



**Facultatea de Electronică, Telecomunicații și  
Tehnologia Informației**

**Margit Antal**

**CONTRIBUȚII LA RECUNOAȘTEREA VORBIRII  
ȘI A VORBITORULUI**

**Teză de doctorat**

PREȘEDINTE: Prof.dr.ing. Dan Pitică -decan, Facultatea de Electronică,  
Telecomunicații și Tehnologia Informației

MEMBRI:

1. Prof.dr.ing. Gavril Todorean -conducător științific, Universitatea Tehnică din Cluj-Napoca
2. Prof.dr.ing. Inge Gavăt -referent, Universitatea Politehnică din București
3. Prof.dr.ing. Corneliu Toma -referent, Universitatea Politehnică din Timișoara
4. Prof.dr.ing. Corneliu Rusu -referent, Universitatea Tehnică din Cluj-Napoca

**2006**

## Mulțumiri

Realizarea acestei lucrări a fost posibilă datorită suportului și sprijinului a mai multor persoane.

Pentru posibilitatea cercetării realizate și a scrierii acestei lucrări mulțumesc domnului profesor Dr. Gavril Toderean, care m-a sprijinit pe tot parcursul perioadei de doctorat. El a fost cel care a avut încredere în mine și a făcut să-mi pară munca de cercetare o experiență plăcută.

Doresc să-i mulțumesc soțului meu Győző Nemes pentru primele lecții de prelucrarea semnalelor și pentru suportul și răbdarea de-a lungul anilor de studii de doctorat. Fără el nu aş fi putut realiza această lucrare.

Dintre cei care m-au sprijinit în special prin discuții științifice doresc să-i mulțumesc membrilor grupului de cercetare de inteligență artificială de la universitatea Szeged, în special domnului profesor János Csirik, cercetătorilor András Kocsor și László Tóth.

În final, dar nu în ultimul rând doresc să mulțumesc familiei pentru sprijinul uman, fără care n-aș fi avut posibilitatea de a elabora această lucrare.



# Cuprins

<b>1</b>	<b>Introducere</b>	<b>1</b>
1.1	Prezentarea capitolelor . . . . .	1
1.2	Contribuții originale . . . . .	3
<b>2</b>	<b>Bazele teoretice</b>	<b>5</b>
2.1	Introducere . . . . .	5
2.2	Elemente de fonetică . . . . .	7
2.2.1	Semnalul vocal . . . . .	7
2.2.2	Producere și percepere . . . . .	8
2.2.3	Transcriere fonetică . . . . .	10
2.2.4	Clasificarea fonemelor . . . . .	10
2.3	Extragerea caracteristicilor . . . . .	14
2.3.1	Preprocesare . . . . .	14
2.3.2	Analiză cepstrală . . . . .	15
2.3.2.1	Coeficienții Mel-cepstrali (MFCC) . . . . .	15
2.3.2.2	Coeficienții cepstrali de predicție liniară (LPCC) .	16
2.3.3	Postprocesare . . . . .	17
2.4	Tehnici de grupare . . . . .	17
2.4.1	Introducere . . . . .	17
2.4.2	Algoritmul K-means . . . . .	19
2.4.3	Algoritmul Fuzzy K-means . . . . .	22
2.4.4	Cuantizare vectorială . . . . .	23
2.4.5	Algoritmul LVQ . . . . .	25
2.4.6	Mixturi gaussiene . . . . .	27
2.4.7	Grupare ierarhică . . . . .	30
2.5	Modele Markov ascunse . . . . .	32
2.5.1	Elemente ale unui model Markov ascuns cu observații discrete	33
2.5.2	Cele trei probleme de bază ale modelelor Markov ascunse .	34

2.5.2.1	Soluții pentru Problema 1 . . . . .	35
2.5.2.2	Soluție pentru Problema 2 . . . . .	40
2.5.2.3	Soluție pentru Problema 3 . . . . .	42
2.5.3	Topologii de modele Markov ascunse . . . . .	44
2.5.4	Tipuri de modele Markov ascunse . . . . .	46
2.5.4.1	Modele continue . . . . .	46
2.5.4.2	Modele semicontinuе . . . . .	47
2.5.5	Detalii de implementare . . . . .	48
2.5.5.1	Scalarea . . . . .	48
2.5.5.2	Secvențe de observații multiple . . . . .	51
2.5.5.3	Estimarea parametrilor inițiali . . . . .	52
2.6	Recunoașterea vorbirii . . . . .	53
2.6.1	Introducere . . . . .	53
2.6.2	Modele de limbaj . . . . .	55
2.6.3	Măsurarea erorii în sisteme de recunoaștere a vorbirii . . . . .	56
2.6.4	Tehnici de căutare . . . . .	57
2.6.4.1	Căutare cadru-sincron cu o singură trecere . . . . .	59
2.6.4.2	Decodarea cu stivă . . . . .	64
2.7	Baze de date . . . . .	66
<b>3</b>	<b>Recunoașterea fonetică a vorbitorului</b>	<b>67</b>
3.1	Introducere . . . . .	67
3.2	Modelarea vorbitorului . . . . .	70
3.3	Distribuția conținutului fonetic . . . . .	72
3.3.1	Modele VQ . . . . .	72
3.3.2	Modele GMM . . . . .	73
3.4	Modele GMM pure . . . . .	75
3.5	Modele GMM structurate fonetic . . . . .	79
3.6	Concluzii . . . . .	83
<b>4</b>	<b>Clasificarea fonemelor cu metode statistice</b>	<b>87</b>
4.1	Introducere . . . . .	87
4.2	Analiza acustică . . . . .	88
4.3	Evaluarea clasificării . . . . .	88
4.4	Rezultate experimentale . . . . .	89
4.4.1	TIMIT . . . . .	89
4.4.2	OASIS Numbers . . . . .	96
4.5	Concluzii . . . . .	101

<b>5 Recunoașterea fonemelor cu modele GMM</b>	<b>103</b>
5.1 Introducere . . . . .	103
5.2 Rezultate experimentale . . . . .	105
5.2.1 TIMIT . . . . .	105
5.2.2 OASIS Numbers . . . . .	106
5.3 Concluzii . . . . .	109
<b>6 Modelarea duratei fonemelor</b>	<b>111</b>
6.1 Introducere . . . . .	111
6.2 Rezultate experimentale . . . . .	112
6.2.1 OASIS Numbers . . . . .	112
6.2.2 TIMIT . . . . .	114
6.3 Concluzii . . . . .	116
<b>7 Recunoașterea cuvintelor prin modelare fonetică</b>	<b>117</b>
7.1 Introducere . . . . .	117
7.2 Modele de cuvinte . . . . .	118
7.3 Modele de foneme . . . . .	119
7.4 Concluzii . . . . .	122
<b>8 Concluzii finale</b>	<b>123</b>
8.1 Contribuții la recunoașterea vorbitorului . . . . .	123
8.2 Contribuții la recunoașterea vorbirii . . . . .	124
8.3 Direcții de cercetare . . . . .	126
<b>A Duratele fonemelor din baza de date TIMIT</b>	<b>127</b>
<b>B Duratele fonemelor din baza de date OASIS</b>	<b>137</b>
<b>C Foneme TIMIT</b>	<b>143</b>
<b>D Vocabular OASIS Numbers</b>	<b>147</b>



# Listă de figuri

2.1	Schema producerii și perceperei sunetului - varianta 1 (după [67]) . . . . .	8
2.2	Schema producerii și a perceperei sunetului - varianta 2 (după [67]) . . . . .	9
2.3	Reprezentarea modulară a extragerii coeficienților Mel cepstrali . . . . .	15
2.4	Algoritm de grupare și VQ (după [65]) . . . . .	25
2.5	Rețea neuronală LVQ . . . . .	26
2.6	Calculul funcției de probabilitate . . . . .	28
2.7	Dendogramă reprezentând gruparea datelor . . . . .	30
2.8	a) Secvențe de operații necesare în vederea calculării variabilei $\alpha_{t+1}(j)$ . b) Implementarea calculării lui $\alpha_t(i)$ . . . . .	37
2.9	Operații necesare pentru calculul variabilei $\beta_t(i)$ . . . . .	39
2.10	Operații necesare pentru calculul probabilității evenimentului, ca sistemul să se afle în momentul $t$ în starea $i$ , respectiv în momentul $t + 1$ în starea $j$ . . . . .	43
2.11	Trei tipuri diferite de MMA . . . . .	45
2.12	Regula lui Bayes pentru un sistem de recunoașterea vorbirii . . . . .	54
2.13	Un sistem ASR cu structură integrată . . . . .	58
2.14	Un sistem ASR cu structură modulară . . . . .	59
2.15	Aliniere timp-perechi stări cuvinte . . . . .	61
2.16	Tipuri de tranziții . . . . .	64
2.17	Arbore de căutare pentru decodorul cu stivă . . . . .	65
3.1	Structura unui sistem de recunoaștere a vorbitorului . . . . .	68
3.2	Distribuția categoriilor fonetice în clustere . . . . .	74
3.3	Selectia cadrelor din mijlocul fonemelor . . . . .	76
3.4	felc0-Distribuția fonemelor în componentele gaussiene ale unui model de vorbitor . . . . .	77
3.5	mdab0-Distribuția fonemelor în componentele gaussiene ale unui model de vorbitor . . . . .	78
3.6	Calculul similarității într-un model de vorbitor structurat fonetic .	82

3.7	Probabilități de identificare a vorbitorilor pentru primii 20 de vorbitori . . . . .	84
4.1	Clasificare foneme - TIMIT 61 - GMM 8 . . . . .	91
4.2	Clasificare foneme - TIMIT 61 - GMM 16 . . . . .	92
4.3	Clasificare foneme - TIMIT 61 - GMM 32 . . . . .	93
4.4	Matricea de confuzie TIMIT 61 . . . . .	94
4.5	Clasificator cu două niveluri . . . . .	95
4.6	Clasificare categorii fonetice TIMIT . . . . .	95
4.7	Matricea de confuzie OASIS 31 . . . . .	98
4.8	Ratele de clasificare OASIS 31, modele cu 16 componente gaussiene	99
5.1	Algoritmul Viterbi - Căutare în spațiul stărilor . . . . .	104
5.2	TIMIT-Variația ratei de recunoaștere în funcție de variația parametrului $\beta$ . . . . .	107
5.3	OASIS-Variația ratei de recunoaștere în funcție de variația parametrului $\beta$ . . . . .	108
6.1	Model Markov ascuns cu trei stări . . . . .	112
6.2	Durata minimă, maximă și medie a fonemelor OASIS . . . . .	113
6.3	Durata minimă, maximă și medie a fonemelor TIMIT . . . . .	115
7.1	Recunoașterea cuvintelor izolate (după [67]) . . . . .	119
7.2	OASIS-Vocabular organizat sub formă de trie . . . . .	120
7.3	Recunoașterea cuvintelor cu modele de foneme . . . . .	121
A.1	Histograme ale duratelor fonemelor: /aa/, /ae/, /ah/, /ao/, /aw/, /ax/, /ax-h/, /axr/ . . . . .	128
A.2	Histograme ale duratelor fonemelor: /ay/, /b/, /bcl/, /ch/, /d/, /dcl/, /dh/, /dx/ . . . . .	129
A.3	Histograme ale duratelor fonemelor: /eh/, /el/, /em/, /en/, /eng/, /epi/, /er/, /ey/ . . . . .	130
A.4	Histograme ale duratelor fonemelor: /f/, /g/, /gcl/, /h#//, /hh/, /hv/, /ih/, /ix/ . . . . .	131
A.5	Histograme ale duratelor fonemelor: /iy/, /jh/, /k/, /kcl/, /l/, /m/, /n/, /ng/ . . . . .	132
A.6	Histograme ale duratelor fonemelor: /nx/, /ow/, /oy/, /p/, /pau/, /pcl/, /q/, /r/ . . . . .	133

A.7	Histograme ale duratelor fonemelor: /s/, /sh/, /t/, /tcl/, /th/, /uh/, /uw/, /ux/ . . . . .	134
A.8	Histograme ale duratelor fonemelor: /v/, /w/, /y/, /z/, /hz/ . . .	135
B.1	Histograme ale duratelor fonemelor: /2/, /2:/, /A/, /E/, /J/, /O/, /'d'/, /e:/ . . . . .	138
B.2	Histograme ale duratelor fonemelor: /h/, /i/, /i:/, /'k/, /j/, /o/, /l/, /l:/ . . . . .	139
B.3	Histograme ale duratelor fonemelor: /m/, /n/, /o:/, /r/, /s/, /t/, /'ts/, /u/ . . . . .	140
B.4	Histograme ale duratelor fonemelor: /u:/, /v/, /z/ . . . . .	141



# Listă de tabele

1	Alfabetul fonetic ARPAbet . . . . .	11
2	Clasificarea fonemelor limbii engleze (după [26]) . . . . .	12
3	Variante pe categorii fonetice . . . . .	75
4	Foneme TIMIT . . . . .	80
5	Categorii de foneme TIMIT . . . . .	81
6	Rata de identificare pentru 630 de vorbitori utilizând GMM-uri fonetic pure . . . . .	81
7	Rate de recunoaștere pentru modele de vorbitori structurate fonetic - TIMIT 630 de vorbitori . . . . .	82
8	Lista identificatorilor vorbitorilor identificați în mod incorrect . . . . .	83
9	Frecvența fonemelor în mulțimea de test și cea de antrenare . . . . .	90
10	Rata de recunoaștere pentru setul TIMIT-61 . . . . .	90
11	Rata de recunoaștere pentru setul TIMIT-39 . . . . .	90
12	Cele mai frecvente erori - TIMIT 61 . . . . .	91
13	Clasificare categorii fonetice TIMIT . . . . .	92
14	Clasificare foneme TIMIT 61 cu clasificatorul pe două niveluri . . . . .	96
15	Clasificare cadru cu cadru-TIMIT 61 . . . . .	96
16	Frecvența fonemelor în mulțimea de antrenare și test . . . . .	97
17	Clasificare foneme OASIS 31 . . . . .	98
18	Cele mai frecvente erori - TIMIT 61 . . . . .	100
19	Clasificare cadru cu cadru - OASIS 31 . . . . .	101
20	Rezultate comparative pentru clasificarea fonemelor TIMIT 39 . . . . .	101
21	Rezultate comparative pentru clasificarea fonemelor OASIS 31 . . . . .	102
22	Recunoașterea fonemelor - <b>timit_test_core</b> . . . . .	106
23	Recunoașterea fonemelor - <b>timit_test</b> . . . . .	106
24	Recunoașterea fonemelor - <b>oasis_test</b> . . . . .	106
25	Recunoașterea fonemelor - OASIS - 31 . . . . .	114
26	Durata medie de fonem și valoare $\beta$ . . . . .	116

27	Recunoașterea cuvintelor cu modele GMM - OASIS-10 cifre . . . . .	118
28	Recunoaștere cuvinte - OASIS 26 de cuvinte - $\beta = 20$ . . . . .	121
29	Foneme TIMIT-1 . . . . .	144
30	Foneme TIMIT-2 . . . . .	145
31	Vocabular alofonic OASIS Numbers . . . . .	148

# Listă de algoritmi

1	Algoritmul K-means . . . . .	21
2	Algoritmul Fuzzy K-means . . . . .	23
3	Algoritmul LVQ-1 . . . . .	26
4	Algoritmul de grupare ierarhică aglomerativă . . . . .	31
5	Algoritmul de cuantizare vectorială binară . . . . .	32
6	Algoritmul forward . . . . .	36
7	Algoritmul backward . . . . .	39
8	Algoritmul lui Viterbi . . . . .	41
9	Algoritmul Viterbi modificat . . . . .	42
10	Algoritmul WER . . . . .	57
11	Algoritmul cu o singură trecere Single-best . . . . .	63



# **Capitolul 1**

## **Introducere**

Această lucrare investighează modele statistice pentru două probleme, și anume recunoașterea vorbitorului și recunoașterea vorbirii. Modelele utilizate au fost tehnice de grupare și modelele Markov ascunse.

În cazul recunoașterii vorbitorului obiectivul principal este studierea conținutului modelelor de vorbitori din punct de vedere fonetic. Un al doilea obiectiv este utilizarea altor informații, decât cele acustice pentru recunoașterea vorbitorului. În acest caz se dorește investigarea categoriilor fonetice, ca informație de ordin superior.

În cazul recunoașterii vorbirii obiectivul principal este investigarea modelelor mixturielor Gassiene pentru modelarea fonemelor. Deoarece foarte multe articole [19, 40, 52, 57, 59, 66, 80] subliniază efectul neglijabil al probabilităților de tranziție în cazul modelelor Markov ascunse, obiectivul este de a studia comportamentul sistemelor de recunoaștere fără acești parametri. Un al doilea obiectiv este modificarea algoritmilor de decodare pentru integrarea modelării duratei fonemelor în acestea.

### **1.1 Prezentarea capituloelor**

Capitolul 1 este introducerea în care se prezintă motivațiile acestei teze, structura pe capitulo și contribuțiile originale.

Capitolul 2 este dedicat prezentării tuturor noțiunilor teoretice care au fost utilizate în această lucrare. Sunt prezentate noțiuni de fonetică, metodele convenționale de extragerea caracteristicilor, tehnice de grupare, modelele Markov ascunse și tehnice de decodare.

Capitolul 3 prezintă problema recunoașterii vorbitorului precum și tehnice

principale de modelare a vorbitorului. Capitolul se limitează la prezentarea metodelor de recunoaștere independente de text. Ca element original am studiat distribuția conținutului fonetic în modele de vorbitori de tip mixturi gaussiene. Au fost studiate două aspecte, distribuția fonemelor individuale precum și distribuția categoriilor fonetice. Un alt element original introdus în acest capitol este modelarea vorbitorilor utilizând datele specifice ale unei singure categorii fonetice, cum ar fi vocalele sau nazalele. Categoriile fonetice au fost ordonate din punctul de vedere al puterii lor de discriminare a vorbitorului. În sfârșit a fost propus un model original, care grupează datele acustice pe categorii fonetice. Capitolul prezintă pentru toate modelele propuse și evaluarea acestora pe baze de date publice.

Capitolul 4 prezintă problema clasificării fonemelor, unitatea acustică utilizată pentru recunoașterea vorbirii. Fonemele sunt modelate cu mixturi gaussiene. Capitolul prezintă rezultate experimentale pentru clasificarea fonemelor, grupurilor, respectiv categoriilor de foneme. Pe baza rezultatelor experimentale, care indică faptul, că fonemele sunt confundate cu alte foneme din aceeași categorie, se propune un nou tip de clasificator, clasificatorul cu două niveluri. Sunt prezentate rezultatele experimentale pentru acest nou tip de clasificator. În final sunt prezentate rezultatele pentru clasificarea cadru cu cadru, care se utilizează tot pentru evaluarea modelelor. Măsurătorile sunt prezentate pentru două baze de date publice, TIMIT și OASIS Numbers.

