILCE, UNESP, São José do Rio Preto (Trabalho de Conclusão de Graduação em Ciência da Computação), 2014.
- 16 LI, X.; YANG, Y.; PANG, Z.; WU, X. A Comparative Study on Selecting Acoustic Modeling Units in Deep Neural Networks Based large Vocabulary Chinese Speech Recognition. Neurocomputing, v. 170, p. 251-256, 2015
- 17 KARPOV, A.; MARKOV, K.; KIPYATKOVA, I.; VAZHENINA, D.; RONZHIN, A.; Large Vocabulary Russian Speech Recognition Using Syntactico-statistical language Modeling. Speech Communication, v.56, p.213-228, 2014.
- 18 Figura extraída de: <<http://penta3.ufrgs.br/RNP/cap3/3.2>
- 19 Figura extraída de: <<http://www.koreapost.com.br/wp-content/uploads/2016/01/fisiologia-da-voz.jpg>>. Acesso em: 18 Abril 2017.





## **Apêndice I - Gráficos das características extraídas**