

Universidade Estadual de Campinas
Faculdade de Engenharia Elétrica e de Computação
Departamento de Comunicações

Reconhecimento Automático de Fala Contínua
Empregando
Modelos Híbridos ANN + HMM

por

Edmilson S. Morais
Eng. Eletricista (Universidade de Brasília) 1995

Orientador : Prof. Dr. Fábio Violaro
DECOM - FEEC - UNICAMP

Este exemplar corresponde a redação final da tese
defendida por Edmilson S. Morais -
e aprovada pela Comissão
Julgada em 17.1.12.1997
Fábio Violaro
Orientador

Dissertação apresentada à Faculdade de Engenharia
Elétrica e de Computação da UNICAMP como re-
quisito parcial para a obtenção do título de Mestre
em Engenharia Elétrica.

Campinas, dezembro de 1997

M792r
00507100

UNICAMP
SERVIDOR CENTRAL

Aos meus pais Espedito e Dirce.

Reconhecimento Automático de Fala Contínua
Empregando
Modelos Híbridos ANN + HMM

por

Edmilson S. Morais

Dissertação apresentada à Faculdade de
Engenharia Elétrica e de Computação da UNICAMP
em Dezembro de 1997, como requisito parcial para a obtenção do título de
Mestre em Engenharia Elétrica

Resumo

Atualmente, os sistemas que representam o estado-da-arte em reconhecimento de fala contínua baseiam-se em modelos ocultos de Markov - HMM (“Hidden Markov Models”), uma estrutura duplamente estocástica capaz de modelar tanto as *variabilidades acústicas* como *temporais* do sinal de fala. Porém, para viabilizar o modelamento matemático de um HMM, são realizadas inúmeras suposições simplificadoras que limitam o seu potencial efetivo. Redes neurais artificiais - ANN (“Artificial Neural Networks”) não necessitam fazer uso de muitas destas suposições, podem aprender e generalizar superfícies complexas de decisão, tolerar ruídos e suportar paralelismo. Todas estas vantagens tornam as ANNs extremamente poderosas para modelar as variabilidades acústicas da fala. Entretanto, ao contrário dos HMMs, as ANNs não têm se mostrado eficientes para o modelamento das variabilidades temporais. Com o objetivo de unir em uma única estrutura o que há de melhor nas tecnologias de redes neurais artificiais e de modelos ocultos de Markov, têm sido estudados e avaliados nos últimos sete anos [36, 14, 3, 33, 21, 8, 1], modelos híbridos ANN-HMM nos quais o modelamento das variabilidades acústicas é confiado à ANN enquanto o HMM responsabiliza-se pela absorção das variabilidades temporais.

Os objetivos desta Tese foram o estudo e a implementação dos principais módulos de um sistema para reconhecimento de fala contínua baseado em modelos híbridos ANN-HMM. As análises desenvolvidas concentram-se em três tópicos fundamentais : (1) Capacidade das ANNs, treinadas a partir dos critérios MSE (“Minimum Square Error”) ou entropia relativa, de estimarem verossimilhanças de emissão de símbolos de HMMs contínuos, (2) Algoritmo REMAP (“Recursive Estimation and Maximization of A Posteriori Probabilities”) para a reestimação dos parâmetros dos modelos híbridos, (3) Algoritmos de busca para a determinação da seqüência mais provável de palavras que compõem uma determinada sentença a ser reconhecida.

Para análise de desempenho são apresentados resultados da avaliação do sistema no reconhecimento de 100 sentenças (constituídas de um total de 319 palavras distintas) dependentes de locutor e gravadas sob condições de estúdio. A taxa de acerto de palavras do sistema sem o uso de restrições gramaticais foi de 83,26% e, com o uso de restrições gramaticais do tipo *pares de palavras*, foi de 99,47%.

Sumário

1	Introdução	11
1.1	Considerações Iniciais	11
1.2	Objetivos	12
1.3	Estrutura da Tese	13
2	Reconhecimento de Fala Contínua	15
2.1	Introdução	15
2.1.1	Pré-Processamento	16
2.1.2	Modelamento Acústico	16
2.1.3	Modelo da língua	18
2.1.4	Decodificação	19
2.2	Definições e Notações	19
2.3	Modelos Ocultos de Markov - HMM	21
2.3.1	Introdução	21
2.3.2	Definição	21
2.3.3	Algoritmos Básicos	22
2.4	Discussão	23
2.4.1	Crítérios de Treinamento - Problemas com o Critério ML	23
2.4.2	Mais Problemas com o Critério ML	25
2.4.3	Alternativa - Modelos Híbridos ANN-HMM	26
3	Modelos Híbridos ANN-HMM	27
3.1	Redes Neurais Artificiais - ANNs	27
3.1.1	Multilayer Perceptrons - MLPs	27
3.1.2	Outras Arquiteturas Conexionistas	29

3.2	ANNs como Estimadores Estatísticos	30
3.2.1	O Somatório das Saídas da ANN Deve Ser Igual a Um	34
3.3	Estimação de Verossimilhanças para HMMs através de ANN	34
3.4	Vantagens da ANN	35
3.5	Tipos de Modelos Híbridos ANN-HMM	36
3.5.1	Modelo Híbrido ANN-HMM Padrão	36
3.5.2	Modelo Híbrido ANN-HMM Discriminativo	38
4	REMAP	41
4.1	Introdução	41
4.2	Estimação de Parâmetros - Pré-REMAP	42
4.3	Reestimação de Parâmetros - REMAP	44
4.3.1	Alvos Abruptos	44
4.3.2	Alvos Suaves	44
4.3.3	Algoritmo para Reestimação	48
4.3.4	Verossimilhanças de Emissão de Símbolos	50
4.3.5	Probabilidades de Transição	50
5	Algoritmos de Busca	52
5.1	Introdução	52
5.2	Level Building	53
5.2.1	Motivação	53
5.2.2	Algoritmo Level Building	55
5.2.3	Incorporando Restrições de Duração ao Algoritmo Level Building	58
5.2.4	Exemplo	59
5.3	One Stage	63
6	Sistema Implementado	65
6.1	Base de Dados	65
6.1.1	Introdução	65
6.1.2	Unidades Sub-Lexicais (fones)	66
6.1.3	Sentenças	66
6.1.4	Dicionário de Pronúncias	66

6.2	Pré-Processamento	66
6.3	Modelos Híbridos ANN-HMM	70
6.3.1	Rede Neural Artificial	70
6.3.2	Modelos Ocultos de Markov	71
6.4	Treinamento	72
6.4.1	Estimação - Pré-REMAP	72
6.4.2	Reestimação - REMAP	73
6.5	Decodificação	73
6.5.1	Algoritmo Level Building	73
6.5.2	Modelo de Duração de Palavras	74
6.5.3	Restrições Gramaticais	74
7	Resultados Experimentais	76
7.1	Introdução	76
7.2	Reconhecimento a Nível de Unidades Sub-Lexicais (fones)	77
7.2.1	Taxas de Acertos da ANN ao Longo das Épocas de Treinamento	77
7.2.2	Taxas de Acertos Finais da ANN	77
7.3	Reconhecimento de Sentenças com uso de Modelos de Sentenças	80
7.3.1	Exemplo	80
7.3.2	Discussão	82
7.4	Desempenho do Algoritmo de Reestimação - REMAP	83
7.4.1	Reestimação dos parâmetros Ψ da ANN	83
7.4.2	Reestimação das Probabilidades de Transição	86
7.4.3	Reestimação das Probabilidades a Priori das Classes	86
7.4.4	Discussão	88
7.5	Reconhecimento de Sentenças com uso de Modelos de Palavras	89
8	Considerações Finais	92
8.1	Análise de Desempenho	92
8.1.1	Reconhecimento de Fones	92
8.1.2	Reconhecimento de Sentenças com o uso de Modelos de Sentenças	93
8.1.3	Reconhecimento de Sentenças com o uso de Modelos de Palavras	93
8.2	Contribuições	94

<i>SUMÁRIO</i>	6
8.3 Sugestões para Trabalhos Futuros	94
A Critério MMI	96
B Dicionário de Pronúncias	98
C Lista com as Sentenças Reconhecidas	105

Lista de Figuras

2-1	Diagrama de blocos de um sistema para reconhecimento estatístico de fala contínua. . .	17
2-2	Estimativa da seqüência de atributos (vetores acústicos) associados a uma seqüência de amostras de um sinal de fala	17
2-3	Trajectoria de Viterbi : define a seqüência de estados com maior verossimilhança de ter gerado a seqüência de símbolos X.	23
3-1	Rede MLP para estimativa das probabilidades das classes dado o vetor acústico de entrada e mais alguma informação contextual.	31
3-2	Rede neural recursiva para estimativa das “probabilidades de transição condicionadas” a serem utilizadas em um modelo híbrido ANN-HMM discriminativo.	40
4-1	Sentença : “ É suficiente ” segmentada em termos de unidades sub-lexicais (fones). . .	42
4-2	Ilustração do processo de retirada dos exemplos de treinamento da ANN, centrados entre as marcas de segmentação das unidades sub-lexicais.	43
4-3	Exemplos utilizados durante a etapa de avaliação - exemplos tomados a cada quadro de análise.	44
4-4	Alvos suaves obtidos segundo as probabilidades a posteriori dos estados dada toda a seqüência de vetores acústicos “probabilidades a posteriori globais”.	45
4-5	Diagrama de blocos do processo de reestimação de parâmetros REMAP.	49
5-1	Avaliação de todas as palavras do léxico (eixo i) no primeiro nível de procura.	56
5-2	Avaliação de todas as palavras do léxico (eixo i) no segundo nível de procura.	58
5-3	Avaliação de todas as sentenças constituídas por apenas uma única palavra.	60
5-4	Três das nove sentenças possíveis de serem contruídas a partir da concatenação de duas palavras.	61

5-5	Avaliação das palavras candidatas ao segundo nível da sentença, após a redução do primeiro nível (utilizando apenas as informações realmente úteis no primeiro nível).	62
5-6	Ilustração do procedimento de “backtracing”. Caso em que a sentença reconhecida é composta por três palavras e consiste na concatenação dos dígitos : Sete - Seis - Três.	63
6-1	Histograma da amplitude dos componentes de todos os vetores acústicos que formam os exemplos de treinamento X_e	69
6-2	Histograma da amplitude dos componentes de todos os vetores acústicos utilizados no treinamento, após a normalização.	69
6-3	HMM “left-right” da palavra “casa” com um único estado por fone.	71
7-1	Taxas de acertos da ANN para os exemplos de validação cruzada em função do número de épocas de treinamento.	77
7-2	Taxa de acertos da ANN tanto para os exemplos de treinamento $T_{train.}$, como para os exemplos de validação T_{VC} , em função do número de épocas de treinamento - da época 171 a 249.	78
7-3	Taxa de acertos global, que é dada por uma média ponderada entre as taxas de acerto para os exemplos de validação cruzada e para os exemplos de treinamento, $T_G = 0,2 \cdot T_{VC} + 0,8 \cdot T_{train.}$ - da época 171 a 249.	78
7-4	Saída da ANN convertida em verossimilhança de emissão de símbolos, para os vetores acústicos de entrada correspondentes à sentença, “O saldo é suficiente”.	81
7-5	Modelos ocultos de Markov associados às sentenças : (a) “É suficiente”, (b) “Isto é suficiente”, (c) “O saldo é suficiente” e (d) “O saldo de sua conta é suficiente”.	81
7-6	Log’s das verossimilhanças normalizadas para os caminhos de Viterbi apresentados na Figura 7-5.	82
7-7	Resultados da apresentação da seqüência de vetores acústicos associados à sétima sentença da Tabela 6.2, aos modelos das vinte sentenças desta mesma tabela.	83
7-8	Alvos suaves $P(q_k X, M, \Theta)$ para a sentença “É suficiente”.	84
7-9	Valores máximos dos alvos $P(q_k X, M, \Theta)$ para todos os estados k , para cada um dos vetores acústicos $x_n \in X$	84
7-10	Saídas da ANN correspondentes aos fones presentes na sentença “É suficiente”, antes da reestimação, Pré-REMAP.	85

7-11 Saídas da ANN correspondentes aos fones presentes na sentença “É suficiente”, após duas reestimações, REMAP-2.	85
7-12 Probabilidades de permanência nos estados $P(q_k^n q_k^{n-1}, \Theta)$	87
7-13 Probabilidade de transição de estados $P(q_k^n q_j^{n-1}, \Theta)$	87
7-14 Probabilidades das 36 classes (fones) calculadas segundo a frequência relativa e segundo a Equação 6.4 para os casos : Pré-REMAP e REMAP-2.	88
7-15 Log's das verossimilhanças normalizadas correspondentes a apresentação da seqüência de vetores acústicos associada à sétima sentença da Tabela 6.2, aos 20 modelos de sentenças desta mesma tabela.	89

Lista de Tabelas

6.1	Lista com os tipos de fones utilizados seguidos por exemplos, e suas respectivas quantidades e durações médias.	67
6.2	Lista com 20 das 100 sentenças a serem utilizadas durante o treinamento e avaliação do sistema.	68
7.1	Taxa de acertos da ANN após o Pré-REMAP, para cada um dos 36 tipos de fones. . .	79
7.2	Resultados da avaliação do sistema híbrido ANN-HMM implementado, no reconhecimento das 100 sentenças do Apêndice C.	90
7.3	Resultados da avaliação de um sistema em desenvolvimento no LPDF UNICAMP, baseado apenas em HMM discreto, no reconhecimento das 100 sentenças apresentadas no Apêndice C.	90

Capítulo 1

Introdução

1.1 Considerações Iniciais

Muitos problemas de reconhecimento de padrões, de fundamental importância hoje em dia, são de natureza inerentemente seqüencial. Alguns exemplos incluem o reconhecimento de movimentos humanos através de uma seqüência de imagens de vídeo, ou o reconhecimento de fala contínua a partir de uma seqüência de amostras de um sinal de voz. Nos últimos anos, vários pesquisadores têm apresentado soluções baseadas em modelos ocultos de Markov - HMM ("Hidden Markov Models") para o reconhecimento de fala contínua [16, 27]. Os modelos ocultos de Markov são estruturas duplamente estocásticas que têm se mostrado eficientes para modelar, tanto as *variabilidades acústicas* como *temporais do sinal de fala*. O critério de treinamento mais usual de um HMM baseia-se no método da máxima verossimilhança - ML ("Maximum Likelihood") [27], porém este método apresenta um baixo poder discriminativo, uma vez que procura maximizar a verossimilhança de um determinado modelo gerar uma dada seqüência observada, mas não se preocupa em minimizar a verossimilhança dos outros modelos gerarem esta mesma seqüência. Um critério mais adequado para o treinamento de um HMM seria o de maximização da informação mútua - MMI ("Maximum Mutual Information") o qual ressalta a capacidade de discriminação dos modelos que competem entre si [14]. Infelizmente o critério MMI não pode ser solucionado por análise direta ou por reestimação [14]. Por outro lado, pode-se mostrar que o critério MMI é equivalente à maximização a posteriori - MAP ("Maximum A Posteriori Probabilities"), que é o mesmo método utilizado por uma rede neural artificial - ANN ("Artificial Neural Network") otimizada a partir do critério de minimização do erro quadrático médio - MSE ("Minimum Square Error") [33, 30, 4].

Além de ser não-discriminativo, o critério ML necessita fazer uso de muitas suposições simplifi-

adoras para viabilizar o modelamento matemático de um HMM. ANNs treinadas segundo o critério MSE não precisam fazer uso de várias destas suposições, e além disso são capazes de aprender e generalizar funções complexas, tolerar ruídos e suportar paralelismo. Estas vantagens tornam as ANNs extremamente poderosas para modelar as variabilidades acústicas da fala. Entretanto, ao contrário dos HMMs, as ANNs não têm se mostrado eficientes para o modelamento das variabilidades temporais. Com o objetivo de unir em uma única estrutura o que há de melhor na tecnologia de redes neurais artificiais e dos modelos ocultos de Markov, têm sido estudados e avaliados nos últimos sete anos [36, 14, 3, 33, 21, 8, 1], sistemas para reconhecimento de fala contínua baseados em modelos híbridos ANN-HMM, que utilizam ANNs para estimar os parâmetros, responsáveis pelo modelamento das variabilidades acústicas, de HMMs contínuos.

Os primeiros sistemas para reconhecimento de fala contínua baseados em modelos híbridos ANN-HMM treinavam a ANN uma única vez, utilizando uma base de dados (conjunto de sentenças) previamente segmentada em termos de unidades sub-lexicais (sílabas, polifones, fones, ou unidades sub-fônicas) para estimar $P(q_k|\mathbf{x}_n)$, a probabilidade a posteriori de um determinado estado q_k de um HMM contínuo, dado um vetor acústico \mathbf{x}_n (segmento do sinal de fala parametrizado em termos de alguma transformação específica, por exemplo : coeficientes Mel Cepstrais [5]). Porém em 1996, König [14] estabeleceu as bases matemáticas de um algoritmo recursivo para a reestimação dos parâmetros desta ANN. Este algoritmo foi denominado “Recursive Estimation and Maximization of A Posteriori Probabilities” - REMAP e possui a propriedade de melhorar progressivamente a estimativa de $P(q_k|\mathbf{x}_n)$. O algoritmo REMAP é motivado por estudos perceptuais que mostram a elevada importância das informações transicionais entre unidades sub-lexicais para o reconhecimento da fala.

1.2 Objetivos

Os objetivos desta Tese foram o estudo e implementação dos principais módulos de um sistema para reconhecimento de fala contínua baseado em modelos híbridos ANN-HMM. As análises realizadas concentraram-se em três tópicos fundamentais :

- Estimação de verossimilhanças de emissão de símbolos de HMMs contínuos empregando ANN.
- Reestimação dos parâmetros dos modelos híbridos utilizando o algoritmo REMAP.
- Utilização do algoritmo “Level Building” [29] na determinação da seqüência mais provável de palavras que compõem uma determinada sentença a ser reconhecida.

A avaliação do sistema híbrido ANN-HMM implementado foi realizada com o uso de uma base de dados composta por 100 sentenças (construídas a partir de 319 palavras distintas) dependentes de locutor e gravadas sob condições de estúdio (sem ruído ambiente).

1.3 Estrutura da Tese

Capítulo 2 Apresenta a estrutura básica dos sistemas que representam o estado da arte em reconhecimento estatístico de fala contínua. Estabelece as definições e notações que serão utilizadas ao longo de toda a Tese. Faz uma breve revisão sobre modelos ocultos de Markov. Discute alguns problemas básicos do critério de ML como por exemplo : treinamento independente e suposição de independência estatística dos vetores acústicos. Anuncia os sistemas baseados em modelos híbridos ANN-HMM como uma possível alternativa aos sistemas que empregam apenas os modelos ocultos de Markov.

Capítulo 3 Apresenta uma breve revisão sobre as redes neurais artificiais do tipo multilayer perceptron - MLP (“Multilayer Perceptron”) [10], redes recorrentes [31] e “Time Delay Neural Networks” - TDNN [40]. Introduz a ANN como um estimador estatístico e discute a estimação de verossimilhanças de emissão de símbolos de HMMs contínuos com o uso de ANNs. Lista algumas das vantagens das redes neurais artificiais na estimativa dos parâmetros que modelam a variabilidade acústica da fala. Apresenta os dois tipos básicos de modelos híbridos ANN-HMM segundo Konig, 1996 [14] : (1) Modelo híbrido ANN-HMM *padrão*, treinado a partir do critério MSE, que é equivalente ao critério MAP, e avaliado segundo o critério ML, (2) Modelo híbrido ANN-HMM *discriminativo*, treinado e avaliado segundo o critério MAP.

Capítulo 4 Relata a importância das informações transicionais entre unidades sub-lexicais no reconhecimento da fala. Descreve os procedimentos para a estimação (etapa Pré-REMAP) dos parâmetros de um modelo híbrido ANN-HMM a partir de uma base de dados previamente segmentada em termos de unidades sub-lexicais. Apresenta o algoritmo de reestimação de parâmetros - REMAP discutindo duas alternativas para a estimação dos alvos a serem utilizados no retreinamento da ANN: (1) *Alvos abruptos* - define transições instantâneas entre as unidades sub-lexicais, (2) *Alvos suaves* - define transições suaves entre as unidades sub-lexicais fornecendo informações sobre o grau de confusão acústica que existe nas fronteiras entre estas unidades.

Capítulo 5 Apresenta um breve relato sobre a evolução histórica dos algoritmos de busca baseados em programação dinâmica. Introduz formalmente o algoritmo “Level Building”. Apresenta um modelo de duração para as palavras a ser empregado durante a realização da busca pela seqüência “ótima”

de palavras. Discorre alguns comentários sobre o algoritmo “One Stage”.

Capítulo 6 Relata o sistema híbrido ANN-HMM implementado. Descreve detalhes sobre o conjunto de exemplos de treinamento utilizado e sobre o pré-processamento das amostras do sinal de fala para o cálculo de atributos. Apresenta a topologia e o treinamento da ANN utilizada assim como a topologia dos HMMs empregados. Descreve o processo de Estimação (Pré-REMAP) e reestimação (REMAP) dos parâmetros do sistema. Para finalizar, discute aspectos de implementação do algoritmo de busca incluindo a gramática do tipo *pares de palavras*, utilizada na avaliação final do sistema.

Capítulo 7 Apresenta uma análise dos resultados obtidos, tanto a nível de unidades sub-lexicais (fones) como a nível de sentenças. Uma ênfase especial é dada à evolução das taxas de acerto do sistema após cada etapa de reestimação.

Capítulo 8 Apresenta as considerações finais, concluindo a Tese e apresentando sugestões para trabalhos futuros.

Capítulo 2

Reconhecimento de Fala Contínua

2.1 Introdução

Atualmente os sistemas para reconhecimento de fala contínua baseiam-se fundamentalmente em princípios de reconhecimento estatístico de padrões. A tarefa destes sistemas é construir uma função probabilística capaz de mapear amostras de um sinal de fala - representadas através de uma seqüência de vetores acústicos - no espaço das possíveis sentenças (palavras ou frases). Em geral esta função probabilística é parametrizada em função das principais fontes de variabilidade da fala : a *variabilidade acústica* e a *variabilidade temporal*.

Sendo $s(n)$ a seqüência de amostras de uma sinal de fala desconhecido, este sinal pode ser submetido a um **pré-processamento** e representado por um conjunto de vetores acústicos $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$. Seja M_i ($i = 1, 2, \dots, I_M$) o conjunto de todas as possíveis sentenças e Θ o conjunto de parâmetros responsáveis pelo modelamento das variabilidades *acústicas* e *temporais*. O objetivo de um sistema para reconhecimento de fala contínua consiste em estimar durante o treinamento, e utilizar durante o reconhecimento, a seguinte função probabilística, $P(M|\mathbf{X},\Theta)$, isto é, a probabilidade da sentença M dado o conjunto de vetores acústicos \mathbf{X} e os parâmetros Θ .

Uma vez que uma sentença M pode ser construída a partir da concatenação de palavras $M = \{W_1, W_2, \dots, W_{N_M}\}$, a tarefa de um sistema para reconhecimento de fala contínua também pode ser interpretada como a determinação da seqüência de palavras mais prováveis \widehat{M} , dada a seqüência de vetores acústicos \mathbf{X} e o conjunto de parâmetros Θ . Se a regra de Bayes for utilizada para decompor $P(M|\mathbf{X},\Theta)$, então \widehat{M} pode ser determinada a partir da seguinte expressão,

$$\widehat{M} = \arg \max_M P(M|\mathbf{X},\Theta) = \arg \max_M \frac{p(\mathbf{X}|M, \Theta) \cdot P(M|\Theta)}{p(\mathbf{X}|\Theta)} \quad (2.1)$$

Esta equação mostra que encontrar \widehat{M} , é equivalente a encontrar a seqüência de palavras que maximiza o produto entre $p(\mathbf{X}|M, \Theta)$ e $P(M|\Theta)$, uma vez que durante o reconhecimento o termo $p(\mathbf{X}|\Theta)$ será o mesmo para todas as sentenças avaliadas. O primeiro termo $p(\mathbf{X}|M, \Theta)$, representa a verossimilhança da seqüência de vetores acústicos \mathbf{X} dada uma seqüência de palavras específica M e o conjunto de parâmetros Θ , e esta verossimilhança pode ser determinada a partir de um **modelo acústico** (modelo estatístico) da sentença M . O segundo termo $P(M|\Theta)$, representa a probabilidade da sentença M dado o conjunto de parâmetros Θ , e esta probabilidade pode ser determinada por um **modelo da língua**.

O processo de determinação de \widehat{M} é denominado **decodificação** e projetos de decodificadores (algoritmos de busca) eficientes são cruciais para a realização prática de sistemas para reconhecimento de fala contínua. Portanto, um sistema estatístico para reconhecimento de fala contínua pode ser dividido em quatro módulos principais : (1) **Pré-processamento**, (2) **Modelamento Acústico**, (3) **Modelo da língua**, (4) **Decodificação (Busca)**. A Figura 2-1 mostra um diagrama de blocos com estes quatro módulos principais e as ligações entre eles.

2.1.1 Pré-Processamento

O pré-processamento pode ser visto como um mapeamento da seqüência de amostras $s(n)$ (sinal de fala digitalizado) em uma seqüência de vetores acústicos $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_N\}$. Este mapeamento possui dois objetivos básicos : (1) Ressaltar características, (2) Reduzir dimensionalidade sem perdas de informações cruciais ao reconhecimento. A forma mais usual para realização deste mapeamento consiste em janelar o sinal $s(n)$ com um certo grau de superposição, conforme Figura 2-3, e em seguida calcular um conjunto de atributos (através de uma transformação específica) para cada uma das janelas. Os atributos mais comumente utilizados, são : coeficientes Mel cepstrais [5] , log da energia normalizada [26], e as primeira e segunda derivadas destes coeficientes [6].

2.1.2 Modelamento Acústico

O objetivo deste módulo é a construção de um **modelo acústico** (modelo estatístico), para cada uma das possíveis sentenças M_i ($i = 1, 2, \dots, I_M$), que permita a avaliação da verossimilhança $p(\mathbf{X}|M_i, \Theta)$ durante a etapa de reconhecimento. Em geral estes modelos são construídos através de um procedimento hierárquico, que utiliza como base a decomposição de sentenças em palavras e estas em unidades sub-lexicais como sílabas, polifones, fones ou unidades sub-fônicas. A construção destes modelos de sentenças consiste de três etapas básicas :

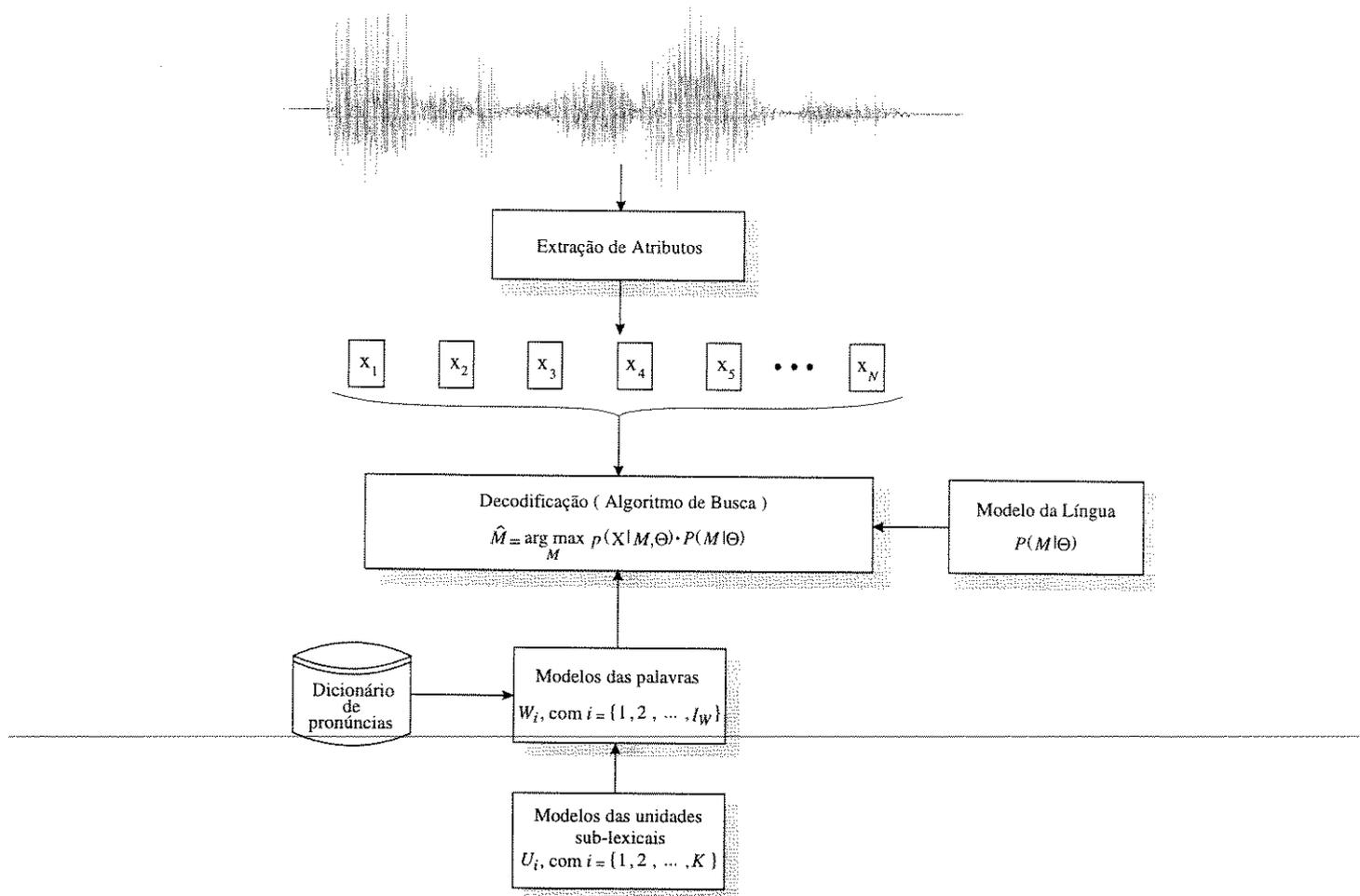


Figura 2-1: Diagrama de blocos de um sistema para reconhecimento estatístico de fala contínua.

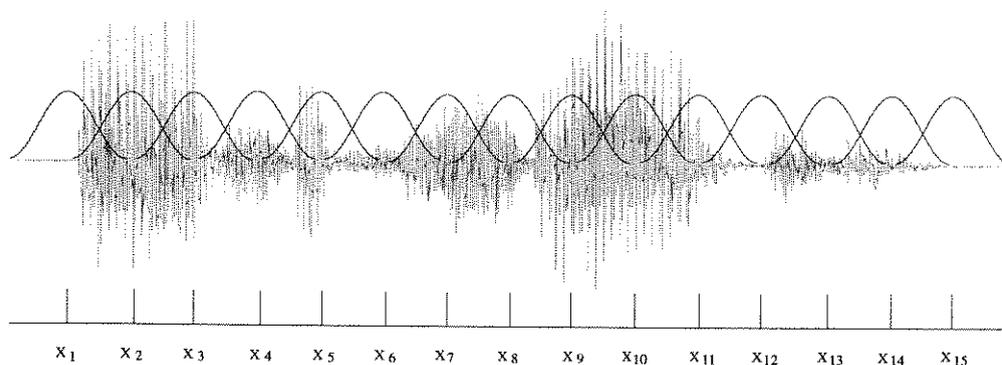


Figura 2-2: Estimativa da seqüência de atributos (vetores acústicos) associados a uma seqüência de amostras de um sinal de fala.

1. Construir um modelo estatístico para cada uma das unidades sub-lexicais.
2. De posse de um **dicionário de pronúncias**, montar os modelos estatísticos para cada uma das palavras do léxico através da concatenação dos correspondentes modelos das unidades sub-lexicais. Por exemplo, o modelo da palavra “*casa*” pode ser obtido pela concatenação dos modelos dos fones $/k/ + /a/ + /z/ + /a/$.
3. Por último, através de um procedimento de concatenação, sujeito ou não a restrições gramaticais, montar os modelos das sentenças a partir dos modelos de palavras.

O modelamento acústico mais comumente utilizado baseia-se em modelos ocultos de Markov - HMMs, e as explicações para este fato são :

- Por serem máquinas de estado finito¹, os HMMs permitem facilmente a concatenação de modelos, possibilitando a construção hierárquica de modelos acústicos para as sentenças e também a aplicação de restrições gramaticais tanto a nível de palavras como de unidades sub-lexicais.
- Por serem estruturas duplamente estocásticas os HMMs são capazes de modelar, tanto a *variabilidade acústica* como *temporal* da fala, permitindo a avaliação de $p(\mathbf{X}|M, \Theta)$.

2.1.3 Modelo da língua

A tarefa do modelo da língua é estimar $P(M|\Theta)$, isto é, a probabilidade da sentença dado o conjunto de parâmetros Θ . O Modelo da língua também pode ser utilizado para determinar a probabilidade de uma palavra W_i , em uma sentença, dadas todas as palavras que a precedem, W_1, W_2, \dots, W_{i-1} . Isto pode ser obtido, verificando que $P(M|\Theta)$ pode ser reescrito como

$$P(M|\Theta) = P(W_1|\Theta) \cdot \prod_{i=2}^{N_M} P(W_i|W_{i-1}, W_{i-2}, \dots, W_1, \Theta) \quad (2.2)$$

sendo N_M o número total de palavras presentes na sentença M .

Em geral a probabilidade de um modelo M é considerada independente dos parâmetros Θ , logo,

$$P(M|\Theta) \approx P(M) = P(W_1|\Theta) \cdot \prod_{i=2}^{N_M} P(W_i|W_{i-1}, W_{i-2}, \dots, W_1) \quad (2.3)$$

¹Máquina de estados com um número finito de estados.

Uma maneira simples mas efetiva para calcular (2.3) é assumir que a probabilidade das palavras em uma sentença dependem apenas das $N - 1$ palavras anteriores,

$$P(M) \approx P(W_1|\Theta) \cdot \prod_{i=2}^{N_M} P(W_i|W_{i-1}, W_{i-2}, \dots, W_{i-N+1}) \quad (2.4)$$

Esta solução é denominada modelamento estatístico N -gram (ou gramática N -gram). Os modelos mais utilizados são os Bi-gram, $P(W_i|W_{i-1})$ e Tri-gram, $P(W_i|W_{i-1}, W_{i-2})$. Estas probabilidades podem ser calculadas através da leitura automática² de vários textos³ e do levantamento das frequências relativas de pares e triplas de palavras.

2.1.4 Decodificação

Dado um modelo da língua e os modelos acústicos de todas as palavras do léxico (construídos através da concatenação dos modelos acústicos de unidades sub-lexicais), a tarefa do decodificador (algoritmo de Busca) é estimar de maneira otimizada a sequência de palavras \widehat{M} que maximiza a Equação (2.1), isto é, deve realizar uma busca da sentença \widehat{M} que maximiza o produto entre $p(\mathbf{X}|M, \Theta)$ e $P(M|\Theta)$. Os algoritmos de busca podem ser classificados, quanto à forma de implementação, em duas categorias:

Algoritmos síncronos por palavras, destacando-se o “Level Building” [29] e os algoritmos síncronos por quadro de análise (síncrono por vetor acústico), destacando-se o “One Stage”, [23, 17]. Estes algoritmos serão discutidos em detalhe no Capítulo 6.

2.2 Definições e Notações

A seguir são apresentadas algumas definições e notações que serão utilizadas nas seções seguintes e em todos os Capítulos subsequentes.

- $Q = \{q_1, q_2, \dots, q_K\}$ - Conjunto que contém todos os estados com os quais serão construídos os modelos de unidades sub-lexicais, palavras e sentenças. Por exemplo, considere que o conjunto de unidade sub-lexicais seja igual a 36 fones, então no caso de se modelar cada fone por um estado, o número de estados será $K = 36$. Se for desejado modelar a produção acústica de início, meio e final dos diferentes fones, cada fone deverá ser modelado por 3 estados e portanto $K = 3 \times 36 = 108$ estados.

²Detecção automática de palavras ao longo de um texto escrito.

³Escolhidos criteriosamente segundo critérios lingüísticos.

- $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ - Seqüência de vetores acústicos associados a uma sentença específica.
 - $\mathbf{X}_{n-c}^{n+d} = \{\mathbf{x}_{n-c}, \mathbf{x}_{n-c+1}, \dots, \mathbf{x}_n, \dots, \mathbf{x}_{n+d}\}$ - Sub-seqüência de \mathbf{X} com $c + d + 1$ vetores acústicos.
 - M_i - Modelo oculto de Markov da i - ésima sentença. Com $i \in \mathcal{I}_M = \{1, 2, \dots, I_M\}$, sendo I_M o número total de seqüências possíveis de serem construídas a partir das palavras presentes no léxico e permitidas pelo modelamento da língua. M_i também pode ser visto como um grafo direcionado com C_i estados, todos eles pertencentes a Q .
 - W_i - Modelo oculto de Markov da i - ésima palavra. Com $i \in \mathcal{I}_W = \{1, 2, \dots, I_W\}$, sendo I_W o número total de palavras presentes no léxico.
 - $\mathbf{X}_{M_j} = \left\{ \mathbf{X}_{M_j}^{(1)}, \mathbf{X}_{M_j}^{(2)}, \dots, \mathbf{X}_{M_j}^{(N_{e_j})} \right\}$ - Conjunto de seqüências de vetores acústicos associadas ao modelo M_j (seqüências de vetores acústicos utilizadas no treinamento de M_j).
 - q^n - Denota um estado no instante n .
 - q_k^n - Indica que o estado k ocorreu no instante n . Rigorosamente falando, q_k^n denota a associação do valor k à variável aleatória q^n .
-
- Γ - Seqüência de estados de comprimento N . Cada Γ é a realização de um processo estocástico, onde os valores associados a este processo em cada instante de tempo pertencem a Q . Ao longo desta dissertação também será usual escrever explicitamente $\Gamma_j = \{q_{\Gamma_j}^1, \dots, q_{\Gamma_j}^n, \dots, q_{\Gamma_j}^N\}$. Uma seqüência de estados permitida por um determinado modelo M_i também será denominada de *caminho* ao longo de M_i .
 - τ_i - Conjunto de todas as seqüências de estado Γ permitidas pelo modelo M_i .
 - \mathfrak{S} - Conjunto de todas as seqüências de estado Γ permitidas por todos os modelos M_i , com $i \in \mathcal{I}_M$.
 - $\Theta = \{\theta_1, \theta_2, \dots, \theta_I\}$ - Conjunto de parâmetros que descrevem todos os modelos M_i com $i \in \mathcal{I}_M$. θ_i representa apenas os parâmetros presentes em M_i . Nos modelos híbridos discutidos nesta dissertação, $\Theta = \{\Psi \cup \mathbf{A} \cup \boldsymbol{\pi}\}$, sendo Ψ o conjunto de parâmetros de uma rede neural artificial (conjunto de pesos sinápticos), \mathbf{A} o conjunto com as probabilidades de transição de cada um dos estados $q \in Q$, e $\boldsymbol{\pi}$ as probabilidades dos estados iniciais q^1 de cada um dos modelos M_i com $i \in \mathcal{I}_M$.
 - $P(\cdot)$ - Será utilizado para representar probabilidade ou distribuição de probabilidade.

- $p(\cdot)$ - Será utilizado para representar função densidade de probabilidade e verossimilhança.

2.3 Modelos Ocultos de Markov - HMM

2.3.1 Introdução

Atualmente os sistemas que representam o estado-da-arte em reconhecimento de fala contínua constroem seus *modelos acústicos* baseados em modelos ocultos de Markov. O sucesso destas estruturas deve-se, principalmente, à sua capacidade de modelar tanto as *variabilidades acústicas* como *temporais* do sinal fala e também por permitir a construção hierárquica dos modelos acústicos das sentenças.

A teoria de modelos ocultos de Markov apóia-se na suposição de que um processo contínuo pode ser aproximado por uma sucessão de curtos estados estacionários e modela uma seqüência de vetores acústicos como um processo estacionário por partes. Isto é, um HMM modela uma seqüência de vetores acústicos $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ como se estes tivessem sido emitidos (gerados) por uma sucessão de L *estados estacionários discretos* $q_k \in Q$, com *transições instantâneas* entre eles. Neste caso, um HMM é definido (e representado) como uma máquina de estados finito com uma topologia particular (usualmente *left-right*, uma vez que a fala é estritamente seqüencial). Esta abordagem define dois processos estocásticos concorrentes : A seqüência dos L estados estacionários discretos (responsável pelo modelamento da estrutura temporal da fala), e a seqüência de símbolos observados (seqüência de vetores acústicos) formada pela emissão de um símbolo em cada um dos L estados (responsável pelo modelamento acústico da fala). O HMM é chamado de “oculto” porque existe um processo estocástico escondido, a seqüência de estados, que não é observável mas que afeta a seqüência de símbolos emitidos. O modelo é chamado de “Markov” porque as estatísticas do estado corrente q^n são modeladas como sendo dependentes apenas do símbolo atual \mathbf{x}_n e do estado prévio q^{n-1} (no caso de modelo de Markov de primeira ordem).

2.3.2 Definição

Um HMM pode ser definido formalmente por :

- $\mathbf{A} = \{a_{ij} | a_{ij} = P(q^{n+1} = j | q^n = i) = P(q_j^{n+1} | q_i^n)\}$ - Função distribuição de probabilidade de transição, sendo que a_{ij} denota a probabilidade de transição do estado i para o estado j . Geralmente assume-se que esta distribuição de probabilidade é a mesma para todo instante de tempo.