Capitolul 5 este o continuare naturală a capitolului precedent, se trece de la problema clasificării la problema recunoașterii fonemelor. Ca algoritm de decodare se utilizează algoritmul Viterbi. Deoarece modelele de foneme utilizate sunt modele independente de context și se utilizează modele GMM, algoritmul Viterbi produce un număr mare de inserări. Din acest motiv se propune o modificare a acestui algoritm, introducând un parametru empiric pentru a putea controla erorile de tip inserare. Acest parametru empiric poate fi privit în mod indirect ca și un model simplist pentru modelarea duratei fonemelor. Rezultatele noastre obținute cu modele mai simple (cu un număr redus de parametri) se pot compara cu cele obținute cu modele HMM. Rezultatele experimentale sunt prezentate atât pentru determinarea parametrului empiric, cât și pentru recunoașterea fonemelor.

Capitolul 6 investighează problema modelării duratei fonemelor. Sunt prezentate rezultate statistice obținute pentru durata fonemelor pentru cele două baze de date și sunt propuse o serie de modele simple pentru modelarea duratei.

Capitolul 7 prezintă problema recunoașterii cuvintelor izolate și sunt prezentate două modalități pentru rezolvarea acestei probleme. Prima modalitate, cea

a modelării cuvintelor doar cu mixturi gaussiene, abordare originală și simplistă, dar care dă rezultate excepționale pentru un număr redus de cuvinte suficient de diferite, cum ar fi cifrele. A doua modalitate este cea generală, modelele de cuvinte sunt formate din modele de foneme. În acest caz trebuie rezolvată organizarea vocabularului, pentru care am utilizat varianta arborescentă, structura de date trie. Capitolul se termină cu prezentarea rezultatelor experimentale.

Ultimul capitol al acestei teze, capitolul 8, prezintă concluziile generale și direcțiile de dezvoltare.

## 1.2 Contribuții originale

- Am studiat conținutul fonetic al modelelor de vorbitori. Am studiat distribuția categoriilor fonetice, respectiv a fonemelor în componentele acestor modele. Am arătat că aceste componente nu sunt pure din punct de vedere fonetic.
- Ordonarea categoriilor fonetice pe baza puterii lor de discriminare a vorbitorilor. Am arătat că vocalele și nazalele au puterea discriminativă cea mai mare și că utilizând doar vocalele unui vorbitor, acesta poate fi identificat cu mare acuratețe.
- Construirea modelelor de vorbitori structurate fonetic. Am arătat că structurând datele unui vorbitor pe bază de categorii fonetice și creând submodelele acestor categorii fonetice, modelul obținut are performanțe similare cu modelele obișnuite.
- Am demonstrat în mod experimental că modelele GMM au performanțe similară cu modelele HMM pentru problema clasificării fonemelor. Aceasta demonstrează faptul că probabilitățile de tranziție din cadrul modelelor HMM pot fi neglijate.
- Am propus un clasificator cu două niveluri pentru problema clasificării fonemelor, care este mult mai eficient din punctul de vedere al calculelor cu o acuratețe ușor scăzută față de clasificarea obișnuită.
- Am propus o modificare a algoritmului Viterbi, pentru a putea fi folosit și pentru modele GMM. Modelele de foneme GMM cu algoritmul Viterbi modificat furnizează rezultate la fel de bune ca și modelele HMM, dar cu reducerea semnificativă a calculelor.

- Am prezentat un studiu referitor la duratele fonemelor pentru două baze de date. Pe baza acestor studii am elaborat o nouă modificare a algoritmului Viterbi pentru a putea include și modelarea duratei fonemelor.
- În final am realizat un sistem cu componente software necesare pentru recunoașterea vorbitorului și a vorbirii.

## Capitolul 2

# Bazele teoretice

### 2.1 Introducere

În acest capitol prezentăm bazele teoretice necesare atât pentru recunoașterea vorbitorului cât și pentru recunoașterea vorbirii. Vom începe cu descrierea extragerii caracteristicilor utilizată în ambele sisteme, specificând particularitățile fiecărui sistem. Continuăm cu prezentarea noțiunilor fonetice strict necesare pentru elaborarea acestei lucrări. Următoarele două secțiuni prezintă metode matematice ca, tehnica grupărilor, respectiv teoria modelelor Markov ascunse utilizate în implementarea celor două sisteme. Penultima secțiune prezintă trei noțiuni de bază ale sistemelor pentru recunoașterea vorbirii continue: modelele de limbaj, măsurarea erorii, precum și tehniciile de decodare. În final sunt prezentate cele două baze de date de sunet utilizate în experimente.

Scopul unui sistem de recunoașterea vorbirii este extragerea conținutului lexical din semnalul vocal. Aceasta se realizează în mai mulți pași. Detalii cu privire la acest subiect se găsesc în [26, 40, 67, 79]. Componentele unui asemenea sistem sunt următoarele:

- *Extragerea caracteristicilor:* În această parte a sistemului din forma de undă a semnalului vocal se extrag vectorii acustici. Acești vectori trebuie să conțină informații pe baza cărora se pot distinge unitățile lingvistice cum ar fi fonemele. Aceste informații ar trebui să fie robuste, adică pentru aceeași unitate lingvistică să fie asemănătoare indiferent de vorbitor sau de condițiile de înregistrare. Desigur, asemenea caracteristici nu sunt ușor de găsit, semnalul fiind influențat atât de diferite surse de zgomot, cât și de vorbitor.

- *Modelarea acustică:* Vectorii acustici se utilizează pentru a crea câte un model pentru fiecare unitate lingvistică (de obicei foneme). Acest pas se numește și pasul de antrenare. Scopul antrenării este estimarea parametrilor modelelor astfel încât aceste modele să reprezinte cât mai adecvat datele de antrenare și de a minimiza eroarea de recunoaștere în cazul datelor neprezentate sistemului în faza de antrenare. Există două abordări de bază pentru modelarea acustică: una este utilizarea modelelor Markov ascunse utilizând mixturi gaussiene pentru estimarea probabilităților de emisie a stărilor [68](HMM/GMM) și cealaltă este utilizarea rețelelor neuronale în loc de mixturi gaussiene în modelul Markov ascuns [18](HMM/ANN).
- *Decodare:* În această fază se determină secvența de cuvinte cea mai probabilă pentru secvența de vectori acustici extrași din semnalul vocal. Decodarea este de fapt o căutare într-un spațiu acustic utilizând mai multe surse de cunoștințe, cum ar fi probabilități ale secvențelor de vectori acustici în modelele fonetice, constrângeri impuse de vocabular, probabilități de secvențe de cuvinte, care se numește model limbaj. În cazul sistemelor HMM/GMM decodarea se realizează de obicei cu algoritmul Viterbi sau o formă modificată a acestuia.

Recunoașterea vorbitorului este procesul în care se recunoaște vorbitorul pe baza unui material vocal furnizat de către vorbitor. Scopul este extragerea acelor informații din semnalul vocal, care caracterizează vorbitorul și nu conținutul lingvistic al semnalului. Fiecare vorbitor este caracterizat prin rata vorbirii, prin modul de pronunțare a diferențelor unități fonetice etc. Aceste particularități caracteristice vorbitorului sunt dificil de extras. Componentele standard a unui sistem de recunoașterea vorbitorului sunt:

- *Extragerea caracteristicilor:* Se extrag informații pe baza căror vorbitorii pot fi identificați, discriminați față de ceilalți vorbitori ai sistemului.
- *Calculul scorului:* În acest modul se verifică ipoteza că un semnal vocal este produs de un anumit vorbitor. În funcție de tipul problemei (identificare sau verificare) se va folosi un model adițional numit și model fundal [16].
- *Decizie:* În cazul sistemelor de verificarea vorbitorului vorbitorul este acceptat sau respins, iar în cazul sistemelor de identificare se va determina identitatea vorbitorului pe baza similarităților față de vorbitorii cunoscuți

de către sistem. Problema de verificare este de fapt o decizie binară, iar problema de identificare este o problemă de clasificare.

Dacă vorbitorul poate rosti orice text, când solicită identificarea sau verificarea, sistemul se numește independent de text, iar dacă conținutul lingvistic ce trebuie reprodus este determinat în prealabil, sistemul se numește dependent de text. În cel de-al doilea caz se utilizează atât recunoașterea vorbitorului cât și recunoașterea materialului vorbit. Următoarele articole furnizează detalii despre acest subiect [16, 21, 36].

## 2.2 Elemente de fonetică

### 2.2.1 Semnalul vocal

Semnalul vocal este compus dintr-o succesiune de sunete, al căror aranjament este guvernat de legi fonetice [75]. Unda acustică vocală este generată prin mișcări fiziolegice voluntare ale următoarelor elemente anatomiche: coarde vocale, tract vocal, limbă, cavitatea nazală, buze, văl palatin.

Aerul expulzat din plămâni trece prin trahee, iar apoi printre coardele vocale. La generarea sunetelor sonore aerul împins spre buze duce la deschiderea și închiderea coardelor vocale într-un ritm dependent de presiunea aerului în trahee și de ajustările fiziolegice, ce au loc. Cu cât tensiunea în coarde este mai mare, cu atât este mai mare și frecvența fundamentală a semnalului vocal. Deschiderea dintre coardele vocale se numește glotă.

Tractul vocal este un tub acustic neuniform, ce se întinde de la glotă la buze, forma sa variind în timp. Componentele anatomiche, ce produc aceste variații, sunt buzele, maxilarele, limba și vălul palatin. În timpul generării sunetelor nenazale vălul palatin separă tractul vocal de cavitatea nazală, care constituie un tub acustic adițional pentru generarea și transmiterea sunetelor nazale: aerul trece prin glotă, iar tensiunea în coardele vocale este astfel ajustată, încât acestea vor vibra într-o oscilație de relaxare, generând pulsuri cvasiperiodice de aer, care excita tractul vocal.

Sunetele nesonore sunt generate prin crearea unei constrictii pe tractul vocal și forțarea aerului prin ea cu o viteză suficient de mare pentru a se genera un zgomot cu spectru larg, care va excita tractul vocal. Constrictia, împreună cu vibrația coardelor vocale, produce fricativele sonore.

Sunetele plozive se obțin prin creșterea presiunii aerului în cavitatea bucală, urmată de eliberarea sa bruscă printre buze.

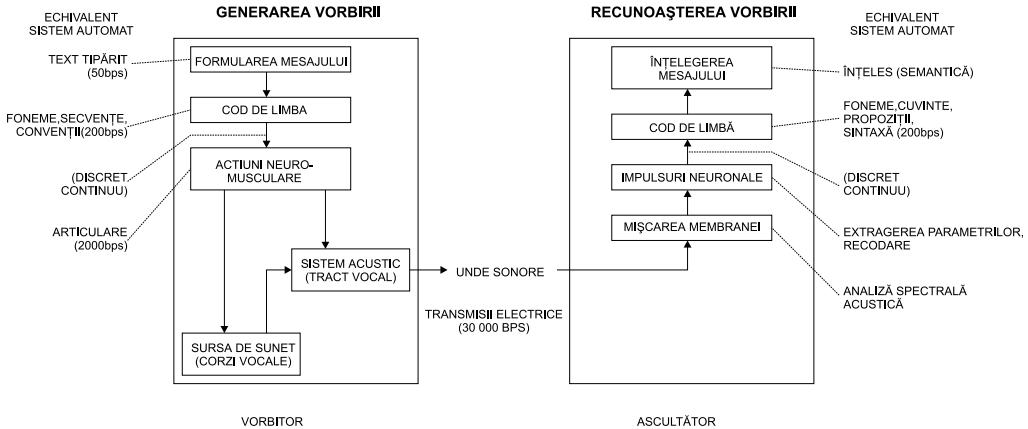


Figura 2.1: Schema producerii și perceprii sunetului - varianta 1 (după [67])

### 2.2.2 Producere și percepere

Figura 2.1 ne arată o diagramă schematică privind producerea și perceperea vorbirii la ființe umane. Procesul de producere (generarea vorbirii) începe atunci, când vorbitorul formulează un mesaj (în minte - memorie) pe care dorește să-l transmită cuiva care îl ascultă, prin vorbire. În cazul mașinii, acestui pas îi corespunde crearea textului, care conține cuvintele mesajului. Pasul următor este conversia mesajului într-un cod al limbajului. Acestei pas îi corespunde conversia textului scris al mesajului într-o mulțime de foneme corespunzătoare sunetelor care alcătuiesc cuvintele. Odată ce a fost selectat limbajul, vorbitorul trebuie să execute o serie de comenzi neuromusculare pentru a pune în vibrație coardele vocale și pentru a forma tractul vocal pentru a produce secvența de sunete corespunzătoare, astfel obținându-se rezultatul: semnalul vocal.

Odată ce semnalul a fost generat și propagat către ascultător, începe procesul de percepere a semnalului. Primul pas în perceperea semnalului este la nivelul urechilor: partea interioară a acesteia, membrana, realizează o analiză spectrală a semnalului intrat. Un proces neuronal realizează conversia semnalului spectral obținut la ieșirea membranei într-un semnal de activare al nervului auzului. Acest pas se mai numește și pasul extragerii trăsăturilor. După acest pas activitatea neuronală de-a lungul nervului auzului, cu ajutorul creierului, este convertită într-un cod al limbajului, obținându-se în final semnificația mesajului.

O altă reprezentare a procesului producere-percepere este dată de figura 2.2. Aici putem urmări stadiile procesului, precum și viteza de propagare a semnalului între stadii. Viteza formulării mesajului este cel mai lent, aproximativ 50 bps

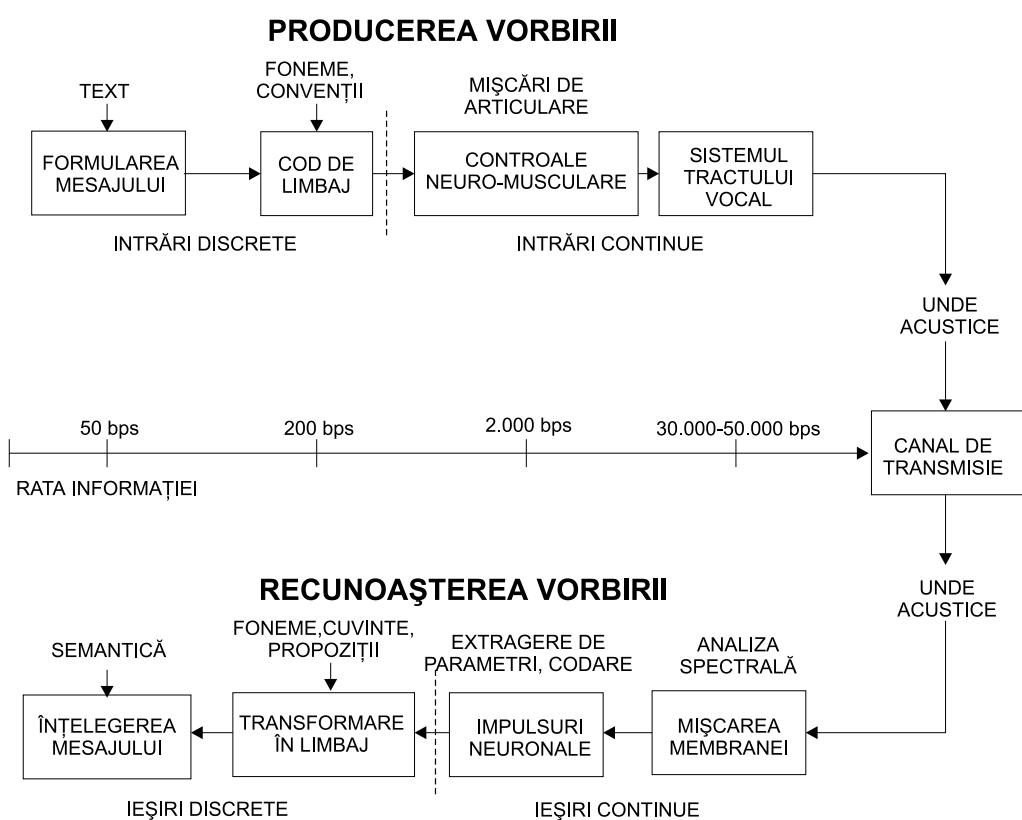


Figura 2.2: Schema producerii și a percepției sunetului - varianta 2 (după [67])

corespunzând cu 8 sunete pe secundă, unde fiecare sunet este unul dintre cele 50 de simboluri. După această codificare viteza crește la 200 bps. În stadiul următor la conversia informației în semnal viteza atinge rata de 2000 bps la nivelul neuromuscular și 30000-50000 bps la nivelul semnalului acustic.

### 2.2.3 Transcriere fonetică

Procesul de translatare a unei rostiri într-o secvență de simboluri, numite foneme, se numește transcriere fonetică. Stabilirea unui alfabet fonetic a stârnit multe discuții. În cazul limbii engleze utilizarea reprezentării ortografice cu ajutorul alfabetului roman nu este corespunzătoare. În anul 1888, un grup de foneticieni din Europa, au publicat alfabetul fonetic internațional (**IPA** - *International Phonetic Alphabet*), pentru a facilita și standardiza transcrițiile fonetice. Acest alfabet fonetic are suficiente simboluri pentru a putea descrie toate limbile lumii.

IPA este foarte comod pentru transcriere cu scris de mână, dar simbolurile nu pot fi scrise cu tastatura unui calculator. Din această cauză s-a dezvoltat un nou alfabet fonetic, alfabetul **ARPAbet**. Acest alfabet fonetic s-a dezvoltat sub auspiciul ARPA (Advanced Research Projects Agency) din SUA. Există două versiuni ale alfabetului ARPAbet, prima versiune utilizând simboluri de o literă și o altă versiune utilizând doar secvențe de litere. Tabelul 1 prezintă a doua versiune a alfabetului.

### 2.2.4 Clasificarea fonemelor

Există mai multe criterii, pe baza cărora fonemele pot fi clasificate. Fonemele pot fi grupate pe baza următoarelor proprietăți:

- Forma de undă
- Caracteristica de frecvență
- Modul de articulare
- Locul de articulare
- Tipul excitației
- Staționaritate

Simbol	Exemplu	Simbol	Exemplu
iy	heed	v	vice
ih	hid	th	thing
ey	hayed	dh	then
eh	head	s	so
ae	had	z	zebra
aa	hod	sh	show
ao	hawed	zh	measure
ow	hoed	hh	help
uh	hood	m	mom
uw	who'd	n	noon
er	heard	nx	sing
ax	ago	l	love
ah	mud	el	cattle
ay	hide	em	some
aw	how'd	en	son
oy	boy	dx	batter
ix	roses	q	'glottal stop'
p	pea	w	want
b	bat	y	yard
t	tea	r	race
d	deep	ch	church
k	kick	jh	just
g	go	wh	when
f	five		

Tabelul 1: Alfabetul fonetic ARPAbet

Un fonem se numește staționar sau continuu în funcție de configurația tractului vocal pe durata pronunțării. Un fonem va fi necontinuu dacă este necesară schimbarea configurației tractului vocal. Tabelul 2 prezintă acest tip de clasificare.

Continue		Necontinue	
Vocale	iy, ih, e y, eh, ae, er, ax, ah, uw, uh, ow, ao, aa	Vocale (diftongi)	ay, aw, oy
Semivocale	hh	Semivocale	r, l, w, y
Fricative	v, dh, z, zh, f, th, s, sh	Plozive	b, d, g, p, t, k
Africate	jh, ch		
Nazale	m, n, nx		

Tabelul 2: Clasificarea fonemelor limbii engleze (după [26])

Vocalele, fricativele, africatele și nazalele aparțin categoriei de foneme continue, iar diftongii, semivocalele și plozivele formează clasa fonemelor necontinue. Vocalele se află printre fonemele cu cele mai mari amplitudini. Durata vocalelor este variabilă, în mod tipic între 40 și 400 milisecunde. În funcție de poziția limbii pe durata rostirii, fonemele se pot clasifica în vocale față (/iy/, /ih/, /ey/, /eh/, /ae/), centrale (/er/, /ax/, /ah/) și spate (/uw/, /uh/, /ow/, /ao/, /aa/).

Formele de undă ale vocalelor ne arată și structura cvaziperiodică ale acestora. Vocalele se disting de celelalte foneme prin locația formanților. De obicei primii trei formanți sunt suficienți. Pentru o voce bărbătească formanții apar în jurul frecvențelor 500, 1500, 2500, 3500 Hz. Primii doi formanți  $F_1$  și  $F_2$  sunt definiți de forma tractului vocal, care este specific fiecărui vorbitor în parte.  $F_3$  este specific doar câtorva sunete. Începând cu formantul  $F_4$  locul acestor formanți nu este influențat de articulațiile produse în timpul vorbirii.

Vocalele sunt sunete sonore, teoretic având forma tractului vocal constantă pe durata pronunțării. O definiție rezonabilă pentru un diftong ar putea fi următoarea: o silabă vocalică conținând două voci. Totuși un diftong nu este pronunția consecutivă a două voci, deoarece prima vocală se pronunță mai lung decât cea de-a doua vocală. Diftongii universal acceptați ar fi următoarele: /ay/ - pie ; /aw/ - out; /oy/ - toy. Deci sunetele vocalice, care nu pot fi pronunțate fără mișcări articulatorii se numesc diftongi.

Semivocalele sunt similare cu vocalele. În limba engleză avem următoarele semivocale: /l/ - lawn; /r/ - run; /w/-wet; /y/-yam.

Pronunțarea consoanelor prezintă mult mai multe constrângeri, decât cea a vocalelor. Anumite consoane pot necesita mișcări articulatorii foarte precise. Sunetele, a căror pronunție necesită trecerea fluxului de aer prin cavitatea orală

și suspendarea totală a fluxului pentru o perioadă scurtă, se numesc consoane întrerupte sau plozive.

Fricativele sunt produse prin excitarea tractului vocal cu un flux de aer constant, care devine turbulent într-un anumit punct. Constrângerea în tractul vocal va rezulta o excitație nesonoră. Totuși anumite fricative pot avea și o componentă sonoră, caz în care obținem excitație mixtă. Astfel fricativele, care nu au această componentă sonoră, se numesc fricative nesonore, iar cele care au excitație mixtă, se numesc fricative sonore. Fricativele sonore sunt următoarele: /f/ - free; /th/ - thick; /s/ - cease; /sh/ - mesh. Fricativele nesonore corespunzătoare celor sonore sunt următoarele: /v/ - vice; /dh/- then; /z/ - zephyr; /zh/ - measure.

În mod similar cu diftongii și africatele se formează prin trecerea de la o plozivă spre o fricativă. În limba engleză sunt două africate: /ch/ - change; /jh/ - jam.

Africata nesonoră /ch/ se produce prin pronunțarea plozivei nesonore /t/, urmată de o tranziție spre o fricativă nesonoră /sh/. Africata sonoră /jh/ se produce prin pronunțarea plozivei sonore /d/, urmată de fricativa sonoră /zh/.

Plozivele /b/, /d/, /g/, /p/, /t/, /k/ se produc prin acumularea unei presiuni undeva în tractul vocal și eliberarea bruscă a acesteia printre buze. În cazul plozivelor nesonore /p/, /t/, /k/ eliberarea bruscă a presiunii este urmată de o fricație nesonoră. Această fricație nesonoră este urmată de o aspirație. Plozivele nesonore sunt caracterizate prin perioade mai lungi de fricație. Intervalul de timp până la eliberarea presiunii se numește intervalul de închidere.

Plozivele sonore sunt similare plozivelor nesonore, cu diferența că acestea din urmă includ și vibrația corzii vocale. Caracteristicile plozivelor variază în funcție de poziția lor în cuvânt sau propoziție. Intervalul de închidere poate să lipsească, dacă apare la sfârșitul silabelor. Plozivele se modifică și în cazul, când apar între vocale. Deoarece plozivele sunt foneme tranzitorii, proprietățile lor sunt extrem de puternic influențate de fonemele care le preced, respectiv le urmează. Tocmai din această cauză forma de undă nu este suficientă pentru distingerea lor. Spectrogramele ne oferă mai multe caracteristici, prin care putem distinge plozivele.

Consoanele nazale /m/, /n/, /nx/ sunt produse prin excitarea cu o undă glotală a cavității nazale deschise și a cavității orale închise. Această categorie de foneme prezintă anumite similarități cu vocalele, dar sunt caracterizate de un nivel energetic mai scăzut datorită faptului că, cavitarea nazală nu poate emite sunetul la fel de puternic ca și cavitarea orală. În funcție de locul unde se produce închiderea cavității orale distingem 3 tipuri de nazale: /m/ - more (nazal frontal); /n/ - noon (nazal central); /nx/ - sing (nazal posterior)

Cu toate că nazalele sunt singurele care incorporează cavitatea nazală în scopul producerii structurii lor de frecvențe rezonante, există și alte foneme care devin nazalizate - fonemele care le preced sau urmează după nazale. În primul rând categoria vocalelor este cea mai afectată de acest fenomen. Nazalizarea produce foneme cu o lățime de bandă mai largă a primului formant ( $F_1$ ) [39]. Gradul în care nazalizarea afectează spectrul, depinde de modul de cuplare între cele două cavități.

## 2.3 Extragerea caracteristicilor

Din semnalul vocal putem extrage atât caracteristici acustice de nivel scăzut cum ar fi ampliitudinea spectrală cât și caracteristici de nivel înalt cum ar fi accentul, dialectul, stilul de vorbire sau modul de pronunțare a unor cuvinte specifice. Caracteristicile de nivel înalt sunt utilizate de oameni pentru recunoașterea vorbitorului în cauză. Extragerea caracteristicilor de nivel înalt sunt greu de cuantificat, de aceea majoritatea sistemelor pentru recunoașterea vorbitorilor utilizează doar caracteristicile de nivel scăzut. În continuare prezentăm extragerea caracteristicilor care se utilizează atât în sisteme de recunoașterea vorbirii cât și în sisteme de recunoașterea vorbitorilor.

Extragerea caracteristicilor se realizează în următorii pași: preprocesare, analiză cepstrală și postprocesare.

### 2.3.1 Preprocesare

Preprocesarea are ca scop amplificarea frecvențelor înalte din spectru. Aceasta se realizează prin utilizarea unui filtru de ordin I. (FIR), care are următoarea formă

$$s_p(n) = s(n) - as(n-1) \quad (2.1)$$

unde  $a$  se numește coeficient de preaccentuare și are valori cuprinse între  $0.9 \leq a \leq 1$ . În sistemul nostru această valoarea utilizată este 0.97. Semnalul vocal nu este un semnal staționar, dar pe perioade scurte de timp este cvazistaționar. Ferestruirea semnalului are ca scop descompunerea semnalului pe bucăți pe care acesta este aproape staționar. Lungimea ferestrei este de obicei între 10 și 30 de milisecunde. Pentru recunoașterea vorbitorului am utilizat ferestre de 32 de milisecunde iar pentru recunoașterea vorbirii ferestre de 20 de milisecunde. Aceste ferestre se aplică la fiecare 10 milisecunde până la sfârșitul semnalului. Rezultatul ferestruirii este obținerea unei secvențe de bucăți de semnal numite și cadre. Cea

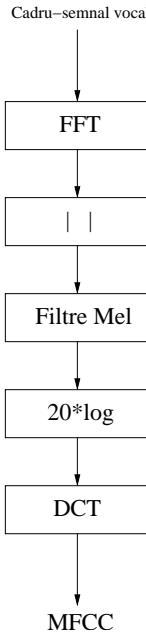


Figura 2.3: Reprezentarea modulară a extragerii coeficienților Mel cepstrali

mai utilizată funcție de fereastră este fereastra de tip Hamming dată de ecuația următoare

$$w(n) = \begin{cases} 0.54 - 0.46 \cos\left(\frac{2\pi n}{N-1}\right) & n = 0, 1, \dots, N-1 \\ 0 & altfel \end{cases} \quad (2.2)$$

unde  $N$  este numărul eșantioanelor din cadru.

### 2.3.2 Analiză cepstrală

#### 2.3.2.1 Coeficienții Mel-cepstrali (MFCC)

Figura 2.3 ne arată pașii necesari pentru extragerea coeficienților Mel cepstrali.

După faza de preprocesare pentru fiecare cadru se aplică transformata Fourier rapidă (FFT) pe un număr de puncte  $N'$ , putere a lui 2 și  $N' \geq N$ . Din spectrul de putere obținut se utilizează doar jumătate, deoarece spectrul este simetric. Spectrul astfel obținut trebuie netezit pe de o parte fiindcă nu ne interesează fluctuațiile mici ale spectrului, pe de altă parte din considerante de reducere a dimensionalității. Pentru a realiza această netezire și a obține anvelopa spectrală, se aplică un banc de filtre asupra spectrului. Un banc de filtre este o serie de filtre trece-bandă cu care se multiplică spectrul, unul după altul, pentru a obține câte

o valoare medie pentru fiecare bandă de frecvență. Un banc de filtre se definește prin forma filtrelor precum și prin așezarea acestora pe scara de frecvențe. În practică așezarea acestor filtre se face pe baza scării Mel, scară determinată a.î. să fie cât mai apropiat de modul de percepere ale frecvențelor de către urechea umană. Localizarea frecvențelor centrale este dată de ecuația

$$f_{MEL} = 1000 \cdot \frac{\log(1 + f_{LIN}/1000)}{\log 2} \quad (2.3)$$

În final se consideră logaritmul acestei anvelope spectrale și se multiplică cu 20 pentru a obține anvelopa spectrală în decibeli (dB).

Ultima transformare aplicată este transformata cosinus discretă (DCT).

$$c_n = \sum_{k=1}^K S_k \cdot \cos \left[ n \left( k - \frac{1}{2} \right) \frac{\pi}{K} \right], \quad n = 0, 1, 2, \dots, D,$$

unde  $K$  este numărul filtrelor din bancul de filtre,  $S_k$  sunt coeficienții spectrali obținuți după aplicarea filtrelor, iar  $D$  este numărul coeficienților cepstrali calculați ( $D \leq K$ ).

### 2.3.2.2 Coeficienții cepstrali de predicție liniară (LPCC)

Analiza prin predicție liniară (LPC) este bazată pe un model liniar de producerea vorbirii. Este un model de autoregresie care este descrisă foarte detaliat în cărțile [26, 67].

Producerea vorbirii poate fi descrisă prin combinarea următoarelor patru elemente:

- sursa glotală - văzută ca un tren de impulsuri pentru sunete sonore sau ca un zgomot alb pentru sunete nesonore
- tractul vocal
- cavitatea nazală
- buzele

Fiecare dintre acestea poate fi reprezentat printr-un filtru, un filtru trece-jos pentru sursa glotală, un filtru autoregresiv pentru tractul vocal, un filtru ARMA (Auto Regressive Moving Average) pentru tractul nazal și un filtru (MA -moving average) pentru buze. În final tot sistemul de producere a vorbirii poate fi privit ca

un filtru ARMA, iar modelele uzuale modelează vorbirea cu o variantă simplificată a filtrului ARMA, cu un filtru AR.

Principiul predicției liniare este estimarea parametrului unui filtru AR pe o porțiune ferestruată a unui semnal vocal. Fereastra este mutată periodic și se calculează parametru filtrului pe următorul segment vocal. Pentru fiecare fereastră coeficienții estimări se numesc coeficienți de predicție liniară (LPC). Coeficienții cepstrali se calculează în mod direct din coeficienții de predicție liniară folosind următoarea formulă

$$\begin{cases} c_0 = \ln \sigma^2 \\ c_m = a_m + \sum_{k=1}^{m-1} \left( \frac{k}{m} \right) c_k a_{m-k}, \quad 1 \leq m \leq p \\ c_m = \sum_{k=1}^{m-1} \left( \frac{k}{m} \right) c_k a_{m-k}, \quad p < m \end{cases} \quad (2.4)$$

unde  $\sigma^2$  este termenul de câștig din modelul LPC,  $a_m$  sunt coeficienții de predicție liniară iar  $p$  este ordinea analizei de predicție liniară.

### 2.3.3 Postprocesare

În faza de postprocesare se poate scădea din fiecare vector cepstral vectorul medie a vectorilor cepstrali (CMS). Această fază se aplică obligatoriu în cazul când înregistrările au fost făcute de pe linie telefonică.

Pentru a încorpora și variația coeficienților de-a lungul timpului, se calculează derivatele de ordin întâi și doi. În mod clasic acești parametri se numesc coeficienți  $\Delta$  și  $\Delta\Delta$  și se calculează cu ajutorul următoarei formule

$$\Delta C_m(t) = \frac{\sum_{k=-K}^K k C_m(t+k)}{\sum_{k=-K}^K k^2} \quad (2.5)$$

$K$  este numărul de vectori vecini care se consideră, în aplicație am utilizat  $K=2$ .

## 2.4 Tehnici de grupare

### 2.4.1 Introducere

În multe aplicații de recunoașterea formelor este foarte dificilă sau costisitoare, chiar imposibilă etichetarea datelor de antrenare cu eticheta categoriei corespunzătoare. Clasificarea nesupervizată se referă la situații în care obiectivul este

determinarea regiunilor de decizie folosind date neetichetate. Etichetele categoriilor și alte informații legate de sursa datelor influențează doar interpretarea rezultatelor grupării nu și formarea acestora.

Clasificarea nesupervizată sau gruparea este o problemă foarte dificilă, deoarece datele pot forma grupe cu mărimi și forme diferite. Mai mult, și numărul grupelor depinde de rezoluția, cu care vizualizăm datele. În literatură au fost propuse următoarele definiții pentru noțiunea de grupă:

- forme care aparțin unei grupe sunt mai apropiate de forme din aceeași grupă și mai îndepărte de formele aflate în alte grupe (*grupe omogene*)
- grupele sunt formate din noruri de puncte cu o densitate ridicată și sunt separate de alte grupe prin puncte izolate sau zone având o densitate scăzută (*grupele sunt bine separate*)

Simplitatea, consistența și naturalețea cu care algoritmii de grupare organizează datele au condus la utilizarea acestor tehnici în aplicații ca mineritul datelor (data mining), regăsirea informației, segmentarea imaginilor, compresia și codarea semnalelor și machine learning. O prezentare amănunțită a algoritmilor de grupare găsim în cărțile [29, 30, 35]. Majoritatea algoritmilor se bazează pe două tehnici de bază populare:

- grupare partițională iterativă bazată pe eroarea pătratică (square error partitional algorithm)
- grupare ierarhică aglomerativă

Tehnicile ierarhice organizează datele în grupe suprapuse, care pot fi afișate în forma unei dendograme sau arbori. Cele partiționale bazate pe eroarea pătratică încearcă să găsească acea partiție, care minimizează dispersia intra-grupă și maximizează dispersia inter-grupe. Algoritmul, care își propune găsirea soluției optime, trebuie să examineze toate partițiile posibile, care se pot obține prin gruparea a  $n$  vectori  $d$ -dimensionale în  $K$  grupe. Un asemenea algoritm ar avea o complexitate exponențială, total inutilizabil în practică. Pentru evitarea complexității exponențiale au fost propuse o serie de euristici care reduc timpul de căutare, negarantând însă optimalitatea soluției.

Tehnicile de grupare partiționale se utilizează mai frecvent, decât cele ierarhice, de aceea vom începe prezentarea cu această tehnică de grupare. Studiile recente au adus la iveală următoarele concluzii:

- orice algoritm de grupare va produce un set de grupe ca și rezultat, indiferent dacă în date există sau nu tendință de grupare
- nu există algoritmul „cel mai bun” de grupare; se recomandă testarea a mai multor algoritmi de grupare
- extragerea caracteristicilor dintr-o colecție de date precum și tehniciile de normalizare sunt la fel de importante ca și alegerea strategiei de grupare

Gruparea partitioană se poate formula în modul următor [41]:

Se dau  $n$  puncte într-un spațiu  $d$ -dimensional. Să se determine gruparea acestor puncte în  $K$  grupe astfel, încât punctele aparținând aceleiași grupe să fie mai apropiate de restul punctelor din aceeași grupă, decât de restul punctelor din afara grupei. Specificarea numărului de grupe este optională, însă este necesară adoptarea unui criteriu de grupare. Criteriul de grupare poate fi un criteriu global sau unul local. Un criteriu global, cum ar fi eroarea minimă pătratică, reprezintă grupele cu câte un vector prototip și clasifică punctele pe baza similarității cu aceste prototipuri. Un criteriu local formează grupele prin identificarea regiunilor de densitate mare sau prin adăugarea punctelor împreună cu cei  $k$ -vecini cei mai apropiati la aceeași grupă.

Tabelul 2.1 prezintă algoritmii cei mai importanți de grupare. La alegerea unui algoritm de grupare trebuie să considerăm următoarele aspecte:

1. Ce definește o grupă? Există un prototip reprezentativ pentru fiecare grupă?
2. Cum definim apartenența la o grupă? Ce metrică utilizăm?
3. Câte grupe există? Este necesară fixarea numărului de grupe?
4. Care este capacitatea de generalizare a modelului obținut? Prin ce metodă se determină apartenența unui punct nou la grupele deja existente?