- $\mathbf{B} = \{b_j(\mathbf{x}_i) | b_j(\mathbf{x}_i) = P(\mathbf{x}_i | q = j) = P(\mathbf{x}_i | q_j)\}$ ou $\mathbf{B} = \{b_j(\mathbf{x}_i) | b_j(\mathbf{x}_i) = p(\mathbf{x}_i | q = j) = p(\mathbf{x}_i | q_j)\}$
- Para cada estado existe uma correspondente distribuição de probabilidade de saída no caso de um HMM discreto ou uma função densidade de probabilidade contínua para o caso de um HMM contínuo. $b_j(\mathbf{x}_i)$ se refere à probabilidade (no caso discreto) ou verossimilhança (no caso contínuo) do símbolo \mathbf{x}_i ser gerado pelo estado q_j . Usualmente $b_j(\mathbf{x}_i)$ é chamado de probabilidade (ou verossimilhança) de emissão de símbolos.
- $\boldsymbol{\pi} = \{\pi_i | \pi_i = P(q^1 = i) = P(q_i^1)\}$ - Distribuição de probabilidade do estado inicial. Nos modelos *left-right*, normalmente é assumido $\pi_1 = 1$ e $\pi_i = 0$ para todo $i \neq 1$.

De posse destas definições pode-se representar um modelo oculto de Markov de forma compacta por $M = (\mathbf{A}, \mathbf{B}, \boldsymbol{\pi})$.

2.3.3 Algoritmos Básicos

Em geral [27], três problemas são levantados quando utilizam-se HMMs para modelar seqüências de observações :

-
1. **Avaliação** Qual a verossimilhança de um HMM gerar uma dada seqüência de observações ?
 2. **Decodificação** Dada uma seqüência de observações e um HMM, qual é a seqüência de estados, ao longo deste HMM, com maior verossimilhança de gerar estas observações ?
 3. **Estimação de Parâmetros** Dado um HMM e um conjunto de seqüências de observações (seqüências de vetores acústicos) a serem modeladas por este HMM, como adaptar seus parâmetros (\mathbf{A} , \mathbf{B} e $\boldsymbol{\pi}$) para que este modelo maximize a verossimilhança de geração da seqüência de observações ?

Todos estes três problemas apresentam soluções bastante eficientes baseadas em casos particulares de algoritmos de programação dinâmica. Por exemplo, o problema de **avaliação** ocorre no reconhecimento de palavras isoladas em que se deseja avaliar diferentes modelos ocultos de Markov de palavras para verificar qual deles apresenta a maior verossimilhança de ter gerado a palavra a ser reconhecida. Este problema pode ser solucionado utilizando-se o algoritmo *Forward* [27]. O problema de **decodificação** é utilizado no reconhecimento de fala contínua em que se deseja descobrir qual o caminho (seqüência de estados) com maior verossimilhança, ao longo de um HMM bastante extenso (HMM de uma sentença), de ter gerado uma determinada seqüência de vetores acústicos. Este problema de

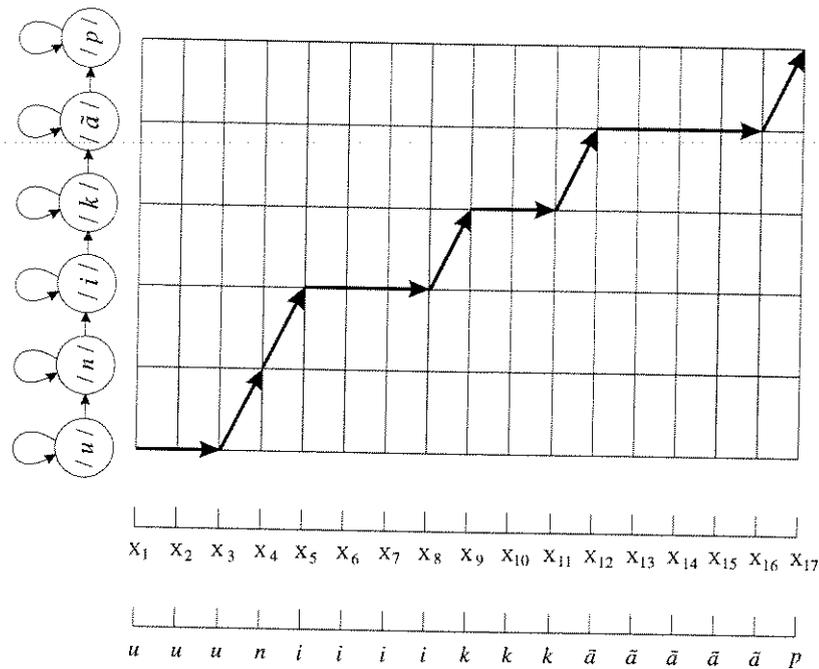


Figura 2-3: Trajetória de Viterbi : define a seqüência de estados com maior verossimilhança de ter gerado a seqüência de símbolos X .

decodificação pode ser solucionado utilizando-se o algoritmo de *Viterbi* [27]. A Figura 2-3 mostra um diagrama de treliça típico⁴ com o caminho ótimo produzido pelo algoritmo de Viterbi.

O último problema, também chamado de **treinamento** (estimação de parâmetros), pode ser solucionado pelo algoritmo *Forward-Backward* [27] também conhecido como *Baum-Welch*, que é essencialmente uma versão do algoritmo EM (“Expectation-Maximization”) [14]. No caso de HMMs estritamente *left-right* com transições de estado constantes no tempo, todos estes três algoritmos têm uma complexidade da ordem de apenas $(C_i \cdot N)$, sendo C_i o número de estados do HMM em questão e N o número de símbolos observados.

2.4 Discussão

2.4.1 Critérios de Treinamento - Problemas com o Critério ML

Como foi discutido no início deste Capítulo, o objetivo de um sistema para reconhecimento de fala contínua é, estimar durante o treinamento e utilizar durante o reconhecimento, a seguinte função

⁴Trajetória de Viterbi produzida pelo modelo da palavra “UNICAMP”. Neste modelo foi utilizado o fone como unidade sub-lexical e estabelecida uma *correspondência um-a-um* entre fones e estados.

probabilística, $P(M|\mathbf{X},\Theta)$, que pode ser reescrita através da regra de Bayes como,

$$P(M|\mathbf{X},\Theta) = \frac{p(\mathbf{X}|M, \Theta) \cdot P(M|\Theta)}{p(\mathbf{X}|\Theta)} \quad (2.5)$$

Na fase de treinamento, a probabilidade $P(M|\mathbf{X},\Theta)$ do modelo M dado a seqüência de vetores acústicos \mathbf{X} , deve ser maximizada. O espaço de parâmetros onde esta otimização é realizada pode fazer a diferença entre *modelos treinados independentemente* (não-discriminativos) e *modelos discriminativos*.

Durante a fase de reconhecimento, $p(\mathbf{X}|\Theta)$ é uma constante pois os parâmetros do modelo são fixados. Entretanto, durante o treinamento, esta verossimilhança depende dos parâmetros de todos os possíveis modelos. Escrevendo $p(\mathbf{X}|\Theta)$ em termos de probabilidade marginal,

$$p(\mathbf{X}|\Theta) = \sum_i p(\mathbf{X}, M|\Theta) \quad (2.6)$$

$$= \sum_i p(\mathbf{X}|M, \Theta) \cdot P(M|\Theta) \quad (2.7)$$

$$\approx \sum_i p(\mathbf{X}|M, \Theta) \cdot P(M) \quad (2.8)$$

A Equação (2.7) estende o somatório em i sobre todos os possíveis modelos de sentenças e assume a independência das probabilidades a priori dos modelos versus os parâmetros ($P(M|\Theta) \approx P(M)$). Substituindo (2.8) em (2.5), e tratando do caso de um modelo particular M_j , tem-se :

$$P(M_j|\mathbf{X}_{M_j}^{(s)}, \Theta) \approx \frac{p(\mathbf{X}_{M_j}^{(s)}|M_j, \Theta) \cdot P(M_j)}{p(\mathbf{X}_{M_j}^{(s)}|M_j, \Theta) \cdot P(M_j) + \sum_{i \neq j} p(\mathbf{X}_{M_j}^{(s)}|M_i, \Theta) \cdot P(M_i)} \quad (2.9)$$

sendo $\mathbf{X}_{M_j}^{(s)}$ com $s \in \{1, 2, \dots, N_{e_j}\}$ as seqüências de treinamento associadas ao modelo M_j .

A maximização de $P(M_j|\mathbf{X}_{M_j}^{(s)}, \Theta)$ como dado por (2.9) é usualmente simplificada restringindo a otimização ao subespaço de parâmetros de M_j . Esta restrição dá origem ao critério de máxima verossimilhança - ML ("Maximum Likelihood"). Neste critério mantém-se o somatório no denominador constante sobre todo o espaço de parâmetros de M_j . Como $P(M_j)$ é uma constante determinada pelo modelo da língua, então maximizar $P(M_j|\mathbf{X}_{M_j}^{(s)}, \Theta)$ segundo o critério ML é semelhante a maximizar $p(\mathbf{X}_{M_j}^{(s)}|M_j, \Theta)$. Portanto, o treinamento de HMMs de acordo com o critério ML procura encontrar o

melhor conjunto de parâmetros $\hat{\Theta}$, tal que⁵,

$$\hat{\Theta} = \arg \max_{\Theta} \prod_{j=1}^{I_M} \prod_{s=1}^{N_{e_j}} p(\mathbf{X}_{M_j}^{(s)} | M_j, \Theta) \quad (2.10)$$

sendo I_M o número total de sentenças.

Esta otimização “*modelo – por – modelo*” permite importantes simplificações no algoritmo de treinamento, porém pagando o preço de um baixo poder de discriminação entre os modelos.

Uma outra forma de maximizar $P(M_j | \mathbf{X}_{M_j}^{(s)}, \Theta)$ é considerar todo o espaço de parâmetros (isto é, os parâmetros de todos os modelos possíveis). Este critério é, de fato, discriminativo, uma vez que a contribuição de $p(\mathbf{X}_{M_j}^{(s)} | M_j, \Theta) \cdot P(M_j)$ deve ser maximizada enquanto que a de todos os modelos rivais, representada por,

$$\sum_{i \neq j} p(\mathbf{X}_{M_j}^{(s)} | M_i, \Theta) \cdot P(M_i) \quad (2.11)$$

deve ser minimizada. O Apêndice A mostra que um treinamento baseado no critério de maximização da informação mútua - MMI (“Maximum Mutual Information”) é capaz de maximizar a probabilidade a posteriori $P(M_j | \mathbf{X}_{M_j}^{(s)}, \Theta)$ com relação a todo o espaço de parâmetros Θ .

2.4.2 Mais Problemas com o Critério ML

Para tornar a estimação de $p(\mathbf{X}_{M_j}^{(s)} | M_j, \Theta)$ segundo o critério ML, tratável matematicamente, são realizadas algumas suposições simplificadoras que reduzem ainda mais a eficiência do critério de máxima verossimilhança - ML :

- **Suposição de independência das saídas** Os vetores acústicos não são correlacionados (independência das observações). O vetor acústico corrente \mathbf{x}_n é assumido ser condicionalmente independente dos vetores acústicos prévios ($\mathbf{X}_1^{n-1} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{n-1}\}$). Esta suposição assume que observações pertencentes a um mesmo segmento da fala (geradas pelo mesmo estado) são independentes, uma suposição não realista dada a natureza não estacionária do sinal de fala. Para reduzir os problemas causados por esta suposição, o vetor acústico no instante n é geralmente complementado por suas primeira e segunda derivadas [6] calculadas sobre uma janela com alguns quadros, permitindo modelar um limitado contexto acústico. Por exemplo, a primeira derivada de \mathbf{x}_n , também conhecida como coeficiente Delta, $\Delta \mathbf{x}_n$, é comumente representada por

⁵Na Equação (2.10) também foi assumido que cada uma das seqüências de vetores acústicos é independente das demais.

$\Delta \mathbf{x}_n = a_c \cdot \mathbf{x}_{n-c} + \dots + a_n \cdot \mathbf{x}_n + \dots + a_d \cdot \mathbf{x}_{n+d}$, sendo que estes coeficientes a_i 's, com $c \leq i \leq d$ são, em geral, fixados segundo critérios heurísticos. Da mesma forma a segunda derivada, coeficiente Delta-Delta, é geralmente definida por : $\Delta^2 \mathbf{x}_n = a_c \cdot \Delta \mathbf{x}_{n-c} + \dots + a_n \cdot \Delta \mathbf{x}_n + \dots + a_d \cdot \Delta \mathbf{x}_{n+d}$.

- **Suposição de Markov** Modelos de Markov são geralmente cadeias de Markov de primeira ordem. Explicitamente, a probabilidade da cadeia de Markov estar no estado q_k no instante n depende apenas do estado da cadeia de Markov no instante $n - 1$, e é condicionalmente independente de todos os outros parâmetros passados (tanto do vetor acústico passado, \mathbf{x}_{n-1} , como dos estados anteriores ao estado q^{n-1}).
- **Suposição de superposição de f.d.p's gaussianas** A aplicação do critério ML a HMMs contínuos assume que as funções densidade de probabilidade de emissão de símbolos podem ser expressas pela superposição de um número finito de gaussianas, necessitando fixar a priori quantas gaussianas devem ser superpostas.

2.4.3 Alternativa - Modelos Híbridos ANN-HMM

Como alternativa aos sistemas baseados apenas em modelos ocultos de Markov, alguns autores, [36, 14, 3, 33, 21, 8, 1] têm apresentado nos últimos sete anos sistemas híbridos baseados tanto em redes neurais artificiais-ANN como em modelos ocultos de Markov. Nestes sistemas as ANNs são utilizadas para estimar, de forma mais eficiente que o tradicional algoritmo de reestimação de *Baum-Welch*, baseado no critério ML, parâmetros de HMMs contínuos. O Capítulo 3 apresenta a formulação matemática dos modelos híbridos assim como as suas principais vantagens em relação aos sistemas baseados apenas em HMM.

Capítulo 3

Modelos Híbridos ANN-HMM

3.1 Redes Neurais Artificiais - ANNs

3.1.1 Multilayer Perceptrons - MLPs

Nesta dissertação, a discussão sobre redes neurais será focalizada principalmente sobre o Perceptron multicamadas - MLP (“Multilayer Perceptron”), uma forma de ANN que é comumente utilizada em reconhecimento de fala. Entretanto, todas as conclusões básicas sobre a utilidade destas estruturas para estimar parâmetros (densidade de probabilidade de emissão de símbolos) para um HMM são perfeitamente estendidas para outros tipos de ANN, como por exemplo redes neurais recorrentes [31], ou “Time Delay Neural Network” (TDNN) [40].

Tipicamente, MLPs possuem uma arquitetura do tipo “*feed – forward*” com uma camada de entrada (consistindo do vetor de entrada), zero ou mais camadas escondidas, e uma camada de saída. Cada camada calcula um conjunto de funções discriminantes lineares [10] (via uma matriz de pesos sinápticos) e em seguida aplica à saída de cada uma destas funções uma não-linearidade, freqüentemente uma função do tipo sigmóide. As mais comuns são :

$$\varphi_l(v) = \frac{a}{1 + \exp(-b \cdot v)} \quad (3.1)$$

assim como,

$$\varphi_t(v) = \frac{2 \cdot a}{1 + \exp(-b \cdot v)} - a \quad (3.2)$$

Se $a = b = 1$, $\varphi_l(v)$ e $\varphi_t(v)$ se transformam, respectivamente, nas funções logística e tangente

hiperbólica.

Como discutido em [10], estas não-linearidades executam diferentes funções para as unidades escondidas e de saída. Nas unidades escondidas, suas funções são gerar momentos de ordem elevada do vetor de entrada; isto pode ser feito efetivamente com o uso de vários tipos de funções não lineares, não apenas por sigmóides. Nas unidades de saída, as não-linearidades podem ser vistas como uma aproximação diferencial de um limiar de decisão. Por este motivo, as funções não lineares das saídas devem ser sigmóides.

Pode ser provado que MLPs com um número elevado de unidades escondidas podem (em princípio) fornecer um mapeamento arbitrário $g(\mathbf{x})$ entre a entrada e a saída. O conjunto de parâmetros (elementos das matrizes de pesos) de uma MLP será denominado Ψ . Estes parâmetros são treinados para associar o vetor de saída “desejado” ao vetor de entrada. Isto é geralmente alcançado via algoritmo de treinamento EBP - “Error Back-Propagation” [10] que usa um procedimento baseado em gradientes para iterativamente minimizar a função de custo em seu espaço de parâmetros.

As funções de custo mais utilizadas, entre outras, são o critério de minimização do erro quadrático médio - MSE “Minimum Square Error”,

$$\mathcal{E} = \sum_{n=1}^{N_e} \|g(\mathbf{x}_n, \Psi) - \mathbf{d}(\mathbf{x}_n)\|^2 \quad (3.3)$$

e o critério de entropia relativa. Uma versão discreta segundo König [14] da versão clássica de entropia relativa é,

$$\mathcal{E}_e = \sum_{n=1}^{N_e} \sum_{k=1}^K d_k(\mathbf{x}_n) \ln \frac{d_k(\mathbf{x}_n)}{g_k(\mathbf{x}_n, \Psi)} \quad (3.4)$$

onde $\mathbf{g}(\mathbf{x}_n, \Psi) = (g_1(\mathbf{x}_n, \Psi), \dots, g_k(\mathbf{x}_n, \Psi), \dots, g_K(\mathbf{x}_n, \Psi))^t$ representa o vetor de saída atual da MLP (em função do vetor de entrada \mathbf{x}_n e dos parâmetros Ψ), $\mathbf{d}(\mathbf{x}_n) = (d_1(\mathbf{x}_n), \dots, d_k(\mathbf{x}_n), \dots, d_K(\mathbf{x}_n))^t$ representa o vetor de alvos desejados na saída (obtidos através de uma base de treinamento devidamente rotulada em termos de unidades sub-lexicais), K é o número total de classes e N_e o total de exemplos de treinamento.

MLPs, e outras estruturas conexionistas, têm sido utilizadas para uma grande variedade de tarefas relacionadas com o reconhecimento de fala [33]. O caso mais comum é o de reconhecimento de palavras isoladas, em que uma seqüência temporal de vetores acústicos é tratada como um padrão espacial. A função da rede neste caso é mapear um vetor de entrada (palavra a ser reconhecida) em uma das N_W unidades de saída, sendo N_W o número de palavras possíveis. A aplicação deste tipo de estrutura

no reconhecimento de fala contínua exigiria inicialmente uma segmentação da sentença em palavras, porém isto é extremamente complexo, uma vez que as palavras encontram-se, em geral, altamente coarticuladas.

3.1.2 Outras Arquiteturas Conexionistas

Como será visto na próxima seção, um modelo híbrido ANN-HMM é baseado em uma perspectiva estatística que é válida para qualquer arquitetura conexionista, dependendo apenas de algumas suposições sobre o procedimento de treinamento. Por este motivo vários trabalhos têm sido realizados utilizando outros tipos de estruturas conexionistas em sistemas híbridos ANN-HMM, destacando-se :

- **Redes Neurais Recorrentes - RNN** Vários trabalhos têm sido publicados pelo grupo de processamento digital de fala da Universidade de Cambridge descrevendo sistemas híbridos, cujas verossimilhanças de emissão de símbolos de HMMs contínuos são geradas por redes neurais recorrentes - RNN ("Recurrent Neural Network") [31]. Estas redes utilizam um conjunto de unidades de entrada para receber vetores acústicos e têm um conexionismo recorrente da saída para algumas unidades de entrada. A rede é treinada utilizando o algoritmo de retropropagação do erro ao longo do tempo - EBPTT ("Error Back Propagation Through Time") [10]. Resultados experimentais [31], mostram que este tipo de rede pode apresentar resultados similares aos de uma MLP, com o uso de um número significativamente menor de parâmetros (conexões sinápticas). Uma das principais limitações deste tipo de estrutura refere-se a alguns problemas inerentes à instabilidade.
- **Time-Delay-Neural Network - TDNN** Redes recorrentes podem ser aproximadas em um intervalo de tempo finito por uma rede do tipo feed-forward, onde os laços (realimentações) são substituídos pelo uso explícito de valores de ativações precedentes. Em 1989, Waibel et al [40] utilizaram, com sucesso, um rede como esta para o reconhecimento de fones, e em 1993 Waibel et al [39], aplicaram-na em um sistema híbrido denominado MS-TDNN ("Multi-State-TDNN") [39], que utilizava além de uma TDNN um procedimento de alinhamento temporal dinâmico - DTW ("Dynamic Time Warping"). A TDNN é uma estrutura bastante flexível e na prática, muitos laboratórios de pesquisa (inclusive o Laboratório de Processamento Digital de Fala da Universidade Estadual de Campinas - LPDF-UNICAMP) estão realizando pesquisas sobre este tipo de ANN para estimar verossimilhanças de emissão de símbolos para modelos ocultos de Markov em sistemas híbridos ANN-HMM.

Assim como em uma MLP, tanto a RNN como o TDNN podem utilizar as funções do tipo sigmóide em suas unidades escondidas e de saída e também podem fazer uso de treinamentos baseados nos critérios MSE e em entropia relativa.

3.2 ANNs como Estimadores Estatísticos

Redes neurais artificiais e modelos estatísticos para a análise de dados não constituem, em absoluto, metodologias disjuntas. Existe uma considerável superposição entre estas duas áreas do conhecimento. A metodologia estatística é diretamente aplicável a muitos modelos de redes neurais, resultando em uma maior eficiência tanto na estimação de parâmetros como dos algoritmos de otimização (aprendizado). Adicionalmente, métodos estatísticos providenciam ferramentas tais como intervalos de confiança e testes de hipótese que são ausentes no campo das redes neurais.

Recentemente, muitos trabalhos estatísticos têm sido publicados estabelecendo conexões entre modelos estatísticos e redes neurais, muitos deles mostrando não só a complementariedade mas também a equivalência entre as duas metodologias [9].

Foi descoberto recentemente [4], que se uma ANN for projetada como um classificador de padrões e otimizada segundo o critério da minimização do erro quadrático médio ou entropia relativa, então suas saídas podem aproximar a probabilidade a posteriori da classe dada a entrada, isto é, $P(\text{classe}|\text{entrada})$. Os modelos denominados híbridos ANN-HMM têm utilizado com sucesso esta propriedade para estimar probabilidades a posteriori de um estado q_k de um HMM dado um vetor acústico \mathbf{x}_n , $P(q_k|\mathbf{x}_n)$. Na verdade, um parâmetro útil para um HMM seria a verossimilhança de um vetor acústico \mathbf{x}_n ser emitido por um determinado estado q_k , isto é, $p(\mathbf{x}_n|q_k)$. Felizmente, $p(\mathbf{x}_n|q_k)$ pode ser facilmente obtido a partir de $P(q_k|\mathbf{x}_n)$ com o uso da regra de Bayes, como será visto na seção 3.3.

Como a ANN será utilizada para estimar probabilidades associadas a cada estado q_k de um HMM, então deve existir uma *equivalência um-a-um* entre as classes c_k 's com $k = 1, \dots, K$, e os estados q_k 's $\in Q$ (conjunto de estados correspondentes a cada uma das unidades sub-lexicais).

O treinamento desta ANN deve ser realizado de forma supervisionada, e os exemplos de treinamento (vetor de entrada, vetor de saída (alvo desejado)) podem ser obtidos através da segmentação da base de dados (sentenças de treinamento) em termos de unidades sub-lexicais. Existem duas formas básicas para a realização desta segmentação: (1) segmentação manual, realizada por especialistas (pode ser bastante precisa, mas é muito demorada); (2) segmentação automática, em geral baseada em HMMs e alinhamentos de Viterbi (não é tão precisa quanto a manual, mas é muito mais rápida). Uma ANN

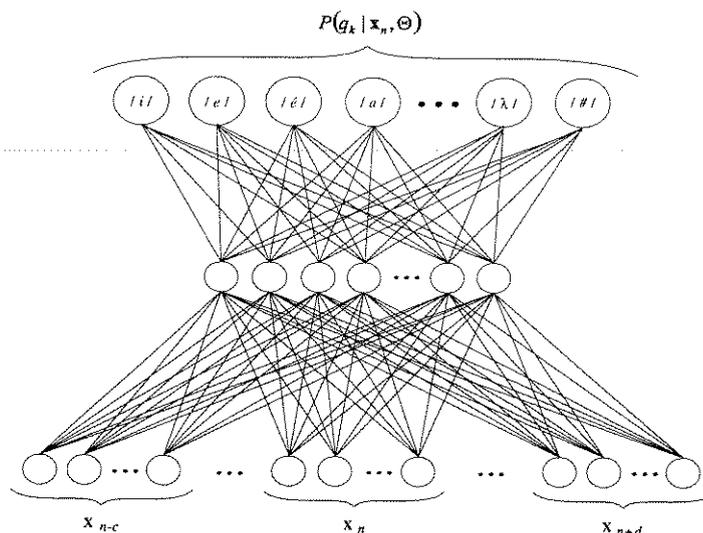


Figura 3-1: Rede MLP para estimativa das probabilidades das classes dado o vetor acústico de entrada e mais alguma informação contextual.

típica para estimar $P(q_k | \mathbf{x}_n)$ é mostrada na Figura 3-1.

A seguir é apresentada a demonstração formal, realizada por Bourlard, 1995 [22], de que uma ANN pode ser utilizada como um estimador da probabilidade da classe de saída dada a entrada, $P(\text{classe} | \text{entrada})$.

Demonstração Seja $\mathbf{X}_e = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{N_e}\}$ o conjunto total de vetores acústicos para treinamento. Os pares entrada e saída para treinamento serão expressos por (\mathbf{x}_i, q_c) com $\mathbf{x}_i \in \mathbf{X}_e$ e $q_c =$ classe correta associada a \mathbf{x}_i , com $c \in \{1, 2, \dots, K\}$. A k -ésima saída da ANN dado o conjunto de parâmetros Θ e a entrada \mathbf{x}_i será denotada por $g_k(\mathbf{x}_i, \Psi)$. O vetor de saída desejado (vetor alvo) para a entrada \mathbf{x}_i será expresso por $\mathbf{d}(\mathbf{x}_i) = (d_1(\mathbf{x}_i), \dots, d_k(\mathbf{x}_i), \dots, d_K(\mathbf{x}_i))^t = \delta_{kc}(\mathbf{x}_i)$, sendo que $\delta_{kc}(\mathbf{x}_i) = 1$ para $k = c$ e $\delta_{kc}(\mathbf{x}_i) = 0$ para $k \neq c, k \in \{1, 2, \dots, K\}$. A partir destas definições o valor médio do erro quadrático na saída da rede neural pode ser expresso por :

$$\mathcal{E} = \sum_{i=1}^{N_e} \sum_{c=1}^K \sum_{k=1}^K [\delta_{kc}(\mathbf{x}_i) - g_k(\mathbf{x}_i, \Psi)]^2 \cdot P(q_c, \mathbf{x}_i) \quad (3.5)$$

$$\mathcal{E} = \sum_{i=1}^{N_e} \sum_{k=1}^K \underbrace{\left\{ \sum_{c=1}^K [\delta_{kc}(\mathbf{x}_i) - g_k(\mathbf{x}_i, \Psi)]^2 \cdot P(q_c, \mathbf{x}_i) \right\}}_{\mathcal{E}_{kx}} \quad (3.6)$$

$$\mathcal{E} = \sum_{i=1}^{N_e} \sum_{k=1}^K \mathcal{E}_{kx} \quad (3.7)$$

De (3.5) e (3.9) pode ser verificado que minimizar \mathcal{E} é equivalente a minimizar \mathcal{E}_{kx} .

Dividindo \mathcal{E}_{kx} em duas partes, $k = c$ e $k \neq c$, resulta :

$$\begin{aligned} \mathcal{E}_{kx} &= \sum_{c=1}^K P(q_c, \mathbf{x}_i) \cdot \delta_{kc}^2(\mathbf{x}_i) - \\ & 2 \cdot P(q_c, \mathbf{x}_i) \cdot \delta_{kc}(\mathbf{x}_i) \cdot g_k(\mathbf{x}_i, \Psi) + P(q_c, \mathbf{x}_i) \cdot g_k^2(\mathbf{x}_i, \Psi) \end{aligned} \quad (3.8)$$

$$\begin{aligned} &= P(q_k, \mathbf{x}_i) - 2 \cdot P(q_k, \mathbf{x}_i) \cdot g_k(\mathbf{x}_i, \Psi) + \\ & P(q_k, \mathbf{x}_i) \cdot g_k^2(\mathbf{x}_i, \Psi) + \sum_{\substack{c=1 \\ c \neq k}}^K P(q_c, \mathbf{x}_i) \cdot g_k^2(\mathbf{x}_i, \Psi) \end{aligned} \quad (3.9)$$

$$\begin{aligned} &= P(q_k, \mathbf{x}_i) - 2 \cdot P(q_k, \mathbf{x}_i) \cdot g_k(\mathbf{x}_i, \Psi) + \\ & P(q_k, \mathbf{x}_i) \cdot g_k^2(\mathbf{x}_i, \Psi) + (P(\mathbf{x}_i) - P(q_k, \mathbf{x}_i)) \cdot g_k^2(\mathbf{x}_i, \Psi) \end{aligned} \quad (3.10)$$

$$= P(q_k, \mathbf{x}_i) - 2 \cdot P(q_k, \mathbf{x}_i) \cdot g_k(\mathbf{x}_i, \Psi) + P(\mathbf{x}_i) \cdot g_k^2(\mathbf{x}_i, \Psi) \quad (3.11)$$

Utilizando o fato $P(q_k, \mathbf{x}_i) = P(q_k|\mathbf{x}_i) \cdot P(\mathbf{x}_i)$ e somando e subtraindo $P(q_k|\mathbf{x}_i) \cdot P(q_k, \mathbf{x}_i) = P(\mathbf{x}_i) \cdot P^2(q_k|\mathbf{x}_i)$ na expressão (3.11) tem-se,

$$\mathcal{E}_{kx} = P(\mathbf{x}_i) \cdot P^2(q_k|\mathbf{x}_i) - 2 \cdot P(\mathbf{x}_i) \cdot P(q_k|\mathbf{x}_i) \cdot g_k(\mathbf{x}_i, \Psi) + \quad (3.12)$$

$$P(\mathbf{x}_i) \cdot g_k^2(\mathbf{x}_i, \Psi) + P(q_k, \mathbf{x}_i) - P(q_k, \mathbf{x}_i) \cdot P(q_k|\mathbf{x}_i) \quad (3.13)$$

que resulta em,

$$\mathcal{E}_{kx} = P(\mathbf{x}_i) \cdot (P(q_k|\mathbf{x}_i) - g_k(\mathbf{x}_i, \Psi))^2 + P(q_k, \mathbf{x}_i) \cdot (1 - P(q_k|\mathbf{x}_i)) \quad (3.14)$$

A partir de (3.14) pode ser verificado que \mathcal{E}_{kx} será mínimo quando os parâmetros Ψ forem ajustados para que $g_k(\mathbf{x}_i, \Psi) = P(q_k|\mathbf{x}_i)$. Portanto, se uma ANN for otimizada segundo o critério MSE (ou entropia relativa, veja demonstração em [10]) suas **saídas ótimas** serão estimativas de probabilidades de classes (estados de uma HMM) condicionadas à entrada (vetor acústico), $\hat{P}(q_k|\mathbf{x}_n)$:

$$g_k(\mathbf{x}_n, \Psi^{\text{ótimo}}) = \hat{P}(q_k|\mathbf{x}_n) \quad (3.15)$$

Na Equação (3.15), $\Psi^{\text{ótimo}}$ representa o conjunto de parâmetros que minimiza tanto (3.3) quanto (3.4).

Porém, para que a ANN alcance estes valores ótimos, algumas condições de treinamento devem ser satisfeitas :

- A ANN deve possuir parâmetros “suficientes” para ser capaz de estimar com razoável aproximação a função de mapeamento *entrada* \rightarrow *saída*.
- A rede não deve ser sobre-treinada. Isto pode ser evitado parando o treinamento antes do declínio da curva que expressa a capacidade de generalização da rede, [27].

Tem sido observado experimentalmente que, para sistemas treinados com grandes bases de dados, a saída de uma MLP treinada apropriadamente, de fato aproxima a probabilidade a posteriori, [22]. Esta conclusão pode facilmente ser estendida para outros casos. Por exemplo, se for fornecida à entrada de uma ANN o vetor acústico \mathbf{x}_n no instante n , mais alguma informação contextual, expressa através do vetor, $\mathbf{X}_{n-c}^{n+d} = \{\mathbf{x}_{n-c}, \dots, \mathbf{x}_n, \dots, \mathbf{x}_{n+d}\}$, os valores de saída da MLP irão estimar

$$g_k(\mathbf{x}_n, \Psi^{\acute{o}timo}) = \hat{P}(q_k | \mathbf{X}_{n-c}^{n+d}), \forall k = 1, \dots, K \quad (3.16)$$

Este janelamento no tempo tem sido utilizado nos sistemas híbridos ANN-HMM para representar a correlação entre os vetores acústicos. Se a classe anterior é também fornecida à entrada da ANN, fornecendo uma rede recorrente, os valores na saída da MLP serão estimativas de

$$g_k(x_n, \Psi^{\acute{o}timo}) = \hat{P}(q_k^n | \mathbf{X}_{n-c}^{n+d}, q_j^{n-1}), \forall k = 1, \dots, K \quad (3.17)$$

sendo que q_j^{n-1} indica a classe vencedora no instante $n - 1$ (saída da rede no instante anterior).

Será mostrado na Seção 3.5.2 que esta é uma forma de probabilidade local que pode ser utilizada na teoria de modelos híbridos ANN-HMM. Esta probabilidade será denotada como “*probabilidade de transição condicionada*”.

Existe uma outra generalização importante desta propriedade que será essencial nesta tese. Se ANNs são treinadas para aproximar as probabilidades a posteriori dos estados de saída dada toda a seqüência de vetores acústicos $P(q_k | \mathbf{X}, \Theta)$ (em oposição aos alvos binários 0 e 1 impostos à K 's saídas da rede durante o período de treinamento), então (3.17) se mantém válido [14]. Em outras palavras, se algum sistema especialista independente fornecer alvos para uma ANN, a rede será capaz de aprender a produzir as probabilidades a posteriori destes novos alvos dado o vetor de entrada. Esta propriedade será utilizada no procedimento de reestimação de parâmetros dos modelos híbridos ANN-HMM, a ser discutido no Capítulo 4.

3.2.1 O Somatório das Saídas da ANN Deve Ser Igual a Um

Para que o vetor de saída da ANN $\mathbf{g}(\mathbf{x}_n, \Psi)$, para um determinado vetor acústico de entrada \mathbf{x}_n , possa ser interpretado como um função distribuição de probabilidade, então $\sum_k g_k(\mathbf{x}_n, \Psi) = 1$. Se a função logística for utilizada na camada de saída e se a rede convergir para um mínimo global então esta propriedade estará assegurada. Porém se a ANN convergir para um mínimo local então não haverá garantias que o somatório de suas saídas seja igual a 1. Para contornar este problema a literatura relata duas alternativas :

1. Após a otimização dos parâmetros Ψ realizar a seguinte normalização : $\bar{g}_k(\mathbf{x}_n, \Psi) = \frac{g_k(\mathbf{x}_n, \Psi)}{\sum_i g_i(\mathbf{x}_n, \Psi)}$, $\forall k \in \{1, 2, \dots, K\}$.
2. Durante a otimização dos parâmetros, e também durante a avaliação das saídas da ANN, substituir a função logística da camada de saída pela função “*softmax*” [4] que possui a propriedade de garantir $\sum_k g_k(\mathbf{x}_n, \Psi) = 1$. Para a k -ésima unidade (neurônio) da camada de saída a função “*softmax*” é definida como :

$$g_k(\mathbf{x}_n) = \frac{\exp(v_k(\mathbf{x}_n))}{\sum_{k=1}^K \exp(v_k(\mathbf{x}_n))} \quad (3.18)$$

e sua derivada pode ser calculada de forma extremamente simples, por :

$$g_k'(\mathbf{x}_n) = g_k(\mathbf{x}_n)(1 - g_k(\mathbf{x}_n)) \quad (3.19)$$

3.3 Estimação de Verossimilhanças para HMMs através de ANN

Uma vez que a saída da rede neural aproxima probabilidades Bayesianas, $g_k(\mathbf{x}_n, \Theta)$ é uma estimativa de

$$P(q_k|\mathbf{x}_n) = \frac{p(\mathbf{x}_n|q_k)}{p(\mathbf{x}_n)} \cdot P(q_k) \quad (3.20)$$

As verossimilhanças de emissão de símbolos $p(\mathbf{x}_n|q_k)$ a ser utilizadas pelos HMMs podem ser obtidas dividindo-se as saídas da rede, $g_k(\mathbf{x}_n)$, pela frequência relativa da classe q_k no conjunto de treinamento $P(q_k)$, e multiplicando pela densidade de probabilidade dos símbolos, $p(\mathbf{x}_n)$.

$$p(\mathbf{x}_n|q_k) = \frac{P(q_k|\mathbf{x}_n)}{P(q_k)} \cdot p(\mathbf{x}_n) \quad (3.21)$$

Richard e Lippmann, [30], afirmam que como as probabilidades a priori das classes $P(q_k)$, que

aparecem em (3.20), representam apenas um termo multiplicativo, então estas probabilidades podem ser alteradas durante a fase de reconhecimento com o objetivo de compensar dados de treinamento com probabilidades de classes que não sejam representativas das condições de teste. No sistema implementado fez-se uso desta propriedade, porque a base de dados utilizada para treinamento do sistema não era balanceada, e o valor de $P(q_k)$ foi estimado segundo um outro procedimento conforme será discutido na Seção 6.3.2.

Durante o reconhecimento o valor de $p(\mathbf{x}_n)$ é constante para todas as classes (estados do HMM) e não exercerá qualquer influência sobre o processo de classificação. Portanto é usual utilizar em sistemas híbridos ANN-HMM a verossimilhança de emissão de símbolos normalizada

$$\frac{p(\mathbf{x}_n|q_k)}{p(\mathbf{x}_n)} = \frac{P(q_k|\mathbf{x}_n)}{P(q_k)} \quad (3.22)$$

A Equação (3.22) mostra que se forem conhecidas boas estimativas de $P(q_k)$ então o treinamento discriminativo da ANN (utilizando o critério MSE ou entropia relativa) irá garantir, indiretamente, boas estimativas de $\frac{p(\mathbf{x}_n|q_k)}{p(\mathbf{x}_n)}$, isto é, boas estimativas das verossimilhanças de emissão de símbolos normalizadas $\frac{p(\mathbf{x}_n|q_k)}{p(\mathbf{x}_n)}$ estão condicionadas a bons estimadores tanto de probabilidades a posteriori $P(q_k|\mathbf{x}_n)$ quanto de probabilidades a priori $P(q_k)$.

3.4 Vantagens da ANN

ANNs possuem várias vantagens que as fazem particularmente atrativas para reconhecimento automático de fala, por exemplo :

- Podem providenciar aprendizado discriminativo entre unidades da fala ou estados de HMM que são representados por saídas de ANN treinadas como classificadoras de padrões. Isto é, quando treinadas para classificação (utilizando funções de custo como MSE ou entropia relativa), os parâmetros são ajustados para que as saídas desta ANN minimizem a taxa de erros e ao mesmo tempo maximizem a discriminação entre a classe de saída correta e as demais classes de saída. Em outras palavras, ANNs não se preocupam apenas em otimizar os parâmetros de cada classe a partir dos seus respectivos exemplos de treinamento, mas também em rejeitar os exemplos de treinamento não associados a ela. O critério de ML não apresenta este caráter discriminativo uma vez que não se preocupa em minimizar a taxa de erro.

- Uma vez que as ANNs podem incorporar múltiplas restrições e encontrar uma combinação ótima destas restrições para classificação, os vetores acústicos não necessitam ser assumidos independentes. Mais genericamente, não existe necessidade da forte suposição sobre as distribuições estatísticas dos dados de entrada, como usualmente é requerido por um HMM padrão.
- Ao contrário de HMMs treinados a partir do critério ML, as ANN não necessitam supor que as funções densidade de probabilidade de emissão de símbolos são formadas pela combinação de funções densidade de probabilidade gaussianas.
- Possuem arquiteturas bastante flexíveis que podem facilmente acomodar informações contextuais e realimentações, e também permitem o uso de entradas binárias ou contínuas.
- ANNs são estruturas altamente paralelas e regulares, o que fazem delas especialmente receptivas para implementações de arquitetura e hardware de alto-desempenho.

3.5 Tipos de Modelos Híbridos ANN-HMM

H. Bourlard e N. Morgan [22] classificam os sistemas híbridos em dois tipos básicos : (1) Modelos híbridos ANN-HMM padrão; (2) Modelos híbridos ANN-HMM discriminativo. Os modelos híbridos que têm sido discutidos até este momento e que foram simulados nesta Tese são do tipo padrão, mas com o objetivo de motivar trabalhos futuros com modelos híbridos discriminativos, também será apresentada uma breve discussão sobre este modelo. Para maiores informações sobre modelos híbridos ANN-HMM discriminativos, veja König [14].

3.5.1 Modelo Híbrido ANN-HMM Padrão

Consiste de um HMM contínuo com verossimilhanças de emissão de símbolos normalizadas $\frac{p(\mathbf{x}_n|q_k)}{p(\mathbf{x}_n)}$, estimadas por uma ANN. Neste caso o treinamento é realizado com um critério discriminativo, em geral a minimização do erro quadrático médio - MSE e estes HMMs permitem a avaliação de $\frac{p(\mathbf{X}|M_i,\Theta)}{p(\mathbf{X}|\Theta)}$, da verossimilhança normalizada de uma sentença \mathbf{X} dado um determinado modelo M_i e o conjunto de parâmetros Θ . Este tipo de modelo é perfeitamente adequado ao modelamento acústico dos sistemas para reconhecimento de fala contínua discutidos no Capítulo 2.