#### 2.4.2 Algoritmul K-means

Se dau punctele  $X = \{x_1, x_2, \dots, x_n\}$ , unde  $x_i \in R^d$  formează datele de intrare ale algoritmului de grupare. Obiectivul este să se găsească acea partiție, care pentru un  $K$  fixat (număr de partiții) minimizează eroarea pătratică. Notăm grupele cu  $\{C_1, C_2, \dots, C_K\}$ , fiecare având  $n_k$  puncte, a. i. fiecare  $x_i$  să aparțină unei singure grupe  $\sum_{k=1}^K n_k = n$ .

Algoritm	Descriere	Observații
K-means	Identifică grupele hipersferice; Grupele hiperelipsoidale se pot determina prin utilizarea distanței Mahalanobis	K se fixează dinante; se fixează și centrele inițiale ale grupelor
Fuzzy K-means	Similar cu K-means Fie căru punct i se atașează un grad de apartenență	Se fixează: K, centrele inițiale ale grupelor și funcția de apartenență
Arborele de acoperire minim	Grupele se formează prin ștergerea muchiilor inconsistente	Se fixează definiția muchiei inconsistente
Mutual Neighbourhood	Se calculează valoarea mutuală de vecinătate pentru fiecare pereche de puncte	Trebuie specificată adâncimea vecinătății
Single-Link	Un algoritm ierarhic de grupare	Detalii: 2.4.7
Complete-Link	Un algoritm ierarhic de grupare	Detalii: 2.4.7
Decompoziție în mixturi	Se consideră că punctele sunt generate de $K$ distribuții	K se fixează, însă se pot defini criterii pentru estimarea acestui parametru

Tabela 2.1: Algoritmi de grupare (După [41])

---

**Algoritm 1** Algoritmul K-means

---

Pas 1. Se selectează o partiție inițială cu  $K$  grupe

Pas 2. Se generează o nouă partiție atașând fiecare punct  $x_i$  grupei celei mai apropiate

Pas 3. Se recalculează centrele de greutate ale grupelor cu formula 2.6.

Pas 4. Se repetă Pas 2 și Pas 3. până când se atinge valoarea optimală a funcției criteriu

---

Centrul de greutate al grupei  $C_k$  se definește prin

$$\mu^{(k)} = \left( \frac{1}{n_k} \sum_{i=1}^{n_k} x_i^{(k)} \right) \quad (2.6)$$

unde  $x_i^{(k)}$  este al  $i$ -lea punct aparținând grupei  $C_k$ . Definim eroarea pătratică pentru  $C_k$  în modul următor

$$e_k^2 = \sum_{i=1}^{n_k} d^2 \left( x_i^{(k)}, \mu^{(k)} \right) \quad (2.7)$$

unde  $d \left( x_i^{(k)}, \mu^{(k)} \right)$  este o funcție distanță. Această eroare se numește și dispersia intra-grupă. Pentru toate grupele eroarea se definește prin

$$E_K^2 = \sum_{k=1}^K e_k^2. \quad (2.8)$$

Obiectivul grupării este găsirea acelei partiții conținând  $K$  grupe, pentru care 2.8 este minimă.

Algoritmul K-means se poate formula în modul arătat în (1).

Pentru oprirea algoritmului (1) se pot utiliza diverse criterii, dintre care amintim următoarele:

- Nu se modifică conținutul grupelor (stabilizare)
- S-a depășit numărul maxim de iterări
- Modificarea erorii totale este sub o valoare de prag  $1 - \frac{E_K^{(n+1)}}{E_K^{(n)}} < T$

Funcția de distanță utilizată în formula 2.7 influențează puternic forma grupelor. Ca și funcție de distanță se utilizează cel mai des distanța Euclideană și distanța Euclideană ponderată definite cu formulele 2.9, respectiv 2.10. Dacă în

formula 2.10 se utilizează  $W = \Sigma^{-1}$ , unde  $\Sigma$  este matricea de covarianță, vorbim despre distanță Mahalanobis. Utilizând distanța euclideană se obțin grupe hipersferice, iar cu distanța Mahalanobis aceste grupe vor fi hiperelipsoidale.

$$d^2(x, y) = (x - y)^T (x - y), \quad x, y \in R^d \quad (2.9)$$

$$d^2(x, y) = (x - y)^T W (x - y), \quad x, y \in R^d \quad (2.10)$$

Algoritmul K-means formulat în (alg.1) este un algoritm foarte eficient d.p.v. computațional și produce rezultate surprinzătoare de bune dacă grupele sunt compacte, cu forme hipersferice și bine separate. Utilizând distanța Mahalanobis algoritmul este capabil să detecteze chiar și grupele de forme hiperelipsoidale. Au fost multe încercări în direcția îmbunătățirii performanțelor algoritmului K-means. Incorporarea unei funcții criteriu fuzzy a rezultat în algoritmul fuzzy K-means. O altă metodă frecvent utilizată este execuția de mai multe ori a algoritmului și reținerea acelei soluții pentru care funcția criteriu este optimă [4, 45].

#### 2.4.3 Algoritmul Fuzzy K-means

În fiecare iterație a algoritmului K-means clasic fiecare punct este atașat la o singură grupă ( $\arg \min_{i=1, K} d(x, \mu_i)$ ). Putem relaxa această condiție definind un grad de apartenență (apartenență fuzzy) la o grupă pentru fiecare punct. La origine aceste grade de apartenență corespund probabilităților  $\hat{P}(C_i|x_j, \hat{\Theta})$ , unde  $\hat{\Theta}$  este vectorul de parametri pentru funcția de apartenență. Funcția criteriu care trebuie optimizată devine

$$J_{fuzz} = \sum_{i=1}^K \sum_{j=1}^n \left[ \hat{P}(C_i|x_j, \hat{\Theta}) \right]^b d^2(x_j, \mu_i) \quad (2.11)$$

unde  $b$  este un parametru liber care se utilizează pentru gradul de amestecare al grupelor. Dacă  $b = 0$ , ecuația 2.11 este echivalentă cu 2.8. Dacă  $b > 1$ , criteriul permite ca fiecare punct să aparțină mai multor grupe. Suma probabilităților de apartenență al fiecărui punct la diferite grupe este 1

$$\sum_{i=1}^K \hat{P}(C_i|x_j) = 1, \quad j = 1, \dots, n, \quad (2.12)$$

**Algoritm 2** Algoritmul Fuzzy K-means

- 
- Pas 1. Se inițializează  $K, b, \mu_1, \mu_2, \dots, \mu_K, P(C_i|x_j), i = 1, \dots, K; j = 1, \dots, n$ ;  
Se normalizează  $\hat{P}(C_i|x_j)$  cu ecuația 2.12  
Pas 2. Se recalculează  $\mu_i$  cu ecuația 2.14  
Pas 3. Se recalculează  $P(\hat{C}_i|x_j)$  cu ecuația 2.15  
Pas 4. Se repetă Pas 2 și Pas 3 până când se stabilizează valorile  $\mu_i$  și  $\hat{P}(C_i|x_j)$
- 

unde din considerente de simplitate nu am notat dependența de  $\hat{\Theta}$ . Notând cu  $\hat{P}_j$  probabilitatea apriori  $\hat{P}(C_j)$ , soluția optimă a ecuației 2.11 se obține prin

$$\partial J_{fuz} / \partial \mu_i = 0 \quad \partial J_{fuz} / \partial \hat{P}_j = 0 \quad (2.13)$$

Acesta ne conduce la soluțiile

$$\mu_j = \frac{\sum_{j=1}^n [\hat{P}(C_i|x_j)]^b x_j}{\sum_{j=1}^n [\hat{P}(C_i|x_j)]^b} \quad (2.14)$$

și

$$\hat{P}(C_i|x_j) = \frac{(1/d_{ij})^{1/(b-1)}}{\sum_{r=1}^K (1/d_{rj})^{1/(b-1)}}, \quad d_{ij} = d^2(x_j, \mu_i) \quad (2.15)$$

În general  $J_{fuz}$  este minim când centrele grupelor  $\mu_j$  sunt aproape de punctele cu un grad mare de apartenență la grupa  $j$ . Algoritmul fuzzy K-means este dat în (Alg. 2).

În mod clar algoritmul clasic K-means formează un caz special al algoritmului fuzzy K-means, caz în care gradul de apartenență se definește în modul următor

$$P(C_i|x_j) = \begin{cases} 1 & \text{daca } d(x_j, \mu_i) < d(x_j, \mu_{i'}) \forall i' \neq i \\ 0 & \text{altfel} \end{cases} \quad (2.16)$$

Încorporarea probabilităților conduce căteodată la îmbunătățirea convergenței față de varianta clasică K-means.

#### 2.4.4 Cuantizare vectorială

Este interesant de notat că, concepte aparent diferite folosite pentru grupări parționale conduc în final la același algoritm. Se poate verifica ușor că algoritmul generalizat Lloyd pentru cuantizare vectorială dezvoltat de către inginerii din domeniul comunicațiilor este echivalent cu algoritmul K-means.

Obiectivul unui cuantizor vectorial [65] (Vector Quantisation - **VQ**) este repre-

zentarea unui set de puncte  $X \subseteq R^d$  printr-un set  $Y \subseteq R^d$ ,  $|Y| = K$ . Elementele mulțimii  $Y$  se numesc vectori de cod sau cuvinte de cod, iar mulțimea  $Y$  se numește dicționar de vectori de cod. Un cuantizor vectorial se poate reprezenta printr-o funcție  $q : X \rightarrow Y$ . Cunoașterea lui  $q$  ne permite să obținem o partiționare  $S$  a mulțimii de puncte  $X$  în  $K$  partiții ( $S_i, i = 1 \dots K$ ) sau celule.

$$S_i = \{x \in X : q(x) = y_i\}, \quad i = 1, \dots, K. \quad (2.17)$$

Eroarea de cuantizare este valoarea definită de  $d(x, q(x))$ , unde  $d$  este o funcție de distanță. Pentru  $d$  putem considera eroarea pătratică

$$d(x, y) = \sum_{i=1}^d (x_i - y_i)^2 \quad (2.18)$$

În acest caz eroare pătratică medie se definește prin

$$MSE = D(\{Y, S\}) = \frac{1}{K} \sum_{p=1}^K d(x_p, q(x_p)) = \frac{1}{K} \sum_{i=1}^K D_i \quad (2.19)$$

unde  $D_i$  este distorsiunea totală pentru celula  $i$ .

$$D_i = \sum_{n:x_n \in S_i} d(x_n, y_i).$$

Obiectivul cuantizării este determinarea unui dicționar de vectori de cod pentru care formula 2.19 are valoare minimă. Acest lucru este practic echivalent cu minimizarea ecuației 2.8 în cazul algoritmului K-means. În cazul algoritmului VQ, la intrarea algoritmului trebuie furnizat numărul vectorilor de cod  $K$ . Versiunea supervizată a algoritmului  $VQ$  se numește  $LVQ$  [49, 79] (Learning Vector Quantisation).

După ce am prezentat ideea de bază a algoritmilor de grupare, precum și cea a cuantizării, putem trage concluzia că acestea sunt concepte similare. Singura deosebire ar fi că, un algoritm de grupare are ca și scop identificarea grupărilor într-o mulțime de date multidimensionale, iar algoritmul de cuantizare nu are ca și scop principal această identificare, ci reprezentarea mulțimii de date printr-o mulțime restrânsă de date. Practic ambele algoritmi ne pot conduce la același rezultat în cazul când minimizează aceeași funcție criteriu. Diferența dintre algoritmul de grupare și VQ este arătată pe figura 2.4. Imaginea din colțul stânga-sus reprezintă mulțimea originală de date. Următoare imagine, din colțul dreapta-sus, ne

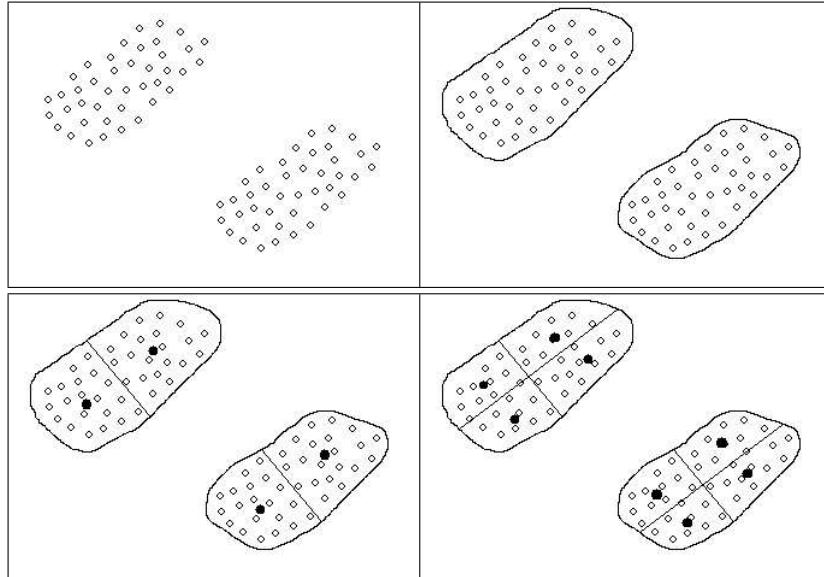


Figura 2.4: Algoritm de grupare și VQ (după [65])

arată cele două grupe găsite de un algoritm de grupare. A treia și a patra imagine ne arată rezultatul unui algoritm VQ cu 4 vectori de cod respectiv cu 8 vectori de cod.

#### 2.4.5 Algoritmul LVQ

LVQ este un algoritm adaptiv de clasificare bazat pe datele de antrenare și informații despre clasele acestora [42, 49]. Cu toate că aparține algoritmilor de antrenare supervizate, LVQ utilizează tehnici nesupervizate pentru obținerea centroizilor grupelor. Arhitectura rețelelor neuronale LVQ este foarte asemănătoare cu cea a rețelelor neuronale competitive, cu excepția faptului că, fiecare ieșire a rețelei este asociată cu o clasă. Figura 2.5 prezintă un exemplu, unde dimensiunea vectorului de intrare este 2, iar spațiul vectorilor de intrare este împărțit în şase grupe făcând parte din două clase. Primele două grupe aparțin clasei 1, iar celelalte patru clasei 2. Algoritmul LVQ are doi pași. În primul pas se utilizează un algoritm nesupervizat de grupare pentru localizarea centroizilor, fără a utiliza informațiile de clasă. În pasul următor, utilizând informațiile de clasă se îmbunătățesc aceste centre de grupe. Pentru pasul întâi se poate utiliza orice algoritm de grupare nesupervizată. După ce s-au obținut grupele, acestea vor fi etichetate cu etichetele claselor. O asemenea etichetare se poate obține prin metoda votării, adică o grupă va fi etichetată cu clasa  $k$ , dacă majoritatea punctelor din acea

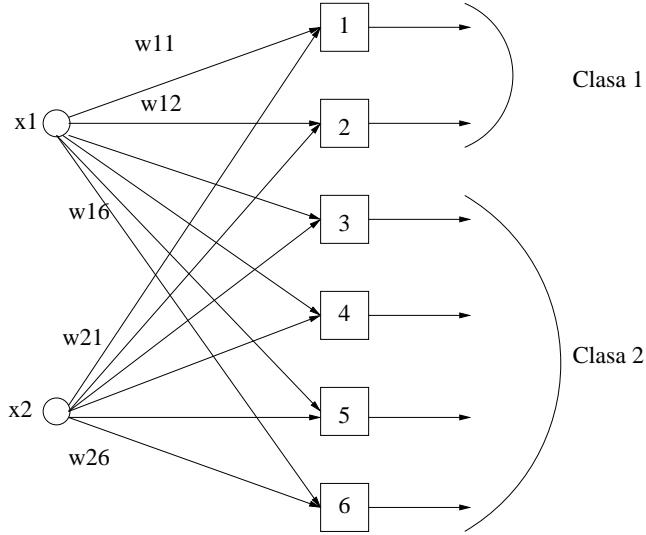


Figura 2.5: Rețea neuronală LVQ

**Algoritm 3** Algoritmul LVQ-1

---

Pas 1. Se inițializează centroizii cu un algoritm de grupare

Pas 2. Se etichetează grupele cu informațiile de clasă

Pas 3. Se selectează un punct  $x$  în mod aleator și se determină centroidul cel mai apropiat  $k = \arg \min_{i=1, M} d(x, \mu_i)$ ,  $M$  reprezintă numărul centroizilor

Pas 4. Dacă  $x$  și  $\mu_i$  aparțin aceleiași clase, se modifică  $\mu$  conform formulei

$$\Delta\mu_k = \eta(x - \mu_k)$$

altfel conform formulei

$$\Delta\mu_k = -\eta(x - \mu_k)$$

---

$\eta$  este rata de învățare, o constantă pozitivă mică care descrește la fiecare iterație.

grupă aparțin clasei  $k$ . În pasul 2 centroizii sunt îmbunătățiți prin următoarea metodă:

- Pentru orice punct de intrare  $x$  se determină centroidul cel mai apropiat  $\mu$ .
- Dacă  $x$  și  $\mu$  aparțin aceleiași clase, centroidul  $\mu$  se apropie la punctul  $x$ , altfel se depărtează.

Algoritmul  $LVQ - 1$  este descris în (alg. 3).  $LVQ - 2$  și  $LVQ - 3$  sunt două îmbunătățiri ale algoritmului, ambele încearcă să utilizeze datele de antrenare într-un mod mai eficient [49].

#### 2.4.6 Mixturi gaussiene

Mixturile finite reprezintă un model probabilistic puternic, des utilizat în domeniul recunoașterii formelor [41]. Mixturile sunt capabile să modeleze situații în care fiecare formă este generată de către o sursă, aceasta din urmă făcând parte dintr-o mulțime de surse. Mixturile pot fi privite ca o clasă de modele capabile de reprezentarea unor densități probabilistice complexe.

Considerăm  $K$  surse, fiecare dintre acestea fiind caracterizată de o funcție de probabilitate (măsură sau densitate)  $p_m(x|\Theta)$ , parametrizată prin  $\Theta_m$ ,  $m = 1, \dots, K$ . La generarea unei forme prima dată se alege o sursă cu probabilitatea  $\{w_1, w_2, \dots, w_K\}$  din cele  $K$  surse și o observație generată de sursa respectivă cu probabilitatea  $p_m(x|\Theta_m)$ . Variabila aleatoare definită de acest mecanism de generare este caracterizată de următoarea funcție de probabilitate

$$p(x|\Theta_{(K)}) = \sum_{m=1}^K w_m p_m(x|\Theta_m) \quad (2.20)$$

unde fiecare  $p_m(y|\Theta_m)$  se numește o componentă iar

$$\Theta_{(K)} = \{\Theta_1, \Theta_2, \dots, \Theta_K, w_1, w_2, \dots, w_K\}. \quad (2.21)$$

De obicei se presupune că, toate componentele au forme funcționale identice, de exemplu sunt gaussiene. Potrivirea unui model cu mixturi la un set de observații  $X = \{x_1, x_2, \dots, x_n\}$  înseamnă estimarea acelor parametri, care descriu cel mai potrivit datele. Cu toate că se pot construi mixturi din mai multe tipuri de componente, în majoritatea literaturii de specialitate sunt descrise mixturile gaussiene (GMM-Gaussian Mixture Model).

În continuare prezentăm mai detaliat metoda mixturilor gaussiene. Această metodă este utilizată atât pentru modelarea vorbitorilor în cazul sistemelor de recunoaștere a vorbitorilor [16, 36, 55, 70] cât și pentru estimarea densităților asociate stărilor unui model Markov ascuns.

Cele mai bune rezultate în verificarea vorbitorului prin metode dependente de text s-au obținut cu modelul GMM. Metoda s-a dovedit a fi metoda statistică cea mai robustă în acest sens. În acest caz distribuția vectorilor spectrali extrași dintr-un semnal vocal specific unui vorbitor este modelată printr-o mixtură gaussiană. Funcția de probabilitate a unei mixturi gaussiene putem defini prin

$$p(x|\Theta_{(K)}) = \sum_{i=1}^K w_i b_i(x|\Theta_i) \quad (2.22)$$

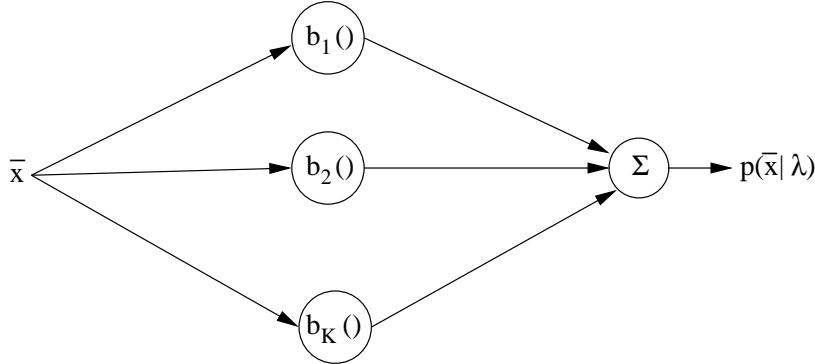


Figura 2.6: Calculul funcției de probabilitate

unde fiecare termen  $b_i$  este o funcție gaussiană de  $d$  variabile ( $x \in R^d$ ) cu parametri  $\Theta_i = (\mu_i, \Sigma_i)$  de forma:

$$b_i(x) = \frac{1}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} e^{-\frac{1}{2}(x-\mu_i)^T \Sigma_i^{-1}(x-\mu_i)} \quad (2.23)$$

unde  $\mu_i$  și  $\Sigma_i$  reprezintă valoarea medie, respectiv matricea de covarianță a densității de ordin  $i$ . Calculul funcției de probabilitate a mixturii este prezentată de figura 2.6. Ponderile satisfac constrângerile  $\sum_{i=1}^K w_i = 1$ , care garantează că 2.25 este o funcție de densitate adevărată. Matricea de covarianță  $\Sigma_i$  poate fi completă sau diagonală. Pentru a reduce numărul parametrilor, se poate utiliza o singură matrice de covarianță comună pentru toate componentele densității, numită și matrice de covarianță globală (care la rândul ei poate fi și diagonală). Aceste constrângerile se utilizează pentru a limita numărul parametrilor liberi care trebuie estimati din datele de antrenare. Astfel o mixtură gaussiană este parametrizată prin vectori de medie, matrici de covarianță și ponderile mixturii:

$$\lambda = \{w_i, \mu_i, \Sigma_i\} \quad i = 1, \dots, K. \quad (2.24)$$

Există două avantaje majore ale aplicării metodei GMM în aplicații de modelare a vorbirii. Primul avantaj constă în posibilitatea de a reprezenta prin densitățile componentelor o mulțime de clase acustice. Aceste clase acustice vor reprezenta diferite configurații de tract vocal care sunt utile atât în caracterizarea identității vorbitorilor cât și pentru caracterizarea unităților de vorbire. Clasa acustică de ordin  $i$  va fi reprezentată prin media  $\mu_i$  și matricea de covarianță  $\Sigma_i$ .

Al doilea avantaj al utilizării metodei este observația empirică că, o combinație

liniară de funcții de bază gaussiene poate reprezenta o clasă largă de distribuții. Una dintre caracteristicile cele mai puternice ale metodei este capacitatea acestia de a forma aproximări netede pentru densități arbitrate.

**Estimarea parametrilor modelului** Există multe tehnici pentru estimarea parametrilor modelului. Cea mai populară metodă rămâne estimarea conform criteriului probabilității maxime (Maximum Likelihood-ML). Dându-se datele de antrenare, obiectivul este de a estima parametrii modelului 2.24, adică găsirea parametrilor care maximizează probabilitatea datelor de antrenare. Probabilitatea unei secvențe de  $n$  vectori de antrenare  $X = \{x_1, x_2, \dots, x_n\}$  se poate calcula cu formula 2.25 în cazul în care datele de antrenare sunt observații independente.

$$p(X|\lambda) = \prod_{i=1}^n p(x_i|\lambda) \quad (2.25)$$

Din considerante computaționale, în loc de a calcula  $p(X|\lambda)$ , se calculează logarithmul acestuia

$$\log P(X|\lambda) = \sum_{i=1}^n \log p(x_i|\lambda). \quad (2.26)$$

Utilizând metoda multiplicatorilor Lagrange, funcția criteriu care trebuie maximizată devine

$$J = \sum_{j=1}^n \log p(x_j|\lambda) - \mu \left( \sum_{i=1}^K w_i - 1 \right),$$

unde  $\mu$  este multiplicatorul Lagrange. Calculând derivatele partiale  $\frac{\partial J}{\partial w_i}$ ,  $\frac{\partial J}{\partial \mu_i}$ ,  $\frac{\partial J}{\partial \Sigma_i}$  și egalându-le cu zero, obținem următoarele formule:

$$w_i = \frac{1}{n} \sum_{j=1}^n q_i(x_j) \quad (2.27)$$

$$\mu_i = \frac{\sum_{j=1}^n q_i(x_j) x_j}{\sum_{j=1}^n q_i(x_j)} \quad (2.28)$$

$$\Sigma_i = \frac{\sum_{j=1}^n q_i(x_j) (x_j - \mu_i)(x_j - \mu_i)^T}{\sum_{j=1}^n q_i(x_j)}, \quad (2.29)$$

unde

$$q_i(x) = \frac{w_i b_i(x)}{\sum_{j=1}^K w_j b_j(x)} \quad (2.30)$$

este probabilitatea aposteriori a componentei de ordin  $i$ .

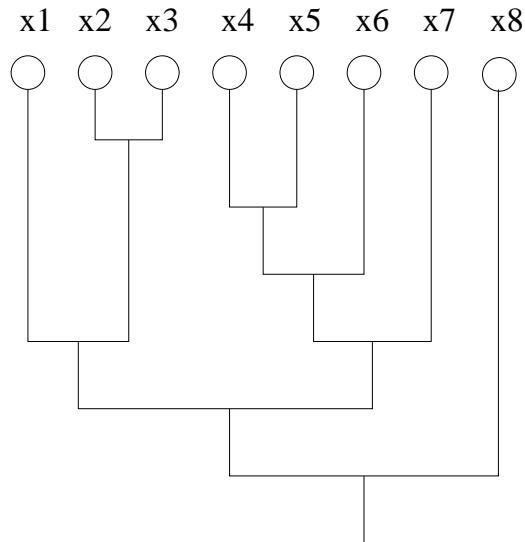


Figura 2.7: Dendogramă reprezentând gruparea datelor

#### 2.4.7 Grupare ierarhică

În algoritmii de grupare prezentați înainte, grupele create de algoritmi au fost netede, în sensul că algoritmii erau incapabili să detecteze cazurile în care anumite grupe erau formate din alte subgrupe. Totuși există mulțimi de date în care datele au o structură ierarhică, grupele fiind formate din alte subgrupe.

O ierarhie de clase se poate construi în două moduri [30]:

- aglomerativ - se pornește cu clase formate dintr-un singur obiect și se unesc clasele apropiate până se obține numărul de clase dorite
- diviziv - inițial obiectele formează o singură clasă, iar clasele unui nivel se obțin prin divizarea claselor nivelului precedent

Să considerăm mulțimea de puncte  $X = \{x_1, x_2, \dots, x_n\}$  și o partitioare a acestor puncte în  $K$  grupe. Se pornește cu  $n$  clase, fiecare conținând exact un punct  $D_i = \{x_i\}$ ,  $i = 1, \dots, n$ . La fiecare pas se aleg acele două clase care sunt cele mai asemănătoare și se unesc. Algoritmul se repetă până când se obține numărul de clase dorite. Acest algoritm se poate reprezenta în mod grafic cu ajutorul unui arbore, care se numește dendogramă și ne arată modul de grupare a punctelor (vezi figura 2.7). Această metodă de grupare este dată de algoritm (alg. 4).

Algoritmul (alg. 4) se termină când numărul specificat de clase este atins. Pentru a măsura asemănarea între două clase se pot utiliza următoarele funcții

**Algoritm 4** Algoritmul de grupare ierarhică aglomerativă

---

Pas 1. Se inițializează  $K$ -numărul claselor; Fie  $K' = n$   $D_i = \{x_i\}$ ,  $i = 1, \dots, n$ .

Pas 2. Se caută grupele cele mai asemănătoare  $D_i$  și  $D_j$  și se fuzionează,  $K' = K' - 1$

Pas 3. Se repetă Pas 2 până când  $K' = K$

---

de distanță:

$$d_{min}(D_i, D_j) = \min_{x \in D_i, x' \in D_j} d(x, x') \quad (2.31)$$

$$d_{max}(D_i, D_j) = \max_{x \in D_i, x' \in D_j} d(x, x') \quad (2.32)$$

$$d_{avg}(D_i, D_j) = \frac{1}{n_i n_j} \sum_{x \in D_i} \sum_{x' \in D_j} d(x, x') \quad (2.33)$$

$$d_{mean}(D_i, D_j) = d(\mu_i, \mu_j) \quad (2.34)$$

Dacă se utilizează funcția de distanță definită prin 2.31, algoritmul se numește algoritmul de grupare pe criteriul vecinului celui mai apropiat. Dacă algoritmul se termină, când distanța dintre clasele cele mai apropiate depășește o valoare de prag, se numește algoritmul legăturii unice (Single Linkage Algorithm). Privim punctele, care trebuie grupate, ca nodurile unui graf. Acestea, împreună cu arcele care leagă nodurile din aceeași grupă  $D_i$ , formează un drum. Dacă  $d_{min}(., .)$  este utilizat pentru calculul distanței între submulțimi, nodurile cele mai apropiate vor determina submulțimile cele mai apropiate. Reuniunea grupelor  $D_i$  cu  $D_j$  corespunde adăugării unui arc la perechea de noduri corespunzătoare celor mai apropiate submulțimi. Deoarece aceste arce care leagă grupele întotdeauna leagă două grupe distincte, graful rezultat nu conține cicluri, adică formează un arbore. Dacă algoritmul se continuă până când toate grupele vor fi legate, arborele astfel obținut se numește arborele minim de acoperire.

Dacă se utilizează funcția  $d_{max}(., .)$  definită de ecuația 2.32, algoritmul se numește grupare cu vecinii cei mai îndepărtați (Farthest-Neighbor Clustering). Dacă algoritmul se termină, când distanța dintre clasele cele mai apropiate depășește o valoare de prag, se numește algoritmul cu legături complete (Complete Linkage Algorithm). Acest algoritm produce un graf în care toate nodurile dintr-o grupă sunt legate prin arce. Exprimându-ne în terminologia grafurilor, fiecare grupă formează un subgraf complet. Distanța dintre grupe este determinată de perechea de noduri cele mai distante ale celor două grupe. La reuniunea a două grupe, se

---

**Algoritm 5** Algoritmul de cuantizare vectorială binară

---

- Pas 1. Se formează un dicționar de vectori de cod având un singur vector de cod, centroidul tuturor punctelor
- Pas 2. Se dublează numărul vectorilor de cod din dicționar conform ecuațiilor 2.35, 2.36
- Pas 3. Se utilizează algoritmul K-means clasic pentru a obține un dicționar de vectori de cod optim conform noii împărțiri
- Pas 4. Se repetă pașii 2 și 3 până când se obține dicționarul de vectori de cod cu dimensiunea dorită.
- 

vor adăuga arce între toate perechile de noduri ale celor două grupe.

Un algoritm ierarhic diviziv este algoritmul de cuantizare vectorială binară [26, 38, 67]. Acest algoritm este cunoscut și sub numele de algoritmul Linde, Buzo, Gray (LBG). Acest algoritm se utilizează pentru crearea unui dicționar de coduri, format din  $K = 2^M$  vectori de cod. Se pornește cu o singură grupă și se împart iterativ grupele în câte două grupe. Divizarea grupelor se realizează în mod uniform conform următoarelor ecuații:

$$\mu_i^+ = (1 + \epsilon)\mu_i \quad (2.35)$$

$$\mu_i^- = (1 - \epsilon)\mu_i \quad (2.36)$$

Valoarea  $\epsilon$  se alege între 0.01 și 0.05. Pentru determinarea noilor centroizi de după împărțire, se utilizează algoritmul K-means clasic (alg. 1). O versiune mai eficientă aplică algoritmul K-means pentru fiecare grupă separat și se reține și structura de arbore pentru o căutare eficientă în cadrul dicționarului de vectori de cod.

Algoritmul 5. este un algoritm eficient din punct de vedere computațional, care produce întotdeauna același rezultat, fiind inițializat în mod deterministic. Algoritmul produce întotdeauna un dicționar de vectori de cod având dimensiunea  $K = 2^M$ .

## 2.5 Modele Markov ascunse

Începând cu anii '80 cercetătorii din domeniul recunoașterii vorbirii au început să aplice metode stohastice pentru modelarea vorbirii. Expresia stochastică se utilizează pentru a indica faptul că aceste modele încearcă să caracterizeze variabilitatea prezentă în pronunțarea sunetelor. Un model stochastic este modelul Markov ascuns (Hidden Markov Model - HMM), care stă la baza mai multor sis-

teme comerciale de recunoașterea vorbirii. Aceste modele au fost introduse în domeniul recunoașterii vorbirii în mod independent de către Baker la universitatea Carnegie-Mellon și Jelinek la IBM în anul 1975. Vorbirea a fost modelată ca un proces stochastic dublu, în care rezultatele primului proces stochastic sunt ascunse, doar cel de-al doilea proces este observabil. Cele două procese trebuiau caracterizate prin prisma celui de-al doilea, cel care era observabil. Algoritmul, care permitea identificarea parametrilor modelului, este astăzi cunoscut sub numele de Expectation Maximisation (**EM**) și a fost prima dată publicat de către Dempster în 1977 [27]. Baum și colegii săi au elaborat algoritmul pentru estimarea parametrilor modelului HMM [13], care poate fi considerat un caz special al algoritmului **EM** dat de către Dempster. Acest algoritm poartă numele de algoritmul de reestimare **Baum-Welch**. Întârzierea aplicării acestor algoritmi pentru recunoașterea vorbirii se datorează faptului că rezultatele lui Dempster precum și rezultatele obținute de către Baum și colegii săi au fost publicate în reviste care sunt citite de prea puțini ingineri. Descrierea amănunțită a acestor modele o găsim în [14, 26, 29, 38, 40, 43, 44, 67, 68].

### 2.5.1 Elemente ale unui model Markov ascuns cu observații discrete

Un model Markov ascuns cu observații discrete este caracterizat de următoarele elemente:

$N$ -numărul de stări ascunse ale modelului, care în anumite aplicații sunt asociate cu semnificații fizice; vom eticheta stările individuale cu  $\{1, 2, \dots, N\}$ , și vom nota starea la momentul  $t$  cu  $q_t$ .

$M$ -numărul de simboluri distincte, adică dimensiunea alfabetului discret generat de model; vom nota simbolurile individuale astfel:  $\{v_1, v_2, \dots, v_M\}$ .

$A = \{a_{ij}\}$  – matricea probabilităților de tranziție între stări, unde:

$$a_{ij} = P[q_{t+1} = j | q_t = i], \quad 1 \leq i, j \leq N \quad (2.37)$$

În cazul special, când se poate trece din orice stare în orice altă stare într-un singur pas (cazul modelului ergodic)  $a_{ij} > 0 \quad \forall i, j$

$B = \{b_j(k)\}$  – matricea probabilității de emisie a simbolurilor, unde:

$$b_j(k) = P[o_t = v_k | q_t = j], \quad 1 \leq k \leq M, \quad (2.38)$$

definește distribuția simbolurilor în starea  $j$ , cu  $1 \leq j \leq N$ .

$\pi = \{\pi_i\}$  – distribuția stării initiale, unde:

$$\pi_i = P[q_1 = i], \quad 1 \leq i \leq N \quad (2.39)$$

Se poate vedea din definițiile de mai sus că, o specificație completă a modelului Markov ascuns  $\lambda$  presupune o stabilire a parametrilor  $N$  și  $M$ , precum și a celor trei seturi de probabilități  $A$ ,  $B$ ,  $\pi$ . Un model HMM va fi notat în modul următor:

$$\lambda = (A, B, \pi) \quad (2.40)$$

### 2.5.2 Cele trei probleme de bază ale modelelor Markov ascunse

Pentru ca modelul Markov ascuns prezentat în secțiunea anterioară să fie utilizabil în aplicații reale, trebuie rezolvate trei probleme de bază. Aceste trei probleme sunt următoarele:

1. **Problema evaluării:** Dându-se secvența de observații  $O = (o_1, o_2, \dots, o_T)$  și un model  $\lambda = (A, B, \pi)$ , să se calculeze în mod eficient probabilitatea secvenței de observații, condiționată de model,  $P(O|\lambda)$ .
2. **Problema determinării secvenței de stări optimale:** Dându-se secvența de observații  $O = (o_1, o_2, \dots, o_T)$ , și un model  $\lambda = (A, B, \pi)$ , să se determine acea secvență de stări  $q = (q_1, q_2, \dots, q_T)$ , care a generat secvența de observații cu probabilitatea maximă.
3. **Problema antrenării:** Dându-se secvența de observații  $O = (o_1, o_2, \dots, o_T)$ , și un model  $\lambda = (A, B, \pi)$ , să se reestimeze parametrii modelului  $\lambda = (A, B, \pi)$  astfel încât  $P(O|\lambda)$  să fie maxim.

*Problema 1* poate fi privită ca o problemă de evaluare a probabilității de potrivire între un model dat și o secvență de observații. În cazul problemelor de recunoaștere avem un set de modele și o secvență de observații. Pentru a decide care model a generat secvența de observații, se va evalua probabilitatea secvenței de observații în fiecare model, și se va alege acel model, care a furnizat probabilitatea maximă.