O processo de avaliação dos modelos híbridos ANN-HMM padrão consiste em estimar o modelo oculto de Markov \widehat{M} , dados os parâmetros acústicos Θ , que apresente a maior verossimilhança de

gerar a seqüência de símbolos observados $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$, isto é :

$$\widehat{M} = \arg \max_{M_i} p(\mathbf{X}|M_i, \Theta) \iff \arg \max_{M_i} \frac{p(\mathbf{X}|M_i, \Theta)}{p(\mathbf{X}|\Theta)} \quad (3.23)$$

Para expressar $\frac{p(\mathbf{X}|M_i, \Theta)}{p(\mathbf{X}|\Theta)}$ em função de parâmetros disponíveis, considere o desenvolvimento a seguir. Escrevendo $\frac{p(\mathbf{X}|M_i, \Theta)}{p(\mathbf{X}|\Theta)}$ em termos de probabilidade marginal, tem-se :

$$\frac{p(\mathbf{X}|M_i, \Theta)}{p(\mathbf{X}|\Theta)} = \frac{1}{p(\mathbf{X}|\Theta)} \cdot p(\mathbf{X}|M_i, \Theta) = \frac{1}{p(\mathbf{X}|\Theta)} \cdot \sum_{\Gamma \in \tau_i} p(\mathbf{X}, \Gamma|M_i, \Theta) \quad (3.24)$$

sendo τ_i o conjunto de todas as seqüências de estados permitidas pelo modelo M_i .

Manipulando (3.24) , tem-se :

$$\frac{p(\mathbf{X}|M_i, \Theta)}{p(\mathbf{X}|\Theta)} = \frac{1}{p(\mathbf{X}|\Theta)} \cdot \sum_{\Gamma \in \tau_i} p(\mathbf{X}|\Gamma, M_i, \Theta) \cdot P(\Gamma|M_i, \Theta) \quad (3.25)$$

$$= \sum_{\Gamma \in \tau_i} \frac{p(\mathbf{X}|\Gamma, M_i, \Theta)}{p(\mathbf{X}|\Theta)} \cdot P(\Gamma|M_i, \Theta) \quad (3.26)$$

$P(\Gamma|M_i, \Theta)$ pode ser determinado a partir das probabilidades de transição,

$$P(\Gamma|M_i, \Theta) = \pi_1 \cdot \prod_{n=1}^{N-1} P(q_{\Gamma}^{n+1}|q_{\Gamma}^n, M_i, \Theta) \quad (3.27)$$

$\frac{p(\mathbf{X}|\Gamma, M_i, \Theta)}{p(\mathbf{X}|\Theta)}$ pode ser reescrito como :

$$\frac{p(\mathbf{X}|\Gamma, M_i, \Theta)}{p(\mathbf{X}|\Theta)} = \frac{p(\mathbf{x}_1|\Gamma, M_i, \Theta) \cdot p(\mathbf{x}_2|\Gamma, M_i, \Theta, \mathbf{x}_1) \cdot \dots \cdot p(\mathbf{x}_N|\Gamma, M_i, \Theta, \mathbf{x}_{N-1}, \dots, \mathbf{x}_1)}{p(\mathbf{x}_1|\Theta) \cdot p(\mathbf{x}_2|\Theta, \mathbf{x}_1) \cdot \dots \cdot p(\mathbf{x}_N|\Theta, \mathbf{x}_{N-1}, \dots, \mathbf{x}_1)} \quad (3.28)$$

Considerando os símbolos não condicionados estatisticamente e lembrando que a verossimilhança de emissão de um símbolo normalizada no instante n depende apenas do estado neste mesmo instante n , então :

$$\frac{p(\mathbf{X}|\Gamma, M_i, \Theta)}{p(\mathbf{X}|\Theta)} = \prod_{n=1}^N \frac{p(\mathbf{x}_n|q_{\Gamma}^n, M_i, \Theta)}{p(\mathbf{x}_n|\Theta)} \quad (3.29)$$

Utilizando as Equações (3.27) e (3.29) $\frac{p(\mathbf{X}|M_i, \Theta)}{p(\mathbf{X}|\Theta)}$ pode ser expresso por :

$$\frac{p(\mathbf{X}|M_i, \Theta)}{p(\mathbf{X}|\Theta)} = \sum_{\Gamma \in \tau_i} \pi_1 \cdot \frac{p(\mathbf{x}_1|q_{\Gamma}^1, M_i, \Theta)}{p(\mathbf{x}_1|\Theta)} \cdot \prod_{n=1}^{N-1} P(q_{\Gamma}^{n+1}|q_{\Gamma}^n, M_i, \Theta) \cdot \frac{p(\mathbf{x}_{n+1}|q_{\Gamma}^{n+1}, M_i, \Theta)}{p(\mathbf{x}_{n+1}|\Theta)} \quad (3.30)$$

3.5.2 Modelo Híbrido ANN-HMM Discriminativo

Consiste de um “HMM” contínuo com *distribuições de probabilidades de transição condicionada* (veja Equação (3.17)) estimadas por uma ANN recursiva, conforme Figura 3.2. Neste caso o treinamento é realizado com um critério discriminativo, em geral a minimização do erro quadrático médio - MSE ou entropia relativa e estes “HMMs”¹ permitem a avaliação da probabilidade a posteriori do modelo M_i dada a sentença \mathbf{X} e os parâmetros do modelo Θ , $P(M_i|\mathbf{X}, \Theta)$. Este modelo é chamado discriminativo porque tanto o treinamento quanto a avaliação são realizados utilizando-se o método de maximização da probabilidade a posteriori MAP. Em sistemas de reconhecimento de fala contínua baseados em modelos híbridos ANN-HMM discriminativos utilizam-se explicitamente *informações da língua L* e neste caso o objetivo do sistema é determinar o modelo M_i com maior probabilidade a posteriori de ter gerado a seqüência de símbolos acústicos \mathbf{X} , dado os parâmetros Θ e L ,

$$\widehat{M} = \arg \max_{M_i} P(M_i|\mathbf{X}, \Theta, L) = \arg \max_{\Gamma \in \tau_i} \sum_{\Gamma} P(\Gamma|\mathbf{X}, L, \Theta) \cdot P(M_i|\Gamma, \mathbf{X}, L, \Theta) \quad (3.31)$$

Assumindo Γ independente da informação da língua L e M_i independente da seqüência acústica \mathbf{X} se Γ for especificado, então,

$$P(M|\mathbf{X}, \Theta, L) = \sum_{\Gamma \in \tau_i} P(\Gamma|\mathbf{X}, \Theta) \cdot P(M_i|\Gamma, L, \Theta) \quad (3.32)$$

Neste caso $P(\Gamma|\mathbf{X}, \Theta)$ é determinado a partir de um *modelamento acústico* de cada sentença e $P(M_i|\Gamma, L, \Theta)$ é determinado por um *modelo da língua*.

Desconsiderando as informações da língua L , o processo de avaliação dos modelos híbridos ANN-HMM discriminativos consiste em estimar o modelo oculto de Markov \widehat{M} que maximiza a probabilidade a posteriori $P(M_i|\mathbf{X}, \Theta)$, isto é, a probabilidade do modelo de Markov M dada a seqüência de símbolos $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ e os parâmetros acústicos Θ ,

$$\widehat{M} = \arg \max_{M_i} P(M_i|\mathbf{X}, \Theta) \quad (3.33)$$

Considere o desenvolvimento proposto por H. Bourlard [4]. Escrevendo $P(M_i|\mathbf{X}, \Theta)$ em termos de

¹O termo HMMs foi colocado entre aspas por não se tratar exatamente de HMMs, pois fazem uso de probabilidades a posteriori e não de verossimilhanças de emissão de símbolos.

probabilidade marginal, resulta :

$$P(M_i|\mathbf{X}, \Theta) = \sum_{\Gamma \in \mathfrak{S}} P(\Gamma, M_i|\mathbf{X}, \Theta) \quad (3.34)$$

$$= \sum_{\Gamma \in \mathfrak{S}} P(M_i|\mathbf{X}, \Gamma, \Theta) \cdot P(\Gamma|\mathbf{X}, \Theta) \quad (3.35)$$

Sendo \mathfrak{S} o conjunto de todas as seqüências de estado permitidas por todos os modelos $M_i = \{1, 2, \dots, I_M\}$.

Assumindo que M_i não precisa ser condicionado a \mathbf{X} se Γ for especificado e que $\tau_i \cap \tau_j = \emptyset \forall i \neq j$, então :

$$P(M_i|\mathbf{X}, \Theta) = \sum_{\Gamma \in \tau_i} P(\Gamma|\mathbf{X}, \Theta) \quad (3.36)$$

porque neste caso $p(M_i|\Gamma, \Theta) = 0 \forall \Gamma \notin \tau_i$ e $P(M_i|\Gamma, \Theta) = 1 \forall \Gamma \in \tau_i$.

Reescrevendo $P(\Gamma|\mathbf{X}, \Theta)$ como

$$P(\Gamma|\mathbf{X}, \Theta) = P(q_\Gamma^2|\mathbf{X}, q_\Gamma^1, \Theta) \cdot P(q_\Gamma^3|\mathbf{X}, q_\Gamma^2, q_\Gamma^1, \Theta) \cdot \dots \quad (3.37)$$

$$= \dots \cdot P(q_\Gamma^N|\mathbf{X}, q_\Gamma^{N-1}, q_\Gamma^{N-2}, \dots, q_\Gamma^1, \Theta) \quad (3.38)$$

resulta,

$$P(\Gamma|\mathbf{X}, \Theta) = \prod_{n=1}^N P(q_\Gamma^n|\mathbf{X}, Q_1^{n-1}, \Theta) \quad (3.39)$$

sendo que q_Γ^n representa o estado observado no tempo n e Q_1^N a seqüência de estados associada com \mathbf{X}_1^N . Portanto as probabilidades $P(\Gamma|\mathbf{X}, \Theta) = P(q_\Gamma^1, q_\Gamma^2, \dots, q_\Gamma^N|\mathbf{X}, \Theta)$ podem ser calculadas em termos de probabilidades "locais" $P(q_\Gamma^n|\mathbf{X}, Q_1^{n-1}, \Theta)$. Estas probabilidades podem ser simplificadas se algumas restrições condicionais forem relaxadas : (1) Suposição de modelo de Markov de primeira ordem (estado atual dependente apenas do estado prévio), (2) Assumir que o estado q_Γ^n está condicionado apenas aos $c + d + 1$ símbolos adjacentes a ele (o símbolo atual, os c anteriores e os d posteriores). Estas suposições permitem aproximar as probabilidades "locais" $P(q_\Gamma^n|\mathbf{X}, Q_1^{n-1}, \Theta)$ por:

$$P(q_\Gamma^n|\mathbf{X}, Q_1^{n-1}, \Theta) \approx P(q_\Gamma^n|\mathbf{X}_{n-c}^{n+d}, q_\Gamma^{n-1}, \Theta) \quad (3.40)$$

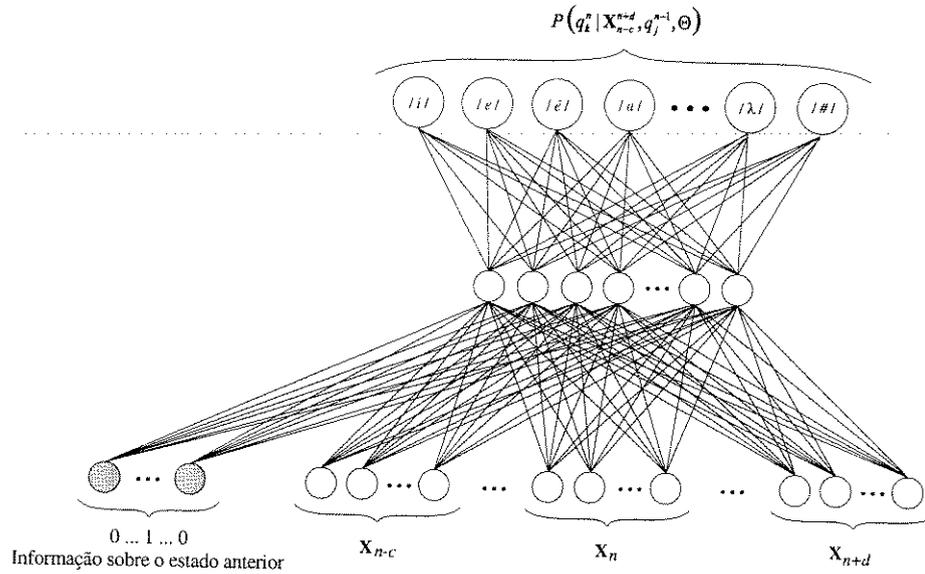


Figura 3-2: Rede neural recursiva para estimativa das “probabilidades de transição condicionadas” a serem utilizadas em um modelo híbrido ANN-HMM discriminativo.

Utilizando (3.39) e (3.40) em (3.36), $P(M_i | \mathbf{X}, \Theta)$ pode ser aproximado por :

$$P(M_i | \mathbf{X}, \Theta) \approx \sum_{\Gamma \in \mathcal{r}_i} \prod_{n=1}^N P(q_{\Gamma}^n | \mathbf{X}_{n-c}^{n+d}, q_{\Gamma}^{n-1}, \Theta) \quad (3.41)$$

Estas probabilidades podem ser estimadas pelas saídas de uma ANN tendo com entrada \mathbf{X}_{n-c}^{n+d} (o símbolo atual e mais informações contextuais) e mais uma realimentação (informações sobre o estado no instante anterior q_{Γ}^{n-1}). A Figura 3-2 mostra um esquema para esta rede neural.

Capítulo 4

REMAP

4.1 Introdução

A questão sobre o reconhecimento de unidades sub-lexicais ser sensível ou não a contexto tem sido estudada por Konig, [14] entre outros. Em seus experimentos Konig verificou o papel das transições dos formantes no reconhecimento de vogais. Especificamente foi testado se o reconhecimento de uma vogal é função, principalmente, da região estacionária dos formantes, ou do contexto acústico a curto-termo - direção e taxa de transição dos formantes adjacentes. Nos experimentos, ouvintes americanos foram submetidos à identificação de monossílabos sem sentido. Cada sílaba *consoante-vogal-consoante* (CVC) era constituída por uma seqüência de três elementos : *transição - vogal alvo - transição*. Os padrões de formantes das vogais alvo foram selecionados de segmentos estacionários das vogais [*i*] e [*u*]. A taxa e a direção das transições adjacentes foram variadas pela escolhas de dois contextos consonantais : [*w - w*] e [*j - j*]. Os resultados destes experimentos mostram que a identificação dos estímulos das vogais são determinados não apenas pelos padrões dos formantes na região estacionária, mas também pela direção das transições dos formantes adjacentes. Por exemplo o mesmo padrão de vogal : F1 = 350 Hz, F2 = 1578 Hz, F3 2604 Hz foi reconhecido por todos os ouvintes como [*i*] no contexto de [*j V j*] e foi reconhecido como [*u*] no contexto [*w V w*]. Em geral, foi mostrado que a categorização do contínuo é ajustada de acordo com os diferentes contextos para poder compensar o efeito de “Undershoot” no estímulo da vogal.

Motivado pelos resultados dos experimentos citados acima, Konig [14] propôs um método para reestimação dos parâmetros de um modelo híbrido ANN-HMM discriminativo baseado tanto em informações sobre as regiões estáveis das unidades sub-lexicais como nas regiões de transição entres estas unidades sub-lexicais. Em 1997 Cole [36] estendeu este método para a reestimação de parâmetros de

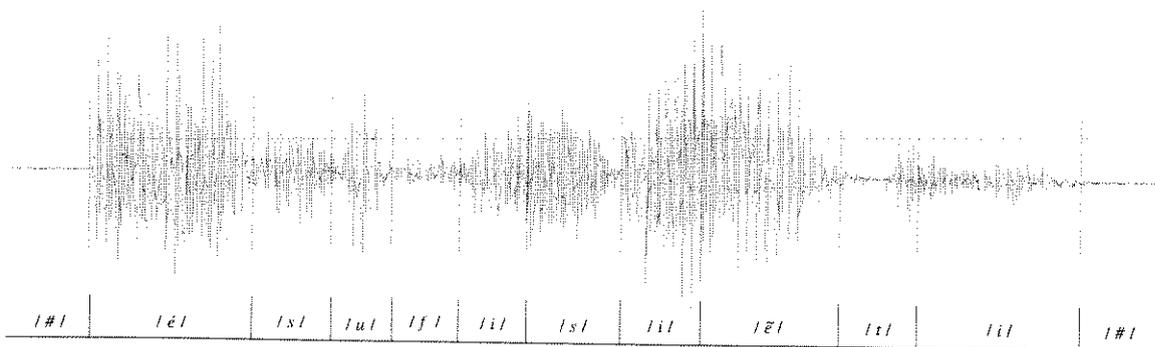


Figura 4-1: Sentença : “ É suficiente ” segmentada em termos de unidades sub-lexicais (fones).

um modelo híbrido ANN-HMM padrão. Nas próximas seções será apresentado e analisado o método proposto por Cole [36].

4.2 Estimação de Parâmetros - Pré-REMAP

Em um modelo híbrido ANN-HMM o treinamento da ANN é realizado de forma supervisionada, necessitando, portanto, de um conjunto de exemplos de treinamento $\mathbf{X}_e = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{N_e}\}$ e de seus correspondentes alvos desejados $\mathbf{d}_e(\mathbf{x}_n) = (d_{e_1}(\mathbf{x}_n), \dots, d_{e_k}(\mathbf{x}_n), \dots, d_{e_K}(\mathbf{x}_n))^t$, isto é, dos pares, $(\mathbf{x}_n, \mathbf{d}_e(\mathbf{x}_n))$. Em geral os exemplos de treinamento são acompanhados de informações contextuais e portanto a n -ésima entrada da ANN será $\{\mathbf{x}_{n-c}, \dots, \mathbf{x}_n, \dots, \mathbf{x}_{n+c}\} = \mathbf{X}_{n-c}^{n+c}$. Conforme já discutido do Capítulo 3 estes exemplos podem ser obtidos através da segmentação das sentenças de treinamento (base de dados) em termos de unidades sub-lexicais. Esta segmentação (manual ou automática) fornecerá marcas ao longo das sentenças indicando as posições inicial e final de cada uma das unidades sub-lexicais, conforme Figura 4-1.

Em geral as unidades sub-lexicais possuem durações distintas, porém o número de entradas da ANN é fixo, o que exige que os exemplos de treinamento $\mathbf{x}_i \in \mathbf{X}_e$ sejam extraídos com o uso de uma janela de tamanho fixo ($\mathbf{x}_i = \mathbf{X}_{n-c}^{n+c}$, com $c = cte, \forall i$). Um alternativa é tomar estas janelas com uma extensão igual à duração médias de todas as unidades sub-lexicais (previamente segmentadas) e centradas entre as marcas de início e fim de cada unidade (esta foi a abordagem utilizada neste trabalho). O conjunto de parâmetros da ANN treinada com o uso de \mathbf{X}_e será denominado Ψ_e . A Figura 4-2 mostra um trecho da forma de onda da frase “É suficiente” e ilustra as posições onde devem ser retirados os exemplo de treinamento, para o caso de em que $c = 1$ (os exemplos de treinamentos correpondem a três janelas de análise).

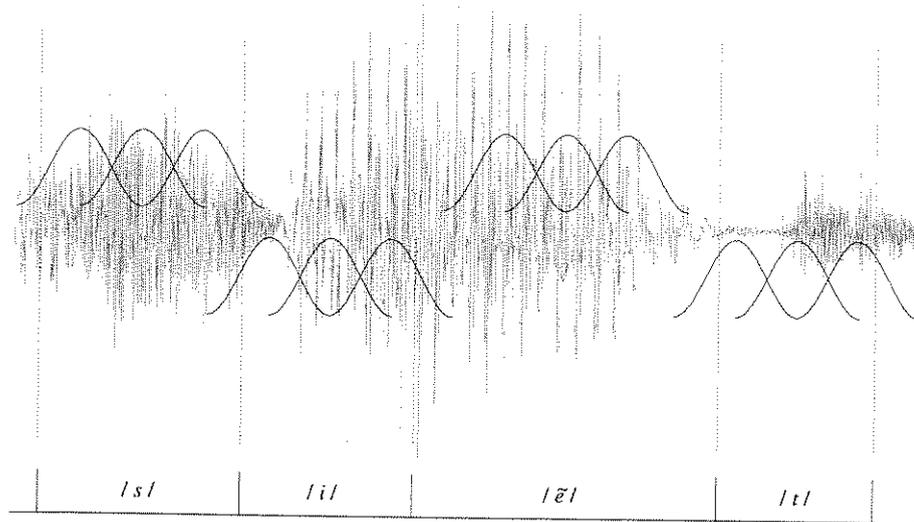


Figura 4-2: Ilustração do processo de retirada dos exemplos de treinamento da ANN, centrados entre as marcas de segmentação das unidades sub-lexicais.

Durante a fase de reconhecimento a ANN treinada com o conjunto de exemplos \mathbf{X}_e descritos acima, terá que ser capaz de reconhecer vetores de entrada extraídos para cada janela de análise (em geral as janelas de análise são de 20 ms, mas tomadas a cada 10 ms). Muito provavelmente esta rede apresentará um bom desempenho para os vetores de entrada razoavelmente centrados em unidades sub-lexicais e um baixo desempenho para os vetores de entrada próximos a transições entre unidades sub-lexicais. Esta dificuldade em reconhecer regiões de transição é óbvia, pois tais informações não fazem parte do conjunto de treinamento. A Figura 4-3 mostra um trecho da sentença “É suficiente” e apresenta as posições dos vetores de entrada a serem utilizados durante a fase de reconhecimento (para o caso de $c = 1$). Para este exemplo a ANN deve apresentar um bom desempenho para as quadros de análise centrados nas unidades sub-lexicais, instantes : $n + 1$, $n + 4$, $n + 7$, $n + 8$, $n + 9$ e $n + 12$, e um desempenho ruim para os quadro de análise localizados nas transições entre estas unidades sub-lexicais, instantes $n + 2$, $n + 3$, $n + 5$, $n + 6$, $n + 10$ e $n + 11$. Uma solução para contornar este problema é retreinar a ANN com o uso de um novo conjunto de exemplos de treinamento \mathbf{X}_e' , com exemplos extraídos a cada quadro de análise, obtendo assim exemplos de treinamento centrados nas marcas de segmentação e exemplos de treinamento localizados em regiões de transição entre as unidades sub-lexicais. A grande dificuldade para a definição deste novo conjunto de treinamento consiste na definição do conjunto dos sinais alvos correspondentes \mathbf{d}_e' .

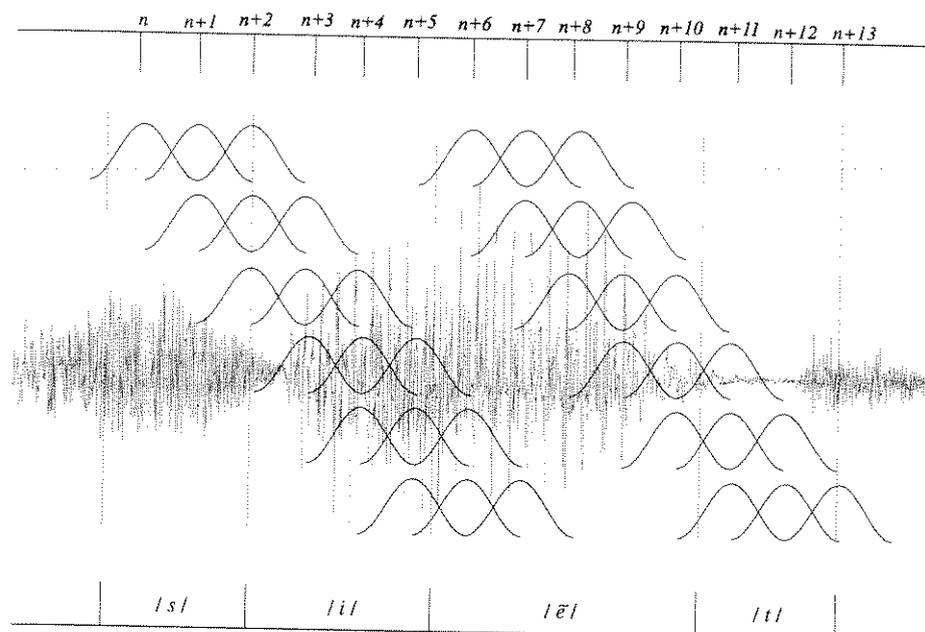


Figura 4-3: Exemplos utilizados durante a etapa de avaliação - exemplos tomados a cada quadro de análise.

4.3 Reestimação de Parâmetros - REMAP

4.3.1 Alvos Abruptos

A forma mais imediata para a determinação dos sinais alvos \mathbf{d}_e' a serem utilizados na reestimação dos parâmetros da ANN consiste em avaliar a ANN (com parâmetros Ψ_e) para cada uma das sentenças de treinamento e em seguida calcular a trajetória de Viterbi para cada uma destas sentenças e utilizar todos os pontos destas trajetórias, um para cada vetor acústico da sentença, como o conjunto de alvos $\mathbf{d}_e'(\mathbf{x}_n) = (d_{1e}'(\mathbf{x}_n), \dots, d_{k_e}'(\mathbf{x}_n), \dots, d_{K_e}'(\mathbf{x}_n))^t$. O problema desta abordagem é que estes alvos representam decisões abruptas, isto é, um exemplo de treinamento $\mathbf{x}_i \in \mathbf{X}_e'$ próximo a uma transição entre duas unidades sub-lexicais seria classificado pela ANN como tendo características pertencentes a apenas umas destas unidades sub-lexicais, desconsiderando quaisquer informações transicionais, que conforme foi discutido na introdução deste Capítulo, são cruciais para o processo de reconhecimento.

4.3.2 Alvos Suaves

Alvos mais adequados para a reestimação dos parâmetros da ANN devem ser capazes de apresentar uma mudança gradual de uma unidade sub-lexical para outra. Uma maneira de obter estes alvos é a proposta por Cole, 1997 [36], e consiste basicamente na determinação das probabilidades a posteriori

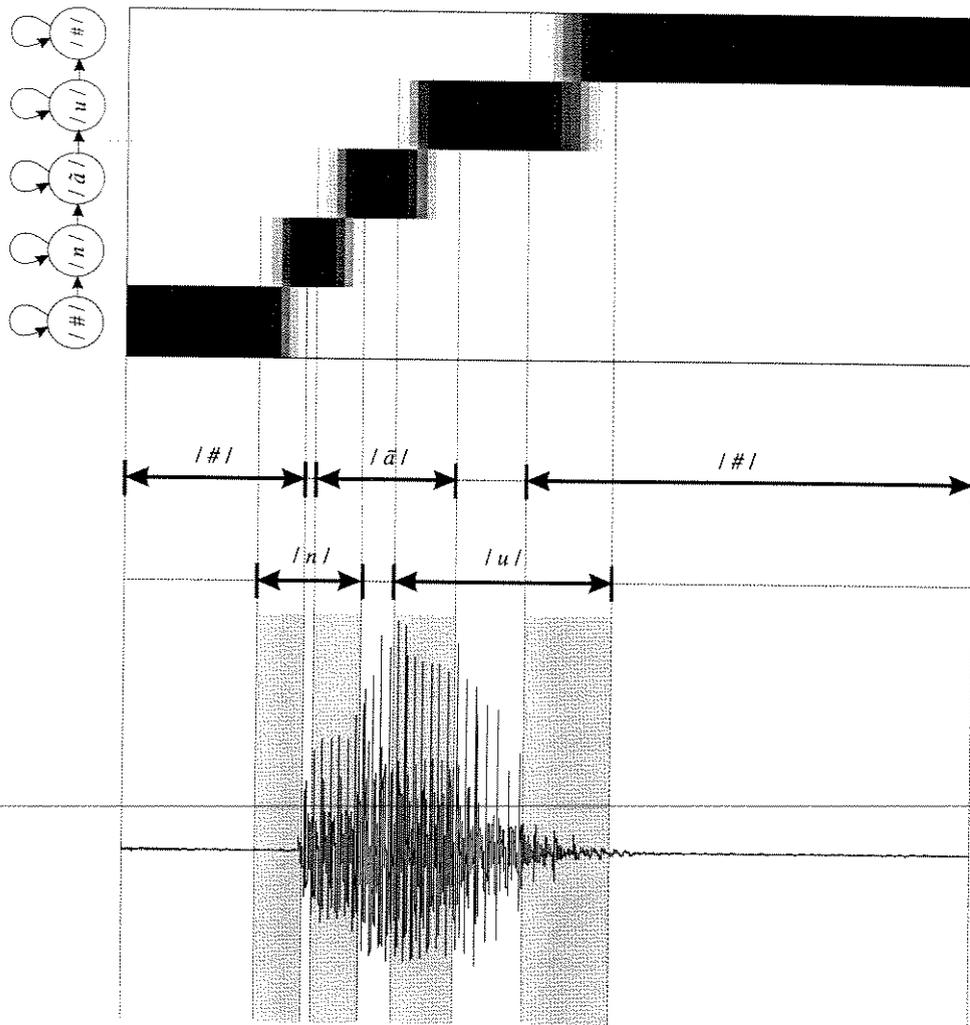


Figura 4-4: Alvos suaves obtidos segundo as probabilidades a posteriori dos estados dada toda a seqüência de vetores acústicos “probabilidades a posteriori globais”.

globais de um estado q_k dada uma seqüência completa de vetores acústicos \mathbf{X} correspondente a uma sentença M , isto é, o k -ésimo alvo da ANN para o exemplo \mathbf{x}_n será dado por $\gamma_n(k) = P(q_k^n | \mathbf{X}, M, \Theta)$. A Figura 4-4 mostra os alvos suaves¹ $\gamma_n(k)$ para o caso da sentença “Não”².

Para a obtenção destas probabilidades a posteriori globais basta avaliar a ANN (com parâmetros Ψ_e) para cada uma das sentenças e, de posse das probabilidades a posteriori locais, $P(q_k^n | \mathbf{x}_n, M, \Theta)$ (saídas da ANN), utilizar o algoritmo *Forward – Backward* para a determinação de $\gamma_n(k)$. A seguir será apresentado o algoritmo *Forward – Backward* para a determinação destas probabilidades a posteriori

¹Na Figura 4-4 os valores dos alvos suaves são dados em termos de tons de cinza : preto = 1; branco = 0.

²O símbolo /#/ corresponde ao silêncio.

globais.

Algoritmo Forward-Backward

O objetivo deste algoritmo é determinar as probabilidades a posteriori globais $\gamma_n(k) = P(q_k^n | \mathbf{X}, M, \Theta)$ em função das verossimilhanças de emissão de símbolos normalizada $\frac{p(\mathbf{x}_n | q_k^n, M, \Theta)}{p(\mathbf{x}_n | \Theta)}$. A seguir será apresentado um desenvolvimento detalhado para o cálculo de $\gamma_n(k)$ utilizando a verossimilhança de emissão de símbolos sem normalização e no final deste desenvolvimento será demonstrado que o uso do fator de normalização não acarreta qualquer alteração no cálculo de $\gamma_n(k)$.

$\gamma_n(k) = P(q_k^n | \mathbf{X}, M, \Theta)$ pode ser escrito como,

$$P(q_k^n | \mathbf{X}, M, \Theta) = \frac{P(q_k^n, \mathbf{X} | M, \Theta)}{p(\mathbf{X} | M, \Theta)} \quad (4.1)$$

$$= \frac{p(\mathbf{X}_1^n, \mathbf{X}_{n+1}^N, q_k^n | M, \Theta)}{p(\mathbf{X} | M, \Theta)} \quad (4.2)$$

$$= \frac{p(\mathbf{X}_1^n, q_k^n | M, \Theta) \cdot p(\mathbf{X}_{n+1}^N | \mathbf{X}_1^n, q_k^n, M, \Theta)}{p(\mathbf{X} | M, \Theta)} \quad (4.3)$$

Assumindo os vetores acústicos condicionalmente independentes, então :

$$P(q_k^n | \mathbf{X}, M, \Theta) = \frac{p(\mathbf{X}_1^n, q_k^n | M, \Theta) \cdot p(\mathbf{X}_{n+1}^N | q_k^n, M, \Theta)}{p(\mathbf{X} | M, \Theta)} \quad (4.4)$$

Definindo

$$\alpha_n(k) = p(\mathbf{X}_1^n, q_k^n | M, \Theta) \text{ e } \beta_n(k) = p(\mathbf{X}_{n+1}^N | q_k^n, M, \Theta) \quad (4.5)$$

então :

$$\gamma_n(k) = P(q_k^n | \mathbf{X}, M, \Theta) = \frac{\alpha_n(k) \cdot \beta_n(k)}{p(\mathbf{X} | M, \Theta)} \quad (4.6)$$

Procedimento recursivo para estimar $\alpha_n(k)$:

$$\alpha_1(k) = \pi_k \cdot b_k(\mathbf{x}_1) \quad (4.7)$$

$$\alpha_n(k) = p(\mathbf{X}_1^n, q_k^n | M, \Theta) \quad (4.8)$$

$$\alpha_{n+1}(k) = p(\mathbf{X}_1^{n+1}, q_k^{n+1} | M, \Theta) \quad (4.9)$$

$$= p(\mathbf{X}_1^n, \mathbf{x}_{n+1}, q_k^{n+1} | M, \Theta) \quad (4.10)$$

$$= \sum_j p(\mathbf{X}_1^n, \mathbf{x}_{n+1}, q_j^n, q_k^{n+1} | M, \Theta) \quad (4.11)$$

$$= \sum_j p(\mathbf{x}_{n+1}, q_k^{n+1} | q_j^n, \mathbf{X}_1^n, M, \Theta) \cdot \overbrace{p(\mathbf{X}_1^n, q_j^n | M, \Theta)}^{\alpha_n(j)} \quad (4.12)$$

Como os vetores acústicos são condicionalmente independentes, tem-se :

$$\alpha_{n+1}(k) = \sum_j \alpha_n(j) \cdot p(\mathbf{x}_{n+1}, q_k^{n+1} | q_j^n, M, \Theta) \quad (4.13)$$

$$= \sum_j \alpha_n(j) \cdot p(\mathbf{x}_{n+1} | q_k^{n+1}, q_j^n, M, \Theta) \cdot P(q_k^{n+1} | q_j^n, M, \Theta) \quad (4.14)$$

$$= \sum_j \alpha_n(j) \cdot \underbrace{p(\mathbf{x}_{n+1} | q_k^{n+1}, M, \Theta)}_{b_k(\mathbf{x}_{n+1})} \cdot \underbrace{P(q_k^{n+1} | q_j^n, M, \Theta)}_{a_{jk}} \quad (4.15)$$

Portanto

$$\alpha_{n+1}(k) = \sum_j \alpha_n(j) \cdot a_{jk} \cdot b_k(\mathbf{x}_{n+1}) \quad (4.16)$$

Procedimento recursivo para estimar $\beta_n(k)$

$$\beta_n(k) = p(\mathbf{X}_{n+1}^N | q_k^n, M) \quad (4.17)$$

$$\beta_n(k) = p(\mathbf{X}_{n+2}^N, \mathbf{x}_{n+1} | q_k^n, M) \quad (4.18)$$

$$= \sum_j p(X_{n+2}^N, \mathbf{x}_{n+1}, q_j^{n+1} | q_k^n, M) \quad (4.19)$$

$$= \sum_j \frac{p(\mathbf{x}_{n+1}, q_k^n | q_j^{n+1}, \mathbf{X}_{n+2}^N, M) \cdot \overbrace{p(\mathbf{X}_{n+2}^N | q_j^{n+1}, M)}^{\beta_{n+1}(j)} \cdot P(q_j^{n+1}, M)}{P(q_k^n, M)} \quad (4.20)$$

Como os símbolos são não condicionados estatisticamente :

$$\beta_n(k) = \sum_j \beta_{n+1}(j) \cdot \frac{p(\mathbf{x}_{n+1}, q_k^n | q_j^{n+1}, M) \cdot p(q_j^{n+1}, M)}{P(q_k^n, M)} \quad (4.21)$$

$$= \sum_j \beta_{n+1}(j) \cdot \underbrace{p(\mathbf{x}_{n+1}, | q_j^{n+1}, M)}_{b_j(\mathbf{x}_{n+1})} \cdot \underbrace{p(q_j^{n+1} | q_k^n, M)}_{a_{kj}} \quad (4.22)$$

Trocando os índices j por k , apenas para melhorar a notação, tem-se :

$$\beta_n(j) = \sum_k \beta_{n+1}(k) \cdot a_{jk} \cdot b_k(\mathbf{x}_{n+1}) \quad (4.23)$$

De (4.16) é fácil observar que $p(\mathbf{X}|M, \Theta) = \sum_{\forall k} \alpha_N(k)$. Logo, resulta :

$$P(q_k^n | \mathbf{X}, M, \Theta) = \frac{\alpha_n(k) \cdot \beta_n(k)}{\sum_k \alpha_N(k)} \quad (4.24)$$

A substituição de $b_k(\mathbf{x}_{n+1})$ (Equações (4.16) e (4.23) pela verossimilhança de emissão de símbolos normalizada $\frac{p(\mathbf{x}_n | q_k^n, M, \Theta)}{p(\mathbf{x}_n | \Theta)}$ transformam a probabilidade a posteriori global em,

$$\hat{P}(q_k^n | \mathbf{X}, M, \Theta) = \frac{\alpha_n(k) \cdot C_n \cdot \beta_n(k) \cdot D_n}{C_N \cdot \sum_k \alpha_N(k)} \quad (4.25)$$

onde : $C_n = \prod_{i=1}^n \frac{1}{p(\mathbf{x}_i)}$ e $D_n = \prod_{i=N}^{n+1} \frac{1}{p(\mathbf{x}_i)}$

Como $C_n \cdot D_n = C_N$, então pode ser verificado que o uso de verossimilhanças de emissão símbolos normalizadas não implica em qualquer alteração na probabilidade a posteriori global, isto é :

$$\hat{P}(q_k^n | \mathbf{X}, M, \Theta) = P(q_k^n | \mathbf{X}, M, \Theta) \quad (4.26)$$

4.3.3 Algoritmo para Reestimação

O processo de reestimação de parâmetros de um modelo híbrido ANN-HMM pelo algoritmo REMAP pode ser sistematizado de acordo com os seguintes passos :

1. **Iteração $i = 0$ (Pré-REMAP)**. Iniciar com os parâmetros $\Psi^i = \Psi_e$, da ANN otimizada apenas com o uso dos exemplos de treinamento \mathbf{X}_e centrados entre as marcas de segmentação das unidades sub-lexicais. Utilizar como probabilidades de transição entre os estados $P^{(i)}(q_k^n | q_j^{n-1})$ e como probabilidades das classes $P^{(i)}(q_k)$, estimativas realizadas apenas em função de \mathbf{X}_e . Propostas para a estimativa destas probabilidades encontram-se no Capítulo 6, Equações (6.1) e (6.3), respectivamente.
2. Através da Equação (4.24) e dos parâmetros Ψ^i , $P^{(i)}(q_k^n | q_j^{n-1})$ e $P^{(i)}(q_k)$ estimar os alvos suaves $P^{(i+1)}(q_k^n | \mathbf{X}_{M_j}^{(s)}, M_j, \Theta^i)$ para cada uma das seqüências de treinamento $\mathbf{X}_{M_j}^{(s)}$ associadas a M_j . Realizar este procedimento para todos os modelos M_j .
3. Para cada uma das seqüências $\mathbf{X}_{M_j}^{(s)} = \{\mathbf{x}_1^{(j,s)}, \mathbf{x}_2^{(j,s)}, \dots, \mathbf{x}_N^{(j,s)}\}$ associadas a todas as sentenças M_j , com $j \in \mathcal{I}_{\mathcal{M}}$, otimizar os parâmetros $\Psi^{(i+1)}$ da ANN para minimizar o erro quadrático médio - MSE (ou entropia relativa), entre a saída da rede $g_k^{(i+1)}(\mathbf{x}_n^{(j,s)}, \Psi^{(i+1)})$ e o correspondente alvo

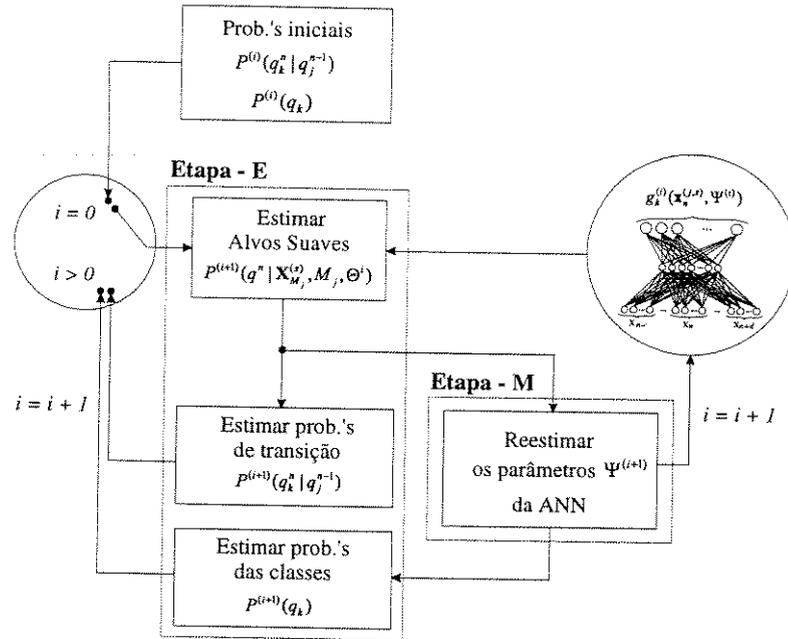


Figura 4-5: Diagrama de blocos do processo de reestimação de parâmetros REMAP.

suave $P^{(i+1)}(q_k^n | \mathbf{X}_{M_j}^{(s)}, M_j, \Theta^i)$. Este procedimento³ fornecerá os novos parâmetros $\Psi^{(i+1)}$ da ANN.

4. A partir dos parâmetros $\Psi^{(i+1)}$ reestimados, calcular as novas probabilidades das classes $P^{(i+1)}(q_k)$, utilizando por exemplo, a Equação (6.4).
5. Ainda utilizando os parâmetros $\Psi^{(i)}$ e as probabilidades $P^{(i)}(q_k)$ e $P^{(i)}(q_k^n | q_j^{n-1})$, estimar os valores das novas probabilidades de transição $P^{(i+1)}(q_k^n | q_j^{n-1})$ conforme a Equação (4.34).
6. Se o sistema ainda “*não convergiu*”, fazer $i = i + 1$ e voltar ao passo 2.

Este procedimento pode ser dividido em duas etapas principais : (1) Estimação - **Etapa-E**, correspondente aos passos 2, 4 e 5; (2) Maximização - **Etapa-M**, correspondente ao passo 3. O algoritmo REMAP é ilustrado na Figura 4-5. O objetivo da **Etapa-M** é aumentar a cada iteração i , o número esperado de estados (saídas da ANN) individualmente “corretos”⁴. Para isto procura minimizar o erro quadrático médio entre a saída da rede, $g_k^{(i)}(\mathbf{x}_n, \Psi^{(i)})$ e a função de probabilidade a posteriori global $P^{(i)}(q^n | \mathbf{X}_{M_j}^{(s)}, M_j, \Theta^{(i-1)})$ que por sua vez maximiza este valor esperado [27].

³Este procedimento de treinamento *seqüência-por-seqüência* assume a independência de cada uma das seqüências de vetores acústicos em relação as demais.

⁴Corretos segundo o critério “*forward - backward*”.

4.3.4 Verossimilhanças de Emissão de Símbolos

A cada iteração do processo de reestimação, um novo conjunto de pesos sinápticos Ψ estará disponível, e as saídas da ANN, $g_k(\mathbf{x}_n, \Psi)$ serão novas estimativas das probabilidades a posteriori do estado q_k dado a entrada \mathbf{x}_n , isto é, $P(q_k|\mathbf{x}_n, \Theta)$. Estas estimativas poderão ser convertidas através da regra de Bayes e das novas probabilidades das classes $P(q_k)$ em verossimilhanças de emissão de símbolos normalizadas $\frac{p(\mathbf{x}_n|q_k, \Theta)}{p(\mathbf{x}_n|\Theta)}$. Porém neste caso as probabilidades das classes $P(q_k)$ não poderão ser obtidas através da frequência relativa, uma vez que os alvos são suaves e não permitem a contagem direta da ocorrência ou não de q_k . Uma alternativa para a determinação de $P(q_k)$ é a proposta por Lippmann, 1991 [30], Seção 6.3.2, Equação (6.4).

4.3.5 Probabilidades de Transição

Considere que para um determinado modelo específico M_i ,

$$\xi_n^{(i,s)}(j, k) = P(q_j^n, q_k^n | \mathbf{X}_{M_i}^{(s)}, M_i, \Theta) \quad (4.27)$$

sendo que $\mathbf{X}_{M_i}^{(s)} = \{\mathbf{x}_1^{(i,s)}, \mathbf{x}_2^{(i,s)}, \dots, \mathbf{x}_{N_i}^{(i,s)}\}$ corresponde à s -ésima seqüência de vetores acústicos associada a M_i . Através de algumas manipulações simples $\xi_n^{(i,s)}(j, k)$ pode ser reescrito como⁵ :

$$\xi_n^{(i,s)}(j, k) = \frac{P(q_j^n, q_k^n, \mathbf{X}_{M_i}^{(s)} | M_i, \Theta)}{p(\mathbf{X}_{M_i}^{(s)} | M_i, \Theta)} = \frac{\alpha_n^{(i,s)}(j) \cdot a_{jk} \cdot b_k(\mathbf{x}_{n+1}^{(i,s)}) \cdot \beta_{n+1}^{(i,s)}(k)}{\sum_{\forall j} \alpha_N^{(i,s)}(j)} \quad (4.28)$$

sendo $\alpha_n^{(i,s)}(j)$ e $\beta_{n+1}^{(i,s)}(k)$ as variáveis *forward* e *backward* para o modelo M_i (conforme (4.5)) quando for apresentada a seqüência $\mathbf{X}_{M_i}^{(s)}$.

Considerando o somatório das probabilidades a posteriori globais, $\gamma_n^{(i,s)}(k) = P(q_k^n | \mathbf{X}_{M_i}^{(s)}, M_i, \Theta)$ ao longo do tempo, pode ser realizada a seguinte interpretação :

$$\sum_{n=1}^{N_i} \gamma_n^{(i,s)}(k) = \text{número esperado de vezes em que o estado } q_k \text{ é visitado.} \quad (4.29)$$

⁵Assim como no caso de $\gamma_n^{(i,s)}$, o uso de verossimilhanças normalizadas no cálculo de $\xi_n^{(i,s)}$ implicará em fatores extras no numerador : $C_n \cdot \frac{1}{p(\mathbf{x}_{n+1})}$ e D_{n+1} e também no denominador : C_{N_i} . Mas pode ser verificado que $C_n \cdot \frac{1}{p(\mathbf{x}_{n+1})} \cdot D_{n+1} = C_{N_i}$.

Se o somatório for até $N - 1$, a Equação (4.29) pode ser interpretada da seguinte maneira,

$$\sum_{n=1}^{N_i-1} \gamma_n^{(i,s)}(k) = \text{número esperado de transições a partir de } q_k. \quad (4.30)$$

Esta afirmação baseia-se no fato que se o estado q_k é visitado, obrigatoriamente tem que haver uma transição a partir dele. Do mesmo modo, o somatório dos $\xi_n^{(i,s)}(j, k)$ pode ser interpretado da seguinte forma :

$$\sum_{n=1}^{N_i-1} \xi_n^{(i,s)}(j, k) = \text{número esperado de transições de } q_j \text{ para } q_k. \quad (4.31)$$

Utilizando-se as equações (4.30) , (4.31) e o conceito de frequência relativa de ocorrência, a equação para reestimação das probabilidades de transição $P(q_k^n | q_j^n, M_i, \Theta)$ é dada por

$$P(q_k^n | q_j^n, M_i, \Theta) = \frac{\text{núm. esperado de transições do estado } q_j \text{ para } q_k}{\text{núm. esperado de transições a partir do estado } q_j} \quad (4.32)$$

Portanto, tem-se

$$P(q_k^n | q_j^n, M_i, \Theta) = \frac{\sum_{s=1}^{N_{e_i}} \left[\sum_{n=1}^{N_i-1} \xi_n^{(i,s)}(j, k) \right]}{\sum_{s=1}^{N_{e_i}} \left[\sum_{n=1}^{N_i-1} \gamma_n^{(i,s)}(j) \right]} \quad (4.33)$$

A Equação 4.33 refere-se apenas a uma única sentença M_i , porém durante a reestimação deseja-se realizar o cálculo das probabilidades de transição levando-se em consideração todas as sentenças de treinamento M_i com $i = \{1, 2, \dots, I_M\}$. Para isto, pode ser utilizado o conceito de treinamento com múltiplas observações [27] que assume a independência entre as seqüências de vetores acústicos associadas aos modelos M_i . Como as fórmulas de reestimação das probabilidades de transição baseiam-se no conceito de frequência relativa de ocorrência, a equação para reestimação a partir de múltiplas observações pode ser obtida realizando-se o quociente entre o número esperado de transições do estado q_j para o estado q_k e o número esperado de transições a partir do estado q_j , para todas as sentenças de treinamento,

$$a_{jk} = \frac{\sum_{i=1}^{I_M} \left\{ \sum_{s=1}^{N_{e_i}} \left[\sum_{n=1}^{N_i-1} \xi_n^{(i,s)}(j, k) \right] \right\}}{\sum_{i=1}^{I_M} \left\{ \sum_{s=1}^{N_{e_i}} \left[\sum_{n=1}^{N_i-1} \gamma_n^{(i,s)}(j) \right] \right\}} \quad (4.34)$$

Capítulo 5

Algoritmos de Busca

5.1 Introdução

O problema de reconhecimento de uma sentença falada de forma contínua baseado na concatenação de modelos de palavras (ou de modelos de unidades sub-lexicais) é extremamente importante para a tarefa de reconhecimento automático de fala. Uma grande variedade de algoritmos para a solução deste problema, em geral baseados em técnicas de programação dinâmica - DP (“Dynamic Program”), tem sido proposta e avaliada nos últimos anos [23, 29, 17]. O primeiro algoritmo para reconhecimento de sentenças a partir da concatenação de modelos de palavras foi proposto por Vintsyuk [35], que mostrou como a técnica de DP poderia ser utilizada para determinar a seqüência ótima de palavras associada a uma determinada seqüência de vetores acústicos observados. O procedimento proposto por Vintsyuk processava os sinais de fala de forma síncrona com os vetores acústicos, e seu trabalho pioneiro estabeleceu os fundamentos para as várias soluções baseadas em programação dinâmica para reconhecimento de fala. Vintsyuk também propôs um esquema rudimentar para incorporar restrições entre as palavras durante a pesquisa (regras gramaticais).

Atualmente existe uma variedade de algoritmos para encontrar a “melhor” seqüência de palavras associada a uma determinada seqüência de vetores acústicos observados. Entre eles destacam-se o “Stack Algorithm” desenvolvido por Jelinek [13]; o “Level Building” de Rabiner [29] e o “One Stage” (também conhecido como “Frame Synchronous Level Building”) de Ney [23] e Lee & Rabiner [17]. Todos estes algoritmos são capazes de utilizar elaboradas regras de restrições gramaticais durante o procedimento de busca pela “melhor” seqüência de palavras.

As diferenças entre estes métodos são principalmente de implementação (por exemplo: síncrono por vetor acústico ou síncrono por palavras), e de características (por exemplo: diferentes formas de impor

restrições de duração, ou diferentes procedimentos para geração das sentenças a serem analisadas a partir dos modelos das palavras ou das unidades sub-lexicais).

5.2 Level Building

5.2.1 Motivação

Sendo $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ uma seqüência de vetores acústicos correspondente a uma sentença desconhecida e sendo M o modelo oculto de Markov associado a uma determinada sentença, então o objetivo do algoritmo “Level Building” - LB, é avaliar de forma eficiente a seguinte expressão

$$\widehat{M} = \arg \max_M p(\mathbf{X}|M, \Theta) \cdot P(M|\Theta) \quad (5.1)$$

Isto corresponde a estimar o modelo \widehat{M} que maximiza o produto entre a verossimilhança da seqüência de símbolos \mathbf{X} dado o modelo M , e o modelo da língua $P(M|\Theta)$.

Considerando o caso em que todas as palavras do léxico podem seguir quaisquer outras palavras com probabilidade igual a um (sem uso de restrições gramaticais),

$$P(M|\Theta) = P(W_1|\Theta) \cdot \prod_{i=2}^{N_M} P(W_i|W_{i-1}, W_{i-2}, \dots, W_1, \Theta) = 1 \quad (5.2)$$

onde N_M é número de palavras presente em M , então o processo de decodificação resume-se a,

$$\widehat{M} = \arg \max_M p(\mathbf{X}|M, \Theta) \quad (5.3)$$

Neste caso será dito que a estimação de \widehat{M} será realizada utilizando-se apenas a *decodificação acústica*.

Utilizando-se o conceito de probabilidade marginal, cada modelo M da Equação (5.3) pode ser avaliado pela seguinte expressão :

$$p(\mathbf{X}|M, \Theta) = \sum_{\Gamma \in \tau} p(\mathbf{X}, \Gamma|M, \Theta) \quad (5.4)$$

sendo τ o conjunto de todas as seqüências de estado de comprimento N permitidas por M .

Porém o algoritmo LB avalia (5.4) fazendo uso da aproximação de Viterbi,

$$p(\mathbf{X}|M, \Theta) \approx p(\mathbf{X}, \Gamma_v|M, \Theta) \quad (5.5)$$

sendo que,

$$\Gamma_v = \arg \max_{\Gamma \in \tau} p(\mathbf{X}, \Gamma | M, \Theta) \quad (5.6)$$

onde Γ_v é denominada seqüência de estados (ou caminho) de Viterbi.

Além de utilizar a aproximação em (5.5) o algoritmo LB explora o fato dos modelos de sentenças poderem ser interpretados como a concatenação de modelos de palavras, $M = \{W_1, W_2, \dots, W_{N_M}\}$. Para entender como isto pode ser feito, acompanhe o seguinte desenvolvimento.

Reescrevendo $p(\mathbf{X}, \Gamma_v | M, \Theta)$ como,

$$p(\mathbf{X}, \Gamma_v | M, \Theta) = p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N, q_{\Gamma_v}^1, q_{\Gamma_v}^2, \dots, q_{\Gamma_v}^N | M, \Theta) \quad (5.7)$$

$$= \frac{p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N, q_{\Gamma_v}^1, q_{\Gamma_v}^2, \dots, q_{\Gamma_v}^N, M | \Theta)}{P(M | \Theta)} \quad (5.8)$$

definindo : $1 < n_1 < n_2 < \dots < n_{N_M} < N$, e realizando algumas manipulações simples com (5.8), tem-se :

$$p(\mathbf{X}, \Gamma_v | M, \Theta) = p(\mathbf{x}_1, \dots, \mathbf{x}_{n_1}, q_{\Gamma_v}^1, \dots, q_{\Gamma_v}^{n_1} | \mathbf{x}_{n_1+1}, \dots, \mathbf{x}_N, q_{\Gamma_v}^{n_1+1}, \dots, q_{\Gamma_v}^N | M, \Theta) \cdot \quad (5.9)$$

$$p(\mathbf{x}_{n_1+1}, \dots, \mathbf{x}_{n_2}, q_{\Gamma_v}^{n_1+1}, \dots, q_{\Gamma_v}^{n_2} | \mathbf{x}_{n_2+1}, \dots, \mathbf{x}_N, q_{\Gamma_v}^{n_2+1}, \dots, q_{\Gamma_v}^N | M, \Theta) \cdot$$

$$\vdots$$

$$p(\mathbf{x}_{n_{N_M}+1}, \dots, \mathbf{x}_N, q_{\Gamma_v}^{n_{N_M}+1}, \dots, q_{\Gamma_v}^N | M, \Theta)$$

Assumindo a independência condicional dos vetores acústicos, e a hipótese de Markov de primeira ordem, então (5.10) pode ser reescrito como :

$$p(\mathbf{X}, \Gamma_v | M, \Theta) = p(\mathbf{x}_1, \dots, \mathbf{x}_{n_1}, q_{\Gamma_v}^1, \dots, q_{\Gamma_v}^{n_1} | M, \Theta) \cdot \quad (5.10)$$

$$p(\mathbf{x}_{n_1+1}, \dots, \mathbf{x}_{n_2}, q_{\Gamma_v}^{n_1+1}, \dots, q_{\Gamma_v}^{n_2} | M, \Theta) \cdot$$

$$\vdots$$

$$p(\mathbf{x}_{n_{N_M}+1}, \dots, \mathbf{x}_N, q_{\Gamma_v}^{n_{N_M}+1}, \dots, q_{\Gamma_v}^N | M, \Theta)$$

Se n_1 for definido tal que a seqüência $q_{\Gamma_v}^1, \dots, q_{\Gamma_v}^{n_1}$ corresponda aos estados percorridos ao longo do modelo W_1 e os demais n_i 's sejam definidos de tal forma que as seqüências $q_{\Gamma_v}^{n_{i-1}+1}, \dots, q_{\Gamma_v}^{n_i}$ correspondam aos estados percorridos ao longo dos modelos W_i 's, $2 \leq i \leq N_M$, então (5.10) pode ser reescrito

como :

$$p(\mathbf{X}, \Gamma_v | M, \Theta) = p(\mathbf{X}_1^{n_1}, \Gamma_{v_1} | W_1, \Theta) \cdot p(\mathbf{X}_{n_1+1}^{n_2}, \Gamma_{v_2} | W_2, \Theta) \cdot \dots \cdot p(\mathbf{X}_{n_{N_M}+1}^N, \Gamma_{v_{N_M}} | W_{N_M}, \Theta) \quad (5.11)$$

sendo Γ_{v_1} o caminho de Viterbi ao longo de W_1 começando em $n = 1$ e finalizando em $n = n_1$; e Γ_{v_i} o caminho de Viterbi ao longo de W_i começando em $n = n_{i-1} + 1$ e finalizando em $n = n_i$, com $2 \leq i \leq N_M$.

O algoritmo LB explora justamente a propriedade expressa pela Equação (5.11), que afirma que $p(\mathbf{X}, \Gamma_v | M, \Theta)$ pode ser calculada como o produto das verossimilhanças, ao longo do caminho de Viterbi, de cada uma das N_M palavras presentes na sentença M .

5.2.2 Algoritmo Level Building

Nesta seção será apresentado o algoritmo LB para realização somente da *decodificação acústica*, (caso em que $P(M|\Theta) = 1$). Durante o processo de decodificação o algoritmo LB desconhece a princípio o número de palavras que compõem a sentença a ser reconhecida e portanto torna-se necessário estabelecer uma quantidade mínima e máxima de palavras a serem procuradas. No desenvolvimento que será apresentado a seguir o número mínimo será de uma palavra e o número máximo será de L palavras. Esta faixa de procura de 1 à L será denominada de *níveis de procura*.

Nível 1 (Primeira palavra da sentença) : Calcula-se $p(\mathbf{X}|W_i, \Theta)$, para todas as palavras W_i do léxico, iniciando-se no instante $n = 1$. Este cálculo é realizado utilizando-se o algoritmo de Viterbi, Figura 5-1, de acordo com os seguintes passos:

1. **Inicialização** - Definir $\delta_n(k)$ como a probabilidade conjunta das seqüências parciais de estado e de observação, $p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n, q_{\Gamma_i}^1, q_{\Gamma_i}^2, \dots, q_{\Gamma_i}^n | W_i, \Theta)$ e fazer

$$\begin{aligned} \delta_1(1) &= p(\mathbf{x}_1, q_1 | W_i, \Theta) \\ \delta_1(k) &= 0, \quad k = 2, 3, \dots, N_{W_i} \\ \alpha_n(N_{W_i}) &= 1, \quad n = 1, 2, \dots, N \end{aligned} \quad (5.12)$$

sendo N_{W_i} o número de estados na palavra W_i , e $\alpha_n(N_{W_i}) = 1$ a indicação de que todas as trajetórias de Viterbi no primeiro nível devem começar no instante $n = 1$.

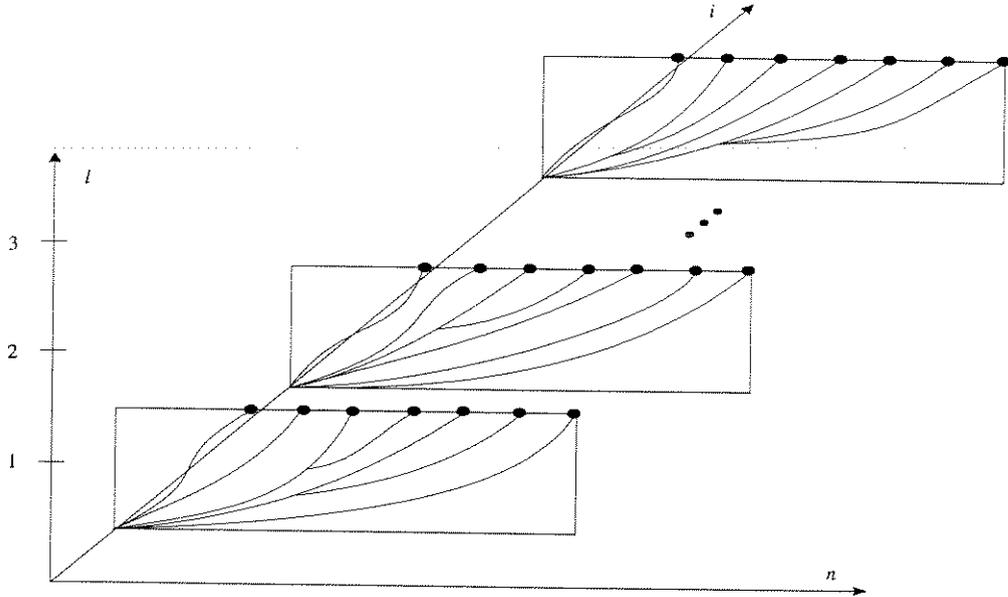


Figura 5-1: Avaliação de todas as palavras do léxico (eixo i) no primeiro nível de procura.

2. **Recursão** - para $2 \leq n \leq N$, e $1 \leq k \leq N_{W_i}$

$$\delta_n(k) = \max_{1 \leq j \leq N_{W_i}} \left[\delta_{n-1}(j) \cdot P(q_k^n | q_j^{n-1}, W_i, \Theta) \right] \cdot p(\mathbf{x}_n | q_k, W_i) \quad (5.13)$$

3. **Término** -

$$\begin{aligned} \mathbf{P}(l, n, i) &= \delta_n(N_{W_i}), \quad 1 \leq n \leq N \\ \mathbf{C}(l, n, i) &= \alpha_n(N_{W_i}), \quad 1 \leq n \leq N \\ \mathbf{T}(l, n, i) &= P(q_k^n | q_{N_{W_i}}^{n-1}, W_i, \Theta), \quad 1 \leq n \leq N \end{aligned} \quad (5.14)$$

De acordo com o passo 3, as verossimilhanças acumuladas (para o último estado de cada palavra) para todos os instantes de tempo, devem ser guardadas em uma matriz tridimensional \mathbf{P} . As coordenadas desta matriz são: l - indicador de nível; n - instante de tempo (no espaço de símbolos) e i - identificador da palavra analisada. A matriz \mathbf{C} guardará os instantes iniciais (do nível anterior) dos caminhos de Viterbi finalizados em cada instante de tempo n (do nível atual) para todas as palavras avaliadas. Finalmente, para cada uma das palavras W_i avaliadas, a matriz \mathbf{T} deve armazenar as probabilidades de transição do estado $q_{N_{W_i}}$ (probabilidade de transição do último estado da palavra W_i).

O nível $l = 1$ deve ser realizado para todas as palavras do léxico conforme descrito pelos 3 passos acima. Após todas as palavras serem analisadas deve ser realizado o que será chamado *redução de*

nível (eliminação do eixo i), formando os seguintes vetores,

$$\widehat{\mathbf{P}}(l, n) = \max_i [\mathbf{P}(l, n, i)] \quad (5.15)$$

$$\widehat{\mathbf{C}}(l, n) = \mathbf{C} \left[l, n, \arg \max_i (\mathbf{P}(l, n, i)) \right] \quad (5.16)$$

$$\widehat{\mathbf{W}}(l, n) = \arg \max_i [\mathbf{P}(l, n, i)] \quad (5.17)$$

$$\widehat{\mathbf{T}}(l, n) = \mathbf{T}(l, b, \widehat{\mathbf{W}}(l, n)) \quad (5.18)$$

sendo que $\widehat{\mathbf{P}}$ selecionará, para cada instante de tempo n , o valor da maior verossimilhança $\delta_n(N_{W_i})$ para todas as palavras W_i analisadas no nível l . $\widehat{\mathbf{W}}$ indicará, para cada instante de tempo n , qual a palavra W_i , que apresentou a maior verossimilhança $\delta_n(N_{W_i})$. $\widehat{\mathbf{C}}$ armazenará o instante n do nível anterior ($l-1$) onde foi iniciado o caminho de Viterbi associado a palavra $\widehat{\mathbf{W}}(l, n)$ (no caso do primeiro nível este valor será igual a 1). $\widehat{\mathbf{T}}$ indicará, para cada instante de tempo n , a probabilidade de transição do último estado da palavra $\widehat{\mathbf{W}}(l, n)$.

Nível 2 (e demais níveis) : Os cálculos envolvidos no nível 2, e em todos os níveis restantes, diferem do nível 1 apenas no primeiro passo, correspondente ao procedimento de inicialização,

$$\delta_1(1) = 0, \quad (5.19)$$

$$\delta_n(1) = \max \left[\widehat{\mathbf{T}}(l-1, n-1) \cdot \widehat{\mathbf{P}}(l-1, n-1), P(q_1^n | q_1^{n-1}, W_i, \Theta) \cdot \delta_{n-1}(1) \right] \cdot p(\mathbf{x}_n | q_k, W_i, \Theta) \quad (5.20)$$

com $2 \leq n \leq N$.

$$\alpha_n(1) = \begin{cases} n-1, & \text{se } \widehat{\mathbf{T}}(l-1, n-1) \cdot \widehat{\mathbf{P}}(l-1, n-1) > P(q_1^n | q_1^{n-1}, W_i, \Theta) \cdot \delta_{n-1}(1) \\ \alpha_{n-1}(1), & \text{caso contrário} \end{cases} \quad (5.21)$$

A Equação (5.19) ajusta $\delta_1(1) = 0$ e a Equação (5.20) faz a ligação do primeiro estado do nível atual com o último estado do nível anterior. O vetor α_n representado pela Equação (5.21) guarda as informações necessárias para o resgate dos instantes do nível anterior onde foram iniciados os caminhos de Viterbi associados a cada instante n do nível atual. Durante o procedimento de recursão α_n deve ser ajustada por,

$$\alpha_n(k) = \alpha_{n-1} \cdot \left(\arg \max_{1 \leq i \leq N} \left(\delta_{n-1}(i) \cdot P(q_k^n | q_j^{n-1}, W_i, \Theta) \right) \right) \quad (5.22)$$

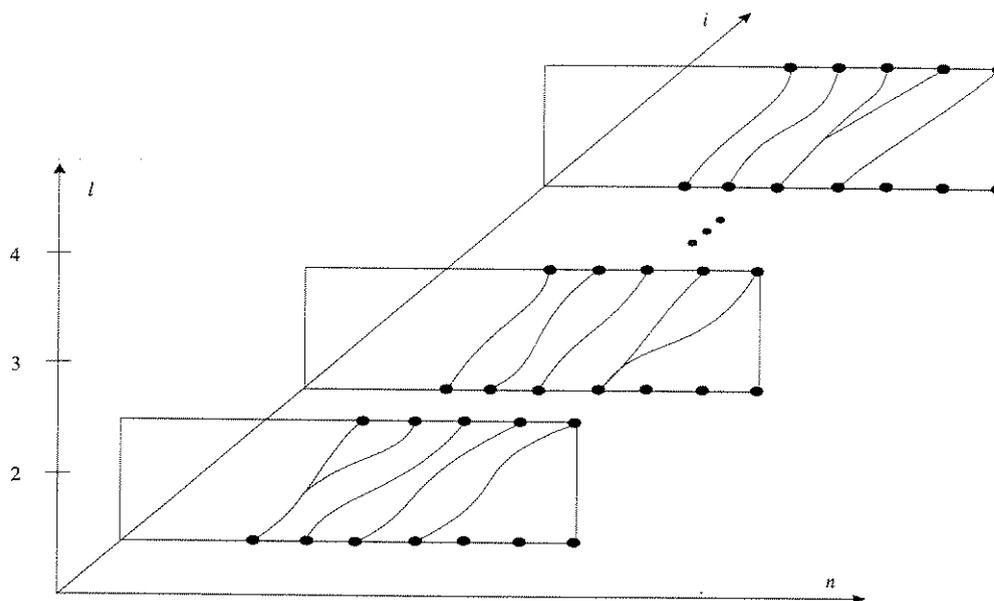


Figura 5-2: Avaliação de todas as palavras do léxico (eixo i) no segundo nível de procura.

Após a aplicação de todas as palavras neste nível, conforme ilustrado na Figura 5-2, os vetores $\hat{\mathbf{P}}$, $\hat{\mathbf{C}}$, $\hat{\mathbf{W}}$ e $\hat{\mathbf{T}}$ devem ser calculados (redução de nível), conforme as equações (5.15), (5.16), (5.17) e (5.18), e então seguir para análise do próximo nível.

Tecnicamente falando, ambos $\delta_n(k)$ e $\alpha_n(k)$ são funções do índice da palavra i . Entretanto, uma vez que as informações relevantes já foram salvas nos vetores $\mathbf{P}(l, n, i)$, $\mathbf{C}(l, n, i)$ e $\mathbf{T}(l, n, i)$, o uso explícito da dependência em relação à palavra i não é necessário nestas funções.

Todo este procedimento termina quando for realizado o último nível L , previamente especificado. Para identificar quantas palavras N_M existem na seqüência de palavras com maior verossimilhança, basta realizar a seguinte análise,

$$N_M = \arg \max_{1 \leq l \leq L} \hat{\mathbf{P}}(l, N) \quad (5.23)$$

Em seguida para encontrar quais são as N_M palavras associadas a esta seqüência mais provável, basta realizar o caminho reverso (“*Backtracking*”) com o uso do vetor $\hat{\mathbf{B}}(l, n)$.

5.2.3 Incorporando Restrições de Duração ao Algoritmo Level Building

Durante a análise de um determinado *nível* o algoritmo LB avalia cada modelo de palavra sob inúmeras condições de compressão e de expansão. Porém, muitas destas compressões e expansões correspondem a durações que diferem demasiadamente das durações médias das palavras presentes nas sentenças de

treinamento do sistema. Desta maneira palavras como “suficiente” podem ser comprimidas excessivamente a ponto de serem confundidas com a palavra “em” e vice-versa. Uma possível solução para este problema é utilizar um modelo de duração para as palavras que leve em consideração a duração média e a variância destas palavras no conjunto de treinamento.

Rabiner, 1985 [29] propõe um modelo de duração bastante simples que associa à i – ésima palavra do léxico, uma função densidade de probabilidade gaussiana $f_i(D)$,

$$f_i(D) = \frac{1}{\sigma_i \sqrt{2\pi}} \cdot \exp \left(-\frac{(D - \bar{D}_i)^2}{2\sigma_i^2} \right)$$

onde, \bar{D}_i e σ_i^2 , representam, respectivamente, a média e a variância, da palavra W_i , obtidas a partir das sentenças de treinamento.

Este modelo de duração pode ser incorporado ao algoritmo LB no final de cada nível e antes do procedimento de redução de nível, de acordo com os seguintes passos :

1. Determinar a duração de cada palavra W_i através do procedimento de “Backtracking”,

$$d_n = n - \mathbf{B}(l, n, i) \text{ para todo } n$$

2. Ponderar a verossimilhança acumulada $\mathbf{P}(l, n, i)$, utilizando a função densidade de probabilidade gaussiana $f_i(\cdot)$ ajustada em função da duração d_n determinada no passo 1,

$$\bar{\mathbf{P}}(l, n, i) = \mathbf{P}(l, n, i) \cdot f_i(d_n) \text{ para todo } n$$

Apesar deste método para incorporar informação duracional ser bastante heurístico, será comprovado, através dos resultados presentes na Tabela 7.2 e no Apêndice C, que na prática ele melhora significativamente a taxa de acertos de palavras do sistema.

5.2.4 Exemplo

Para facilitar a compreensão do algoritmo level building, considere o seguinte exemplo :

1. Número de níveis de busca igual a três (sentenças constituídas por no máximo três palavras).
2. $P(M|\Theta) = 1$ - Não será utilizado nenhum tipo de restrição gramatical.

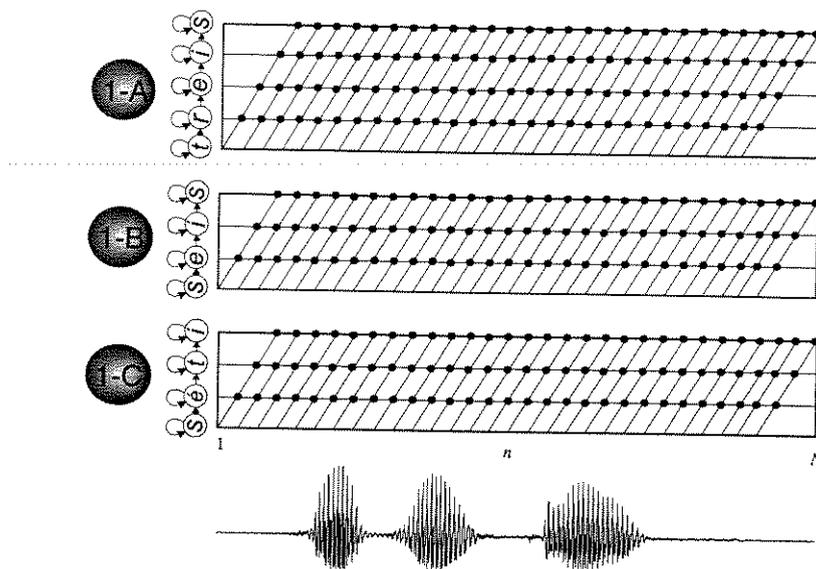


Figura 5-3: Avaliação de todas as sentenças constituídas por apenas uma única palavra.

3. Léxico constituído por apenas três palavras, correspondentes aos dígitos : Três, Seis e Sete. Estes dígitos também serão identificados, respectivamente, pelas letras **A**, **B** e **C**. Serão utilizadas as seguintes conversões ortográfico-fonética.

- Três : $t + r + e + i + s - \mathbf{A}$
- Seis : $s + e + i + s - \mathbf{B}$
- Sete : $s + e + t + i - \mathbf{C}$

A seguir são apresentados os procedimentos que devem ser realizados durante a análise de cada um dos três níveis de procura.

Primerio Nível : Avaliação de todas as sentenças constituídas por apenas uma única palavra As treliças mostradas na Figura 5-3 apresentam todos os possíveis caminhos pelos quais os HMMs das palavras **A**, **B** e **C** podem gerar a seqüência de vetores acústicos associada ao sinal de fala apresentado ao sistema. Em outras palavras estes caminhos descrevem todas as possíveis dilatações temporais (quantizadas em função do tamanho do quadro de análise) das palavras **A**, **B** e **C**.

Segundo Nível : Avaliação de todas as sentenças construídas através da concatenação de duas palavras A Figura 5-4 mostra três das nove sentenças possíveis de serem construídas a partir da combinação de duas palavras.

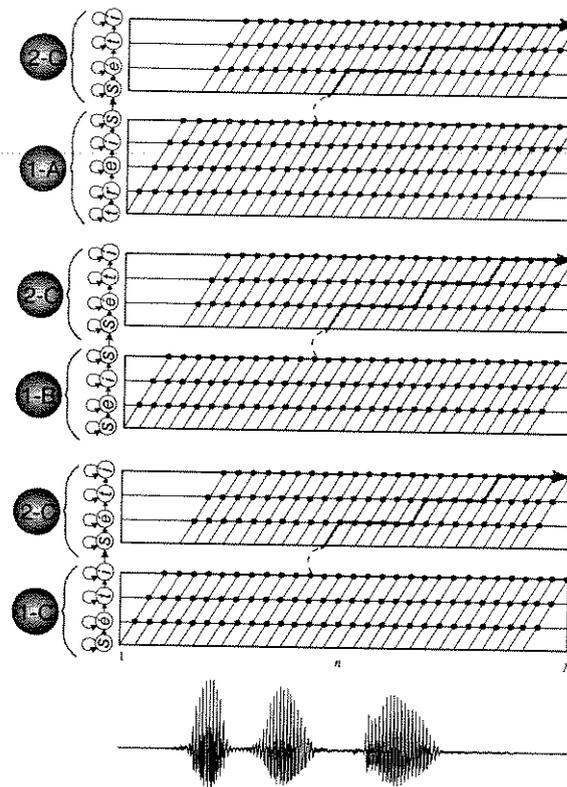


Figura 5-4: Três das nove sentenças possíveis de serem contruídas a partir da concatenação de duas palavras.

A verossimilhança acumulada, para um determinado caminho específico, ao longo da segunda palavra a ser concatenada, será determinada por três parâmetros básicos :

- Verossimilhanças de emissão de símbolos associadas a cada estado.
- Probabilidade de transição de cada estado.
- Probabilidade do estado inicial. Esta probabilidade corresponde a verossimilhança acumulada, no último estado da palavra do nível anterior, para um determinado instante n específico.

A Figura 5-4 mostra a mesma trajetória de Viterbi ao longo da palavra **C** nas sentenças **C-C**, **B-C** e **A-C**. Por se tratar da mesma trajetória de Viterbi, é óbvio que ela compartilha das mesmas verossimilhanças de emissão de símbolos e probabilidades de transição. Portanto, para esta trajetória específica, o parâmetro decisivo na determinação da verossimilhança acumulada com maior valor será a probabilidade do estado inicial. A partir desta análise pode-se concluir que após a realização do

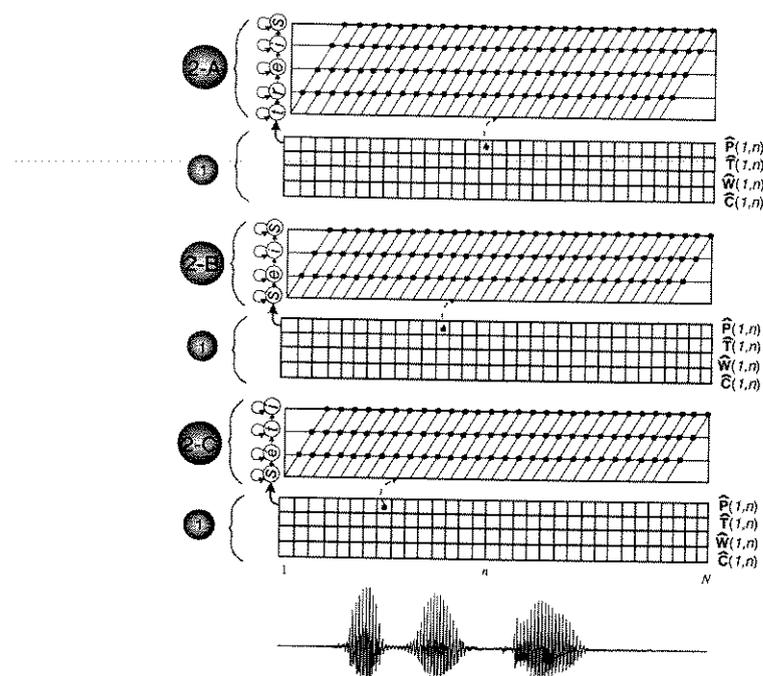


Figura 5-5: Avaliação das palavras candidatas ao segundo nível da sentença, após a redução do primeiro nível (utilizando apenas as informações realmente úteis no primeiro nível).

primeiro nível muitas das informações disponíveis, para cada instante de tempo n são desnecessárias e que a informações realmente úteis, e que portanto devem ser armazenadas, são :

- Entre todas as palavras avaliadas, a verossimilhança de emissão de símbolos acumulada com maior valor - $\hat{P}(1, n)$.
- Indicação da palavra que apresentou esta maior verossimilhança - $\hat{W}(1, n)$.
- Probabilidade de transição associada ao último estado desta palavra vencedora - $\hat{T}(1, n)$.
- Instante inicial onde começou o caminho de Viterbi associado a esta maior verossimilhança acumulada - $\hat{C}(1, n)$.

Este procedimento de análise e armazenamento das informações realmente úteis para cada nível é denominado *redução de nível*.

A Figura 5-5 mostra a avaliação das palavras candidatas ao segundo nível da sentença, porém para o caso em que já foi realizada a *redução do primeiro nível*. Neste caso apenas três combinação são necessárias para avaliar todas as concatenações de duas palavras que podem apresentar os maiores valores de verossimilhança acumulada.

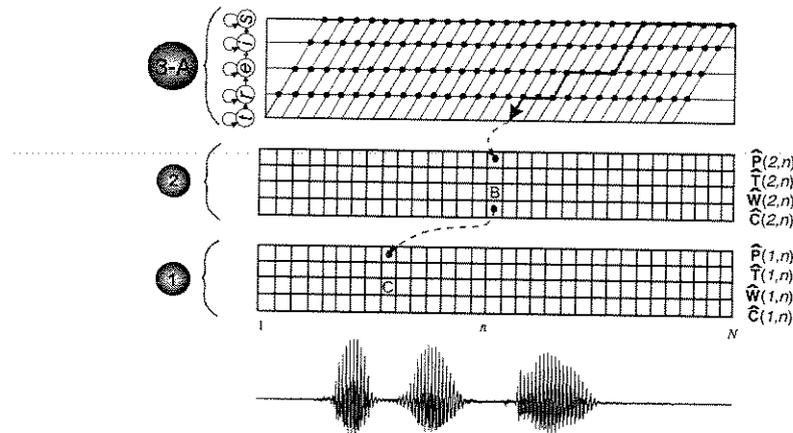


Figura 5-6: Ilustração do procedimento de “backtracing”. Caso em que a sentença reconhecida é composta por três palavras e consiste na concatenação dos dígitos : Sete - Seis - Três.

Terceiro Nível : Avaliação de todas as sentenças construídas através da concatenação de três palavras Após a *redução do primeiro e segundo nível*, avaliar as palavras candidatas ao terceiro nível, verificar a vencedora e em seguida, fazendo uso dos vetores (\hat{P} , \hat{T} , \hat{W} e \hat{C}) com as informações dos níveis já reduzidos, realizar o “backtracing” para identificar a seqüência de palavras com maior verossimilhança de gerar a seqüência de vetores acústicos apresentados ao sistema. A Figura 5-6 ilustra a realização deste procedimento de “backtracing” para o caso em que a sentença com maior verossimilhança é composta por três palavras e consiste na concatenação dos dígitos : **C - B - A** \Leftrightarrow Sete - Seis - Três.