*Problema 2* este o problemă în care se încearcă descoperirea procesului stochastic al modelului prin determinarea succesiunii corecte de stări. De fapt, este foarte greu de găsit această secvență de stări aşa-numită corectă. Pentru a rezolva această problemă cât mai bine posibil, se folosește un criteriu de optimiza-

litate. Unul din algoritmii cei mai cunoscuți pentru rezolvarea acestei probleme este algoritmul Viterbi.

*Problema 3* este o problemă în care se dorește optimizarea parametrilor modelului pentru a descrie cât mai bine o secvență dată de observații. Secvența de observații utilizată la ajustarea parametrilor modelului se numește secvență de antrenare. Problema antrenării este foarte importantă pentru cele mai multe aplicații ale modelelor Markov ascunse, deoarece permite adaptarea optimă a parametrilor modelului la datele de antrenare, prin reestimarea acestora. Cel mai cunoscut algoritm de reestimare a parametrilor este algoritmul Baum-Welch.

### 2.5.2.1 Soluții pentru Problema 1

Dorim să calculăm probabilitatea secvenței de observații,  $O = (o_1, o_2, \dots, o_T)$  în modelul  $\lambda = (A, B, \pi)$ . Metoda cea mai simplă este să enumerăm toate secvențele posibile de lungime  $T$ . Obținem  $N^T$  secvențe de lungime  $T$ . Considerăm secvența fixată:

$$q = (q_1, q_2, \dots, q_T) \quad (2.41)$$

unde  $q_1$  este starea inițială. Probabilitatea secvenței de observații  $O$  este:

$$P(O|q, \lambda) = \prod_{t=1}^T P(o_t|q_t, \lambda) \quad (2.42)$$

Am presupus că observațiile sunt independente, deci:

$$P(O|q, \lambda) = b_{q_1}(o_1) \cdot b_{q_2}(o_2) \dots b_{q_T}(o_T) \quad (2.43)$$

Probabilitatea secvenței  $q$  este dată de formula:

$$P(q|\lambda) = \pi_{q_1} \cdot a_{q_1 q_2} \cdot a_{q_2 q_3} \dots a_{q_{T-1} q_T} \quad (2.44)$$

Probabilitatea ca secvența  $q$  să rezulte secvența de observații  $O$  este:

$$P(O, q|\lambda) = P(O|q, \lambda) \cdot P(q|\lambda) \quad (2.45)$$

Probabilitatea secvenței de observații  $O$  în modelul dat este:

$$\begin{aligned} P(O|\lambda) &= \sum_{\text{toate } q} P(O|q, \lambda) P(q|\lambda) = \\ &= \sum_{q_1, q_2, \dots, q_T} \pi_{q_1} b_{q_1}(o_1) a_{q_1 q_2} b_{q_2}(o_2) \dots a_{q_{T-1} q_T} b_{q_T}(o_T) \end{aligned} \quad (2.46)$$

**Algoritm 6** Algoritmul forward

Pas 1. Inițializare

$$\alpha_1(i) = \pi_i b_i(o_1), \quad 1 \leq i \leq N. \quad (2.48)$$

Pas 2. Inductie

$$\alpha_{t+1}(j) = \left[ \sum_{i=1}^N \alpha_t(i) a_{ij} \right] b_j(o_{t+1}), \quad 1 \leq t \leq T-1, \quad 1 \leq j \leq N \quad (2.49)$$

Pas 3. Terminare

$$P(O|\lambda) = \sum_{i=1}^N \alpha_T(i) \quad (2.50)$$


---

Interpretarea acestei formule este următoarea:

- Inițial ne aflăm în starea  $q_1$  cu probabilitatea  $\pi_{q_1}$  și se generează simbolul  $o_1$  cu probabilitatea  $b_{q_1}(o_1)$ .
- Se trece în starea  $q_2$  cu probabilitatea de tranziție  $a_{q_1 q_2}$  și se generează simbolul  $o_2$  cu probabilitatea  $b_{q_2}(o_2)$ .
- Acest proces se continuă până când se face inclusiv tranziția din starea  $q_{T-1}$  în starea  $q_T$  cu probabilitatea  $a_{q_{T-1} q_T}$  generând simbolul  $o_T$  cu probabilitatea  $b_{q_T}(o_T)$ .

Calculul acestei probabilități presupune efectuarea de aproximativ  $2T \cdot N^T$  de operații. Din cauza ineficienței acestui algoritm se utilizează alți algoritmi, cum ar fi algoritmul înainte (**forward procedure**).

**Algoritmul înainte (The Forward Procedure)**

Considerăm variabila

$$\alpha_t(i) = P(o_1 o_2 \dots o_t, q_t = i | \lambda) \quad (2.47)$$

care este probabilitatea, ca până în momentul  $t$  să se obțină subsecvența de observații  $o_1 o_2 \dots o_t$ , iar sistemul să fie în starea  $i$  la momentul  $t$ . Această probabilitate se poate calcula inductiv cu algoritmul 6.

Primul pas inițializează probabilitățile de a fi în starea  $i$  și de a avea observația  $o_1$ . Pasul inductiv este ilustrat în figura 2.8(a). Figura ne arată cum se poate ajunge din starea  $i$ ,  $1 \leq i \leq N$ , la momentul  $t$ , în starea  $j$  la momentul  $t+1$ .

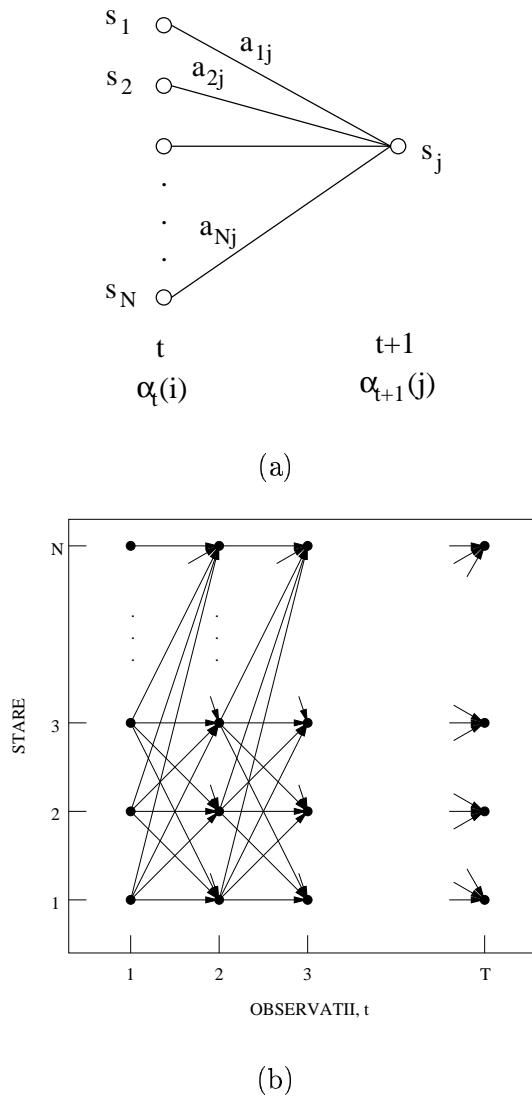


Figura 2.8: a) Secvențe de operații necesare în vederea calculării variabilei  $\alpha_{t+1}(j)$ .  
 b) Implementarea calculării lui  $\alpha_t(i)$

Deoarece  $\alpha_t(i)$  este probabilitatea secvenței de observații  $o_1 o_2 \dots o_t$  la momentul  $t$ , și cu condiția ca sistemul să fie în starea  $i$ , produsul  $\alpha_t(i) a_{ij}$  va fi probabilitatea ca sistemul în momentul  $t + 1$  să fie în starea  $j$ . Însumând acest produs pentru toate stările posibile  $i$ ,  $1 \leq i \leq N$  rezultă probabilitatea de a fi în momentul  $t + 1$  în starea  $j$ . După acest pas se vede că  $\alpha_{t+1}(j)$  se obține prin înmulțirea valorii obținute cu  $b_j(o_{t+1})$ . Calculele din formula 2.49 sunt repetate pentru toate stările  $j$ ,  $1 \leq j \leq N$ , pentru un  $t$  fixat. După aceea se iterează pentru fiecare  $t = 1, 2, \dots, T - 1$ . În pasul 3 obținem  $P(O|\lambda)$  ca suma variabilelor  $\alpha_T(i)$ .

Dacă examinăm complexitatea calculului  $\alpha_t(j)$ ,  $1 \leq t \leq T$ ,  $1 \leq j \leq N$ , observăm că necesită aproximativ  $N^2T$  calcule. Mai precis sunt necesare  $N(N + 1)(T - 1) + N$  multiplicări și  $N(N - 1)(T - 1)$  adunări. Comparând cu calculele directe, obținem raportul:

$$\frac{C_{forma-directa}}{C_{forward}} = \frac{2TN^T}{N^2T} = 2N^{T-2} \quad (2.51)$$

Procedura forward este bazată pe structura de lăție din figura 2.8(b).

### Algoritmul înapoi (The Backward Procedure)

În mod similar cu procedura înapoi, definim variabila  $\beta_t(i)$  definit prin:

$$\beta_t(i) = P(o_{t+1} o_{t+2} \dots o_T | q_t = i, \lambda) \quad (2.52)$$

care este probabilitatea ca în momentul  $t$  sistemul să fie în starea  $i$  și începând cu acest moment să se obțină secvența de observații  $o_{t+1} o_{t+2} \dots o_T$ . Algoritmul descris este dat de (alg. 7).

Pasul 1 definește în mod arbitrar, ca  $\beta_T(i)$  să fie 1 pentru orice  $i$ . Pasul 2 este ilustrat și în figura 2.9.

Pentru a calcula  $\beta_t(i)$  pentru  $t$  și  $i$  fixat, trebuie să însumăm produsele obținute din următoarele probabilități:

- probabilitatea de tranziție din starea  $i$  în starea  $j$  ( $a_{ij}$ )
- probabilitatea de a obține observația  $o_{t+1}$  la momentul  $t + 1$ , sistemul fiind în starea  $j$  ( $b_j(o_{t+1})$ )
- probabilitatea, ca începând cu momentul  $t + 1$  să se obțină secvența de observații  $o_{t+2} o_{t+3} \dots o_T$ . ( $\beta_{t+1}(j)$ )

**Algoritm 7** Algoritmul backward

Pas 1. Initializare

$$\beta_T(i) = 1, \quad 1 \leq i \leq N \quad (2.53)$$

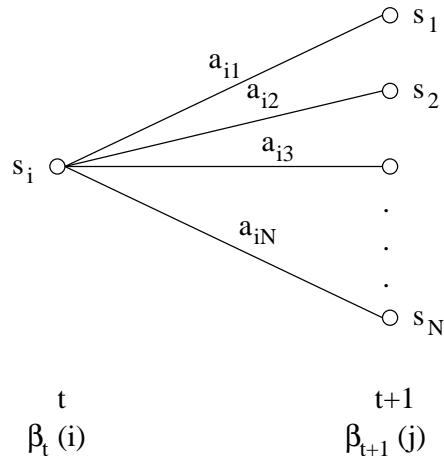
Pas 2. Inducție

$$\beta_t(i) = \sum_{j=1}^N a_{ij} b_j(o_{t+1}) \beta_{t+1}(j), \quad t = T-1, T-2, \dots, 1, \quad 1 \leq i \leq N \quad (2.54)$$

Pas 3. Terminare

$$P(O|\lambda) = \sum_{i=1}^N \pi_i \beta_1(i) b_i(o_1) \quad (2.55)$$


---

Figura 2.9: Operații necesare pentru calculul variabilei  $\beta_t(i)$

Repetând calculele pentru  $1 \leq t \leq T$  și  $1 \leq i \leq N$  algoritmul va fi de complexitate  $O(N^2T)$ .

### 2.5.2.2 Soluție pentru Problema 2

Pentru prima problemă am avut soluții exacte. În cazul problemei 2 va fi mai dificil să dăm o soluție exactă din cauza definiției de secvență optimă de stări, deoarece există mai multe criterii de optimalitate. Un criteriu de optimalitate ar fi să alegem cea mai probabilă stare la fiecare moment  $t$ . Acest criteriu de optimalitate maximizează numărul corect de stări. Pentru a implementa această soluție, definim variabila de probabilitate

$$\gamma_t(i) = P(q_t = i|O, \lambda) \quad (2.56)$$

care este probabilitatea de a fi în starea  $i$  la momentul  $t$  și producând secvența de observații  $O$  în modelul  $\lambda$ . Putem exprima  $\gamma_t(i)$  în modul următor:

$$\begin{aligned} \gamma_t(i) &= P(q_t = i|O, \lambda) \\ &= \frac{P(O, q_t = i|\lambda)}{P(O|\lambda)} \\ &= \frac{P(O, q_t = i|\lambda)}{\sum_{i=1}^N P(O, q_t = i|\lambda)} \end{aligned} \quad (2.57)$$

Deoarece  $P(O, q_t = i|\lambda)$  este egal cu  $\alpha_t(i)\beta_t(i)$ , obținem:

$$\gamma_t(i) = \frac{\alpha_t(i)\beta_t(i)}{\sum_{i=1}^N \alpha_t(i)\beta_t(i)} \quad (2.58)$$

Utilizând  $\gamma_t(i)$  putem obține formula pentru starea cea mai probabilă la momentul  $t$ :

$$q_t^* = \arg \min_{1 \leq i \leq N} [\gamma_t(i)], \quad 1 \leq t \leq T. \quad (2.59)$$

Ecuația 2.59 maximizează numărul stărilor corecte selectând starea cea mai probabilă pentru fiecare moment  $t$ . Totuși, secvența astfel obținută poate fi și una imposibilă, de exemplu dacă avem probabilități de tranziție cu valoarea 0 ( $a_{ij} = 0$ ) pentru o anumită stare  $i$ . O posibilă soluție pentru problema de mai sus este să modificăm criteriul de optimalitate astfel, încât să căutăm perechi de stări corecte sau triplete de stări corecte. Deși acest criteriu poate fi aplicabil pentru anumite aplicații, criteriul cel mai des folosit este de a găsi o singură secvență de stări care maximizează  $P(q|O, \lambda)$ , care este echivalent cu maximizarea  $P(q, O|\lambda)$ . O tehnică

---

**Algoritm 8** Algoritmul lui Viterbi

---

Pas 1. Initializare

$$\delta_1(i) = \pi_i b_i(o_1), 1 \leq i \leq N \quad (2.62)$$

$$\psi_1(i) = 0 \quad (2.63)$$

Pas 2. Inducție

$$\delta_t(j) = \max_{1 \leq i \leq N} [\delta_{t-1}(i)a_{ij}]b_j(o_t), 2 \leq t \leq T, 1 \leq j \leq N \quad (2.64)$$

$$\psi_t(j) = \arg \max_{1 \leq i \leq N} [\delta_{t-1}(i)a_{ij}], 2 \leq t \leq T, 1 \leq j \leq N \quad (2.65)$$

Pas 3. Terminare

$$P^* = \max_{1 \leq i \leq N} [\delta_T(i)] \quad (2.66)$$

$$q_T^* = \arg \max_{1 \leq i \leq N} [\delta_T(i)] \quad (2.67)$$

Pas 4. Reconstituire drum

$$q_t^* = \psi_{t+1}(q_{t+1}^*), t = T - 1, T - 2, \dots, 1 \quad (2.68)$$


---

formală pentru găsirea acestei secvențe există, se bazează pe metoda programării dinamice și se numește algoritmul Viterbi.

### Algoritmul Viterbi

Pentru a găsi secvența de stări cea mai potrivită,  $q = (q_1 q_2 \dots q_T)$ , pentru secvența de observații  $O = (o_1 o_2 \dots o_T)$ , trebuie să definim cantitatea

$$\delta_t(i) = \max_{q_1 q_2 \dots q_{t-1}} P[q_1 q_2 \dots q_{t-1}, q_t = i, o_1 o_2 \dots o_t | \lambda] \quad (2.60)$$

$\delta_t(i)$  este probabilitatea maximă de-a lungul unui singur drum, la momentul  $t$ , definită pentru primele  $t$  observații, terminată în starea  $i$ . Prin inducție obținem

$$\delta_{t+1}(j) = [\max_i \delta_t(i)a_{ij}]b_j(o_{t+1}) \quad (2.61)$$

Pentru a găsi această secvență de stări, trebuie să păstrăm argumentul care maximizează ecuația 2.61 pentru fiecare  $t$  și  $j$ . Pentru acesta vom utiliza tabloul  $\psi_t(j)$ . Algoritmul complet este dat în (alg. 8).

Trebuie să menționăm că algoritmul Viterbi este asemănător cu algoritmul forward (cu excepția pasului 4). Diferența majoră este maximizarea din ecuația

**Algoritm 9** Algoritmul Viterbi modificat

Pas 0. Preprocesare

$$\tilde{\pi}_i = \log(\pi_i), 1 \leq i \leq N$$

$$\tilde{b}_i(o_t) = \log[b_i(o_t)], 1 \leq i \leq N, 1 \leq t \leq T$$

$$\tilde{a}_{ij} = \log(a_{ij}) 1 \leq i, j \leq N$$

Pas 1. Inițializare

$$\tilde{\delta}_1(i) = \log(\delta_1(i)) = \tilde{\pi}_i + \tilde{b}_i(o_1), 1 \leq i \leq N$$

$$\psi_1(i) = 0, 1 \leq i \leq N$$

Pas 2. Recurență

$$\tilde{\delta}_t(j) = \log(\delta_t(j)) = \max_{1 \leq i \leq N} [\tilde{\delta}_{t-1}(i) + \tilde{a}_{ij}] + \tilde{b}_j(o_t)$$

$$\psi_t(j) = \arg \max_{1 \leq i \leq N} [\tilde{\delta}_{t-1}(i) + \tilde{a}_{ij}], \quad 2 \leq t \leq T, 1 \leq j \leq N$$

Pas 3. Terminare

$$\tilde{P}^* = \max_{1 \leq i \leq N} [\tilde{\delta}_T(i)]$$

$$q_t^* = \arg \max_{1 \leq i \leq N} [\tilde{\delta}_T(i)]$$

Pas 4. Reconstituire drum

$$q_t^* = \psi_{t+1}(q_{t+1}^*), t = T - 1, T - 2, \dots, 1$$

2.64, care în cazul algoritmului forward este o sumă.

**O alternativă pentru implementarea algoritmului Viterbi**

Logaritmând parametrii modelului, algoritmul Viterbi prezentat în secțiunea precedentă poate fi implementat fără utilizarea operației de înmulțire (alg. 9).

**2.5.2.3 Soluție pentru Problema 3**

Estimarea parametrilor modelului  $\lambda = (A, B, \pi)$  pentru a satisface un criteriu de optimalitate este cea mai sofisticată problemă. Nu există metode analitice cunoscute pentru a rezolva problema pentru seturi de parametri care maximizează probabilitatea secvenței de observații. Putem de altfel alege  $\lambda = (A, B, \pi)$  astfel, încât probabilitatea  $P(O|\lambda)$  să aibă un maxim local folosind o procedură iterativă,

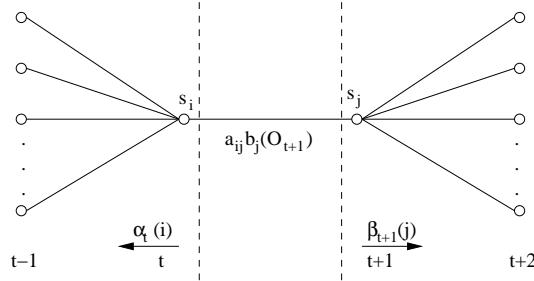


Figura 2.10: Operații necesare pentru calculul probabilității evenimentului, ca sistemul să se afle în momentul  $t$  în starea  $i$ , respectiv în momentul  $t + 1$  în starea  $j$

de exemplu Baum-Welch (caz particular al algoritmului Expectation Maximization - EM) sau folosind tehnici de tip gradient. În această secțiune discutăm o procedură iterativă bazată în principal pe lucrarea clasică a lui Baum și ai altor colegi pentru alegerea parametrilor modelului.

Definim  $\xi_t(i, j)$ , ca fiind probabilitatea, ca sistemul la momentul  $t$  să fie în starea  $i$  și la momentul  $t + 1$  în starea  $j$ , dată fiind secvența de observații  $O$  și modelul  $\lambda$ , ca

$$\xi_t(i, j) = P(q_t = i, q_{t+1} = j | O, \lambda) \quad (2.69)$$

Drumurile care satisfac condițiile cerute de ecuația 2.69 sunt reprezentate în fig. 2.10.

Din definițiile variabilelor directe și inverse putem scrie  $\xi_t(i, j)$  sub forma

$$\begin{aligned} \xi_t(i, j) &= \frac{P(q_t = i, q_{t+1} = j, O | \lambda)}{P(O | \lambda)} \\ &= \frac{\alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)}{P(O | \lambda)} \\ &= \frac{\alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)}{\sum_{i=1}^N \sum_{j=1}^N \alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)} \end{aligned} \quad (2.70)$$

Am definit mai înainte  $\gamma_t(i)$  ca și probabilitatea, ca sistemul la momentul  $t$  să se afle în starea  $i$ , fiind dat întregul set de observații precum și modelul. Astfel putem exprima  $\gamma_t(i)$  prin  $\xi_t(i, j)$  prin formula:

$$\gamma_t(i) = \sum_{j=1}^N \xi_t(i, j) \quad (2.71)$$

Dacă însumăm  $\gamma_t(i)$  după  $t$ , obținem o valoare  $\sum_{t=1}^{T-1} \gamma_t(i)$  care poate fi interpretată ca frecvența trecerii prin starea  $i$ , sau echivalent, frecvența tranzițiilor făcute din starea  $i$ . În mod similar însumarea lui  $\xi_t(i, j)$  după  $t$ ,  $\sum_{t=1}^{T-1} \xi_t(i, j)$ , poate fi interpretată ca fiind numărul de tranziții făcute din starea  $i$  în starea  $j$ .

Utilizând formulele de mai sus și conceptul de numărare a evenimentelor, putem propune o metodă pentru reestimarea parametrilor modelului Markov ascuns.

$\bar{\pi}_j =$  frecvența cu care sistemul este în starea  $i$  la momentul ( $t = 1$ ) =  $\gamma_1(i)$   
 $\bar{a}_{ij} = \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)}$ , fiind raportul dintre numărul de tranziții din starea  $i$  în starea  $j$  și numărul de tranziții din starea  $i$ .

$\bar{b}_j(k) = \frac{\sum_{t=1, o_t=v_k}^T \gamma_t(j)}{\sum_{t=1}^T \gamma_t(j)}$ , fiind raportul dintre frecvența evenimentului ca sistemul să fie în starea  $j$  și să facă observația  $v_k$  și frecvența ca sistemul să fie în starea  $j$ .

Dacă notăm modelul inițial cu  $\lambda = (A, B, \pi)$ , iar modelul reestimat cu  $\bar{\lambda} = (\bar{A}, \bar{B}, \bar{\pi})$ , atunci s-a demonstrat de către echipa lui Baum că, ori modelul inițial  $\lambda$  definește un punct critic al funcției de probabilitate, caz în care  $\bar{\lambda} = \lambda$ ; ori  $\bar{\lambda}$  este mai probabil ca  $\lambda$  în sensul  $P(O|\bar{\lambda}) > P(O|\lambda)$ ; caz în care secvența de observații va fi generată cu o probabilitate mai mare de către modelul  $\bar{\lambda}$ .

Continuând procedura de mai sus, înlocuind la fiecare pas modelul  $\lambda$  cu  $\bar{\lambda}$  și repetând pașii de estimare obținem în final estimata ML (Maximum Likelihood) al modelului Markov ascuns. Trebuie menționat însă că prin algoritmul forward-backward obținem doar un maxim local, iar în problemele reale funcțiile de probabilitate sunt foarte complexe cu mai multe maxime locale.

### 2.5.3 Topologii de modele Markov ascunse

O modalitate de clasificare a modelelor Markov ascunse este aceea, în care se ține cont de structura matricei de tranziție  $A$ . Până acum am discutat doar despre cazul special al modelelor ergodice, în care din fiecare stare se poate ajunge în orice altă stare (eventual chiar în aceeași stare), într-un număr finit de tranziții. Un astfel de model este ilustrat în figura 2.11(a), pentru cazul particular cu 4 stări.

S-a constatat însă, că există situații în care proprietățile anumitor semnale se modeleză mai bine cu altfel de modele decât cel ergodic. Un astfel de model este ilustrat în figura 2.11(b). Acest model se numește model Bakis sau model stânga-dreapta (left-right model); numele se datorează faptului că stările se desfășoară de la stânga la dreapta. Este evident că modelele Markov ascunse de tip stânga-dreapta pot modela cu ușurință semnale ale căror proprietăți variază în timp,

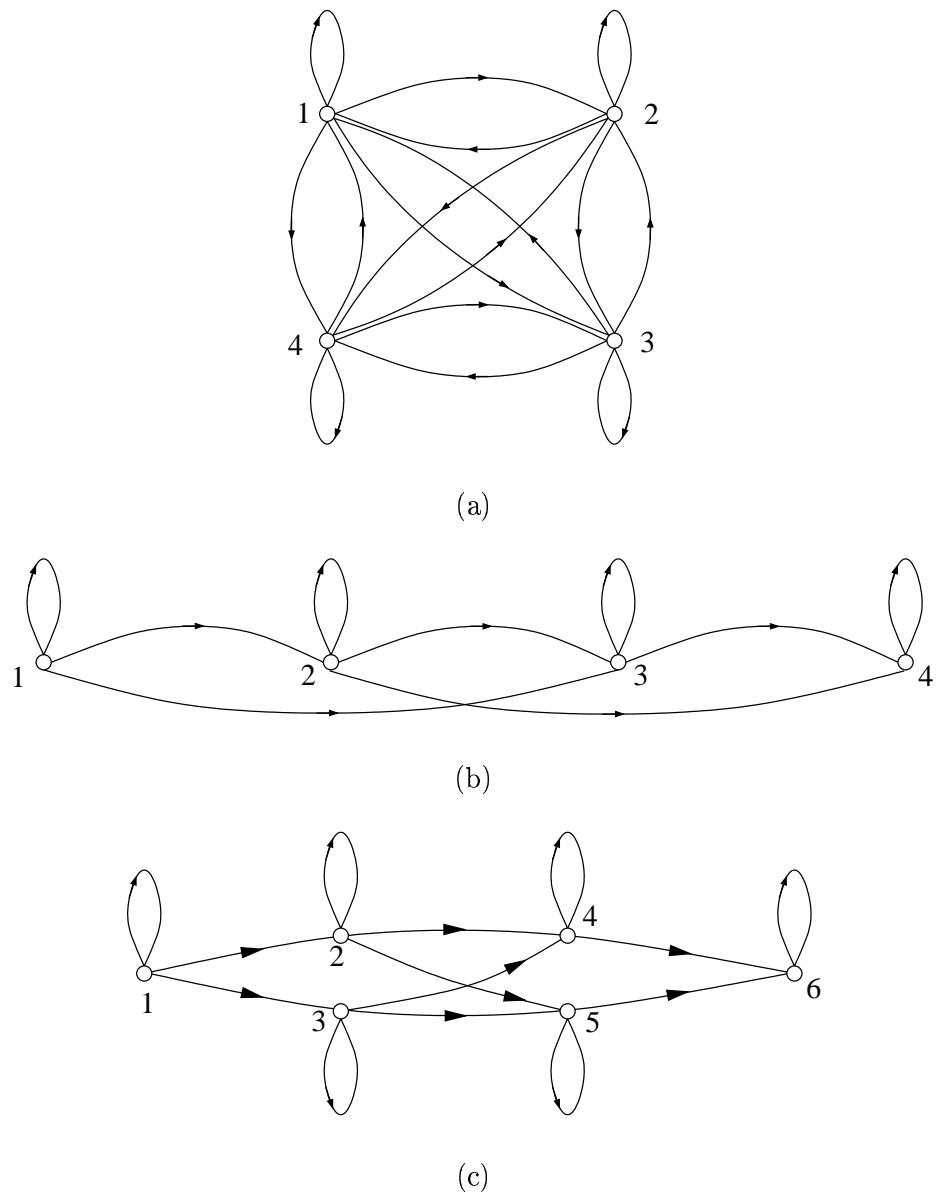


Figura 2.11: Trei tipuri diferite de MMA

ca de exemplu semnalul vocal. Proprietatea modelelor Markov ascunse de tip stânga-dreapta este aceea, că elementele matricei de tranziție respectă condiția:

$$a_{ij} = 0, \forall j < i, \quad (2.72)$$

ceea ce înseamnă că nu sunt permise tranziții în stări cu indici mai mici decât cei ai stărilor curente. Mai mult, probabilitățile stărilor inițiale au proprietatea:

$$\pi_i = \begin{cases} 0, & i \neq 1 \\ 1, & i = 1 \end{cases} \quad (2.73)$$

deoarece secvența de stări trebuie să înceapă cu starea 1 și să se termine în starea  $N$ .

Figura 2.11(c) ne arată încrucișarea a două modele stânga-dreapta. Prin punerea în paralel a două structuri de tip Bakis se poate obține modelul din figură.

#### 2.5.4 Tipuri de modele Markov ascunse

Cele mai utilizate modele Markov ascunse sunt:

- modele discrete
- modele continue
- modele semicontinue

##### 2.5.4.1 Modele continue

În secțiunile precedente am prezentat modelele Markov ascunse de tip discret. În acest caz observațiile aparțin unei mulțimi finite. Dacă observațiile aparțin unei mulțimi continue, trebuie adaptată descrierea modelului discret.

O observație în cazul modelului discret este o valoare discretă  $o_t \in \{v_1, v_2, \dots, v_M\}$ , iar în cazul continuu  $o_t \in R^d$ , un vector  $d$ -dimensional.

Probabilitățile de emisie a observațiilor în stările modelului se descrie cu  $b_j(k) = P[o_t = v_k | q_t = j]$ ,  $1 \leq k \leq M$ ,  $\forall 1 \leq j \leq N$  în cazul discret, iar în cazul continuu va fi modificat la  $b_j(x) = \sum_{k=1}^M w_{jk} N(x, \mu_{jk}, \Sigma_{jk}) = \sum_{k=1}^M w_{jk} b_{jk}(x)$ , unde  $N(x, \mu_{jk}, \Sigma_{jk})$  sau  $b_{jk}(x)$  este funcția de densitate gaussiană cu valoare medie  $\mu_{jk}$  și matricea de covarianță  $\Sigma_{jk}$ .  $M$  este numărul componentelor mixturii, iar  $w_{jk}$ ,  $k = 1, \dots, M$  sunt ponderile componentelor mixturii, care satisfac

constrângerea

$$\sum_{k=1}^M w_{jk} = 1.$$

Ecuatiile de reestimare devin

$$\overline{w_{jk}} = \frac{\sum_{t=1}^T \gamma_t(j, k)}{\sum_{t=1}^T \sum_{k=1}^M \gamma_t(j, k)} \quad (2.74)$$

$$\overline{\mu_{jk}} = \frac{\sum_{t=1}^T \gamma_t(j, k) \cdot o_t}{\sum_{t=1}^T \gamma_t(j, k)} \quad (2.75)$$

$$\overline{\Sigma_{jk}} = \frac{\sum_{t=1}^T \gamma_t(j, k) (o_t - \overline{\mu_{jk}})(o_t - \overline{\mu_{jk}})^T}{\sum_{t=1}^T \gamma_t(j, k)} \quad (2.76)$$

unde

$$\gamma_t(j, k) = \left[ \frac{\alpha_t(j)\beta_t(j)}{\sum_{j=1}^N \alpha_t(j)\beta_t(j)} \right] \left[ \frac{w_{jk} N(o_t, \mu_{jk}, \Sigma_{jk})}{\sum_{m=1}^M w_{jm} N(o_t, \mu_{jm}, \Sigma_{jm})} \right] \quad (2.77)$$

Interpretarea formulelor este simplă, și anume:  $w_{jk}$  se estimează prin raportul a două valori medii, de câte ori este sistemul în starea  $j$  și de câte ori generează o observație folosind componenta  $k$  a mixturii corespunzătoare stării  $j$ . Interpretări similare se pot da și pentru  $\mu_{jk}$  și  $\Sigma_{jk}$ .

#### 2.5.4.2 Modele semicontinuе

Modelele semicontinuе reprezintă un compromis între eficiență și aplicabilitate. Modelele continue au prea mulți parametri, care prezintă probleme atât din punctul de vedere al estimării adecvate cât și din punct de vedere computațional, calculele fiind mult prea lente. În cazul modelor semicontinuе în loc să avem mixturi gaussiene pentru fiecare stare în parte, definim un singur set de mixturi gaussiene, cu un număr adecvat de componente, care vor fi utilizate în mod comun de toate stările. Numărul componentelor se alege astfel încât mixturile gaussiene să fie capabile pentru a reprezenta întregul spațiu acustic.

În cazul modelor HMM discrete se utilizează un dicționar de vectori de cod  $\{o_1, o_2, \dots, o_M\}$  pentru discretizarea vectorilor acustici, adică pentru a cuantiza vectorii acustici în observații discrete. Distribuția vectorilor acustici caracteristici

stărilor se va defini ca

$$b_j(x) = \sum_{k=1}^M b_j(k) f(x|o_k) = \sum_{k=1}^M b_j(k) N(x, \mu_k, \Sigma_k) \quad (2.78)$$

unde  $o_k$  este vectorul cod cu indicele  $k$ ,  $b_j(k)$  este probabilitatea de generare a observației  $k$  în starea  $j$  (cazul modelelor discrete), sau ponderea în cazul modelelor continue, iar  $N(x, \mu_k, \Sigma_k)$  sunt componentele mixturii gaussiene partajate de către toate stările. Dacă privim prin prisma modelelor HMM discrete, dicționarul de vectori de cod va fi generalizat, fiecare vector cod va fi reprezentat de către o valoare medie și o matrice de covarianță. Cuantizarea cu acest dicționar de vectori de cod generalizat produce valoarea  $f(x|o_k)$  pentru fiecare vector de cod  $o_k$ .

Deoarece în cazul modelelor semicontinute  $M$  este de obicei o valoare mare, ecuația 2.78 se poate simplifica utilizând cele mai reprezentative  $L$  valori pentru  $f(x|o_k)$  pentru fiecare  $x$  fără a afecta performanța modelului.

$$b_j(x) = \sum_{k \in \eta(x)} f(x|o_k) b_j(k) \quad (2.79)$$

Deoarece  $\eta(x)$  este mai mic decât  $M$ , ecuația 2.79 reduce în mod semnificativ cantitatea de calcule.

### 2.5.5 Detalii de implementare

În secțiunile precedente am studiat teoria generală a modelelor Markov ascunse și câteva variante ale modelului. În această secțiune vom studia detalii practice legate de implementarea modelelor, cum ar fi scalarea, secvențe de antrenare multiple și estimarea parametrilor inițiali. Pentru o parte dintre aceste detalii putem da soluții analitice exakte, pentru restul doar soluții empirice bazate pe experiență.

#### 2.5.5.1 Scalarea

Pentru a înțelege necesitatea scalării în procesul de reestimare a parametrilor modelului considerăm definiția variabilei  $\alpha_t(i)$  dată de ecuația 2.47. Se poate vedea că,  $\alpha_t(i)$  este o sumă cu termenii având forma

$$\left( \prod_{s=1}^{t-1} a_{q_s q_{s+1}} \prod_{s=1}^t b_{q_s}(o_s) \right)$$

cu  $q_t = i$  și  $b$  fiind o probabilitate discretă definită prin ecuația 2.38. Deoarece elementele matricilor  $A$  și  $B$  sunt probabilități cu valori subunitare (de obicei foarte mici), se poate observa că, dacă  $t$  crește, termenii  $\alpha_t(i)$  tind către 0. Pentru valori mari ale parametrului  $t$  (peste 100) valorile  $\alpha_t(i)$  depășesc și precizia dublă cu care lucrează calculatoarele. Astfel singura modalitate pentru efectuarea calculelor este încorporarea unei scalări la fiecare reestimare a parametrilor.

Procedura de scalare cea mai simplă înmulțește valorile  $\alpha_t(i)$  cu un coeficient, care este independent de  $i$  pentru a păstra valorile scalate  $\alpha_t(i)$  în domeniul de precizie al calculatorului pentru fiecare  $t$ ,  $1 \leq t \leq T$ . O scalare similară se va utiliza și în calculul coeficienților  $\beta_t(i)$ .

Pentru a înțelege mai bine procedura de scalare, considerăm formula de reestimare a coeficienților de tranziție  $a_{ij}$ .

$$\bar{a}_{ij} = \frac{\sum_{t=1}^{T-1} \alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)}{\sum_{t=1}^T \sum_{j=1}^N \alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)} \quad (2.80)$$

Notăm cu  $\hat{\alpha}_t(i)$  valoarea scalată a lui  $\alpha_t(i)$ , iar cu  $\hat{\alpha}_i(t)$  varianta locală a lui  $\alpha_i(t)$ , înainte de scalare.