5.3 One Stage

O algoritmo “One Stage” difere do algoritmo “Level Building” apenas em termos de implementação. Para cada *nível de procura* o algoritmo “Level Building” calcula os caminhos de Viterbi de todas as palavras para todos os instantes de tempo e somente depois de todos estes cálculos realiza a redução de nível e passa a análise do *nível seguinte*. Este tipo de procedimento é denominado síncrono por palavra e está limitado à execução serial dos *níveis de procura*. No caso do algoritmo “One Stage”, para o *nível de procura l* calcula-se a trajetória de Viterbi de todas as palavras para um determinado instante de tempo n e, se os valores das verossimilhanças acumuladas nos últimos estados de todas as palavras W_i 's estiverem disponíveis, então realiza-se a redução de *nível* (apenas para o instante de tempo n) e passa-se a executar simultaneamente o *nível l* e o *nível seguinte l + 1*. O tipo de

procedimento utilizado pelo algoritmo “One Stage” é denominado síncrono por quadro (síncrono por janela de análise) e permite a execução paralela dos *níveis* de procura.

Outra vantagem do algoritmo “One Stage” é que, ao contrário do algoritmo “Level Building”, ele não necessita esperar a chegada de toda a sequência de vetores acústicos para iniciar o processo de decodificação.

Capítulo 6

Sistema Implementado

6.1 Base de Dados

6.1.1 Introdução

A base de dados utilizada para treinamento e avaliação do sistema implementado consiste de um total de 100 sentenças pronunciadas por um único locutor masculino, adquiridas sob condições de estúdio, amostradas a 16 kHz e quantizadas com 16 bits por amostra. Estas sentenças, que totalizam 5,21 minutos de fala, foram gravadas no CPQD TELEBRÁS e segmentadas por especialistas do Laboratório de Fonética Acústica e Psicolinguística Experimental - LAFAPE, do Instituto de Estudos da Linguagem - IEL, da UNICAMP.

Esta base foi gravada inicialmente para análise e levantamento de características prosódicas a serem utilizadas em um sistema para síntese concatenativa de fala a partir de texto. As grandes motivações para a utilização desta base para treinamento e avaliação de um sistema para reconhecimento de fala contínua foram :

- O fato dela se encontrar totalmente rotulada em termos de unidades sub-lexicais - 36 tipos de fones. Esta rotulação foi realizada manualmente por especialistas do LAFAPE.
- Apresentar uma diversidade relativamente grande de duração e de palavras, mostrando-se, portanto, bastante útil para a avaliação do sistema. A menor frase possui apenas uma única palavra e a maior 20 palavras. As 100 frases são constituídas de um total de 319 palavras distintas.

6.1.2 Unidades Sub-Lexicais (fones)

Para treinamento do sistema foram escolhidos aleatoriamente 80% (aproximadamente) dos exemplos de cada fone e os fones restantes foram utilizados para validação cruzada (monitoramento da capacidade de generalização da ANN). Não foi definido um conjunto de teste, porque a avaliação do sistema foi realizada a nível de sentenças e não de fones.

A Tabela 6.1 apresenta todos os tipos de fones que compõem a base de dados, seguidos de exemplos que ilustram a pronúncia de cada um deles. Também são mostradas as quantidades de exemplos de cada tipo de fone utilizado assim como suas durações médias. Na última linha da Tabela 6.1 são apresentados os totais de exemplos utilizados assim como a média das durações médias de todos estes exemplos.

6.1.3 Sentenças

A Tabela 6.2 apresenta uma lista com 20 das 100 sentenças que compõem a base de dados utilizada. Uma lista completa com as 100 encontra-se no Apêndice C.

6.1.4 Dicionário de Pronúncias

O Apêndice B apresenta a pronúncia¹ (conversão ortográfico-fonética) de todas as palavras que compõem o léxico utilizado. A pronúncia de cada uma destas palavras foi obtida a partir da conversão ortográfico-fonética de cada uma das 100 sentenças do Apêndice C. Com isto foi levado em consideração informações sobre o contexto em que cada uma destas palavras se encontrava.

6.2 Pré-Processamento

A parametrização do sinal de fala foi realizada utilizando-se apenas coeficientes Mel cepstrais - MCC (“Mel Cepstral Coeficientes”) baseados em banco de filtros segundo Davis & Mermelstein [5]. Para cada vetor acústico \mathbf{x}_n foram calculados 12 coeficientes Mel Cepstrais $\mathbf{x}_n = \{mcc_{n_1}, mcc_{n_2}, \dots, mcc_{n_{12}}\}$. Estes vetores acústicos foram calculados para janelas de 20 ms tomadas a cada 10 ms (superposição de 50 %).

Para evitar problemas de saturação dos potenciais de ativação da ANN (efeito conhecido como congelamento de sinapses) é comum se realizar uma normalização dos coeficientes $mcc_{n_j}, j = \{1, 2, \dots, 12\}$

¹A definição da pronúncia de cada uma das palavras que compõem o léxico utilizado foi realizada por especialistas do Laboratório de Fonética Acústica e Psicolinguística Experimental - LAFAPE do Instituto de Estudos da Linguagem - IEL da UNICAMP.

Nº	Fones	Exemplos	Nº de fones	Duração Média (ms)
01	<i>i</i>	t <u>i</u> me,cai,mã <u>e</u>	411	57,05
02	<i>e</i>	b <u>e</u> sta	243	75,21
03	<i>é</i>	f <u>e</u> sta	56	129,08
04	<i>á</i>	p <u>a</u> ta	332	91,97
05	<i>ó</i>	p <u>o</u> rta	16	159,62
06	<i>o</i>	b <u>o</u> lo	111	85,27
07	<i>u</i>	p <u>ã</u> o, pau, t <u>u</u> do	482	55,60
08	<i>a</i>	p <u>a</u> ta	137	63,06
09	<i>ã</i>	pl <u>a</u> nta	89	86,17
10	<i>ẽ</i>	t <u>e</u> nta	119	111,90
11	<i>ĩ</i>	pin <u>t</u> a	74	109,82
12	<i>õ</i>	ton <u>o</u>	75	101,94
13	<i>ũ</i>	mun <u>u</u> do	25	90,88
14	<i>p</i>	p <u>p</u> a	101	82,96
15	<i>t</i>	p <u>tt</u> a	217	78,67
16	<i>k</i>	pa <u>kk</u> a	178	85,97
17	<i>T</i>	t <u>tt</u> a	73	97,47
18	<i>b</i>	b <u>bb</u> a	51	65,27
19	<i>d</i>	da <u>ddd</u> o	176	50,44
20	<i>g</i>	ga <u>gg</u> a	29	51,96
21	<i>D</i>	di <u>dd</u> a	55	57,27
22	<i>f</i>	fa <u>ff</u> a	50	98,02
23	<i>s</i>	sa <u>sss</u> o	320	112,59
24	<i>ç</i>	cha <u>ccc</u> o	18	104,61
25	<i>v</i>	va <u>vv</u> a	67	62,35
26	<i>z</i>	ca <u>zz</u> a	145	61,82
27	<i>j</i>	ji <u>jj</u> pe	11	73,09
28	<i>m</i>	ma <u>mm</u> a	112	62,59
29	<i>n</i>	na <u>nn</u> a	113	45,04
30	<i>η</i>	man <u>nnh</u> ã	2	91,50
31	<i>ř</i>	ca <u>rr</u> o	200	36,46
32	<i>r</i>	ca <u>rrr</u> a	32	51,93
33	<i>R</i>	ca <u>RR</u> o	25	66,68
34	<i>l</i>	la <u>ll</u> a	56	44,48
35	<i>λ</i>	ca <u>llh</u> a	6	57,83
36	#	(pausa)	37	314,67
		Total	4244	76.54

Tabela 6.1: Lista com os tipos de fones utilizados seguidos por exemplos, e suas respectivas quantidades e durações médias.

-
-
- 01 A cotação do dólar aumentou, e as bolsas fecharam em baixa.
 - 02 A cotação do dólar aumentou, mas as bolsas fecharam em baixa.
 - 03 A bolsa ficará estável, ou sofrerá uma pequena queda.
 - 04 Não haverá ajustes nem modificações radicais no plano.
 - 05 Foi detectado um problema em seu cartão. Ele deve ser substituído.
 - 06 É necessário que o convênio permita o intercâmbio.
 - 07 Posso afimar-lhes que o convênio permite o intercâmbio.
 - 08 O convênio que foi assinado recentemente permite o intercâmbio.
 - 09 O convênio permite o intercâmbio porque visa a integração entre alunos de culturas diferentes.
 - 10 O convênio que foi assinado recentemente, permite o intercâmbio.
 - 11 O convênio assinado na última reunião é mais interessante do que o anterior.
 - 12 A medida que o tempo passa, mais nos convencemos da eficiência do convênio.
 - 13 É suficiente.
 - 14 Isto é suficiente.
 - 15 O saldo é suficiente.
 - 16 O saldo de sua conta é suficiente.
 - 17 O saldo disponível é insuficiente.
 - 18 O saldo disponível em sua conta é insuficiente.
 - 19 Isto parece insuficiente.
 - 20 O saldo parece insuficiente.
-

Tabela 6.2: Lista com 20 das 100 sentenças a serem utilizadas durante o treinamento e avaliação do sistema.

de todos os vetores acústicos \mathbf{x}_n a serem utilizados tanto no treinamento como na avaliação do sistema. Uma normalização bastante comum é fazer $mcc_{n_j} = \frac{mcc_{n_j}}{f_N}$, sendo f_N o fator de normalização capaz de tornar o maior coeficiente, de todos os vetores acústicos utilizados durante o treinamento, igual a 1. Um problema associado a este tipo de normalização é que ele não leva em consideração a distribuição estatística dos coeficientes dos vetores acústicos de treinamento. Por exemplo, se o maior coeficiente de todos os vetores acústicos for muito maior que todos os coeficientes restantes, então este tipo de normalização reduzirá muito a magnitude da maioria dos dados de entrada da ANN, o que dificultará o processo de treinamento. Para contornar este problema montou-se uma seqüência \mathbf{X}_{MCC} , formada pela concatenação de todos os coeficientes mcc_{n_j} , de todos os vetores acústicos utilizados no treinamento e ajustou-se o valor de f_N para que 95% dos dados (coeficientes mcc_{n_j}) passassem a ter magnitudes entre -1.0 e 1.0 . O valor encontrado para f_N foi equivalente a fazer o desvio padrão da seqüência \mathbf{X}_{MCC} igual a $0,49$, ($\sigma = 0,49$) mostrando que os dados de treinamento apresentavam uma função distribuição de probabilidade quase gaussiana. Este fato pôde ser comprovado através da análise do histograma da seqüência \mathbf{X}_{MCC} , ilustrado na Figura 6-1.

A Figura 6-2 mostra o histograma após a normalização dos coeficientes mcc_{n_j}

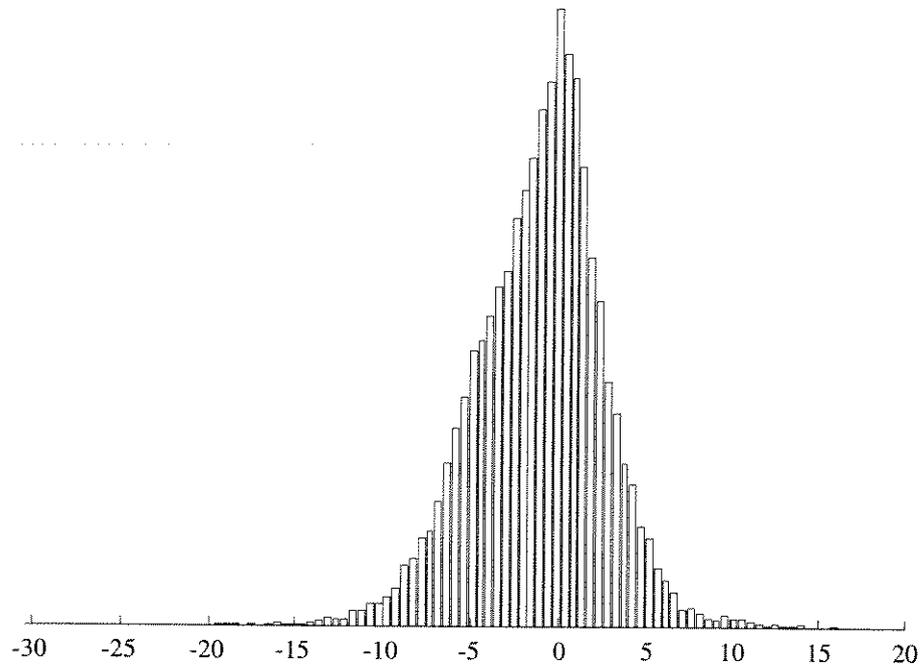


Figura 6-1: Histograma da amplitude dos componentes de todos os vetores acústicos que formam os exemplos de treinamento X_e

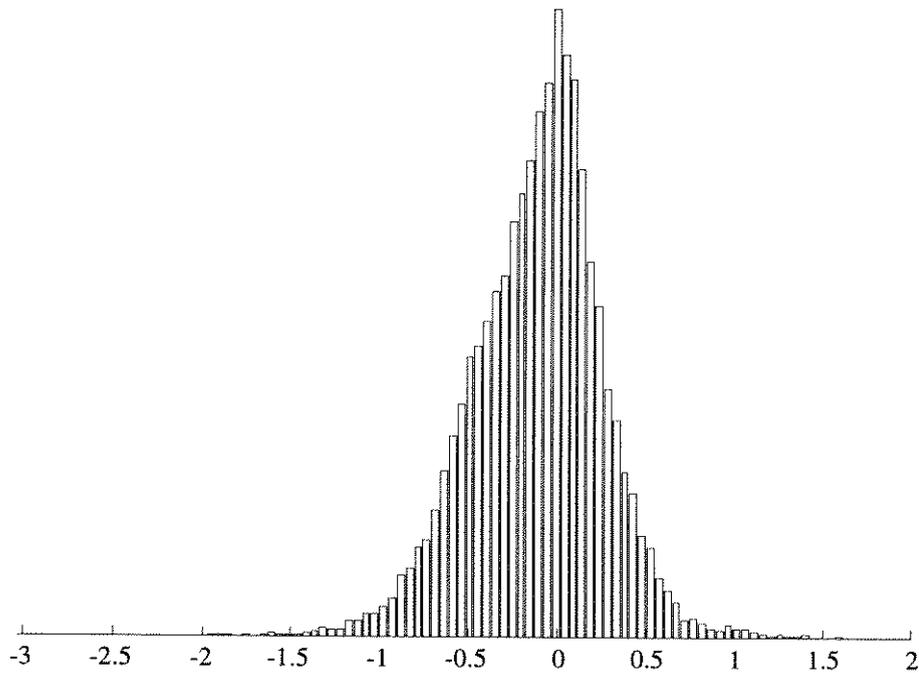


Figura 6-2: Histograma da amplitude dos componentes de todos os vetores acústicos utilizados no treinamento, após a normalização.

6.3 Modelos Híbridos ANN-HMM

6.3.1 Rede Neural Artificial

Foi utilizada uma ANN do tipo perceptron multicamadas - MLP (“Multilayer Perceptron”) com apenas uma camada escondida e com a seguinte configuração :

- **Algoritmo de Treinamento EBP** - “Error Back-Propagation”. Com atualização instantânea dos pesos sinápticos e sem o uso de momento [10].
- **Número de entradas** Igual a 84, correspondente a 7 vetores acústicos. Cada vetor acústico corresponde a 12 coeficientes Mel Cepstrais [5] extraídos a partir de janelas de 20 ms tomadas a cada 10 ms (superposição de 50%). A escolha de 7 vetores acústicos (correspondente a 70 ms de fala) foi motivada pelo fato de a média das durações médias de todos os fones ser igual a 76,57 ms. Também foram testadas MLPs com 108 entradas (9 vetores acústicos). Esta configuração, entretanto, além de ser mais difícil de ser treinada, apresentou piores resultados a nível de reconhecimento de fones.

- **Número de neurônios na camada escondida** Igual a 70, número este escolhido segundo um compromisso entre o desempenho da rede e o tempo de treinamento. Também foram treinadas MLPs com 80, 90 e 100 neurônios na camada escondida, mas para um mesmo número de épocas de treinamento, verificou-se que o aumento na quantidade de neurônios degradava o desempenho da ANN.
- **Número de saídas** Igual a 36, correspondente ao número total de fones necessários para realizar a transcrição ortográfico-fonética de qualquer palavra presente no léxico (cada fone foi modelado por um único estado). Este número de 36 saídas também pode ser interpretado como o conjunto estados $Q = \{q_1, q_2, \dots, q_{36}\}$ necessários para montar o HMM de qualquer palavra presente no léxico. Isto porque existe uma *correspondência um-a-um* entre as classes de saída da MLP e os estados $q_k \in Q$.
- **Funções não lineares** Foi utilizada em todas as camadas, inclusive na camada de saída, a função logística, $\varphi(i) = \frac{1}{(1+\exp(-v_i))}$.
- **Normalização dos dados de entrada** Foram normalizados de modo a resultar um desvio padrão $\sigma = 0,49$, fazendo com que 95% dos dados estivessem entre -1.0 e 1.0 .

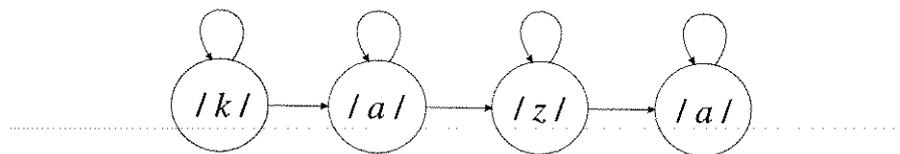


Figura 6-3: HMM “left-right” da palavra “casa” com um único estado por fone.

- **Normalização da saída** Pelo fato de a otimização não ter garantido um mínimo global e como não foi utilizada a função softmax na camada de saída da MLP, então para que a saída da rede $g_k(\mathbf{x}_n, \Psi)$ pudesse ser interpretada como uma medida de probabilidade, foi realizada uma normalização para que $\sum_{k=1}^{36} g_k(x_n, \Psi) = 1$.

6.3.2 Modelos Ocultos de Markov

A definição do HMM de cada uma das palavras do léxico consiste dos seguintes passos

1. **Seqüência de estados que compõem cada HMM** Definir, de acordo como o dicionário de pronúncias Apêndice B, a seqüência de N_{W_i} estados $q_k \in Q$, que devem ser concatenados para a construção do HMM de cada uma das palavras do léxico.
2. **Definir a topologia a ser utilizada** Assumiu-se o modelo *left-right*. Considere, por exemplo, o caso da palavra “casa” cujo correspondente HMM será formado pela concatenação dos estados associados aos fones $/k/ + /a/ + /z/ + /a/$, conforme ilustrado na Figura 6-3.
3. **Determinar as probabilidades de auto-transição de estados a_{jj}** Determinar a probabilidade de transição do estado q_j para ele próprio. Como o modelo é estritamente left-right, então a probabilidade de transição a_{jk} , do estado q_j para um estado q_k , será : $a_{jk} = 1 - a_{jj}$. Durante a estimação (Pré-REMAP) estas probabilidades foram definidas em função das durações médias dos fones, de acordo com a seguinte expressão :

$$a_{jj} = \left[\frac{\left(\frac{D_j}{10} - 1 \right)}{\frac{D_j}{10}} \right] = \left[\frac{D_j - 10}{D_j} \right] \quad (6.1)$$

sendo D_j a duração média em milissegundos do fone associado ao estado q_j e o fator 10 corresponde ao intervalo de tempo em milissegundos em que são tomados as janelas de análise.

Durante a reestimação estas probabilidade de auto-transição (e transição) foram estimadas segundo a Equação (4.34) da Seção 4.3.5.

4. **Definir as verossimilhanças de emissão de símbolos normalizada $\frac{p(\mathbf{x}_n|q_k)}{p(\mathbf{x}_n)}$ associadas a cada estado** Pela regra de Bayes estas verossimilhanças podem ser obtidas a partir das probabilidades a posteriori estimadas pela ANN,

$$\frac{p(\mathbf{x}_n|q_k)}{p(\mathbf{x}_n)} = \frac{P(q_k|\mathbf{x}_n)}{P(q_k)} \quad (6.2)$$

Como já foi discutido na seção 3.3 o cálculo da probabilidade a priori da classe $P(q_k)$ pode ser realizado verificando-se quantas vezes esta classe ocorre no conjunto de treinamento e dividindo pelo número total de exemplos de treinamento (frequência relativa). Porém, devido ao fato do conjunto de exemplos de treinamento, $\mathbf{X}_e = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{N_e}\}$, não ser balanceado (o fone /u/ possui 482 exemplares enquanto os fones /λ/ e /η/ possuem apenas 6 e 2 exemplares, respectivamente), a estimativa das verossimilhanças de emissão de símbolos normalizada $\frac{p(\mathbf{x}_n|q_k)}{p(\mathbf{x}_n)}$ não apresentou bons resultados quando $P(q_k)$ foi calculada segundo o método da frequência relativa. Como foi citado na seção 3.3, Lippmann [30] afirma que durante a avaliação o valor de $P(q_k)$ pode ser ajustado para melhor representar o conjunto de teste, e com este objetivo decidiu-se adotar o seguinte procedimento para o cálculo da probabilidade a priori da classe $P(q_k)$,

$$P(q_k) \cong \sum_{\mathbf{x}_k \in \mathbf{X}_e} P(q_k, \mathbf{x}_k) = \sum_{\mathbf{x}_k \in \mathbf{X}_e} P(q_k|\mathbf{x}_k) \cdot p(\mathbf{x}_k) \quad (6.3)$$

Substituindo (6.3) em (6.2) e assumindo os símbolos \mathbf{x}_k equiprováveis (suposição baseada no fato de que os vetores acústicos \mathbf{x}_k não são quantizados e portanto podem ser considerados distintos entre si), obtém-se :

$$p(\mathbf{x}_n|q_k) = \frac{P(q_k|\mathbf{x}_n)}{\frac{1}{N_e} \cdot \sum_{\mathbf{x}_k \in \mathbf{X}_e} P(q_k|\mathbf{x}_k)} \cdot p(\mathbf{x}_n) \quad (6.4)$$

6.4 Treinamento

6.4.1 Estimação - Pré-REMAP

A rede neural foi treinada utilizando-se os exemplos de treinamento mostrados na Tabela 6.1. As probabilidades de transição $P(q_k^n|q_j^n)$ foram determinadas em função das durações médias dos fones, Equação (6.1) e as probabilidades das classes $P(q_k)$ foram estimadas a partir das saídas da ANN (sem

reestimação) segundo a Equação (6.3).

6.4.2 Reestimação - REMAP

Foram realizadas três reestimações utilizando alvos suaves. A reestimação dos parâmetros $P(q_k^n | q_j^{n-1})$, $P(q_k^n)$ e Ψ foram realizadas de acordo com seguintes passos,

1. Estimou-se para cada uma das M_i sentenças de treinamento, com $i = \{1, 2, \dots, 100\}$, os correspondentes alvos suaves, $P(q_k | \mathbf{X}_{M_i}, M_i, \Theta)$, sendo \mathbf{X}_{M_i} a seqüência de vetores acústicos² correspondente a M_i .
2. Utilizando os 100 alvos suaves do item 1, reestimou-se dos parâmetros Ψ da ANN a partir das correspondentes seqüências de vetores acústicos \mathbf{X}_{M_i} . A apresentação dos pares de treinamento $(\mathbf{X}_{M_i}, P(q_k | \mathbf{X}_{M_i}, M_i, \Theta))$ foi realizada de forma aleatória, segundo uma *função distribuição de probabilidade uniforme*.
3. A partir da ANN com parâmetros Ψ reestimados e da seqüência $\{\mathbf{X}_{M_1}, \mathbf{X}_{M_2}, \dots, \mathbf{X}_{M_{100}}\}$ (concatenação das seqüências de vetores acústicos correspondentes às 100 sentenças de treinamento) calculou-se as probabilidades a priori dos fones $P(q_k)$, com $k \in \{1, 2, \dots, 36\}$, segundo a Equação (6.3).
4. Assumindo que cada uma das 100 sentenças de treinamento são independentes das 99 restantes, realizou-se a reestimação das probabilidades de transição $P(q_k^n | q_j^{n-1})$ para todos os estados q_j , com $j \in \{1, 2, \dots, 36\}$, conforme a Equação (4.34).

6.5 Decodificação

6.5.1 Algoritmo Level Building

Utilizou-se o algoritmo Level Building para realizar a decodificação das sentenças a partir dos modelos de palavras (determinação da seqüência da palavras M , com maior probabilidade dado a seqüência de vetores acústicos \mathbf{X} apresentada ao sistema). Como entre as 100 sentenças a serem avaliadas, existiam sentenças com apenas uma única palavra e outras com até 20 palavras, Apêndice C, decidiu-se ajustar o Level Building para realizar buscas de sentenças compostas pela concatenação de 1 a 22 palavras, (número de níveis máximo, $L = 22$). Este acréscimo de dois níveis ao número máximo de palavras

²A cada sentença M_i , do sistema implementado, foi associado uma única seqüência de vetores acústicos \mathbf{X}_{M_i} .

deve-se ao fato de estar sendo assumido que toda sentença pode ser precedida e seguida por um silêncio (pausa) que por sua vez será modelado como uma palavra adicional.

6.5.2 Modelo de Duração de Palavras

Foi incorporado ao algoritmo Level Building o modelo de duração de palavras discutido na Seção 5.2.3. Para a implementação deste modelo de duração foi realizado um levantamento das durações médias e variâncias de todas as palavras que compõem o léxico utilizado.

6.5.3 Restrições Gramaticais

No sistema implementado foi utilizado uma gramática do tipo Pares-de-palavras, *P-gram* [25]. Uma gramática do tipo *P-gram* faz uso da aproximação mostrada na Equação (6.5) e pode ser interpretada como uma versão determinística de uma gramática estocástica do tipo *Bi-gram*.

$$P(M) \approx G(W_1) \cdot \prod_{n=2}^{N_M} G(W_n|W_{n-1}) \quad (6.5)$$

onde $G(W_1)$ e $G(W_n|W_{n-1})$ são definidos como

$$\begin{cases} G(W_1) = 1, & \text{se for permitido que a palavra } W_1 \text{ inicie uma sentença.} \\ G(W_1) = 0, & \text{se não for permitido que a palavra } W_1 \text{ inicie uma sentença.} \end{cases} \quad (6.6)$$

$$\begin{cases} G(W_n|W_{n-1}) = 1, & \text{se a transição } W_{n-1} \Rightarrow W_n \text{ for permitida.} \\ G(W_n|W_{n-1}) = 0, & \text{se a transição } W_{n-1} \Rightarrow W_n \text{ não for permitida.} \end{cases} \quad (6.7)$$

Com as definições da Equação (6.6) e (6.7) pode ser verificado que, $P(M) = 1$ se for permitido que W_1 inicie a sentença e se todas as transições $W_{n-1} \Rightarrow W_n$ forem permitidas e $P(M) = 0$ se não for permitido que W_1 inicie a sentença ou se alguma das transições $W_{n-1} \Rightarrow W_n$ não forem permitidas.

A determinação da gramática *P-gram*. utilizada no sistema implementado foi realizada levando-se em consideração apenas as 100 sentenças do Apêndice C. Para o estabelecimento destas restrições gramaticais foram necessários os seguintes levantamentos :

- Determinação de todas as possíveis transições de palavras, $W_{n-1} \Rightarrow W_n$.
- Determinação de todas as possíveis palavras que podem iniciar as sentenças.

Além destes levantamentos foram estabelecidas as seguintes regras :

Regra 1 Todas as sentenças podem ser iniciadas por um “silêncio”. Com isto foi definido que o “silêncio” (/ # /) consiste de uma das possíveis “palavras” que podem iniciar as sentenças.

Regra 2 Todas as palavras podem ser precedidas e seguidas por um “silêncio”, isto é, as transições “silêncio” $\Rightarrow W_n$ e $W_n \Rightarrow$ “silêncio”, são permitidas para todas as palavras W_n presentes no léxico.

A incorporação destas restrições gramaticais ao algoritmo Level Building foram realizadas da seguinte forma :

Nível 1 (Primeira palavra da sentença) Realizar este nível apenas para as palavras que puderem iniciar a sentença de acordo com a gramática P -gram.

Nível 2 (E demais níveis) Para cada uma das palavras W_j 's a serem analisadas, inicializar $\delta_n(1)$ e $\alpha_n(1)$ conforme as Equações (6.8) e (??).

$$\delta_n(1) = \max \left[\widehat{\mathbf{T}}(l-1, n-1) \cdot \widehat{\mathbf{P}}(l-1, n-1) \cdot G(W_j|W_i), P(q_1^n|q_1^{n-1}, W_j, \Theta) \cdot \delta_{n-1}(1) \right] \quad (6.8)$$

$$\cdot p(\mathbf{x}_n|q_k, W_j, \Theta)$$

$$\alpha_n(1) = \begin{cases} n-1, & \text{se } \widehat{\mathbf{T}}(l-1, n-1) \cdot \widehat{\mathbf{P}}(l-1, n-1) \cdot G(W_j|W_i) > \\ P(q_1^n|q_1^{n-1}, W_i, \Theta) \cdot \delta_{n-1}(1) & \\ \alpha_{n-1}(1), & \text{caso contrário} \end{cases} \quad (6.9)$$

Capítulo 7

Resultados Experimentais

7.1 Introdução

Neste capítulo são apresentados alguns resultados da avaliação do sistema implementado com o uso das unidades sub-lexicais (fones) definidas na Tabela 6.1 e das 100 sentenças apresentadas no Apêndice C. A partir destes resultados são realizadas quatro análises básicas :

1. **Reconhecimento a nível de unidades sub-lexicais** Desempenho da ANN como *classificadora de unidades sub-lexicais* (fones).
2. **Reconhecimento de sentenças com o uso de modelos de sentenças** Tem como objetivo avaliar o *modelamento acústico* realizado pela ANN e a absorção da *variabilidades temporal* realizada pelos HMMs.
3. **Algoritmo de reestimação de parâmetros - REMAP** Análise das reestimações dos parâmetros Ψ da ANN, das probabilidades de transição de estados $P(q_k^n | q_j^{n-1})$ e das probabilidades a priori das classes (fones) $P(q_k)$.
4. **Reconhecimento de sentenças com o uso de modelos de palavras** Neste caso os objetivos principais são : avaliar o desempenho do *algoritmo de busca* - algoritmo “Level Building” - e comprovar a importância de um bom modelo da língua (restrições gramaticais) no desempenho final de um sistema para reconhecimento de fala contínua.

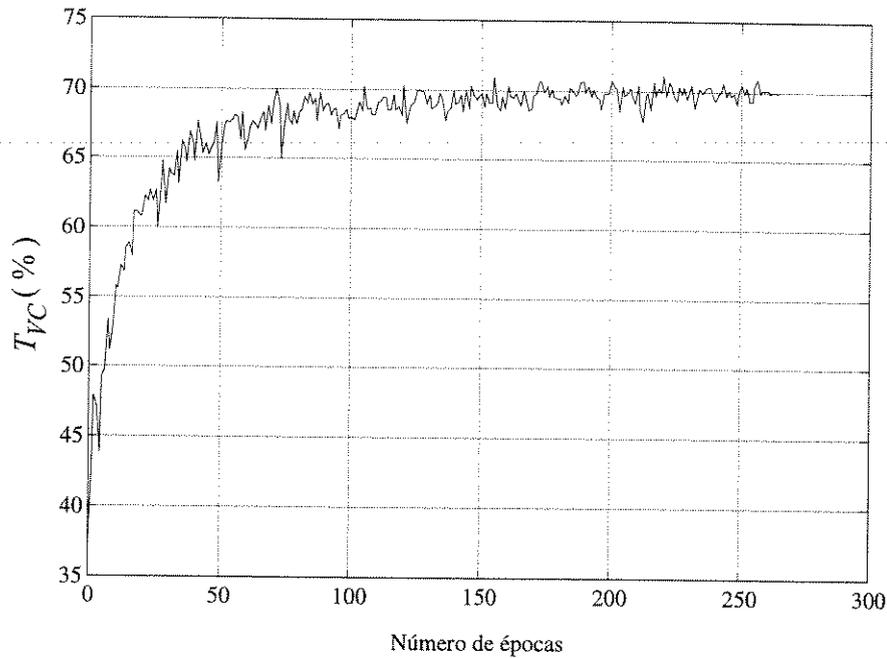


Figura 7-1: Taxas de acertos da ANN para os exemplos de validação cruzada em função do número de épocas de treinamento.

7.2 Reconhecimento a Nível de Unidades Sub-Lexicais (fones)

7.2.1 Taxas de Acertos da ANN ao Longo das Épocas de Treinamento

A Figura 7-1 apresenta as taxas totais de acertos de fones (considerando todos os 36 fones), para os exemplos de validação cruzada T_{VC} , ao longo do número de épocas de treinamento.

A Figura 7-2 apresenta as taxas totais de acertos de fones, tanto para os exemplos de validação cruzada T_{VC} , como para os exemplos de treinamento $T_{trein.}$, ao longo das épocas de 171 a 249.

A taxa de acertos global de fones da ANN T_G , foi definida como a média ponderada dos resultados para o conjunto de validação cruzada T_{VC} e para o conjunto de treinamento $T_{trein.}$, $T_G = 0,2 \cdot T_{VC} + 0,8 \cdot T_{trein.}$

A Figura 7-3 apresenta a taxa de acertos global da ANN T_G , para as épocas de 171 a 249 e mostra que o conjunto de pesos sinápticos Ψ que apresentou a melhor T_G , está associado à época 227.

7.2.2 Taxas de Acertos Finais da ANN

A Tabela 7.1 apresenta as taxas de acertos T_{VC} , $T_{trein.}$ e T_G da ANN, com pesos sinápticos Ψ associados à época número 227, para cada um dos 36 tipos de fones utilizados.

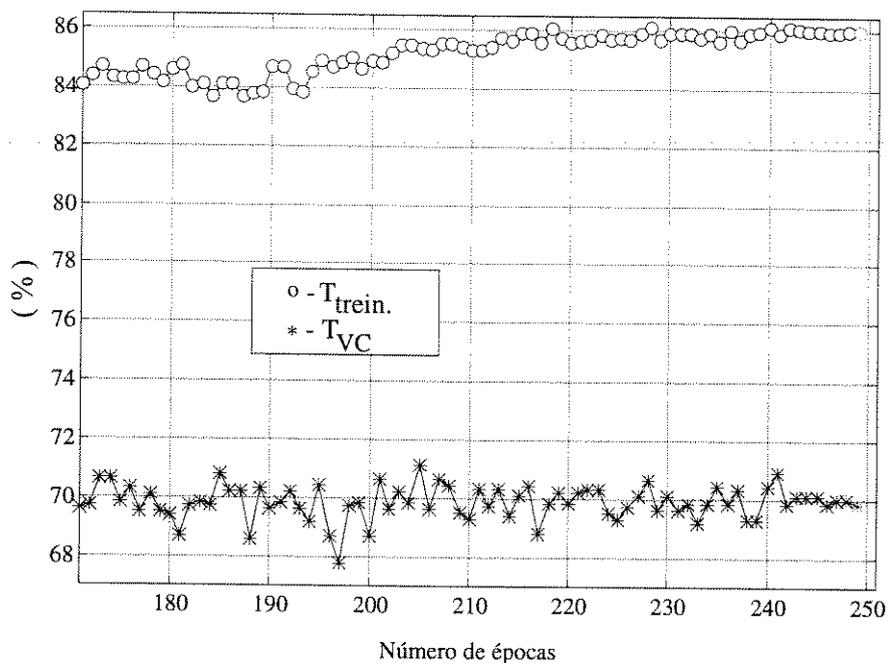


Figura 7-2: Taxa de acertos da ANN tanto para os exemplos de treinamento $T_{\text{trein.}}$, como para os exemplos de validação T_{VC} , em função do número de épocas de treinamento - da época 171 a 249.

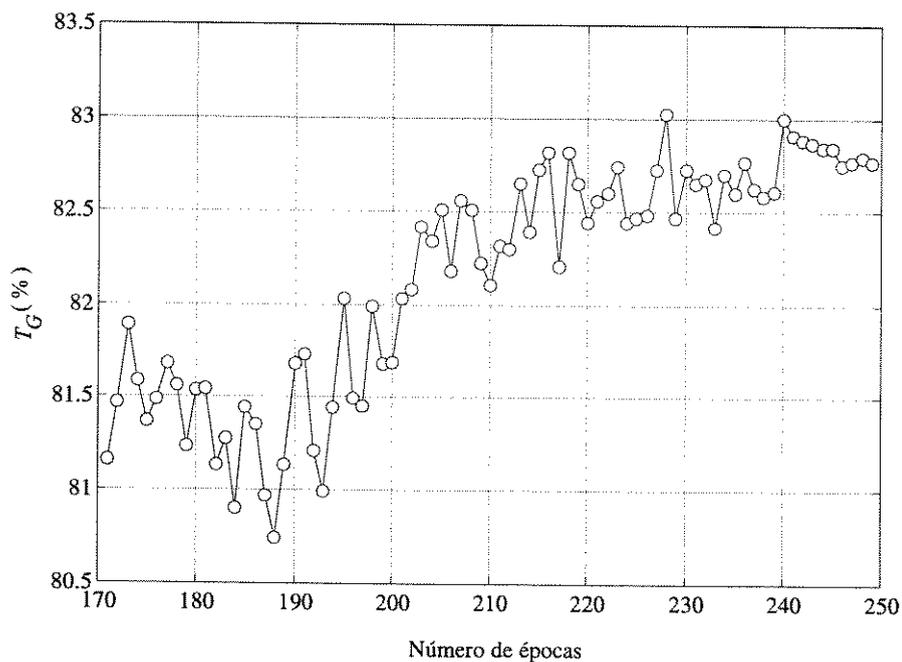


Figura 7-3: Taxa de acertos global, que é dada por uma média ponderada entre as taxas de acerto para os exemplos de validação cruzada e para os exemplos de treinamento, $T_G = 0,2 \cdot T_{VC} + 0,8 \cdot T_{\text{trein.}}$ - da época 171 a 249.

Nº	Fones	Exemplos	Nº de fones	T_{VC} (%)	$T_{trein.}$ (%)	T_G (%)
01	<i>i</i>	t <u>i</u> me,cai,mã <u>e</u>	411	74,70	97,14	92,65
02	<i>e</i>	b <u>e</u> sta	243	81,63	87,11	86,01
03	<i>é</i>	f <u>e</u> sta	56	72,73	95,56	90,99
04	<i>á</i>	p <u>a</u> ta	332	92,54	93,58	93,37
05	<i>ó</i>	p <u>o</u> rta	16	50,00	83,33	76,66
06	<i>o</i>	b <u>o</u> lo	111	65,22	86,36	82,13
07	<i>u</i>	p <u>ã</u> o, p <u>au</u> , t <u>ud</u> o	482	78,35	92,05	89,31
08	<i>a</i>	p <u>a</u> ta	137	46,43	77,98	71,67
09	<i>ã</i>	pl <u>a</u> nta	89	66,67	84,51	80,94
10	<i>ẽ</i>	t <u>e</u> nta	119	70,83	89,47	85,74
11	<i>ĩ</i>	pin <u>t</u> a	74	60,00	93,22	86,57
12	<i>õ</i>	ton <u>o</u> ta	75	46,67	91,67	82,67
13	<i>ũ</i>	mun <u>u</u> do	25	0,00	5,00	4,0
14	<i>p</i>	p <u>a</u> ta	101	75,00	93,82	90,05
15	<i>t</i>	p <u>a</u> ta	217	68,18	95,37	89,93
16	<i>k</i>	pa <u>ç</u> a	178	66,67	94,37	88,83
17	<i>T</i>	t <u>i</u> a	73	66,67	87,93	88,67
18	<i>b</i>	b <u>a</u> ta	51	60,00	85,37	80,29
19	<i>d</i>	d <u>a</u> do	176	55,55	96,43	88,25
20	<i>g</i>	g <u>a</u> ta	29	50,00	47,83	48,26
21	<i>D</i>	d <u>i</u> a	55	45,45	93,18	83,63
22	<i>f</i>	f <u>a</u> ca	50	100,00	92,50	94,00
23	<i>s</i>	s <u>a</u> po	320	93,75	98,44	97,50
24	<i>f</i>	ch <u>a</u> to	18	0,00	42,86	34,29
25	<i>v</i>	y <u>a</u> la	67	64,29	73,58	71,72
26	<i>z</i>	ca <u>s</u> a	145	82,76	92,24	90,34
27	<i>j</i>	ji <u>p</u> e	11	0,00	0,00	0,00
28	<i>m</i>	ma <u>ç</u> a	112	65,22	92,13	86,75
29	<i>n</i>	na <u>d</u> a	113	78,26	84,44	83,20
30	<i>η</i>	man <u>h</u> ã	2	0,00	0,00	0,00
31	<i>ř</i>	ca <u>r</u> o	200	52,00	86,25	79,40
32	<i>r</i>	ca <u>r</u> ta	32	16,67	57,69	49,48
33	<i>R</i>	ca <u>r</u> ro	25	20,00	65,00	56,00
34	<i>l</i>	la <u>t</u> a	56	41,46	90,91	81,02
35	<i>λ</i>	ca <u>l</u> ha	6	0,00	0,00	0,00
36	<i>#</i>	(pausa)	37	87,50	93,10	91,98
		Total	4244	70.67	86.12	83.02

Tabela 7.1: Taxa de acertos da ANN após o Pré-REMAP, para cada um dos 36 tipos de fones.

7.3 Reconhecimento de Sentenças com uso de Modelos de Sentenças

7.3.1 Exemplo

Antes de realizar o reconhecimento de fala contínua propriamente dito, realizou-se um experimento para analisar o desempenho do sistema no reconhecimento de sentenças a partir de modelos de sentenças. Neste experimento montou-se um modelo híbrido ANN-HMM para cada uma das 100 sentenças a serem reconhecidas (uma única ANN e 100 HMMs). O processo de decodificação consistiu basicamente em verificar qual modelo de sentença M_i com $i = \{1, 2, \dots, 100\}$ apresentava a maior verossimilhança de ter gerado a seqüência de vetores acústicos $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ apresentada ao sistema. Para ilustrar o procedimento considere o seguinte exemplo :

- Seqüência de vetores acústicos \mathbf{X} , apresentada ao sistema, correspondente à sentença : “O saldo é suficiente” .
- Modelos ocultos de Markov a serem analisados : (1) “É suficiente”, (2) “Isto é suficiente”, (3) “O saldo é suficiente” e (4) “O saldo de sua conta é suficiente”.

A Figura 7-4 mostra a saída da ANN convertida em verossimilhanças de emissão de símbolos normalizada $\frac{p(\mathbf{x}_n|q_k)}{p(\mathbf{x}_n)}$, com $k = \{1, 2, \dots, 36\}$, para a seqüência de vetores acústicos \mathbf{X} (correspondente a sentença, “O saldo é suficiente”). Nesta figura os tons de cinza indicam o valor de $\frac{p(\mathbf{x}_n|q_k)}{p(\mathbf{x}_n)}$; quanto mais intenso (mais escuro) o tom de cinza, maior o valor da verossimilhança.

Nas Figuras 7.5a, 7.5b, 7.5c e 7.5d são mostrados os modelos e os alinhamentos realizados pelo algoritmo de Viterbi (seqüência de estados com maior verossimilhança acumulada $p(\mathbf{X}|M, \Theta)$, mais precisamente $\frac{p(\mathbf{X}|M, \Theta)}{p(\mathbf{X}|\Theta)}$), relativos, respectivamente, às sentenças “É suficiente” (12 fonos), “Isto é suficiente” (16 fonos), “O saldo é suficiente” (18 fonos) e “O saldo de sua conta é suficiente” (27 fonos), quando a seqüência de vetores acústicos apresentada ao sistema corresponde à sentença, “O saldo é suficiente”.

As Figuras 7.5a, 7.5b, 7.5c e 7.5d, também mostram através das trajetórias de Viterbi, que a inclusão de estados correspondentes a pausas (silêncios) no início e no final dos HMMs tornou o sistema capaz de absorver os silêncios iniciais e finais das sentenças.

A Figura 7-6 apresenta os log's normalizados das verossimilhanças calculadas ao longo dos caminhos de Viterbi apresentados na Figura 7-5, e confirma que o modelo correspondente à Figura 7.5c é o que tem a maior verossimilhança de gerar a seqüência de símbolos correspondentes à sentença : “O saldo é suficiente”.

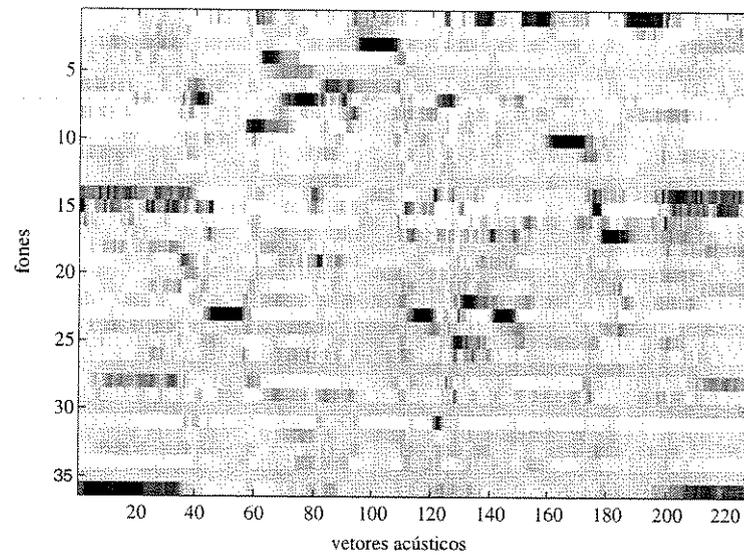


Figura 7-4: Saída da ANN convertida em verossimilhança de emissão de símbolos, para os vetores acústicos de entrada correspondentes à sentença, “O saldo é suficiente”.

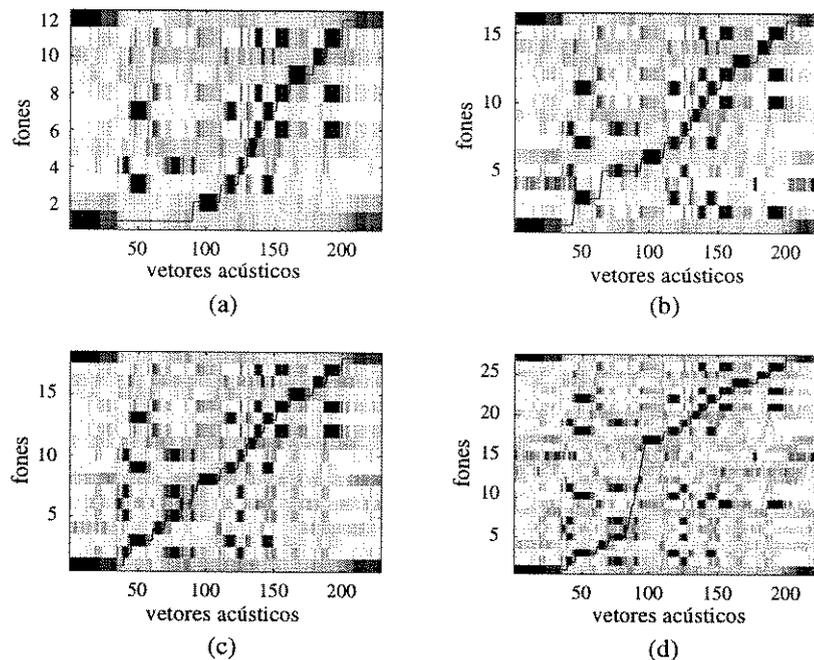


Figura 7-5: Modelos ocultos de Markov associados às sentenças : (a) “É suficiente”, (b) “Isto é suficiente”, (c) “O saldo é suficiente” e (d) “O saldo de sua conta é suficiente”.

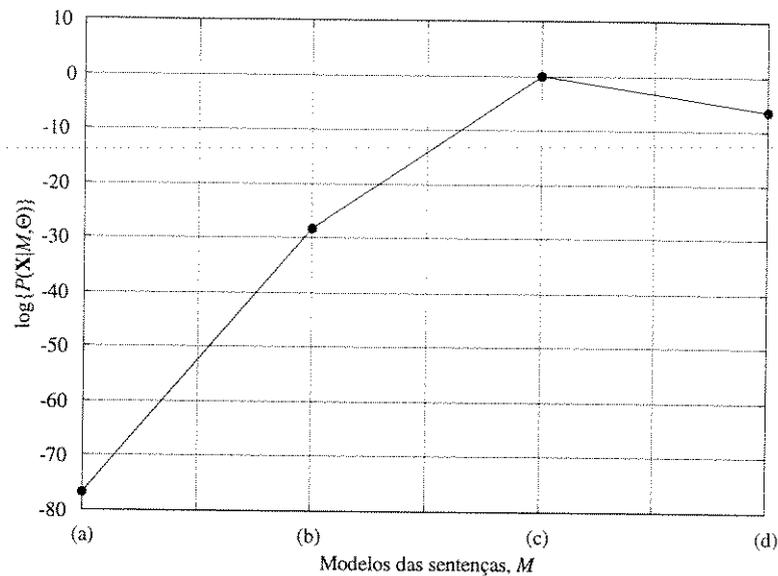


Figura 7-6: Log's das verossimilhanças normalizadas para os caminhos de Viterbi apresentados na Figura 7-5.

7.3.2 Discussão

Para avaliação de desempenho do sistema para reconhecimento de sentenças a partir de modelos de sentenças, utilizou-se as 100 frases presentes no Apêndice C e a taxa de acertos do sistema foi de 100%. Estes resultados comprovam a capacidade das ANNs de realizarem um bom modelamento acústico, uma vez que existem sentenças que diferem entre si em apenas uma única palavra, e também a boa capacidade de absorção das variabilidades temporais por parte dos HMMs, pois existem sentenças com apenas uma única palavra e outras com até vinte palavras.

Alguns autores [33] aplicam o algoritmo de Viterbi diretamente à saída da ANN, $P(q_k|\mathbf{x}_n, \Theta)$. Esta abordagem apresenta o inconveniente de não permitir a interpretação do segundo estágio do modelo híbrido como um HMM, pois este baseia-se em verossimilhanças de emissão de símbolos $p(\mathbf{x}_n|q_k, \Theta)$. Além deste inconveniente a Figura 7-7 mostra que o uso de $p(\mathbf{x}_n|q_k, \Theta)$ apresenta uma maior capacidade de discriminação do modelo que gerou a seqüência de vetores acústicos apresentada ao sistema. A Figura 7-7 mostra os resultados obtidos pelos HMMs correspondentes às 20 sentenças presentes na Tabela 6.2 quando foi apresentado ao sistema a seqüência de vetores acústicos \mathbf{X} correspondente a sétima sentença desta mesma Tabela.

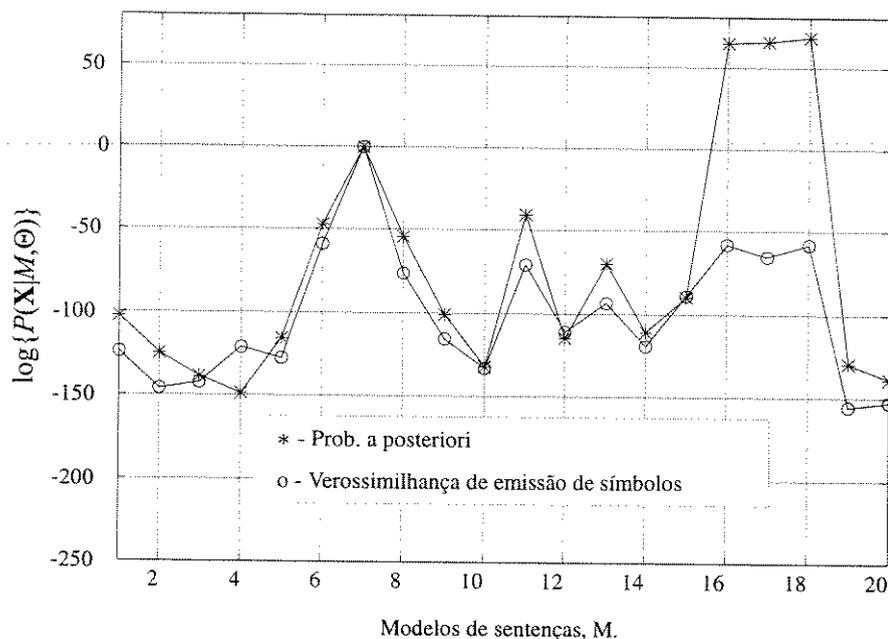


Figura 7-7: Resultados da apresentação da seqüência de vetores acústicos associados à sétima sentença da Tabela 6.2, aos modelos das vinte sentenças desta mesma tabela.

7.4 Desempenho do Algoritmo de Reestimação - REMAP

7.4.1 Reestimação dos parâmetros Ψ da ANN

A Figura 7-8 apresenta os alvos suaves $P(q_k|\mathbf{X}, M, \Theta)$ para a sentença “É suficiente”. A Figura 7-9 apresenta os valores máximos de $P(q_k|\mathbf{X}, M, \Theta)$ para cada um dos vetores acústicos \mathbf{x}_n presentes em \mathbf{X} . Estas duas figuras ilustram o grau de confusão acústica, nas fronteiras entre os fones, estabelecido pelos alvos suaves a serem utilizados na reestimação dos parâmetros Ψ da ANN.

A Figura 7-10 apresenta a saída da ANN após a estimação dos parâmetros Pré-REMAP e a Figura 7-11 apresenta a saída da ANN após o REMAP-2 (duas reestimações). Uma análise destas duas figuras permite as seguintes conclusões :

1. O processo de reestimação aumentou os valores das probabilidades a posteriori $P(q_k|\mathbf{x}_n)$ na saída da ANN.
2. O processo de reestimação reduziu as discontinuidades existentes entre as transições de um estado (fone) para outro.

O acréscimo dos valores das probabilidades a posteriori $P(q_k|\mathbf{x}_n)$ pode ser atribuído ao aumento

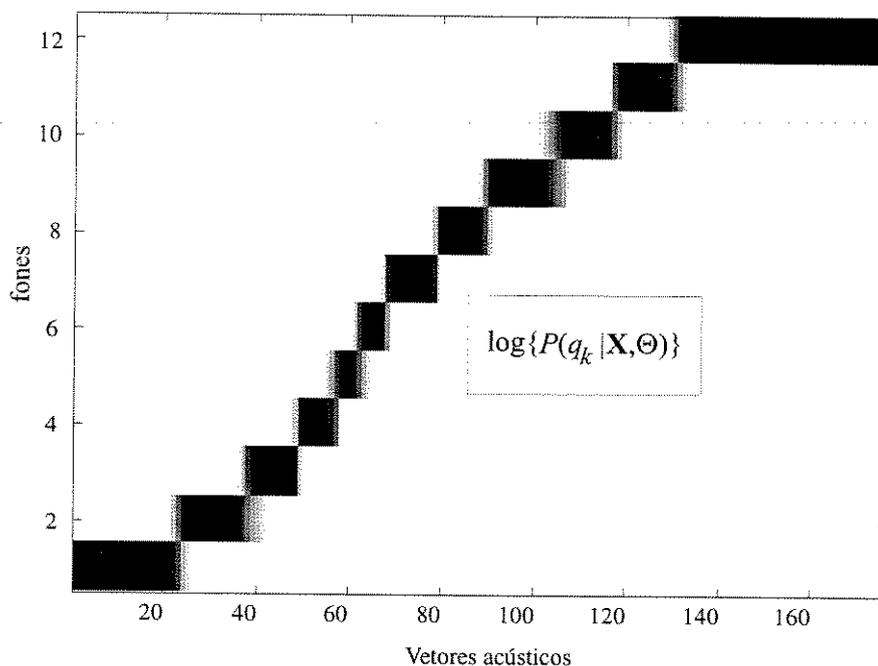


Figura 7-8: Alvos suaves $P(q_k|X, M, \Theta)$ para a sentença “É suficiente”.

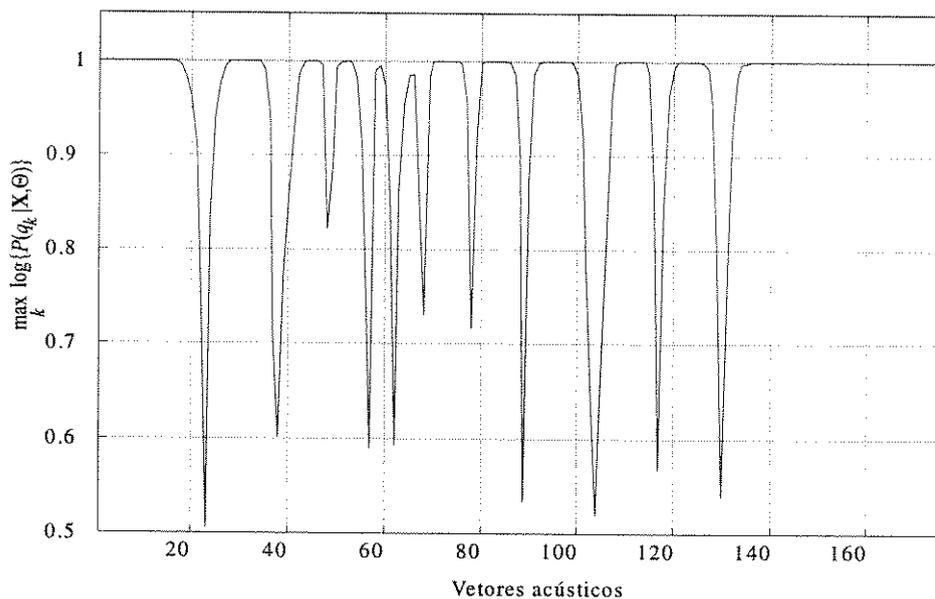


Figura 7-9: Valores máximos dos alvos $P(q_k|X, M, \Theta)$ para todos os estados k , para cada um dos vetores acústicos $x_n \in X$.

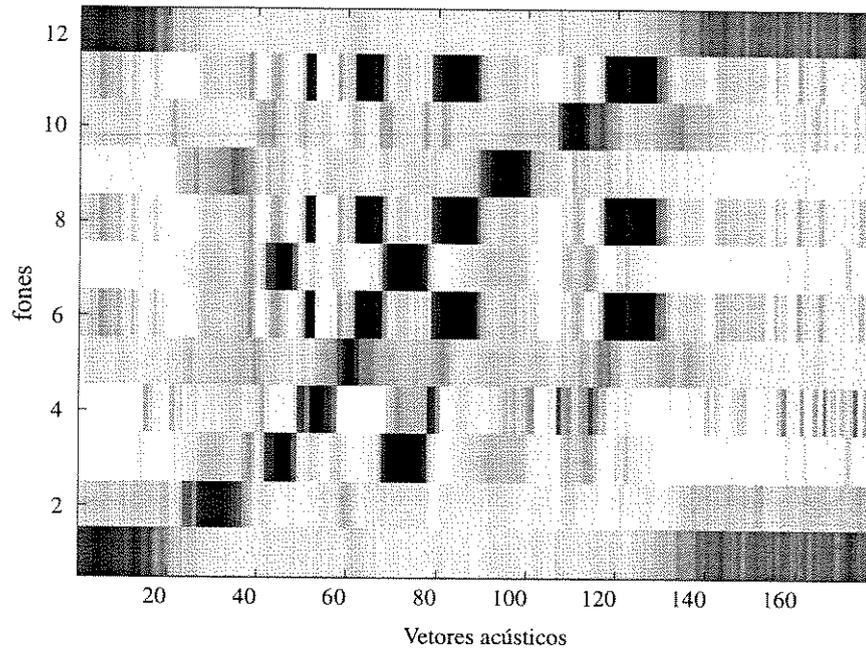


Figura 7-10: Saídas da ANN correspondentes aos fones presentes na sentença “É suficiente”, antes da reestimação, Pré-REMAP.

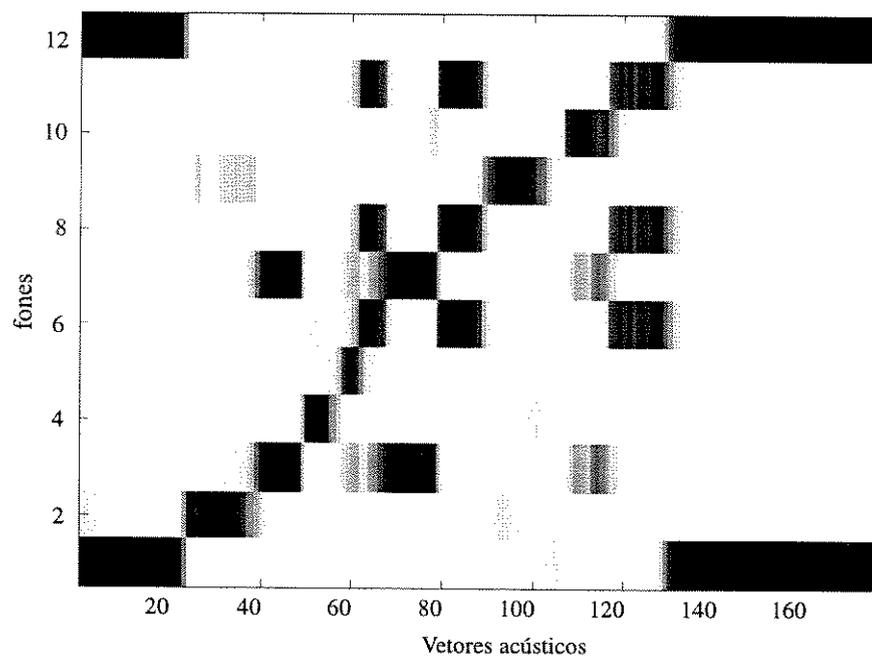


Figura 7-11: Saídas da ANN correspondentes aos fones presentes na sentença “É suficiente”, após duas reestimações, REMAP-2.

na quantidade de exemplos de treinamento utilizados durante a fase de reestimação em relação aos exemplos de treinamento utilizados durante a fase de estimação¹. A redução das descontinuidades deve-se ao fato que durante a fase de reestimação, ao contrário da fase de estimação, a ANN também é treinada com exemplos extraídos das regiões de fronteira entre as unidades sub-lexicais.

7.4.2 Reestimação das Probabilidades de Transição

A Figura 7-12 apresenta as probabilidades de autotransição de estados para três casos :

1. $P^{(0)}(q_k^n | q_k^{n-1})$ - Estimadas em função das durações médias das unidades sub-lexicais, conforme Equação (6.1).
2. $P^{(1)}(q_k^n | q_k^{n-1})$ - Estimadas a partir dos alvos suaves utilizados na 1^a reestimação - REMAP-1.
3. $P^{(2)}(q_k^n | q_k^{n-1})$ - Estimadas a partir dos alvos suaves utilizados na 2^a reestimação - REMAP-2.

A partir da Figura 7-12 pode ser verificado que após o REMAP-2 as probabilidades de autotransição entre estados, $P^{(2)}(q_k^n | q_k^{n-1})$ tornaram-se muito próximas das probabilidades de autotransição estimadas em função das durações médias dos fones $P^{(0)}(q_k^n | q_k^{n-1})$, com exceção dos estados, 30 e 35 (fones /η/ e /λ/). Como as estimativas destas probabilidades de autotransição são baseadas em levantamentos estatísticos, uma possível explicação para as diferenças nos casos dos fones /η/ e /λ/, pode ser associada ao fato destes fones ocorrerem em número extremamente reduzido.

A Figura 7-13 apresenta as probabilidades de transição de cada um dos 36 estados para os mesmos casos apresentados na Figura 7-12

7.4.3 Reestimação das Probabilidades a Priori das Classes

A Figura 7-14 apresenta as probabilidades a priori das classes (estados ou fones) para três casos :

1. $P(q_k)$ - Estimadas a partir das frequências relativas (frequência de ocorrência da classe no conjunto de treinamento \mathbf{X}_e).
2. $P^{(0)}(q_k)$ - Estimadas a partir das saídas da ANN sem reestimação - Pré-REMAP.
3. $P^{(2)}(q_k)$ - Estimadas a partir das saídas da ANN com reestimação - REMAP-2.

¹Durante a etapa de reestimação os exemplos de treinamento são extraídos a cada quadro de análise e durante a etapa de estimação os exemplos são extraídos apenas centrados entre as marcas de segmentação dos fones.

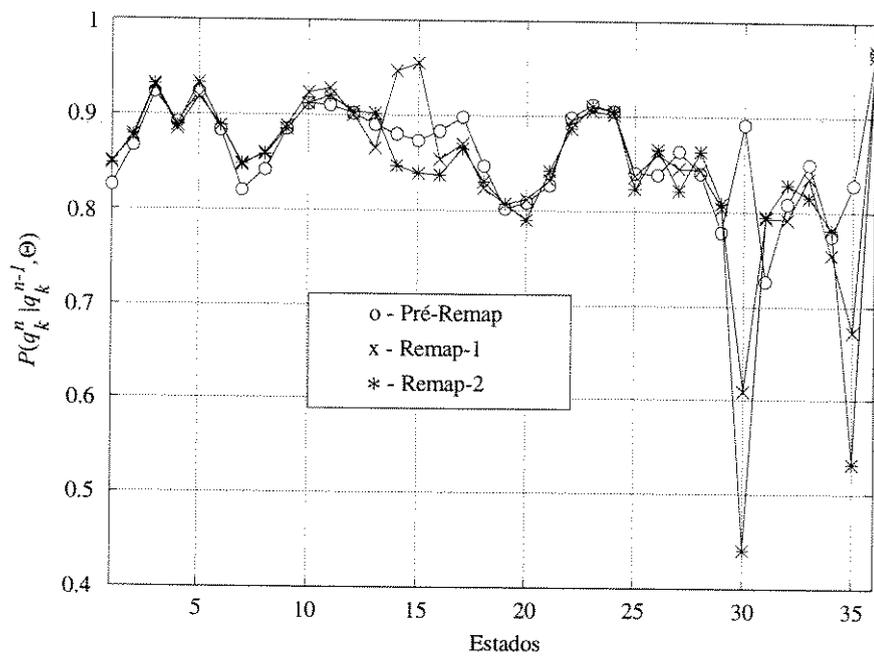


Figura 7-12: Probabilidades de permanência nos estados $P(q_k^n | q_k^{n-1}, \Theta)$

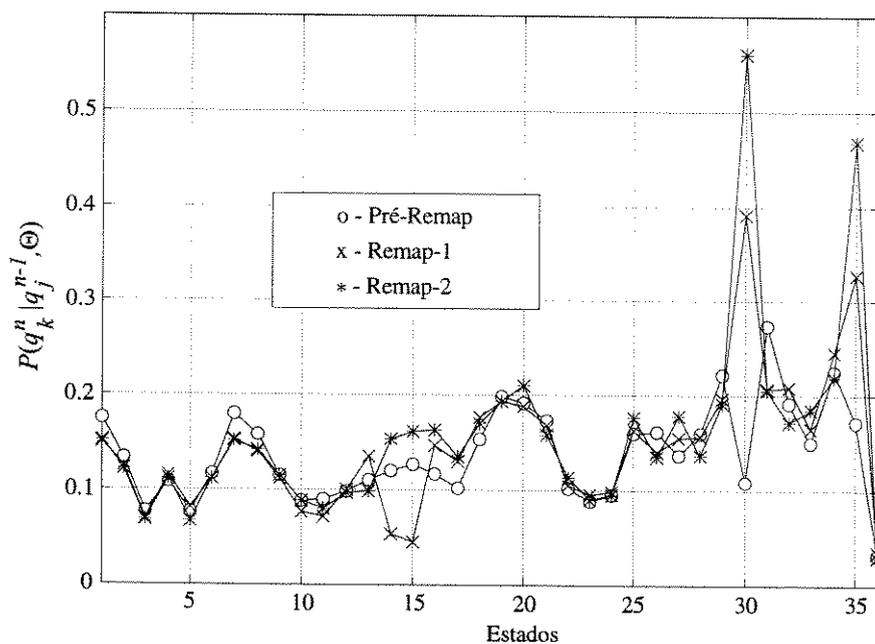


Figura 7-13: Probabilidade de transição de estados $P(q_k^n | q_j^{n-1}, \Theta)$.

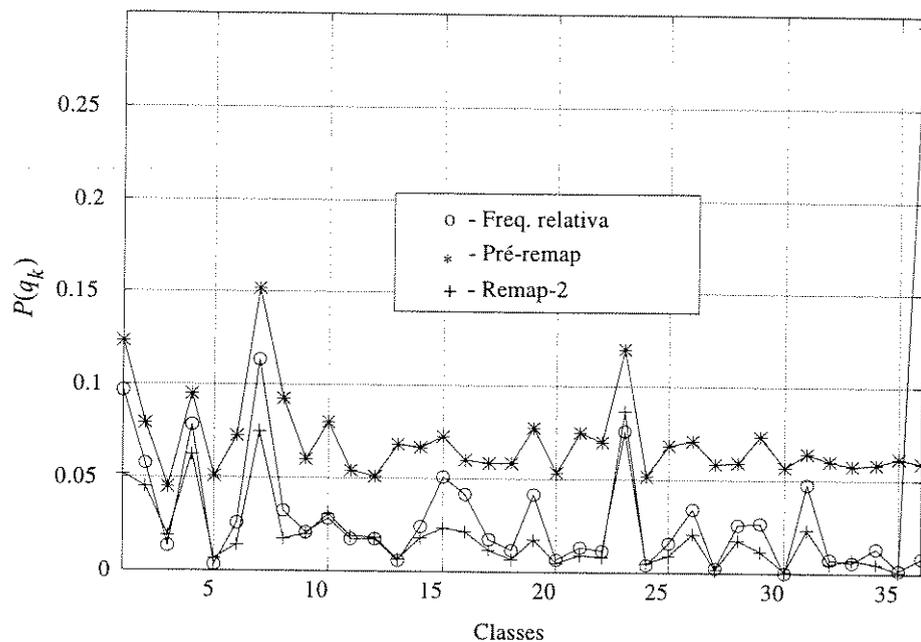


Figura 7-14: Probabilidades das 36 classes (fones) calculadas segundo a frequência relativa e segundo a Equação 6.4 para os casos : Pré-REMAP e REMAP-2.

A partir da Figura 7-14 pode ser verificado que $P^{(2)}(q_k)$ aproxima $P(q_k)$ com uma precisão razoável, com exceção do estado 36 (fone /#/). Esta diferença está associada ao fato da estimativa através da frequência relativa não levar em consideração as durações médias dos fones². Através da Figura 7-14, também pode ser verificado que, a menos de um deslocamento, $P^{(2)}(q_k)$ aproxima com razoável precisão $P^{(1)}(q_k)$. Este deslocamento de $P^{(1)}(q_k)$ deve-se ao fato que durante o Pré-REMAP a ANN ainda apresenta uma elevada superposição entre as classes de saídas (fones).

7.4.4 Discussão

Para verificar a influência do processo de reestimação dos parâmetros dos modelos híbridos ANN-HMM no desempenho final do sistema, considere o exemplo da Figura 7-15 em que os vetores acústicos $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$, correspondentes à sétima sentença da Tabela 6.2, são apresentados aos modelos das 20 sentenças desta mesma Tabela. Neste exemplo são apresentados os resultados (log da verossimilhança normalizada $\log(p(\mathbf{X}|M, \Theta))$, para três casos : (1) Pré-REMAP (sem reestimação), (2) REMAP-1 (primeira reestimação) e (3) REMAP-2 (segunda reestimação). A partir dos resultados

²Em geral as sentenças utilizadas são iniciadas e finalizadas por trechos de silêncio com durações médias em torno de 300 ms.

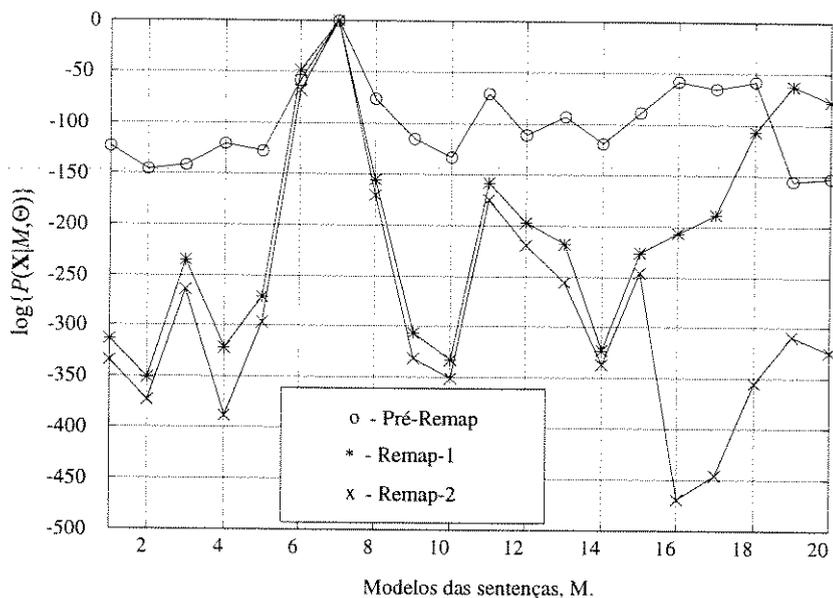


Figura 7-15: Log's das verossimilhanças normalizadas correspondentes a apresentação da seqüência de vetores acústicos associada à sétima sentença da Tabela 6.2, aos 20 modelos de sentenças desta mesma tabela.

da Figura 7-15 pode ser verificado que, em geral, a cada iteração os modelos das sentenças tornam-se mais discriminativos.

7.5 Reconhecimento de Sentenças com uso de Modelos de Palavras

Para avaliar o desempenho do sistema no reconhecimento de sentenças a partir de modelos de palavras foram realizados os seguintes experimentos :

1. Sem o uso de modelo da língua $P(M|\Theta) = 1$, e com o uso do modelo de duração apresentado na seção 5.2.3, foram avaliados os casos : Pré-REMAP, REMAP-1, REMAP-2 e REMAP-3.
2. Sem o uso de modelo da língua $P(M|\Theta) = 1$, e também sem o uso do modelo de duração avaliou-se o REMAP-3.
3. Como o uso de restrições gramaticais do tipo $P - Gram$, e com o uso do modelo de duração avaliou-se o REMAP-3.

Os resultados destes experimentos encontram-se na Tabela 7.2 e uma lista completa com as 100 sentenças reconhecidas para os experimentos 1 e 3 encontra-se no Apêndice C. Nesta tabela tem-se que :

	M. de duração	M. de linguagem	S (%)	D (%)	I (%)	E (%)
Pré-REMAP	Sim	$P(M \Theta) = 1$	11,07	8,50	13,46	33,03
REMAP-1	Sim	$P(M \Theta) = 1$	5,67	5,14	10,27	21,08
REMAP-2	Sim	$P(M \Theta) = 1$	4,69	4,87	8,24	17,80
REMAP-3	Sim	$P(M \Theta) = 1$	4,51	4,96	7,26	16,74
REMAP-3	Não	$P(M \Theta) = 1$	27,54	2,13	12,75	42,42
REMAP-3	Sim	Gramática $P - Gram.$	0	0,35	0,18	0,53

Tabela 7.2: Resultados da avaliação do sistema híbrido ANN-HMM implementado, no reconhecimento das 100 sentenças do Apêndice C.

	S (%)	D (%)	I (%)	E (%)
HMM-1-MCC	24,0	3,10	27,10	54,21
HMM-1-MCC-D	16,82	3,89	20,72	41,45
HMM-3-MCC	10,89	3,80	14,70	29,40
HMM-3-MCC-D	6,64	2,48	9,12	18,2

Tabela 7.3: Resultados da avaliação de um sistema em desenvolvimento no LPDF-UNICAMP, baseado apenas em HMM discreto, no reconhecimento das 100 sentenças apresentadas no Apêndice C.

- $S = 100 \cdot \frac{N_S}{I_W}$ - Taxa de erro de substituição (N_S - número total de erros de substituição e I_W - número total de palavras presentes no léxico).
- $D = 100 \cdot \frac{N_D}{I_W}$ - Taxa de erro de deleção (N_D - número total de erros de deleção).
- $I = 100 \cdot \frac{N_I}{I_W}$ - Taxa de erro de inserção (N_I - número total de erros de inserção).
- $E = 100 \cdot \frac{N_S + N_D + N_I}{I_W}$ - Taxa total de erro de palavras.

Neste momento, está sendo desenvolvido no LPDF - UNICAMP, um sistema para reconhecimento de fala contínua empregando apenas modelos ocultos de Markov discretos com codebook de 256 posições. Este sistema foi avaliado com as 100 sentenças do Apêndice C e os resultados do sistema encontram-se na Tabela 7.3. Nesta tabela tem-se que :

- **HMM-1-MCC** - HMM com um estado por fone e utilizando somente coeficientes Mel Cepstrais (12 coeficientes).

- **HMM-1-MCC-D** - HMM com um estado por fone e utilizando coeficientes Mel Cepstrais e coeficientes Delta Mel Cepstrais.
- **HMM-3-MCC** - HMM com três estados por fone e utilizando somente coeficientes Mel Cepstrais.
- **HMM-3-MCC-D** - HMM com três estados por fone e utilizando coeficientes Mel Cepstrais e coeficientes Delta Mel Cepstrais.

Durante a obtenção dos resultados da Tabela 7.3 fez-se uso do mesmo modelo de duração utilizado no sistema híbrido, e não foi utilizado qualquer tipo de restrição gramatical $P(M|\Theta) = 1$.

Capítulo 8

Considerações Finais

8.1 Análise de Desempenho

8.1.1 Reconhecimento de Fones

- **Bases de Dados Desbalanceada** O desempenho da ANN como classificadora de unidades sub-lexicais poderia ter sido melhor se a base de dados utilizada não fosse tão desbalanceada. Os fones /u/ e /i/ possuem 482 e 411 exemplares, respectivamente, enquanto os fones /η/ e /f/ possuem, respectivamente, apenas 6 e 2 exemplares.
- **Parâmetros de Entrada** No sistema implementado não houve uma preocupação em otimizar o módulo de Pré-processamento. Por este motivo foram utilizados apenas 12 coeficientes Mel Cepstrais [5]. Porém acredita-se que o desempenho do sistema pode ser melhorado se forem utilizados além dos coeficientes Mel Cepstrais, outros parâmetros tais como : o Log de Energia Normalizada [26], coeficientes Delta Mel Cepstrais¹ [6] e Delta Log de Energia Normalizada [26].
- **Dimensões da Rede** As dimensões da ANN foram ajustadas para melhor otimizar o processo de estimação (Pré-REMAP). A rede MLP foi dimensionada com 70 neurônios na camada escondida porque a base de dados inicial (fones segmentados manualmente) possuía poucos exemplares. Mas durante a reestimação a quantidade de exemplos aumentou significativamente uma vez que estes passaram a ser tomados a cada quadro de análise e não apenas centrados nas marcas de segmentação. Neste caso acredita-se que a quantidade de 70 neurônios tenha sido insuficiente durante as etapas de reestimação (REMAP) dos parâmetros da ANN. Duas possíveis soluções

¹Com relação aos coeficientes Delta ainda existem dúvidas quanto a sua eficiência, uma vez que os modelos híbridos já são sensíveis a contexto. A influência destes coeficientes talvez limite-se às bordas do vetor de entrada \mathbf{X}_{n-c}^{n+d} .

para este problema são : (1) dimensionar a ANN de forma a melhor se adequar tanto aos exemplos utilizados durante a etapa de estimação (Pré-REMAP) como aos exemplos utilizados durante as etapas de reestimação (REMAP), (2) Utilizar uma ANN com poucos neurônios durante o Pré-REMAP e a partir desta ANN gerar os alvos suaves a serem utilizados no REMAP-1 e então utilizar uma outra ANN com um número maior de neurônios para as etapas de reestimação. O problema neste caso é que a ANN não poderá utilizar os pesos sinápticos da etapa Pré-REMAP.

8.1.2 Reconhecimento de Sentenças com o uso de Modelos de Sentenças

- **Eficiência do modelamento** Os resultados para o reconhecimento de sentenças mostram a eficiência, tanto da ANN no modelamento das *variabilidades acústicas* como dos HMMs na absorção das *variabilidades temporais* da fala. Os resultados mostram que o sistema foi capaz de reconhecer sentenças de tamanhos significativamente diferentes, frases com apenas 1 palavra e frase com até 20 palavras. Além disto mostra a capacidade do sistema em distinguir frases que diferem entre si em apenas uma única palavra como por exemplo : “A cotação do dólar aumentou, e as bolsas fecharam em baixa” e “A cotação do dólar aumentou, mas as bolsas fecharam em baixa”.

- **Flexibilidade** O reconhecimento com o uso de modelos de sentenças somente faz sentido para uma aplicação específica em que o número de sentenças permitidas pelo sistema seja bastante reduzido (porque devem ser montados explicitamente o modelo de cada uma das sentenças). Apesar desta restrição os modelos híbridos ANN-HMM, conforme definidos neste trabalho, apresentam a grande flexibilidade de permitir ao usuário expandir o número de sentenças permitidas pelo sistema, bastando somente, que ele forneça a conversão ortográfico-fonética desta nova sentença. Isto é possível porque para montar o HMM de uma nova sentença o sistema necessita saber apenas quais HMMs de fones devem ser concatenados.

8.1.3 Reconhecimento de Sentenças com o uso de Modelos de Palavras

- **Modelos dos fones com apenas um único estado** No sistema implementado foi utilizado um único estado para modelar um fone. Em geral, os sistemas baseados em HMM que representam o estado-da-arte em reconhecimento de fala, utilizam 3 estados por fone com o objetivo de modelar melhor as variações estatísticas existentes ao longo de um fone.

- **Modelos de duração de palavras** Os resultados apresentados no Capítulo 7 mostram que o uso da informação sobre as durações médias das palavras durante o procedimento de decodificação realizado pelo algoritmo Level Building é capaz de melhorar, significativamente, a taxa de acertos de palavras do sistema.
- **Algoritmos de busca** O algoritmo de busca implementado, algoritmo “level building”, não é muito eficiente do ponto de vista computacional. Uma versão mais otimizada do “Level Building” é o “Frame Synchronous Level Building”, que realiza a mesma tarefa do “Level Building”, porém, permitindo a realização dos cálculos de todos os L níveis de procura de forma paralela e não seqüencial como é realizado no “Level Building”.

8.2 Contribuições

- **Inovação** Primeiro trabalho na área de sistemas híbridos ANN-HMM aplicados ao reconhecimento de fala contínua realizado no Brasil.
- **Levantamento Bibliográfico** Foi realizado um vasto levantamento bibliográfico sobre reconhecimento de fala contínua, que será utilizado como base, no LPDF (Laboratório de Processamento Digital de Sinais de Fala da UNICAMP), para outros trabalhos sobre modelos híbridos ANN-HMM, técnicas de otimização conjunto de ANNs e HMMs, algoritmos de busca, modelamento de unidades sub-lexicais, independência de locutor e adaptabilidade de locutor.
- **Definição de metas de pesquisa na área de reconhecimento de fala contínua no LPDF - UNICAMP** Um dos grandes méritos deste trabalho foi o estabelecimento de bases para a realização de uma nova linha de pesquisa dentro do LPDF, o reconhecimento de fala contínua. Nessa linha deve-se destacar a utilização de modelos híbridos ANN-HMM para o modelamento acústico de unidades sub-lexicais e o domínio de algoritmos básicos para a decodificação acústica.

8.3 Sugestões para Trabalhos Futuros

- **Pré-Processamento** Realizar um melhor pré-processamento do sinal, utilizando Log de Energia Normalizada, coeficientes Delta e possivelmente análises em multitasas e segmentações não uniformes [42].

- **Construção de um segmentador automático** Para continuação dos trabalhos em reconhecimento de fala contínua é necessário a aquisição de uma nova base de dados com um número maior de sentenças e de palavras, e preferivelmente independente de locutor. No caso de sistemas baseados em modelos híbridos esta nova base deve ser pré-segmentada em termos de unidades sub-lexicais, necessitando, portanto, de um segmentador automático, possivelmente baseado em HMMs e alinhamentos de Viterbi.
- **Explorar melhor as não estacionariedades da fala** Melhorar o modelamento das unidades sub-lexicais utilizando um número maior de estados por fone (possivelmente 3 estados por fone) para modelar melhor as variabilidades estatísticas do sinal ao longo de uma unidade sub-lexical.
- **Unidades sub-lexicais dependentes de contexto** As unidades sub-lexicais utilizadas no sistema implementado (fones) foram considerados independentes de contexto, isto é, o modelo de cada fone utilizado não levava em consideração qual o fone que o precedia ou sucedia. O modelamento de fones dependentes de contexto não é uma tarefa difícil de ser implementada com o uso de um sistema híbrido ANN-HMM [33].

- **Métodos para otimização das ANNs** Desenvolver novos algoritmos para otimização das ANNs com o objetivo de melhor estimar as probabilidades $P(q_k|\mathbf{x}_n)$. Como sugestões podem ser desenvolvidos algoritmos híbridos baseados em algoritmos genéticos e métodos de gradiente [25].
- **Algoritmos de busca mais eficientes computacionalmente** Implementação de algoritmos de busca mais eficientes computacionalmente como é o caso do “Frame Synchronous Level Building” [17].
- **Restrições gramaticais** Realizar um estudo detalhado sobre restrições gramaticais a serem utilizadas durante o procedimento de busca e também como uma etapa de pós-processamento.
- **Uso de bases de dados padrões** Com o objetivo de avaliar melhor o sistema e poder comparar os resultados obtidos com os resultados apresentados em periódicos internacionais, utilizar bases de dados padrão em Inglês, como por exemplo o TIMIT (6300 sentenças, 6299 palavras, previamente rotulada em unidades sub-lexicais e gravada em ambiente de estúdio) e o NTIMIT (semelhante ao TIMIT, porém com a simulação de ruído telefônico).
- **Modelos híbridos ANN-HMM discriminativos** Implementar o sistema híbrido ANN-HMM discriminativo proposto por Konig e H. Bourlard [14].

Apêndice A

Critério MMI

O critério de maximização da informação mútua - MMI (“Maximum Mutual Information”), procura ressaltar a discriminação entre os modelos que competem entre si, na tentativa de utilizar da melhor forma possível as informações disponíveis no limitado conjunto de treinamento. Para uma dada seqüência de vetores acústicos $\mathbf{X}_{M_j}^{(s)}$ o critério MMI treina o modelo correto M_j (modelo associado a $\mathbf{X}_{M_j}^{(s)}$) positivamente e treina negativamente todos os outros modelos M_i , com $i \neq j$, ajudando a diferenciar melhor os modelos que competem entre si, e também melhorando a capacidade de discriminação durante a etapa de avaliação. A informação mútua entre uma seqüência de vetores acústicos $\mathbf{X}_{M_j}^{(s)}$ e o seu modelo associado M_j é definida segundo [27] como :

$$I(\mathbf{X}_{M_j}^{(s)}, M_j | \Theta) = \log \frac{p(\mathbf{X}_{M_j}^{(s)}, M_j | \Theta)}{p(\mathbf{X}_{M_j}^{(s)} | \Theta) \cdot P(M_j | \Theta)} \quad (\text{A.1})$$

$$= \log \frac{p(\mathbf{X}_{M_j}^{(s)} | M_j, \Theta)}{p(\mathbf{X}_{M_j}^{(s)} | \Theta)} \quad (\text{A.2})$$

$$= \log p(\mathbf{X}_{M_j}^{(s)} | M_j, \Theta) - \log p(\mathbf{X}_{M_j}^{(s)} | \Theta) \quad (\text{A.3})$$

utilizando o conceito de probabilidade marginal e realizando algumas manipulações simples, a Equação (A.3) pode ser reescrita, como:

$$I(\mathbf{X}_{M_j}^{(s)}, M_j | \Theta) = \underbrace{\log p(\mathbf{X}_{M_j}^{(s)} | M_j, \Theta)}_{\text{I}} - \underbrace{\sum_i \log p(\mathbf{X}_{M_j}^{(s)} | M_i, \Theta) \cdot P(M_i | \Theta)}_{\text{II}} \quad (\text{A.4})$$

sendo que :

I - Representa o treinamento positivo do modelo correto M_j (justamente como no critério ML).

II - Representa o treinamento negativo de todos os outros modelos M_i com $i \neq j$.

O treinamento com o critério MMI consiste na determinação dos parâmetros Θ que maximizem a informação mútua

$$\Theta_{MMI} = \arg \max I(\mathbf{X}_{M_j}^{(s)}, M_j | \Theta) \quad (\text{A.5})$$

Infelizmente, esta equação não pode ser solucionada por análise direta ou por reestimação [33]. Uma possível solução baseia-se em métodos de gradiente [33]. Entretanto tais procedimentos de otimização estão sujeitos a problemas de implementação numérica [27].

Também pode ser verificado através da Equação (A.2) que o critério MMI é equivalente ao critério de maximização da probabilidade a posteriori MAP ("Maximum A Posteriori"), o qual é expresso por $P(M_j | \mathbf{X}_{M_j}^{(s)}, \Theta)$ conforme Equação (2.9). Para verificar esta equivalência basta reescrever $P(M_j | \mathbf{X}_{M_j}^{(s)}, \Theta)$ em termos das regras de Bayes:

$$P(M_j | \mathbf{X}_{M_j}^{(s)}, \Theta) = \frac{p(\mathbf{X}_{M_j}^{(s)} | M_j, \Theta) \cdot P(M_j | \Theta)}{p(\mathbf{X}_{M_j}^{(s)} | \Theta)} \quad (\text{A.6})$$

Como o termo $P(M_j | \Theta)$ que representa o modelo da língua em geral é considerado independente de Θ , $P(M_j | \Theta) = P(M_j)$ e como a função logaritmo é uma função monotônica, então:

$$\arg \max_{\Theta} P(M_j | \mathbf{X}_{M_j}^{(s)}, \Theta) \Leftrightarrow \arg \max_{\Theta} \left\{ \log \frac{p(\mathbf{X}_{M_j}^{(s)} | M_j, \Theta)}{p(\mathbf{X}_{M_j}^{(s)} | \Theta)} \right\} \quad (\text{A.7})$$

Apêndice B

Dicionário de Pronúncias

A seguir são apresentadas todas as palavras que compõem o léxico utilizado, seguidas de suas respectivas pronúncias (representação em termos de fones). É importante observar que existem palavras com mais de uma pronúncia.

1.	(SILÊNCIO)	#
2a.	A	<i>a</i>
2b.	A	<i>á</i>
3.	ACEITARÃO	<i>á set á r ã u</i>
4.	ACORDO	<i>a k o r ã d u</i>
5.	ADELAIDE	<i>a del á i D i</i>
6.	AFIRMAR	<i>a f i r m á r</i>
7.	AGUARDAREMOS	<i>á g u á r ã d a r ã m u z</i>
8.	AINDA	<i>á i ã d</i>
9.	AJUSTES	<i>a j u s T i s</i>
10.	ALUNOS	<i>a l u n o z</i>
11.	AMANHÃ	<i>á m ã ã</i>
12a.	ANALISTAS	<i>á n á l i s t a z</i>
12b.	ANALISTAS	<i>á n á l i s t a s</i>
13.	ANTECIPAR	<i>ã t e s i p á r</i>
14.	ANTERIOR	<i>ã t e r i o r</i>
15a.	AO	<i>a u</i>
15b.	AO	<i>á u</i>
16a.	AOS	<i>á u s</i>
16b.	AOS	<i>a u s</i>
17.	APLICAÇÕES	<i>á p l i k á s õ i z</i>
18.	APROVEITAR	<i>á p r ã o v e i t á r</i>
19.	AQUI	<i>á k i</i>
20a.	AS	<i>a s</i>
20b.	AS	<i>ás</i>
20c.	AS	<i>á z</i>
21.	ASSIM	<i>á s i</i>
22.	ASSINADO	<i>á s i n á d u</i>
23.	ATÉ	<i>á t é</i>
24a.	ATRASADAS	<i>á t r á z á d a z</i>
24b.	ATRASADAS	<i>á t r á z á d a s</i>
25.	ATRÁS	<i>á t r á z</i>
26.	ATRATIVIDADE	<i>á t r á T i v i d á D i</i>
27.	ATUALIZAÇÃO	<i>á t u á l i z á s ã u</i>

28.	AUMENTO	<i>á u m ě t u</i>
29a.	AUMENTOU	<i>á m ě t o u</i>
29b.	AUMENTOU	<i>á u m ě t o u</i>
30.	BAIXA	<i>b á i f a</i>
31.	BAIXO	<i>b á i f u</i>
32a.	BANCO	<i>b ā k u</i>
32b.	BANCO	<i>b ā k</i>
33.	BANCOS	<i>b ā k u z</i>
34.	BARROSO	<i>b á R o z u</i>
35.	BASTANTE	<i>b á s t ā T i</i>
36.	BENEFICIANDO	<i>b e n e f i s i ā d u</i>
37a.	BOLSA	<i>b o u s a</i>
37b.	BOLSA	<i>b o s a</i>
38a.	BOLSAS	<i>b o u s a s</i>
38b.	BOLSAS	<i>b o s a s</i>
39.	BRASIL	<i>b ř a z i u</i>
40.	BRASILEIRA	<i>b ř a z i l e ř a</i>
41.	CADASTRAL	<i>k á d á s t ř á u</i>
42.	CADERNETA	<i>k á d e ř n e t a</i>
43.	CAFÉ	<i>k á f é</i>
44.	CAIXAS	<i>k á i f a z</i>
45.	CARTÃO	<i>k á r t ā u</i>
46a.	CENTO	<i>s ě t</i>
46b.	CENTO	<i>s ě t u</i>
47.	CENTRAL	<i>s ě t ř á u</i>
48.	CENTRO	<i>s ě t ř u</i>
49a.	CHAMADA	<i>f ā m á d a</i>
49b.	CHAMADA	<i>f á m á d a</i>
50.	CHEGARAM	<i>f e g á ř ā u</i>
51.	CHEQUE	<i>f é k i</i>
52.	CIDADE	<i>s i d á D i</i>
53.	CINCO	<i>s ĩ k u</i>
54.	CINQUENTA	<i>s ĩ k u ě t a</i>
55.	CLIENTE	<i>k l i ě T i</i>
56a.	CLIENTES	<i>k l i ě T i z</i>
56b.	CLIENTES	<i>k l i ě T z</i>
57.	CÓDIGO	<i>k ó D i g u</i>
58.	COLOCARÁ	<i>k o l o k á ř á</i>
59a.	COM	<i>k ō</i>
59b.	COM	<i>k o</i>
60.	COMPARECIMENTO	<i>k ō p á ř e s i m ě t u</i>
61.	CONDOMÍNIO	<i>k ō d o m í n i u</i>
62.	CONFORTO	<i>k ō f o r t u</i>
63.	CONSELHO	<i>k ō s e l u</i>
64a.	CONSIDERADO	<i>k ō s i d e ř á d</i>
64b.	CONSIDERADO	<i>k ō s i d e ř á d u</i>
65a.	CONSIDERAVELMENTE	<i>k ō s i d e ř á v e u m ě T i</i>
65b.	CONSIDERAVELMENTE	<i>k ō s d e ř á v e u m ě T i</i>
65c.	CONSIDERAVELMENTE	<i>k ō s i d e ř á v i u m ě T i</i>
66a.	CONSUMO	<i>k ō s u m u</i>
66b.	CONSUMO	<i>k ō s ũ m u</i>
67.	CONTA	<i>k ō t a</i>
68.	CONTAS	<i>k ō t á s</i>
69a.	CONTINUAM	<i>k ō T i n u ā</i>
69b.	CONTINUAM	<i>k ō T i n u ā u</i>
70a.	CONTRIBUINTES	<i>k ō t ř i b u ĩ T i z</i>
70b.	CONTRIBUINTES	<i>k ō t ř i b u ĩ T i s</i>
71.	CONVENCEMOS	<i>k ō v ě s ě m u z</i>
72a.	CONVÊNIO	<i>k ō v ě n i u</i>
72b.	CONVÊNIO	<i>k ō v e n i u</i>

73.	CORRIGIDAS	<i>koRijidas</i>
74.	COTAÇÃO	<i>kotásãu</i>
75a.	CRÉDITO	<i>křéDitu</i>
75b.	CRÉDITO	<i>křéDtu</i>
76.	CRUZ	<i>křus</i>
77.	CRUZEIROS	<i>křuzeřuz</i>
78.	CULTURAS	<i>kuťuraz</i>
79.	CUMPRIMENTO	<i>kũprimětu</i>
80.	CURVA	<i>kuřva</i>
81a.	DA	<i>d</i>
81b.	DA	<i>da</i>
82.	DADOS	<i>daduz</i>
83a.	DAS	<i>dáz</i>
83b.	DAS	<i>das</i>
83c.	DAS	<i>daz</i>
84a.	DE	<i>de</i>
84b.	DE	<i>Dĩ</i>
84c.	DE	<i>D</i>
85a.	DEPÓSITOS	<i>depózituz</i>
85c.	DEPÓSITOS	<i>depóztuz</i>
86.	DESCONTOS	<i>deskōtuz</i>
87.	DESENVOLVER	<i>dezěvoverě</i>
88.	DESTE	<i>desTi</i>
89.	DESTINO	<i>desTĩn</i>
90.	DETECTADO	<i>detektádu</i>
91.	DETERMINAÇÃO	<i>deteřminásãu</i>
92.	DEVE	<i>děvi</i>
93a.	DEVEM	<i>děvě</i>
93b.	DEVEM	<i>děvěi</i>
94.	DEVEMOS	<i>devěmuz</i>
95a.	DEVIDO	<i>devidu</i>
95b.	DEVIDO	<i>devíd</i>
96.	DEZEMBRO	<i>dezěbřu</i>
97.	DEZENOVE	<i>dezěnóvi</i>
98.	DEZESSEIS	<i>dezesez</i>
99.	DEZESSETE	<i>dezeséTi</i>
100.	DIA	<i>Dia</i>
101.	DIARIAMENTE	<i>DiářiáměTi</i>
102.	DIFERENTES	<i>DifeřěTis</i>
103.	DISPONÍVEIS	<i>Disponives</i>
104a.	DISPONÍVEL	<i>Dsponiviu</i>
104b.	DISPONÍVEL	<i>Disponivi</i>
104c.	DISPONÍVEL	<i>Disponiviu</i>
105.	DISPOSIÇÃO	<i>Dspozisãuũ</i>
106.	DO	<i>du</i>
107.	DOCUMENTO	<i>dokumětu</i>
108.	DOIS	<i>doiz</i>
109a.	DÓLAR	<i>dólař</i>
109b.	DÓLAR	<i>dólá</i>
110.	DUZENTOS	<i>duzětuz</i>
111a.	E	<i>e</i>
111b.	E	<i>i</i>
112.	É	<i>ě</i>
113.	EFICÁCIA	<i>efikásia</i>
114.	EFICIÊNCIA	<i>efisiěsia</i>
115.	ELE	<i>eli</i>
116a.	ELETRÔNICOS	<i>eletřōnikuz</i>
116b.	ELETRÔNICOS	<i>eletřonikus</i>
117a.	EM	<i>ěi</i>
117b.	EM	<i>ě</i>

118.	EMBARQUE	<i>ê b á r k</i>
119a.	EMPRESA	<i>ê p r e z</i>
119b.	EMPRESA	<i>ê i p r e z a</i>
120.	EMPRESÁRIO	<i>ê p r e z á r i u</i>
121.	ENTRE	<i>ê t r i</i>
122.	ENTREGOU	<i>ê t r e g o u</i>
123.	ESTA	<i>é s t</i>
124a.	ESTÁ	<i>í s t á</i>
124b.	ESTÁ	<i>s t á</i>
125.	ESTAÇÃO	<i>e s t á s ã u</i>
126.	ESTARÃO	<i>í s t á r ã u</i>
127.	ESTATAL	<i>e s t á t á u</i>
128.	ESTÃO	<i>í s t ã u</i>
129a.	ESTÁVEL	<i>e s t á v e u</i>
129b.	ESTÁVEL	<i>e s t á v i u</i>
130.	EXPLÍCITA	<i>e s p l i s i t a</i>
131.	EXPRESSO	<i>s p r e s u</i>
132.	FAZER	<i>f a z e r</i>
133.	FECHARAM	<i>f e f á r ã</i>
134.	FICHARÁ	<i>f i k á r á</i>
135.	FINANCIAMENTO	<i>f i n ã s i á m e t u</i>
136.	FOI	<i>f o i</i>
137a.	FORMULÁRIOS	<i>f o r m u l á r i u s</i>
137b.	FORMULÁRIOS	<i>f o r m u l á r i u z</i>
138.	FORTALEZA	<i>f o r t á l e z a</i>
139.	FUNCIONÁRIO	<i>f ũ s õ n á r i u</i>
140.	FUTURO	<i>f u t u r u</i>
141a.	GOVERNO	<i>g o v e r n u</i>
141b.	GOVERNO	<i>g o v e r n u</i>
142.	HAVERÁ	<i>á v e r á</i>
143a.	HORAS	<i>ó r a s</i>
143b.	HORAS	<i>ó r a z</i>
144.	IBGE	<i>i b e j e é</i>
145.	IMEDIATO	<i>i m e D i á t u</i>
146.	IMPORTAÇÃO	<i>ĩ p o r t á s ã u</i>
147.	IMPOSTO	<i>ĩ p o s t u</i>
148.	IMPULSIONADO	<i>ĩ p u s i õ n á d u</i>
149.	IMPULSIONOU	<i>ĩ p u s õ n o u</i>
150a.	INADEQUADO	<i>i n á d e k u á d u</i>
150b.	INADEQUADO	<i>ĩ n á d e k u á d u</i>
151.	INCOMPLETO	<i>ĩ k õ p l é t u</i>
152.	INDEXADORES	<i>ĩ d e k s á d o r i s</i>
153.	INFORMA	<i>ĩ f ó r m a</i>
154a.	INÍCIO	<i>i n i s i u</i>
154b.	INÍCIO	<i>i n i s u</i>
155.	INSTALADOS	<i>ĩ s t á l á d u z</i>
156.	INSTITUIÇÕES	<i>ĩ s T i t u i s õ i z</i>
157.	INSUFICIENTE	<i>ĩ s u f i s i e T i</i>
158.	INTEGRAÇÃO	<i>ĩ t e g r á s ã u</i>
159a.	INTERCÂMBIO	<i>ĩ t e r k ã b i u</i>
159b.	INTERCÂMBIO	<i>ĩ t e r k ã b i</i>
159c.	INTERCÂMBIO	<i>ĩ t e r k ã b i u</i>
160.	INTERESSANTE	<i>ĩ t e r e s ã T i</i>
161.	INTERNO	<i>ĩ t é r n u</i>
162a.	INVESTINDO	<i>ĩ v e s T i d u</i>
162b.	INVESTINDO	<i>ĩ v e s T i n</i>
163.	INVESTIR	<i>ĩ v e s T i r</i>
164.	ISTO	<i>ĩ s t u</i>
165.	JUNHO	<i>j ũ ŋ o</i>
166.	JUROS	<i>j u r u z</i>

167.	-LHES	<i>λ i s</i>
168.	LOCALIZADO	<i>l o k á l i z á d o</i>
169.	MAIORIA	<i>m á i o ř i a</i>
170.	MAIS	<i>m á i z</i>
171a.	MAS	<i>m á s</i>
171b.	MAS	<i>m á z</i>
172.	MEDIDA	<i>m e D i d a</i>
173.	MELHORES	<i>m e λ ó ř i z</i>
174a.	MERCADO	<i>m e ř k á d u</i>
174b.	MERCADO	<i>m e R k á d u</i>
175.	MÊS	<i>m e s</i>
176a.	MIL	<i>m i u</i>
176b.	MIL	<i>m i</i>
177.	MILHÕES	<i>m i λ õ i s</i>
178.	MINUTO	<i>m ĩ n u t u s</i>
179.	MODIFICAÇÕES	<i>m o d i f i k á s õ i z</i>
180.	MOMENTO	<i>m o m ě t u</i>
181.	MONETÁRIO	<i>m õ n e t á ř i u</i>
182.	MUITO	<i>m ũ i t</i>
183.	NA	<i>n a</i>
184.	NAÇIONAL	<i>n á s õ n á u</i>
185.	NÃO	<i>n ã u</i>
186.	NAQUELE	<i>n á k e l i</i>
187.	NAS	<i>n a z</i>
188.	NECESSÁRIO	<i>n e s e s á ř i u</i>
189.	NEM	<i>n ě i</i>
190.	NO	<i>n u</i>
191.	NOME	<i>n õ m i</i>
192.	NORMALMENTE	<i>n o ř m á u m ě T</i>
193.	NOS	<i>n u s</i>
194.	NOTÍCIA	<i>n o T i s i a</i>
195.	NOVE	<i>n ó v i</i>
196.	NOVO	<i>n o v u</i>
196.	NÚMERO	<i>n ũ m i ř u</i>
198a.	O	<i>o</i>
198b.	O	<i>u</i>
199.	OITENTA	<i>o i t ě t a</i>
200.	ONTEM	<i>õ t ě ĩ</i>
201a.	OPERAÇÕES	<i>o p e ř á s õ i s</i>
201b.	OPERAÇÕES	<i>o p e ř á s õ i z</i>
202.	OPINIÃO	<i>o p i n i ã u</i>
203.	OPORTUNIDADE	<i>o p o r t u n i d á D i</i>
204a.	OS	<i>u s</i>
204b.	OS	<i>u z</i>
205a.	OU	<i>o u</i>
205b.	OU	<i>o u</i>
206.	OUTROS	<i>o t ř u z</i>
207.	PÁGINAS	<i>p á j i n a s</i>
208.	PAÍS	<i>p á i s</i>
209a.	PARA	<i>p a ř</i>
209b.	PARA	<i>p á ř a</i>
210.	PARECE	<i>p á ř ě s</i>
211a.	PARTIR	<i>p á r T i r</i>
211b.	PARTIR	<i>p á ř T i ř</i>
212.	PASSA	<i>p á s a</i>
213.	PASSADA	<i>p á s á d a</i>
214.	PASSADO	<i>p á s á d u</i>
215.	PASSEGEIROS	<i>p á s á j e ř u s</i>
216.	PASSARÁ	<i>p á s á r á</i>
217.	PAULO	<i>p á u l u</i>

218.	PELO	<i>p el u</i>
219.	PELOS	<i>p el u z</i>
220.	PEQUENA	<i>p e k e n a</i>
221.	PERDA	<i>p e r d a</i>
222.	PERIGOSA	<i>p e ř i g ó z a</i>
223.	PERMITA	<i>p e ř m i t á</i>
224a.	PERMITE	<i>p e r m i T</i>
224b.	PERMITE	<i>p e ř m i T</i>
224c.	PERMITE	<i>p e ř m ĩ T i</i>
225.	PERMITIDO	<i>p e ř m í T i d u</i>
226.	PESQUISA	<i>p e s k i z a</i>
227.	PLANO	<i>p l ā n u</i>
228.	POR	<i>p u ř</i>
229.	PORCENTO	<i>p u ř s ě t u</i>
230.	PORQUE	<i>p u ř k e</i>
231.	PORTÃO	<i>p o ř t ā u</i>
232.	POSSO	<i>p ó s o</i>
233.	POUPANÇA	<i>p o u p ā s a</i>
234.	PRECISO	<i>p ř e s i z</i>
235.	PREÇO	<i>p ř e s u</i>
236.	PREÇOS	<i>p ř e s u s</i>
237.	PRESTAÇÃO	<i>p ř e s t á s ā u</i>
238.	PREZADO	<i>p ř e z á d u</i>
239.	PROBLEMA	<i>p ř o b l ě m á</i>
240.	PROJETO	<i>p ř o j ě t u</i>
241.	PROVOCANDO	<i>p ř o v o k ā n</i>
242.	PRÓXIMO	<i>p ř ó s i m u</i>
243.	PÚBLICA	<i>p u b l i k a</i>
244.	QUADRA	<i>k u á d ř a</i>
245.	QUANDO	<i>k u ā d u</i>
246.	QUARENTA	<i>k u á ř ě t</i>
247.	QUATRO	<i>k u á t ř u</i>
248.	QUATROCENTOS	<i>k u á t ř u s ě t u z</i>
249.	QUE	<i>k í</i>
250.	QUEDA	<i>k é d a</i>
251.	QUINHENTOS	<i>k ĩ ě t u s</i>
252.	QUINZE	<i>k ĩ z i</i>
253.	RADICAIS	<i>R á d i k á i z</i>
254.	REAIS	<i>R e á i s</i>
255a.	REAL	<i>R e á u</i>
255b.	REAL	<i>R é á u</i>
256a.	RECENTEMENTE	<i>R e s ě T i m ě T i</i>
256b.	RECENTEMENTE	<i>R e s ě T i m ě T</i>
257.	RECIFE	<i>R e s i f i</i>
258.	REGISTRADO	<i>R e j i s t ř á d u</i>
259a.	REGRAS	<i>R é g ř a s</i>
259b.	REGRAS	<i>R é g ř a z</i>
260.	RESOLUÇÃO	<i>R e z o l u s ā u</i>
261.	REUNIÃO	<i>R e u n i ā u</i>
262.	RIO	<i>R i u</i>
263a.	SALDO	<i>s á u d u</i>
263b.	SALDO	<i>s á u d</i>
264.	SANTA	<i>s ā t a</i>
265.	SÃO	<i>s ā u</i>
266.	SAQUES	<i>s á k i s</i>
267.	SE	<i>s i</i>
268.	SEGUINTE	<i>s i g ĩ T i</i>
269.	SEGUIR	<i>s e g i ř</i>
270.	SEGUNDO	<i>s e g ũ d u</i>
271.	SEMANA	<i>s e m ā n a</i>

272.	SEMANAS	<i>s e m ã n a s</i>
273a.	SEMPRE	<i>s ě p ř i</i>
273b.	SEMPRE	<i>s ě p ř e</i>
274.	SER	<i>s e r</i>
275.	SERÁ	<i>s e ř á</i>
276.	SESSENTA	<i>s e s ě t a</i>
277.	SETE	<i>s é T i</i>
278.	SETENTA	<i>s e t ě t</i>
279.	SEU	<i>s e u</i>
280.	SEUS	<i>s e u s</i>
281.	SIM	<i>s ĩ</i>
282.	SOBRE	<i>s o b ř</i>
283.	SOFRERÁ	<i>s o f ř e ř á</i>
284.	SUA	<i>s u a</i>
285.	SUBINDO	<i>s u b ĩ d u</i>
286.	SUBSTITUÍDO	<i>s u b i s T i t u i d u</i>
287.	SUFICIENTE	<i>s u f i s i ě T i</i>
288.	SUL	<i>s u u</i>
289.	TAXAS	<i>t á f a z</i>
290.	TELEBRAS	<i>t e l e b ř á s</i>
291.	TELECOMUNICAÇÕES	<i>t é l i k o m u n i k á s ō i z</i>
292a.	TELEFÔNICA	<i>t e l e f ô n i k a</i>
292b.	TELEFÔNICA	<i>t e l e f ô n i k</i>
293.	TELEFÔNICAS	<i>t e l e f ô n i k a s</i>
294.	TELESP	<i>t e l é s p</i>
295.	TEM	<i>t ě i</i>
296.	TEMPO	<i>t ě p u</i>
297.	TERÁ	<i>t e ř á</i>
298.	TIVEMOS	<i>T i v ě m u z</i>
299.	TODAS	<i>t o d a z</i>
300a.	TODOS	<i>t o d u z</i>
300b.	TODOS	<i>t o d u s</i>
301.	TRABALHO	<i>t ř á b á λ u</i>
302.	TRATA	<i>t ř á t a</i>
303.	TRÊS	<i>t ř e s</i>
304.	TREZENTOS	<i>t ř e z ě t u z</i>
305a.	TRINTA	<i>t ř ĩ t a</i>
305b.	TRINTA	<i>t ř ĩ t</i>
306.	ÚLTIMA	<i>u T i m a</i>
307.	ÚLTIMAS	<i>u T i m a s</i>
308.	UM	<i>ũ</i>
309a.	UMA	<i>ũ m a</i>
309b.	UMA	<i>u m a</i>
310.	UNA	<i>u n a</i>
311.	UNIVERSIDADES	<i>u n i v e r s i d á D i z</i>
312a.	VALOR	<i>v á l o r</i>
312b.	VALOR	<i>v á l o ř</i>
313a.	VENCIMENTO	<i>v ě s i m ě t</i>
313b.	VENCIMENTO	<i>v ě s i m ě t u</i>
314.	VIGOR	<i>v i g o r</i>
315.	VINCULADAS	<i>v ĩ k u l á d a z</i>
316a.	VINTE	<i>v ĩ T</i>
316b.	VINTE	<i>v ĩ T i</i>
317.	VISA	<i>v í z a</i>
318.	VOÇÊ	<i>v o s e</i>
319.	VÔO	<i>v o</i>

Apêndice C

Lista com as Sentenças Reconhecidas

Este Apêndice apresenta uma lista com as 100 sentenças utilizadas para a avaliação do sistema híbrido ANN-HMM implementado. Também são apresentadas as sentenças reconhecidas pelo sistema para os casos : (1) Pré-REMAP (sem reestimação), (2) REMAP-1, (3) REMAP-2 e (4) REMAP-3, estes quatro casos sem o uso de restrições gramaticais e (5) utilizando uma Gramática *P-Gram*. Em todos estes cinco casos foram utilizados o modelo de duração de palavras discutido na Seção 5.2.3.

- 1 Original A cotação do dólar aumentou, e as bolsas fecharam em baixa
Pré-REMAP *para* cotação ** ** aumentou *mas* bolsas *antecipar* em baixa *tempo*
REMAP-1 A cotação ** dólar aumentou *mas* bolsas fecharam em baixa
REMAP-2 A cotação ** dólar aumentou, e as bolsas fecharam em baixa
REMAP-3 A cotação ** dólar aumentou, e as bolsas fecharam em baixa
P-Gram. A cotação do dólar aumentou, e as bolsas fecharam em baixa
- 2 Original A cotação do dólar aumentou, mas as bolsas fecharam em baixa
Pré-REMAP *para* cotação ** dólar aumentou *mas* as bolsas fecharam *o* em baixa *cento*
REMAP-1 A cotação ** dólar aumentou *mais* as bolsas fecharam *uma* baixa
REMAP-2 A cotação do dólar aumentou, *mas* as bolsas fecharam *o* em baixa
REMAP-3 A cotação ** dólar aumentou *mais* as bolsas fecharam *o* em baixa
P-Gram. A cotação do dólar aumentou, *mas* as bolsas fecharam em baixa
- 3 Original A bolsa ficará estável ou sofrerá uma pequena queda
Pré-REMAP *passa* bolsa ficará *e* estável ** sofrerá uma pequena queda *por*
REMAP-1 A bolsa ficará *é* estável ** sofrerá uma pequena queda
REMAP-2 A bolsa ficará *estarão* sofrerá uma pequena queda
REMAP-3 A bolsa ficará *está Rio* sofrerá uma pequena queda
P-Gram A bolsa ficará estável ou sofrerá uma pequena queda

- 4 Original Não haverá ajustes, nem modificações radicais no plano
 Pré-REMAP *plano* haverá ajustes nem modificações *a dia mais* no plano tempo
 REMAP-1 Não haverá *as* ajustes nem modificações *a de queda* no plano
 REMAP-2 Não haverá *as* ajustes nem modificações *a de queda* no plano
 REMAP-2 Não haverá ajustes nem modificações *a de queda* no plano
 P-Gram Não haverá ajustes, nem modificações radicais no plano
- 5 Original Foi detectado um problema em seu cartão; ele deve ser substituído
 Pré-REMAP *para* foi detectado ** problema *São* cartão
 REMAP-1 Foi detectado ** problema *vinte São* cartão *para* , *ele de é ele ser substituído*
 REMAP-2 Foi detectado ** problema *vinte são* cartão, *ele de é ele ser substituído*
 REMAP-3 Foi detectado ** problema *vinte* seu cartão, *ele ele substituído*
 P-Gram Foi detectado um problema em seu cartão, ele deve ser substituído
- 6 Original É necessário que o convênio permita o intercâmbio
 Pré-REMAP *até* necessário ** ** convênio permita o intercâmbio *país*
 REMAP-1 É necessário *do* convênio permita *ao* intercâmbio
 REMAP-2 É necessário *do* convênio permita o intercâmbio
 REMAP-3 É necessário ** ** convênio permita o intercâmbio
 P-Gram É necessário que o convênio permita o intercâmbio
- 7 Original Posso afirmar-lhes que o convênio permite o intercâmbio
 Pré-REMAP *passa* afirmar *início* convênio permite ** intercâmbio *tempo*
-
- REMAP-1 Posso afirmar *início* convênio permite ** intercâmbio
 REMAP-2 Posso afirmar *ele se* o convênio permite ** intercâmbio
 REMAP-3 Posso afirmar *ele se* o convênio permite ** intercâmbio
 P-Gram Posso afirmar-lhes que o convênio permite o intercâmbio
- 8 Original O convênio que foi assinado recentemente permite o intercâmbio
 Pré-REMAP *posso* convênio ** foi e assinado recentemente permite ** intercâmbio *país*
 REMAP-1 o convênio ** foi *é* assinado recentemente permite o intercâmbio
 REMAP-2 O convênio que foi *é* assinado recentemente permite o intercâmbio
 REMAP-3 O convênio que foi *é* assinado recentemente permite o intercâmbio
 P-Gram O convênio que foi assinado recentemente permite o intercâmbio
- 9 Original O convênio permite o intercâmbio porque visa a integração entre alunos de culturas diferentes
 Pré-REMAP *tempo* convênio permite ** intercâmbio porque visa a integração entre *dados* de culturas diferentes *cento*
 REMAP-1 O convênio permite o intercâmbio porque visa *é* integração *em terá* alunos de culturas diferentes
 REMAP-2 O convênio permite o intercâmbio porque visa a *é* integração entre *é* alunos de culturas diferentes
 REMAP-3 O convênio permite o intercâmbio porque visa a integração entre *é* alunos de culturas diferentes
 P-Gram O convênio permite o intercâmbio porque visa a integração entre alunos de culturas diferentes
- 10 Original O convênio permite o intercâmbio quando se trata de universidades vinculadas ao projeto de integração
 Pré-REMAP o convênio permite ** intercâmbio quando se trata de universidades vinculadas ao projeto de integração
 REMAP-1 o convênio permite ** intercâmbio quando se trata de universidades vinculadas ao projeto de integração
 REMAP-2 ** convênio permite ** intercâmbio quando se trata de universidades vinculadas ao projeto de integração
 REMAP-3 ** convênio permite ** intercâmbio *com no* se trata de universidades vinculadas ao projeto de integração
 P-Gram. O convênio permite o intercâmbio quando se trata de universidades vinculadas ao projeto de integração

- 11 Original O convênio, que foi assinado recentemente, permite o intercâmbio
 Pré-REMAP *tempo* convênio que foi e assinado recentemente *permitido* intercâmbio *para*
 REMAP-1 O convênio que foi e assinado *a* recentemente *permitido* intercâmbio
 REMAP-2 O convênio que foi e assinado *a* recentemente *permitido* intercâmbio
 REMAP-3 O convênio *que foi e assinado a* recentemente, permite e intercâmbio
 P-Gram O convênio, que foi assinado recentemente, permite o intercâmbio
- 12 Original O convênio assinado na última reunião é mais interessante do que o anterior
 Pré-REMAP *pelo* convênio assinado *no* última reunião ** mais interessante do que ** anterior *muito*
 REMAP-1 O convênio assinado *no* última reunião ** mais interessante do *Rio* anterior *horas*
 REMAP-2 O convênio assinado *não* última reunião ** mais interessante *do Rio* anterior *horas*
 REMAP-3 O convênio assinado *não* última reunião ** mais interessante do *Rio* anterior
 P-Gram O convênio assinado na última reunião é mais interessante do que o anterior
- 13 Original A medida que o tempo passa, mais nos convencemos da eficácia do convênio
 Pré-REMAP *país dia pelo* tempo passa mais nos convencemos da eficácia do convênio *país*
 REMAP-1 A *mil devido* tempo passa , mais nos convencemos da eficácia do convênio
 REMAP-2 A *mil devido* tempo passa, mais nos convencemos da eficácia do convênio
 REMAP-3 A *mil devido* tempo passa, mais nos convencemos da eficácia do convênio
 P-Gram A medida que o tempo passa, mais nos convencemos da eficácia do convênio
- 14 Original Se o convênio permite o intercâmbio, devemos aproveitar a oportunidade para desenvolver o projeto de integração
 Pré-REMAP *próximo* convênio permite ** intercâmbio devemos aproveitar *ao* oportunidade para desenvolver o projeto de integração
 REMAP-1 Se o convênio permite ** intercâmbio devemos aproveitar *ao* oportunidade para *dezembro convênio* projeto de integração
 REMAP-2 Se o convênio permite ** intercâmbio devemos aproveitar *ao* oportunidade para desenvolver o projeto de integração
 REMAP-3 Se o convênio permite ** intercâmbio devemos aproveitar *ao* oportunidade para desenvolver o projeto de integração
 P-Gram. Se o convênio permite o intercâmbio, devemos aproveitar a oportunidade para desenvolver o projeto de integração
- 15 Original Localizado a uma quadra do centro da cidade, o condomínio permite que você uma trabalho e conforto
 Pré-REMAP *o* localizado ** uma quadra do cento ** cidade *por* condomínio permite que você *uma* trabalho e conforto
 REMAP-1 Localizado ** uma quadra do cento da cidade *por* condomínio permite que você *uma* trabalho e conforto
 REMAP-2 Localizado ** uma quadra do *cento* da cidade *por* condomínio permite que você *um na* trabalho e conforto
 REMAP-3 Localizado ** uma *a* quadra do centro da cidade *por* condomínio permite ** você *um na* trabalho e conforto
 P-Gram Localizado a uma quadra do centro da cidade, o condomínio permite que você uma trabalho e conforto
- 16 Original É suficiente
 Pré-REMAP *até* suficiente *por*
 REMAP-1 É suficiente
 REMAP-2 É suficiente
 REMAP-3 É suficiente
 P-Gram É suficiente
- 17 Original Isto é suficiente
 Pré-REMAP *país com* é suficiente *para*
 REMAP-1 Isto *o* é suficiente
 REMAP-2 Isto *o* é suficiente
 REMAP-3 Isto *o* é suficiente
 P-Gram Isto é suficiente

- 18 Original O saldo é suficiente
 Pré-REMAP *tempo São projeto suficiente para*
 REMAP-1 O saldo *o* é suficiente
 REMAP-2 O saldo *o* é suficiente
 REMAP-3 O saldo *o* é suficiente
 P-Gram O saldo é suficiente
- 19 Original O saldo de sua conta é suficiente
 Pré-REMAP *posso saldo de sua conta até suficiente tempo*
 REMAP-1 O saldo de sua conta *até* suficiente
 REMAP-2 O saldo de sua conta *até* suficiente *em*
 REMAP-3 O saldo de sua conta *até* suficiente *em*
 P-Gram O saldo de sua conta é suficiente
- 20 Original O saldo disponível é insuficiente
 Pré-REMAP *pelo saldo disponível o é insuficiente cento*
 REMAP-1 O saldo disponível *com* é insuficiente *em*
 REMAP-2 O saldo disponível *com* é insuficiente *em*
 REMAP-3 O saldo disponível *com* é insuficiente *em*
 P-Gram O saldo disponível é insuficiente
- 21 Original O saldo disponível em sua conta é insuficiente
 Pré-REMAP *cotação do disponível ** sua conta até insuficiente até*

 REMAP-1 O saldo disponível em sua *o com parece* insuficiente
 REMAP-2 O saldo disponível em sua *o com parece* insuficiente
 REMAP-3 O saldo disponível em sua *o com parece* insuficiente
 P-Gram O saldo disponível em sua conta é insuficiente
- 22 Original Isto parece insuficiente
 Pré-REMAP *país parece insuficiente tempo*
 REMAP-1 *-lhes* parece insuficiente
 REMAP-2 Isto parece insuficiente
 REMAP-3 Isto parece insuficiente
 P-Gram Isto parece insuficiente
- 23 Original O saldo parece ser insuficiente
 Pré-REMAP *passado parece e ser insuficiente tempo*
 REMAP-1 *do* saldo parece *ele* insuficiente
 REMAP-2 O saldo parece *ele* insuficiente
 REMAP-3 O saldo parece *ele* insuficiente
 P-Gram O saldo parece ser insuficiente
- 24 Original O saldo sempre está disponível
 Pré-REMAP *tempo saldo sempre está disponível para*
 REMAP-1 O saldo sempre está disponível
 REMAP-2 O saldo sempre está disponível
 REMAP-3 O saldo sempre está disponível
 P-Gram O saldo sempre está disponível

- 25 Original O saldo está sempre disponível no início do mês
 Pré-REMAP *tempo* saldo está sempre disponível ** início *mil em isto*
 REMAP-1 O saldo está sempre *de se número mil* início *mil em -lhes*
 REMAP-2 O saldo está sempre *de se número mil vencimento as*
 REMAP-3 O saldo está sempre *se número mil vencimento as*
 P-Gram O saldo está sempre disponível no início do mês
- 26 Original No início do mês, o saldo está disponível
 Pré-REMAP *tempo* início *mais* o saldo está disponível *para*
 REMAP-1 *de* início *nem os* o saldo está disponível
 REMAP-2 *de* início do *melhores* o saldo está disponível
 REMAP-3 *de* início *melhores* o saldo está disponível
 P-Gram ** início do mês, o saldo está disponível
- 27 Original Esta é a última chamada para o voo 737 da Rio-Sul
 Pré-REMAP *até as terá últimas* chamada *por os o se esta por e e sete* da Rio São *tempo*
 REMAP-1 Esta é ** última chamada *por os o sete três sete* da Rio Sul *um*
 REMAP-2 Esta é ** última chamada *por os o sete três sete* da Rio Sul *um*
 REMAP-3 Esta é ** última chamada *por os o os esta três sete* da Rio Sul *um*
 P-Gram Esta é a última chamada para o voo sete três sete da Rio-Sul
- 28 Original Isto é uma pesquisa de opinião pública
 Pré-REMAP *país a* uma pesquisa de opinião pública *tempo*
 REMAP-1 *estão com* pesquisa de opinião pública
 REMAP-2 *estão* uma pesquisa de opinião pública
 REMAP-3 *estão* uma pesquisa de opinião pública
 P-Gram Isto é uma pesquisa de opinião pública
- 29 Original O valor de sua conta telefônica é baixo
 Pré-REMAP *portão dois* de sua conta telefônica é baixa *tempo*
 REMAP-1 O valor de sua conta telefônica é baixo
 REMAP-2 O valor de sua conta telefônica é baixo
 REMAP-3 O valor de sua conta telefônica é baixo
 P-Gram O valor de sua conta telefônica é baixo
- 30 Original É de trinta mil cruzeiros, o valor de sua conta telefônica
 Pré-REMAP *parece* de trinta mil cruzeiros o valor de sua conta telefônica *empresa*
 REMAP-1 É de trinta mil cruzeiros o valor de sua conta telefônica
 REMAP-2 É de trinta mil cruzeiros, o valor de sua conta telefônica
 REMAP-3 É de trinta mil cruzeiros o valor de sua conta telefônica
 P-Gram É de trinta mil cruzeiros, o valor de sua conta telefônica
- 31 Original O vencimento de sua prestação será no dia quatro de junho
 Pré-REMAP *todos* vencimento de sua prestação será *última* quatro de *número país*
 REMAP-1 O vencimento de *Sul* o prestação será no dia quatro *visa nome*
 REMAP-2 O vencimento de *Sul* o prestação será no dia quatro *visa nome*
 REMAP-3 O vencimento de *Sul* o prestação será no dia quatro *visa uma um*
 P-Gram O vencimento de sua prestação será no dia quatro de *consumo*

- 32 Original O preço aumentou
 Pré-REMAP *todos* preço o aumentou *para*
 REMAP-1 O preço o aumentou *para*
 REMAP-2 O preço o aumentou *para*
 REMAP-3 O preço o aumentou
 P-Gram O preço aumentou
- 33 Original O preço do café aumentou
 Pré-REMAP ** preço do café aumentou *para*
 REMAP-1 *do* preço do café aumento o
 REMAP-2 O preço do café aumentou
 REMAP-3 O preço do café aumentou
 P-Gram O preço do café aumentou
- 34 Original O preço do café expresso aumentou
 Pré-REMAP ** preço ** café expresso aumentou *para*
 REMAP-1 O preço *do* café expresso o aumentou
 REMAP-2 O preço ** café expresso o aumentou
 REMAP-3 O preço ** café expresso o aumentou
 P-Gram O preço do café expresso aumentou
- 35 Original O preço do café aumentou consideravelmente
 Pré-REMAP *plano* preço ** café ** aumentou consideravelmente *tempo*
-
- REMAP-1 O preço do café ** aumentou consideravelmente
 REMAP-2 O preço do café ** aumentou consideravelmente
 REMAP-3 O preço do café ** aumentou consideravelmente
 P-Gram O preço do café aumentou consideravelmente
- 36 Original O preço do café aumentou consideravelmente na semana passada
 Pré-REMAP *porque São* café aumentou consideravelmente na semana passada *tempo*
 REMAP-1 *do* preço do café aumentou consideravelmente na semana passada
 REMAP-2 *do* preço do café aumentou consideravelmente na semana passada
 REMAP-3 O preço do café aumentou consideravelmente na semana passada
 P-Gram O preço do café aumentou consideravelmente na semana passada
- 37 Original Aumentou o preço do café
 Pré-REMAP *para* aumentou ** preço do café *cento*
 REMAP-1 Aumentou ** preço do café
 REMAP-2 Aumentou ** preço do café
 REMAP-3 Aumentou ** preço do café
 P-Gram Aumentou o preço do café
- 38 Original As taxas de juros no mercado interno estão subindo bastante
 Pré-REMAP *país* taxas *início* mercado interno estão subindo bastante *tempo*
 REMAP-1 As taxas de juros** mercado interno estão subindo bastante
 REMAP-2 As taxas de *Sul* no mercado interno estão subindo bastante *em vinte*
 REMAP-3 As taxas de juros no mercado interno estão subindo bastante *em vinte*
 P-Gram As taxas de juros no mercado interno estão subindo bastante

- 39 Original As contas chegaram atrasadas
 Pré-REMAP *passa* contas chegaram atrasadas *imposto*
 REMAP-1 As contas chegaram atrasadas
 REMAP-2 As contas chegaram atrasadas
 REMAP-3 As contas chegaram atrasadas
 P-Gram As contas chegaram atrasadas
- 40 Original As contas chegaram atrasadas ontem
 Pré-REMAP *passa* contas chegaram *muito* atrasadas *o TELESP*
 REMAP-1 As contas chegaram *muito* atrasadas *com o TELESP*
 REMAP-2 As contas chegaram *muito* atrasadas *com o TELESP*
 REMAP-3 As contas chegaram *muito* atrasadas *com o TELESP*
 P-Gram As contas chegaram atrasadas ontem
- 41 Original As contas telefônicas deste mês chegaram muito atrasadas ao banco
 Pré-REMAP *país* contas telefônicas deste *mil aceitarão* muito atrasadas ao banco *tempo*
 REMAP-1 As contas telefônicas *recentemente* chegaram muito atrasadas ao banco
 REMAP-2 As contas telefônicas *recentemente* chegaram muito atrasadas ao banco
 REMAP-3 As contas telefônicas *recentemente* chegaram muito atrasadas ao banco
 P-Gram As contas telefônicas deste mês chegaram muito atrasadas ao banco
- 42 Original Ontem, as contas chegaram aqui muito atrasadas
 Pré-REMAP *Paulo ele* as contas chegaram *bastante* muito atrasadas *para*
 REMAP-1 *posso ele* as contas chegaram *bastante* muito *a* atrasadas
 REMAP-2 *posso ele* as contas chegaram *partir* muito *a* atrasadas
 REMAP-3 *Com o ele* as contas chegaram *partir* muito *a* atrasadas
 P-Gram Ontem, as contas chegaram aqui muito atrasadas
-
- 43 Original Chegaram atrasadas
 Pré-REMAP Chegaram atrasadas *trinta*
 REMAP-1 Chegaram atrasadas
 REMAP-2 Chegaram atrasadas
 REMAP-3 Chegaram atrasadas
 P-Gram Chegaram atrasadas
- 44 Original Chegaram atrasadas todas as contas telefônicas deste mês
 Pré-REMAP *até* chegaram atrasadas todas as contas telefônicas deste *nem cento*
 REMAP-1 Chegaram *com* atrasadas todas as contas telefônicas deste *nem -lhes*
 REMAP-2 Chegaram *com* atrasadas todas as contas telefônicas deste mês
 REMAP-3 Chegaram *com* atrasadas todas as contas telefônicas deste mês
 P-Gram Chegaram atrasadas todas as contas telefônicas deste mês
- 45 Original O governo aumentou o imposto no mês passado
 Pré-REMAP *tempo* governo aumentou **** imposto **** mês passado *conta*
 REMAP-1 *ao* governo aumentou **** imposto** no mês passado
 REMAP-2 *ao* governo aumentou **** imposto** no mês passado
 REMAP-3 *ao* governo aumentou **** imposto** no mês passado
 P-Gram O governo aumentou o imposto no mês passado

- 46 Original O governo aumentou o imposto sobre importação
 Pré-REMAP *tempo* governo aumentou imposto sobre importação *tempo*
 REMAP-1 O governo aumentou ** imposto sobre importação
 REMAP-2 O governo aumentou ** imposto sobre importação
 REMAP-3 O governo aumentou ** imposto sobre importação
 P-Gram O governo aumentou o imposto sobre importação
- 47 Original O governo entregou os formulários aos contribuintes
 Pré-REMAP *Paulo* governo entregou ** formulários aos contribuintes *para*
 REMAP-1 O governo entregou ** formulários aos contribuintes
 REMAP-2 O governo entregou ** formulários aos contribuintes
 REMAP-3 O governo entregou ** formulários aos contribuintes
 P-Gram O governo entregou os formulários aos contribuintes
- 48 Original O governo entregou aos contribuintes os formulários
 Pré-REMAP *Paulo* governo entregou aos contribuintes os formulários *para*
 REMAP-1 O governo entregou *horas* contribuintes *nos* formulários
 REMAP-2 O governo entregou *horas* contribuintes *nos* formulários
 REMAP-3 O governo entregou aos contribuintes *nos* formulários
 P-Gram O governo entregou aos contribuintes os formulários
- 49 Original O banco colocará a sua disposição o novo cheque
 Pré-REMAP *tempo* banco colocará ** sua disposição ** novo *sete para*

 REMAP-1 *do* banco colocará ** sua disposição ** novo cheque
 REMAP-2 *do* banco colocará ** sua disposição ** novo cheque
 REMAP-3 O banco colocará ** sua disposição ** novo cheque
 P-Gram O banco colocará a sua disposição o novo cheque
- 50 Original A conta telefônica em nome de Adelaide Barroso terá vencimento amanhã
 Pré-REMAP *terá* conta telefônica *e* nome ** Adelaide Barroso *da* , terá vencimento *aumento*
 REMAP-1 A conta telefônica *e* nome *dia* Adelaide Barroso *da* , terá vencimento *aumento da*
 REMAP-2 A conta telefônica *e* nome *dia* Adelaide Barroso *da*, terá vencimento *ao não em a da*
 REMAP-3 A conta telefônica *e* nome *dia* Adelaide Barroso *da*, terá vencimento *a uma sim em a*
 P-Gram A conta telefônica em nome de Adelaide Barroso terá vencimento amanhã
- 51 Original A perda da atratividade das aplicações em caderneta de poupança está provocando um aumento de consumo no país
 Pré-REMAP A *prezado* atratividade das aplicações em caderneta de poupança *estarão quando* aumento de consumo ** país
 REMAP-1 A perda da atratividade das aplicações em caderneta de poupança está provocando *nome entre* de consumo ** país
 REMAP-2 A perda da atratividade das aplicações em caderneta de poupança *de* está provocando *nome entre* de consumo ** país
 REMAP-3 A perda da atratividade das aplicações em caderneta de poupança está provocando *o* aumento de consumo ** país
 P-Gram. A perda da atratividade das aplicações em caderneta de poupança está provocando um aumento de consumo no país
- 52 Original O mercado foi considerado inadequado
 Pré-REMAP *tempo* mercado foi considerado inadequado *tempo*
 REMAP-1 *do* mercado foi considerado inadequado
 REMAP-2 *do* mercado foi considerado inadequado
 REMAP-3 O mercado foi considerado inadequado
 P-Gram O mercado foi considerado inadequado

- 53 Original O mercado foi considerado inadequado pelos analistas
 Pré-REMAP *tempo* mercado foi considerado inadequado pelos analistas *tempo*
 REMAP-1 O mercado foi considerado inadequado pelos analistas
 REMAP-2 O mercado foi considerado inadequado pelos analistas
 REMAP-3 O mercado foi considerado inadequado pelos analistas
 P-Gram O mercado foi considerado inadequado pelos analistas
- 54 Original O mercado foi considerado inadequado naquele momento
 Pré-REMAP *tempo* mercado foi considerado inadequado naquele momento *trinta*
 REMAP-1 *uma é a do* foi considerado inadequado naquele momento
 REMAP-2 *uma é a do* foi considerado inadequado naquele momento
 REMAP-3 *uma é a do* foi considerado inadequado naquele momento
 P-Gram O mercado foi considerado inadequado naquele momento
- 55 Original Naquele momento, o mercado foi considerado inadequado pelos analistas das melhores instituições de pesquisa
 Pré-REMAP Naquele momento *por* mercado foi considerado inadequado pelos analistas das melhores instituições e pesquisa
 REMAP-1 Naquele momento *por* mercado foi considerado inadequado pelos analistas das *ele* melhores instituições de pesquisa
 REMAP-2 Naquele momeno *por* mercado foi considerado inadequado pelos analistas das melhores instituições de pesquisa
 REMAP-3 Naquele momeno *por* mercado foi considerado inadequado pelos analistas das melhores instituições de pesquisa
 P-Gram. Naquele momento, ** mercado foi considerado inadequado pelos analistas das melhores instituições de pesquisa
- 56 Original Diariamente
 Pré-REMAP *perda* diariamente *para*
-
- REMAP-1 *dia é cliente*
 REMAP-2 *dia é cliente*
 REMAP-3 Diariamente
 P-Gram Diariamente
- 57 Original Curva perigosa
 Pré-REMAP Curva *preço das ao*
 REMAP-1 Curva perigosa
 REMAP-2 Curva perigosa
 REMAP-3 Curva perigosa
 P-Gram Curva perigosa
- 58 Original Dia vinte do sete
 Pré-REMAP *pública* vinte do sete *para*
 REMAP-1 Dia vinte do sete *para*
 REMAP-2 Dia vinte do sete
 REMAP-3 Dia vinte do sete
 P-Gram Dia vinte do sete
- 59 Original Sim
 Pré-REMAP sim
 REMAP-1 sim
 REMAP-2 Sim
 REMAP-3 Sim
 P-Gram Sim

- 60 Original Não
 Pré-REMAP *por não tempo*
 REMAP-1 Não
 REMAP-2 Não
 REMAP-3 Não
 P-Gram Não
- 61 Original Saldo: vinte e cinco reais
 Pré-REMAP *passado país vinte ** cinco reais tempo*
 REMAP-1 Saldo , *de vinte ** cinco regras*
 REMAP-2 Saldo, *de vinte ** cinco regras*
 REMAP-3 Saldo, *de vinte ** cinco real -lhes*
 P-Gram Saldo: vinte e cinco reais
- 62 Original Estação Santa Cruz
 Pré-REMAP *prestação Santa Cruz para*
 REMAP-1 Estação Santa Cruz
 REMAP-2 Estação Santa Cruz
 REMAP-3 Estação Santa Cruz
 P-Gram Estação Santa Cruz
- 63 Original Passageiros com destino a São Paulo, Recife e Fortaleza, embarque imediato
 Pré-REMAP Passageiros com deste *na São Paulo Recife ** Fortaleza ainda aqui imediato com é*

 REMAP-1 *as fazer as com deste na São posso Recife ** Fortaleza , vinte Barroso imediato com*
 REMAP-2 *as fazer as com deste na São posso Recife ** Fortaleza, sim banco aqui imediato*
 REMAP-3 *as fazer as com deste na São posso Recife ** Fortaleza, embarque imediato*
 P-Gram Passageiros com destino a São Paulo, Recife e Fortaleza, embarque imediato
- 64 Original Os bancos atrás de mais eficiência
 Pré-REMAP *pelos bancos atrás e mais eficiência tempo*
 REMAP-1 O bancos atrás ** mais eficiência
 REMAP-2 Os bancos atrás ** mais eficiência
 REMAP-3 Os bancos atrás ** mais eficiência
 P-Gram ** bancos atrás de mais eficiência
- 65 Original Descontos de até 50%
 Pré-REMAP *tempo descontos ** até cinquenta por cento para*
 REMAP-1 Descontos e até cinquenta por cento
 REMAP-2 Descontos ** até cinquenta por cento
 REMAP-3 Descontos a até cinquenta por cento
 P-Gram Descontos de até cinquenta por cento
- 66 Original Número incompleto
 Pré-REMAP *por número incompleto conta*
 REMAP-1 da número incompleto
 REMAP-2 Número incompleto
 REMAP-3 Número incompleto
 P-Gram Número incompleto

- 67 Original Vinte e cinco
 Pré-REMAP o vinte ** cinco
 REMAP-1 o vinte ** cinco
 REMAP-2 Vinte ** cinco
 REMAP-3 Vinte ** cinco
 P-Gram Vinte e cinco
- 68 Original Vinte cinco reais
 Pré-REMAP tempo vinte ** cinco real TELESP
 REMAP-1 de vinte cinco regras
 REMAP-2 de vinte cinco regras
 REMAP-3 Vinte cinco regras
 P-Gram Vinte cinco reais
- 69 Original Cento e vinte cinco
 Pré-REMAP até cento ** vinte cinquenta
 REMAP-1 Cento ele vinte ** cinco
 REMAP-2 Cento ele vinte ** cinco
 REMAP-3 Cento ele vinte ** cinco
 P-Gram Cento e vinte cinco
- 70 Original Cento e vinte e cinco reais
 Pré-REMAP posso e entre vinte ** cinco da TELESP
-
- REMAP-1 setenta e vinte ** cinco regras
 REMAP-2 Cento de vinte ** cinco regras
 REMAP-3 Cento e vinte ** cinco regras
 P-Gram Cento e vinte e cinco reais
- 71 Original Quatrocentos e quarenta e nove
 Pré-REMAP para quatrocentos e quarenta estável tempo
 REMAP-1 Quatrocentos e quarenta e nove
 REMAP-2 Quatrocentos e quarenta e nove
 REMAP-3 Quatrocentos e quarenta e nove
 P-Gram Quatrocentos e quarenta e nove
- 72 Original Dois mil, cento e vinte e cinco
 Pré-REMAP todos Rio centro vinte ** cinco tempo
 REMAP-1 Dois mil cento e vinte ** cinco
 REMAP-2 Dois mil cento e vinte ** cinco
 REMAP-3 Dois mil, cento e vinte ** cinco
 P-Gram Dois mil, cento e vinte e cinco
- 73 Original Dezesesseis mil e quinhentos
 Pré-REMAP até dezesesseis mil e quinhentos país
 REMAP-1 Dezesesseis mil e quinhentos
 REMAP-2 Dezesesseis mil e quinhentos
 REMAP-3 Dezesesseis mil e quinhentos
 P-Gram Dezesesseis mil e quinhentos

- 74 Original Oitenta milhões, trezentos e sessenta mil e duzentos e setenta e um
 Pré-REMAP *aproveitar* milhões trezentos ** sessenta mil e duzentos ** setenta e *uma para*
 REMAP-1 Oitenta milhões trezentos ** sessenta mil *mil do cento os* setenta e *uma*
 REMAP-2 Oitenta milhões trezentos ** sessenta mil *mil do cento os* setenta e *uma*
 REMAP-3 *com setenta* milhões trezentos ** sessenta mil *mil* duzentos ** setenta e *uma*
 P-Gram Oitenta milhões, trezentos e sessenta mil e duzentos e setenta e um
- 75 Original A TELEBRÁS, a empresa de telecomunicações brasileira, está investindo em pesquisa
 Pré-REMAP *mas* TELEBRÁS empresa *os* telecomunicações brasileira , está investindo e pesquisa *até*
 REMAP-1 *até* TELEBRÁS *está* empresa *das* telecomunicações brasileira , está investindo e pesquisa
 REMAP-2 *até* TELEBRÁS *está* empresa *das* telecomunicações brasileira, está investindo *de* pesquisa
 REMAP-3 *até* TELEBRÁS, *assim* empresa *das* telecomunicações brasileira, está investindo *de* pesquisa
 P-Gram A TELEBRÁS, a empresa de telecomunicações brasileira, está investindo em pesquisa
- 76 Original A TELEBRÁS, uma empresa estatal, está investindo em pesquisa
 Pré-REMAP *para* TELEBRÁS *com* uma empresa estatal , está investindo ** pesquisa *tempo*
 REMAP-1 *até* TELEBRÁS , uma empresa estatal , está investindo em pesquisa
 REMAP-2 *até* TELEBRÁS, uma empresa estatal, está investindo em pesquisa
 REMAP-3 *até* TELEBRÁS, uma empresa estatal, está investindo em pesquisa
 P-Gram A TELEBRÁS, uma empresa estatal, está investindo em pesquisa
- 77 Original Tivemos recentemente a seguinte notícia : a TELEBRÁS passará a investir mais em pesquisa
 Pré-REMAP *até* tivemos recentemente ** seguinte notícia *país* TELEBRÁS passará ** investir mais em pesquisa
-
- REMAP-1 Tivemos recentemente ** seguinte notícia , a TELEBRÁS passará ** investir mais em pesquisa
 REMAP-2 Tivemos recentemente ** seguinte notícia, a TELEBRÁS passará a investir mais *investir visa*
 REMAP-3 Tivemos recentemente ** seguinte notícia, a TELEBRÁS passará ** investir mais *vinte* pesquisa
 P-Gram Tivemos recentemente a seguinte notícia, a TELEBRÁS passará a investir mais em pesquisa
- 78 Original TELESP informa : dezoito horas e trinta minutos
 Pré-REMAP TELESP informa , dezoito horas e trinta *mil outros sempre*
 REMAP-1 TELESP informa , dezoito horas e trinta minutos
 REMAP-2 TELESP informa, dezoito horas e trinta minutos
 REMAP-3 TELESP informa, dezoito horas e trinta minutos
 P-Gram TELESP informa : dezoito horas e trinta minutos
- 79 Original Empresário, é preciso antecipar o futuro
 Pré-REMAP *para* empresário ** preciso antecipar *os* futuro *pelo*
 REMAP-1 *é* empresário , é preciso antecipar o futuro
 REMAP-2 *é* empresário, é preciso antecipar o futuro
 REMAP-3 *é* Empresário, é preciso antecipar o futuro
 P-Gram Empresário, é preciso antecipar o futuro
- 80 Original Prezado cliente, aguardaremos o seu comparecimento
 Pré-REMAP Prezado cliente *para* aguardaremos o seu comparecimento *para*
 REMAP-1 Prezado cliente , aguardaremos o seu comparecimento
 REMAP-2 Prezado cliente, aguardaremos o seu comparecimento
 REMAP-3 Prezado cliente, aguardaremos o seu comparecimento
 P-Gram Prezado cliente, aguardaremos o seu comparecimento

- 81 Original O código foi registrado pelo funcionário
 Pré-REMAP *Paulo* código foi registrado *de os impulsionado tempo*
 REMAP-1 *no* código foi registrado pelo funcionário
 REMAP-2 *no* código foi registrado pelo funcionário
 REMAP-3 O código foi registrado pelo funcionário
 P-Gram O código foi registrado pelo funcionário
- 82 Original O convênio, um documento de trinta páginas, tem permitido o intercâmbio
 Pré-REMAP *pelo* convênio ** documento de trinta páginas ** permitido ** intercâmbio *com tempo*
 REMAP-1 O convênio ** documento de trinta páginas , tem *na* permitido ** intercâmbio *com*
 REMAP-2 O convênio ** documento de trinta páginas, *trinta* permitido ** intercâmbio *com*
 REMAP-3 O convênio ** documento de trinta páginas, *trinta* permitido ** intercâmbio
 P-Gram O convênio, um documento de trinta páginas, tem permitido o intercâmbio
- 83 Original Os caixas eletrônicos não aceitarão mais depósitos
 Pré-REMAP *poupança caixas* eletrônicos não aceitarão ** depósitos *para*
 REMAP-1 Os *a* caixas eletrônicos não aceitarão ** depósitos
 REMAP-2 Os caixas eletrônicos não aceitarão ** depósitos
 REMAP-3 Os caixas eletrônicos não aceitarão ** depósitos
 P-Gram Os caixas eletrônicos não aceitarão mais depósitos
- 84 Original Os caixas eletrônicos não aceitarão mais depósitos a partir das 15 horas
 Pré-REMAP *bolsas os* eletrônicos não aceitarão mais depósitos a partir das quinze *saques*
-
- REMAP-1 *bolsa* caixas eletrônicos não aceitarão mais depósitos a partir *da sete vinte* horas
 REMAP-2 *bolsa* caixas eletrônicos não aceitarão mais depósitos a partir das *sete vinte* horas
 REMAP-3 Os caixas eletrônicos não aceitarão mais depósitos a partir *da sete vinte* horas
 P-Gram Os caixas eletrônicos não aceitarão mais depósitos a partir das quinze horas
- 85 Original Os caixas eletrônicos não aceitarão mais depósitos a partir das 15 horas do próximo dia vinte e nove
 Pré-REMAP Os *esta os* eletrônicos não aceitarão mais depósitos a partir das quinze horas , do próximo dia vinte ** nove
 REMAP-1 Os *café se os* eletrônicos não aceitarão mais depósitos a partir das quinze horas , do próximo dia vinte ** nove
 REMAP-2 Os caixas eletrônicos não aceitarão mais depósitos a partir das quinze horas, do próximo dia vinte ** nove
 REMAP-3 Os caixas eletrônicos não aceitarão mais depósitos a partir das quinze horas, do próximo dia vinte ** nove
 P-Gram Os caixas eletrônicos não aceitarão mais depósitos a partir das quinze horas do próximo dia vinte e nove
- 86 Original E não estarão disponíveis para saques a partir das 17 horas do dia trinta
 Pré-REMAP *opinião* estarão disponíveis para *a* saques , a partir das dezessete horas do dia trinta *esta*
 REMAP-1 *de* não estarão disponíveis passa saques , a partir das dezessete horas do dia trinta
 REMAP-2 *de* não estarão disponíveis para saques, a partir das dezessete horas do dia trinta
 REMAP-3 E não estarão disponíveis para saques, a partir das dezessete horas do dia trinta
 P-Gram E não estarão disponíveis para saques a partir das dezessete horas do dia trinta
- 87 Original Todos os bancos devem fazer a atualização cadastral até o dia 31 de dezembro
 Pré-REMAP Todos os *o* bancos devem fazer ** atualização cadastral até ** ** trinta *Rio* de dezembro
 REMAP-1 Todos os *o* bancos devem fazer ** atualização cadastral até ** ** trinta *é* e um de dezembro
 REMAP-2 Todos os bancos devem fazer ** atualização cadastral até ** ** trinta *aqui* um de dezembro
 REMAP-3 Todos os bancos devem fazer ** atualização cadastral até ** ** trinta *aqui* um de dezembro
 P-Gram Todos os bancos devem fazer a atualização cadastral até o dia trinta e um de dezembro

- 88 Original Todos os bancos devem fazer a atualização cadastral de seus clientes até 31 de dezembro
 Pré-REMAP Todos os bancos devem fazer ** atualização cadastral de seus clientes até trinta *mil mil* dezembro *com*
 REMAP-1 Todos os bancos devem fazer ** atualização cadastral de seus clientes até trinta *é e um* de dezembro
 REMAP-2 Todos os bancos devem fazer ** atualização cadastral de seus clientes até trinta *é e um* de dezembro
 REMAP-3 Todos os bancos devem fazer ** atualização cadastral de seus clientes até trinta *é e um* de dezembro
 P-Gram Todos os bancos devem fazer a atualização cadastral de seus clientes até trinta e um de dezembro
- 89 Original De acordo com a determinação do Banco Central
 Pré-REMAP De acordo *quadra* determinação ** banco central *conta*
 REMAP-1 De acordo com a determinação ** banco central
 REMAP-2 De acordo com a determinação ** Banco Central
 REMAP-3 De acordo com a determinação ** Banco Central
 P-Gram De acordo com a determinação do Banco Central
- 90 Original Em cumprimento a Resolução 2025 do Conselho Monetário Nacional
 Pré-REMAP Em cumprimento a resolução dois mil ** vinte ** cinco , o conselho monetário nacional *com*
 REMAP-1 Em cumprimento a resolução dois mil ** vinte ** cinco *para* , do *com ser* monetário nacional
 REMAP-2 Em cumprimento *nove nas o São* dois mil ** vinte ** cinco, do *com ser* Monetário Nacional
 REMAP-3 Em cumprimento *nove* resolução dois mil ** vinte ** cinco, do *com ser* Monetário Nacional
 P-Gram Em cumprimento a Resolução 2025 do Conselho Monetário Nacional
- 91 Original De acordo com a determinação explícita do Banco Central
 Pré-REMAP De acordo com ** determinação explícita do banco central *para*

 REMAP-1 *a* , de acordo com ** determinação explícita do banco central
 REMAP-2 *a*, de acordo com ** determinação explícita do Banco Central
 REMAP-3 De acordo com a determinação explícita do Banco Central
 P-Gram De acordo com a determinação explícita do Banco Central
- 92 Original As operações continuam
 Pré-REMAP *atrás* operações continuam *país*
 REMAP-1 As operações continuam
 REMAP-2 As operações continuam
 REMAP-3 As operações continuam
 P-Gram As operações continuam
- 93 Original As operações de crédito continuam
 Pré-REMAP *atrás* operações de crédito continuam *tempo*
 REMAP-1 As operações de crédito continuam
 REMAP-2 As operações de crédito continuam
 REMAP-3 As operações de crédito continuam
 P-Gram As operações de crédito continuam
- 94 Original As operações de crédito e financiamento continuam a seguir as regras
 Pré-REMAP *para* as operações de crédito ** financiamento continuam ** seguir as regras *tempo*
 REMAP-1 As operações de crédito e financiamento continuam ** seguir as regras
 REMAP-2 As operações de crédito e financiamento continuam ** seguir as regras
 REMAP-3 As operações de crédito e financiamento continuam ** seguir as regras
 P-Gram As operações de crédito e financiamento continuam a seguir as regras

- 95 Original As operações de crédito e financiamento continuam a seguir as regras do banco central
 Pré-REMAP *mas* operações de crédito ** financiamento continuam ** seguir as regras do banco central *com*
 REMAP-1 As operações *dia* crédito *mil* financiamento continuam ** seguir *nas* regras do banco central
 REMAP-2 As operações de crédito *mil* financiamento continuam ** seguir as regras do banco central
 REMAP-3 As operações de crédito *mil* financiamento continuam ** seguir *nas* regras do banco central
 P-Gram As operações de crédito e financiamento continuam a seguir as regras do banco central
- 96 Original As operações de crédito e financiamento corrigidas por outros indexadores
 Pré-REMAP *mas* operações de crédito ** financiamento corrigidas por outros indexadores *que*
 REMAP-1 As operações de crédito e financiamento corrigidas por outros indexadores *que*
 REMAP-2 As operações de crédito e financiamento corrigidas por outros indexadores *com*
 REMAP-3 As operações de crédito e financiamento corrigidas por outros indexadores
 P-Gram As operações de crédito e financiamento corrigidas por outros indexadores
- 97 Original Continuam normalmente a seguir as regras do banco central ainda em vigor
 Pré-REMAP Continuam ** ** seguir as regras do banco central , ainda ** vigor
 REMAP-1 Continuam ** ** seguir as regras do banco central , ainda *e e* vigor
 REMAP-2 Continuam ** ** seguir as regras do banco central, ainda *e e* vigor
 REMAP-3 Continuam ** ** seguir as regras do banco central, ainda *e e* vigor
 P-Gram Continuam ** a seguir as regras do banco central ainda em vigor
- 98 Original O aumento no consumo devido ao Plano Real impulsionou os preços
 Pré-REMAP *momento* consumo devido ** plano real impulsionou os preços *o*
-
- REMAP-1 O aumento no consumo devido ao plano real impulsionou os preços
 REMAP-2 O aumento ** consumo devido ao Plano Real impulsionou os preços
 REMAP-3 O aumento ** consumo devido ao Plano Real impulsionou os preços
 P-Gram O aumento no consumo devido ao Plano Real impulsionou os preços
- 99 Original Segundo dados do IBGE, o aumento no consumo devido ao plano real impulsionou os preços
 Pré-REMAP *para seu* dados do IBGE , ** *documento* consumo devido ** plano real impulsionou os preços
 REMAP-1 Segundo dados do IBGE , *documento* no consumo devido plano real impulsionou os preço *cheque*
 REMAP-2 Segundo dados do IBGE, *o documento* no consumo devido ** plano real impulsionou os preços *cheque*
 REMAP-3 Segundo dados do IBGE, *o documento* no consumo devido ** plano real impulsionou os preços *cheque*
 P-Gram Segundo dados do IBGE, o aumento no consumo devido ao plano real impulsionou os preços
- 100 Original O aumento no consumo devido ao Plano Real tem impulsionado consideravelmente os preços nas últimas semanas
 Pré-REMAP
 REMAP-1 * aumento no consumo devido * Plano Real *para em* impulsionado consideravelmente os preço *bancos última* semanas.
 REMAP-2 * aumento no consumo devido * Plano Real *para em* impulsionado consideravelmente os preço *não os última* semanas.
 REMAP-3 * aumento no consumo devido * Plano Real *por em* impulsionado consideravelmente os preço *não os última* semanas.
 P-Gram. * aumento no consumo devido ao Plano Real tem impulsionado consideravelmente os preços nas últimas semanas

Bibliografia

- [1] BENGIO, Y., MORI, R. D., FLAMMIA, G. AND KOMPE, R., "Global Optimization of a Neural Network-Hidden Markov Model Hybrid", IEEE Trans. on Neural Networks., vol. 3, pp. 252-259, **1992**.
- [2] BISHOP, C. M., *Neural Networks for Patterns Recognition*, Oxford University Press, **1996**.
- [3] BOURLARD, H., HERMANSKY, H. AND MORGAN, N., "Towards increasing speech recognition error rates", Speech Communication., vol. 12, pp. 205-231, **1996**.

- [4] BOURLARD, H. AND WELLEKENS, C. J., "Links Between Markov Models and Multilayer Perceptrons", IEEE Trans. Patt. Anal. and Mach. Intell., vol. 12, pp. 1167-1178, **1990**.
- [5] DAVIS, S. B. AND MELMELSTEIN, P., "Comparision of Parametric Representation for Monosyllabic Word Recognition in Continuously Spoken Sentences", IEEE Transactions on Acoustic, Speech and Signal Processing, ASSP - 28, No. 4, August **1980**.
- [6] DELLER, J., PROAKIS, J. AND HANSEN, J., *Discrete-Time Processing of Speech Signals*, Macmillan, New York, **1993**.
- [7] DEVILLERS, L. AND DUGAST, C., "Hybrid System Combining Expert-TDNN and HMMs for Continuous Speech Recognition" , IEEE ICASSP94, , pp 165-168, **1994**.
- [8] DUGAST, C., DEVILLERS, L. AND AUBERT, X., "Combining TDNN and HMM in a Hybrid System for Improved Continuous-Speech Recognition" , IEEE Trans. Acoustic, Speech and Audio Processing, vol. 2, pp 217-223, jan. **1994**.
- [9] FRITSCH, J., "Modular Neural Networks for Speech Recognition," Diploma Thesis, Interactive Systems Laboratories, Carnegie Mellon University, **1996**.
- [10] HAYKIN, S., *Neural Networks : A Comprehensive Foundation*, Prentice Hall, New Jersey, **1994**.

- [11] HUANG, X. C., HON, H. W., HWANG H. W. AND LEE K. F., "A Comparative Study of Discrete, Semicontinuous, and Continuous Hidden Markov Models", *Computer Speech and Language*, vol(7), pp. 359-368, **1993**.
- [12] HUANG, X. D., ARIKI, Y. AND JACK, M. A., "Hidden Markov Models for Speech Recognition", Edinburgh University Press, Edinburgh, **1990**.
- [13] JELINEK, F., "A Fast Sequential Decoding Algorithm Using A Stack," *IBM Journal, Res. Develop.*, Vol. 13, pp. 675-685, Nov. **1969**.
- [14] KONIG, Y., "REMAP : Recursive Estimation and Maximization of A Posteriori Probabilities in Transition-Based Speech Recognition", PhD Thesis, University of California at Berkeley, **1996**.
- [15] KONIG, Y., BOURLARD, H. AND MORGAN, N., "REMAP : Experiments with Speech Recognition", *IEEE ICASSP96*, Atlanta, pp. 7-10, **1996**.
- [16] LEE, K. F., *Automatic Speech Recognition*, Kluwer Academic Publishers, 2^a edition, **1992**.
- [17] LEE, C. H. AND RABINER, L. R., "A Frame-Synchronous Network Search Algorithm for Connected Word Recognition", *IEEE Transaction on Acoustics, Speech, and Signal Processing*, Vol. 37, No. 11, November **1989**.
- [18] MORAIS E. S. E VIOLARO, F., "Sistema Híbrido ANN-HMM para Reconhecimento de Fala Contínua", *Anais do XV SBT - Simpósio Brasileiro de Telecomunicações*, Recife, pp. 117-120, **1997**.
- [19] MORAIS E. S., VIOLARO, F., YNOGUTI, C. A. E NETTO, M. A., "Sistemas Híbridos ANN-HMM Baseados nos Critérios ML e MAP para Reconhecimento de Séries Temporais", *Anais do 3^o SBAi - Simpósio Brasileiro de Automação Inteligente*, Vitória, pp. 406-411, **1997**.
- [20] MORAIS E. S., VIOLARO, F. E NETTO, M. A., "Modelo Oculto de Markov Parametrizado através de uma ANN e sua Aplicação no Reconhecimento de Fala Contínua", *Anais do III CBRN - Congresso Brasileiro de Redes Neurais*, Florianópolis, pp. 324-329, **1997**.
- [21] MORGAN, N. AND BOURLARD, H., "Continuous Speech Recognition : An introduction to the Hybrid HMM/Connectionist Approach", *IEEE Signal Processing Magazine*, pp. 25-49, may **1995**.
- [22] MORGAN, N. AND BOURLARD, H., "Neural Networks for Statistical Recognition of Continuous Speech", *Proceedings of the IEEE.*, vol. 83, pp. 741-770, may **1995**.

- [23] NEY, H., "The Use of a One-Stage Dynamic Programming Algorithm for Connected Word Recognition", IEEE Transactions on Acoustics, Speech, and Signal Processing, Vol. ASSP-32, pp 253-271, April 1984.
-
- [24] NUNES, H. F., "Reconhecimento de Fala Baseado em HMM", Tese de Mestrado, UNICAMP, Campinas, 1996.
- [25] PACHECO, S. S. AND THOMÉ, A. G., "Turbo-Shake", Anais do III Congresso Brasileiro de Redes Neurais, pp. 51-55, Julho 1997.
- [26] PICONE, J. W., "Signal Modeling Techiques in Speech Recognition", Proceedings of the IEEE, Vol. 81, No 9, 1215-1247, September 1993.
- [27] RABINER, L. AND JUANG, B. H., *Fundamentals of Speech Recognition*, Prentice Hall 1993.
- [28] RABINER, L. R., "A Tutorial on Hidden Markov Models an Selected Applications in Speech Recognition", Proceedings of the IEEE, vol. 77, No. 2, pp. 257-286, 1989.
- [29] RABINER, L. R. AND LEVINSON, S. E., "A Speaker-Independent, Syntax-Directed, Connected Word Recognition System Based on Hidden Markov Models and Level Building," IEEE Transactions on Acoustics, Speech, and Signal Processing, Vol. ASSP - 33, No. 3, pp. 561-573, June 1985.
- [30] RICHARD, M. D. AND LIPPMANN, R., "Neural Network Classifers Estimate Bayesian A-Posteriori Probabilities", Neural Computation 3(4), pp. 461-483, 1991.
- [31] ROBINSON, A. J., "An Application of Recurrent Nets to Phone Probability Estimation", IEEE Trans. on Neural Networks., vol. 5, pp. 298-305, 1994.
- [32] TATMAN, G. AND JANNARONE, R., "Real Time Neural Networks, III: Alternative Neural Networks For Speech Aplications", IEEE Computer Society Press., pp. 591-596, 1991.
- [33] TEBELSKIS, J., "Speech Recognition using Neural Networks", PhD Thesis, Carnegie Mellon University, 1995.
- [34] TEBELSKIS, J., "Performance Through Consistency : Connectionist Large Vocabulary Contiunuous Speech Recognition", IEICE on Neuro-Computing, pp. 391-398, 1990.
- [35] VINTSYUK, T. K., "Element-Wise Recognition of Continuous Speech Composed of Words from a Specified Dictionary", Kibernetika, Vol 7, pp. 133-143, March-April. 1971.

- [36] YAN, Y., FANTY, M. AND COLE, R., "Speech Recognition Using Neural Networks with Forward-Backward Probability Generated Targets", IEEE ICASSP97, Munich, April **1997**.
- [37] YOMA, N. B., "Reconhecimento Automático de Palavras Isoladas : Estudo e Aplicação dos Métodos Determinístico e Estocástico", Tese de Mestrado, UNICAMP, Campinas, **1993**.
- [38] YOUNG, S., "A Review of Large-Vocabulary Continuous-Speech Recognition", IEEE Signal Processing Magazine, pp. 45-57, september **1996**.
- [39] WAIBEL, A., ZEPPENFELD, T. AND HOUGHTON, R., "Improving The MS-TDNN for Word Spotting", International Acoustic Speech and Signal Processing, pp. 475-478, **1993**.
- [40] WAIBEL, A., HANAZAWA, T., HINTON, G., SHIKANO, K. AND LANG, K. J., "Phoneme Recognition Using Time-Delay Neural Networks", IEEE Trans. on Acoustics Speech and Signal Processing., vol. 37, pp. 328-338, march **1989**.
- [41] WAN, E. A., "Times Series Prediction by Using a Connectionist Network with Internal Delay Lines", *Forecasting the Future and Undertanding the Past*, A. S. Weigend and N. A. Gershenfeld, pp. 195-217, Addison-Wesley, **1993**.
- [42] OSTENDORF, M., AND S. ROUKOS., "A Stochastic Segment Model for Phoneme-Based Contiuous Speech Recognition", IEEE Transactions on Acoustics, Speech, and Signal Processing, Vol. ASSP-37, pp 1857-1869, **1989**.