Pentru  $t = 1$  calculăm  $\alpha_1(i)$  conform ecuației 2.48 și inițializăm  $\hat{\alpha}_1(i) = \alpha_1(i)$  cu  $c_1 = \frac{1}{\sum_{i=1}^N \alpha_1(i)}$  și  $\hat{\alpha}_1(i) = c_1 \alpha_1(i)$ . Pentru fiecare  $t$ ,  $2 \leq t \leq T$ , calculăm  $\hat{\alpha}_1(i)$  conform formulei 2.49 și obținem

$$\hat{\alpha}_t(i) = \sum_{j=1}^N \hat{\alpha}_{t-1}(j) a_{ji} b_i(o_t). \quad (2.81)$$

Calculăm coeficientul de scalare

$$c_t = \frac{1}{\sum_{i=1}^N \hat{\alpha}_t(i)} \quad (2.82)$$

având

$$\hat{\alpha}_t(i) = c_t \hat{\alpha}_t(i) \quad (2.83)$$

Din ecuațiile 2.81-2.83 obținem

$$\hat{\alpha}_t(i) = \frac{\sum_{j=1}^N \hat{\alpha}_{t-1}(j) a_{ji} b_i(o_t)}{\sum_{i=1}^N \sum_{j=1}^N \hat{\alpha}_{t-1}(j) a_{ji} b_i(o_t)} \quad (2.84)$$

Prin inducție putem scrie  $\hat{\alpha}_{t-1}(j)$  ca

$$\hat{\alpha}_{t-1}(j) = \left( \prod_{\tau=1}^{t-1} c_{\tau} \right) \alpha_{t-1}(j) \quad (2.85)$$

Din care rezultă că  $\hat{\alpha}_t(i)$  poate fi scrisă ca

$$\hat{\alpha}_t(i) = \frac{\sum_{j=1}^N \alpha_{t-1}(j) \left( \prod_{r=1}^{t-1} c_r \right) a_{ji} b_i(o_t)}{\sum_{i=1}^N \sum_{j=1}^N \alpha_{t-1}(j) \left( \prod_{r=1}^{t-1} c_r \right) a_{ji} b_i(o_t)} = \frac{\alpha_t(i)}{\sum_{i=1}^N \alpha_t(i)} \quad (2.86)$$

Similar calculăm și  $\beta_t(i)$  utilizând aceiași coeficienți de scalare ca și pentru sirul  $\alpha$ . Astfel obținem

$$\hat{\beta}_t(i) = c_t \beta_t(i). \quad (2.87)$$

Ecuația 2.80 devine

$$\bar{a}_{ij} = \frac{\sum_{t=1}^{T-1} \hat{\alpha}_t(i) a_{ij} b_j(o_{t+1}) \hat{\beta}_{t+1}(j)}{\sum_{t=1}^{T-1} \sum_{j=1}^N \hat{\alpha}_t(i) a_{ij} b_j(o_{t+1}) \hat{\beta}_{t+1}(j)} \quad (2.88)$$

dar fiecare  $\hat{\alpha}_t(i)$  poate fi scris ca

$$\hat{\alpha}_t(i) = \left[ \prod_{s=1}^t c_s \right] \alpha_t(i) = C_t \alpha_t(i) \quad (2.89)$$

și fiecare  $\hat{\beta}_{t+1}(j)$  poate fi scris ca

$$\hat{\beta}_{t+1}(j) = \left[ \prod_{s=t+1}^T c_s \right] \beta_{t+1}(j) = D_{t+1} \beta_{t+1}(j). \quad (2.90)$$

Astfel ecuația 2.80 devine

$$\bar{a}_{ij} = \frac{\sum_{t=1}^{T-1} C_t \alpha_t(i) a_{ij} b_j(o_{t+1}) D_{t+1} \beta_{t+1}(j)}{\sum_{t=1}^{T-1} \sum_{j=1}^N C_t \alpha_t(i) a_{ij} b_j(o_{t+1}) D_{t+1} \beta_{t+1}(j)} \quad (2.91)$$

În final termenul  $C_t D_{t+1}$  poate fi scris în forma

$$C_t D_{t+1} = \prod_{s=1}^t c_s \prod_{s=t+1}^T c_s = \prod_{s=1}^T c_s = C_T \quad (2.92)$$

care este independent de  $t$ .

Este clar că, procedura de reestimare aplicată pentru  $a_{ij}$  se poate utiliza atât

pentru coeficienții  $B$ , cât și pentru  $\pi$ . Singura schimbare datorată scalării se produce în procedura de calcul al probabilității  $P(O|\lambda)$ . Nu putem însuma termenii  $\hat{\alpha}_T(i)$ , deoarece acestea sunt deja scalate, dar putem utiliza următoarea proprietate

$$\prod_{t=1}^T c_t \sum_{i=1}^N \alpha_T(i) = C_T \sum_{i=1}^N \alpha_T(i) = 1 \quad (2.93)$$

Astfel obținem

$$\prod_{t=1}^T c_t P(O|\lambda) = 1 \quad (2.94)$$

sau

$$\log[P(O|\lambda)] = \sum_{t=1}^T \log c_t. \quad (2.95)$$

Astfel se poate calcula logaritmul lui  $P$ , care oricum ar fi mult prea mic ca să-l putem reprezenta exact într-un calculator.

Menționăm că, în cazul utilizării algoritmului Viterbi pentru determinarea secvenței de stări cu probabilitate maximă ( $ML$ ), nu este necesară scalare, dacă utilizăm varianta cu logaritmi a algoritmului.

### 2.5.5.2 Secvențe de observații multiple

În secțiunea 2.5.3 am discutat despre modelul Left-Right unde stările pot fi vizitate într-o ordine crescătoare. Acest model impune o serie de constrângeri pentru elementele matricii de tranziție  $A$  și pentru vectorul probabilităților inițiale  $\pi$ . Marele neajuns al acestor modele este că o singură secvență de observații nu ajunge pentru antrenare. Astfel pentru a putea utiliza aceste modele la estimări reale, trebuie utilizate mai multe secvențe de observații la antrenare.

Modificarea procesului de reestimare este simplă și se face în modul următor. Notăm mulțimea secvențelor de observații

$$O = [O^{(1)}, O^{(2)}, \dots, O^{(k)}] \quad (2.96)$$

unde  $O^{(k)} = (o_1^{(k)}, o_2^{(k)}, \dots, o_{T_k}^{(k)})$  este secvența de observații cu numărul de ordine  $k$ . Presupunem că fiecare secvență de observații este independentă de restul secvențelor de observații și obiectivul nostru este ajustarea parametrilor modelului  $\lambda$  astfel, încât să maximizăm  $P(O|\lambda)$ .

$$P(O|\lambda) = \prod_{k=1}^K P(O^{(k)}|\lambda) = \prod_{k=1}^K P_k.$$

Deoarece formulele de reestimare se bazează pe frecvențele producerii anumitor evenimente, formulele de reestimare pentru secvențe de observații multiple se modifică prin adăugarea acestor frecvențe. Astfel formulele modificate pentru  $\bar{a}_{ij}$  și  $\bar{b}_j(l)$  vor fi

$$\bar{a}_{ij} = \frac{\sum_{k=1}^K \frac{1}{P_k} \sum_{t=1}^{T_k-1} \alpha_t^k(i) a_{ij} b_j(o_{t+1}^{(k)}) \beta_{t+1}^k(j)}{\sum_{k=1}^K \frac{1}{P_k} \sum_{t=1}^{T_k-1} \alpha_t^k(i) \beta_t^k(i)} \quad (2.97)$$

și

$$\bar{b}_j(l) = \frac{\sum_{k=1}^K \frac{1}{P_k} \sum_{t=1}^{T_k-1} \alpha_t^k(i) \beta_t^k(i)}{\sum_{k=1}^K \frac{1}{P_k} \sum_{t=1}^{T_k-1} \alpha_t^k(i) \beta_t^k(i)} \quad (2.98)$$

iar  $\pi$  nu necesită reestimare deoarece  $\pi_1 = 1$  și  $\pi_i = 0$ ,  $i \neq 1$ .

Scalarea corespunzătoare ecuațiilor 2.97 și 2.98 este evidentă, deoarece orice secvență de observații are factorul de scalare propriu. Ideea cheie este eliminarea acestui factor din fiecare termen înainte de însumare. Aceasta se poate realiza prin scrierea ecuației de reestimare cu termenii variabilelor scalate,

$$\bar{a}_{ij} = \frac{\sum_{k=1}^K \frac{1}{P_k} \sum_{t=1}^{T_k-1} \hat{\alpha}_t^k(i) a_{ij} b_j(o_{t+1}^{(k)}) \hat{\beta}_{t+1}^k(j)}{\sum_{k=1}^K \frac{1}{P_k} \sum_{t=1}^{T_k-1} \hat{\alpha}_t^k(i) \hat{\beta}_t^k(i)} \quad (2.99)$$

Astfel, pentru fiecare secvență  $O^{(k)}$  aceiași factori de scalare vor apărea în fiecare termen al sumei pentru fiecare  $t$ , care apar și în termenul  $P_k$  și astfel se simplifică exact. Astfel, utilizând valorile scalate atât pentru termenii  $\alpha$ , cât și pentru termenii  $\beta$ , ne conduce la o valoare nescalată  $\bar{a}_{ij}$ . Un rezultat similar se obține și pentru termenul  $\bar{b}_j(l)$ .

#### 2.5.5.3 Estimarea parametrilor inițiali

În teorie ecuațiile de reestimare ale parametrilor unui MMA trebuie să ne conducă la niște valori ai parametrilor, care corespund maximului local al funcției de probabilitate. O întrebare cheie este, cum să alegem valorile inițiale ale parametrilor, încât maximul local să fie egal sau cât mai aproape de maximul global al funcției de probabilitate.

Pentru această întrebare nu există un răspuns evident și simplu. Experiența însă a demonstrat că valori inițiale aleatoare și uniform distribuite ale parametrilor  $\pi$  și  $A$  sunt adecvate pentru cele mai multe cazuri. Cu toate acestea, tot experiența ne-a demonstrat că, pentru parametrul  $B$  sunt necesare estimări inițiale adecvate. Există mai multe căi pentru obținerea acestor estimări inițiale și anume:

- segmentarea manuală a secvenței de observații în stări, și medierea obser-

vațiilor pentru fiecare stare în parte

- segmentarea Maximum Likelihood și mediere
- segmentare K-means segmentală cu grupare

## 2.6 Recunoașterea vorbirii

### 2.6.1 Introducere

Fie  $X = \{x_1, x_2, \dots, x_t\}$  o secvență de vectori acustici corespunzătoare unei rostiri. Problema recunoașterii înseamnă determinarea acelei secvențe de cuvinte  $\hat{W}$  care satisface

$$\hat{W} = \max_W P(W|X) = \max_W P(X|W) \cdot P(W) \quad (2.100)$$

unde  $W = w_1, w_2, \dots, w_n$ ,  $w_i \in V$  este o secvență de  $n$  cuvinte aparținând unui vocabular  $V$  fixat. În formula 2.100  $P(W)$  reprezintă probabilitatea secvenței de cuvinte  $W$  și este furnizată de componenta de model de limbaj al unui sistem pentru recunoașterea vorbirii.

Modelarea limbajului este o problemă care apare nu doar în domeniul recunoașterii vorbirii ci și în domeniul prelucrării limbajului natural. Există o serie de cărți, care tratează această problemă foarte detaliat. O parte din cărți sunt din domeniul recunoașterii vorbirii [14, 40, 43, 79], altele din domeniul procesării limbajului natural (Natural Language Processing-NLP) [58], iar altele pot fi privite ca și cărți interdisciplinare, deoarece tratează ambele domenii [44].

În cazul recunoașterii vorbirii, cuvintele vor însemna forma cuvintelor pronunțate. Dacă același cuvânt are două pronunții diferite, atunci se consideră două cuvinte. Utilizând regula Bayes  $P(W)$  se poate descompune în modul următor

$$P(W) = \prod_{i=1}^n P(w_i|w_1 w_2 \dots w_{i-1}) \quad (2.101)$$

unde  $P(w_i|w_1 w_2 \dots w_{i-1})$  este probabilitatea ca  $w_i$  să fie pronunțat după ce  $w_1, w_2, \dots, w_{i-1}$  au fost pronunțate. Secvența  $w_1, w_2, \dots, w_{i-1}$  se numește istoric și se notează cu  $h_i$ . Rolul modelului de limbaj este să furnizeze sistemului de recunoaștere estimări corespunzătoare ale probabilităților  $P(w_i|w_1 w_2 \dots w_{i-1})$ , ceea ce înseamnă practic predicția cuvântului ce urmează după secvența de cuvinte  $w_1 w_2 \dots w_{i-1}$ . Schema unui sistem, care utilizează regula lui Bayes pentru decizie este prezentată în figura 2.12.

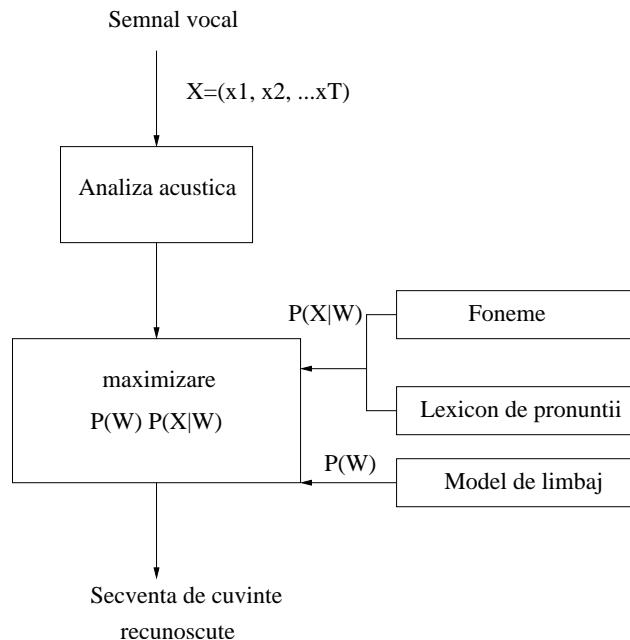


Figura 2.12: Regula lui Bayes pentru un sistem de recunoașterea vorbirii

Relativ la acest sistem putem nota următoarele observații:

- Modelul de limbaj  $P(W)$  este independent de observațiile acustice și are rolul de a ghida sistemul în concatenarea cuvintelor limbajului pentru a forma propoziții.
- Modelul acustico-fonetic,  $P(X|W)$ , este probabilitatea condiționată de a observa vectorii acustici  $X$ , când vorbitorul pronunță secvența de cuvinte  $W$ . Atât aceste probabilități, cât și probabilitățile corespunzătoare modelului de limbaj sunt estimate în faza de antrenare a sistemului de recunoaștere. În cazul unui sistem de recunoaștere având un vocabular mare, unitatea de recunoaștere este o parte a cuvântului, de obicei fonemul, silaba sau demisilaba. Modelele cuvintelor se obțin prin concatenarea acestor modele conform transcripțiilor fonetice ale cuvintelor.
- Decizia, care are ca și rezultat secvența de cuvinte pronunțate, trebuie să fie luată cu ajutorul unei proceduri de optimizare, care combină informații din diferite surse de cunoaștere: modelul de limbaj, modelele acustico-fonetice ale fonemelor și lexicul de pronunții sau dicționar. Procedura de optimizare este de fapt o căutare în spațiul stărilor definite de aceste surse de cunoaștere.

### 2.6.2 Modele de limbaj

Există mai multe metode pentru modelarea limbajului, însă nu toate s-au dovedit viabile. O primă metodă pentru reprezentarea limbajului ar fi cu ajutorul unui automat finit (Finite State Automaton-FSA). Automatele finite acceptă limbajele regulate, care sunt descrise de gramatici regulate [26, 40, 44]. Un prim dezavantaj ar fi că, gramaticile regulate nu sunt suficient de puternice pentru generarea limbajelor naturale. O a doua metodă mai puternică pentru modelarea limbajelor este oferită tot de teoria limbajelor formale prin gramaticile independente de context, care sunt capabile să genereze limbaje mai puternice. În practică se utilizează varianta stocastică a acestor gramatici. Cărțile [40, 44] tratează detaliat această metodă de reprezentare a limbajului. O metodă simplă, dar utilizată în practică este cea de a reprezenta spațiul limbajului prin perechi de cuvinte corecte. Această metodă se mai numește și model  $n$ -gram și se intergrează cu ușurință în căutarea Viterbi. Modelul limbaj se utilizează de către componenta de căutare, care de obicei este o căutare de la stânga-dreapta. Notând cu  $h_i = (w_1, w_2, \dots, w_{i-1})$  istoricul cuvântului  $w_i$ , ecuația 2.101 devine

$$P(W) = \prod_{i=1}^n P(w_i|h_i). \quad (2.102)$$

De obicei se utilizează clase de echivalențe pentru  $h_i$ . Notând cu  $\phi(h_i)$  clasa de echivalență, căreia aparține  $h_i$ , formula 2.102 se scrie

$$P(W) = \prod_{i=1}^n P(w_i|\phi(h_i)). \quad (2.103)$$

Modelele  $n$ -gram utilizează cele  $n - 1$  cuvinte ca și istoricul cuvântului  $w_i$ , adică  $\phi(h_i) = \{w_{i-1}, w_{i-2}, \dots, w_{i-(n-1)}\}$ . Pentru cazul modelelor trigram în calcule se utilizează formula

$$f(w_3|w_1w_2) = \frac{c(w_1w_2w_3)}{c(w_1w_2)} \quad (2.104)$$

unde  $c()$  reprezintă frecvența de apariție a unei anumite secvențe de cuvinte, iar  $f()$  reprezintă frecvența relativă. Din nefericire formula 2.104 nu este una foarte adecvată, deoarece foarte multe triplete  $w_1w_2w_3$  nu apar nici în cazul textelor de antrenare de dimensiuni uriașe. În anii 1970 cercetătorii de la IBM au făcut un experiment pe un corpus format din descrierea unor inovații având un vocabular de 1000 de cuvinte. Acest text a fost împărțit într-un set de antrenare și într-

unul de testare format din 300000 repectiv 150000 cuvinte. În mod surprinzător 23% din trigramele din setul de testare nu au apărut deloc în setul de antrenare. Utilizarea unui model de limbaj definit cu formula 2.104 ar fi cauzat o eroare cel puțin de 23%.

De aceea este necesară netezirea frecvențelor relative. Acest lucru se poate realiza simplu prin interpolarea frecvențelor relative trigram, bigram și unigram dateă de următoarea formulă

$$P(w_3|w_1w_2) = \lambda_3 f(w_3|w_1w_2) + \lambda_2 f(w_3|w_2) + \lambda_1 f(w_3)$$

unde  $\lambda_1 + \lambda_2 + \lambda_3 = 1$ ,  $\lambda_i \geq 0$ .

### 2.6.3 Măsurarea erorii în sisteme de recunoaștere a vorbirii

Evaluarea performanței unui sistem de recunoașterea vorbirii este o problemă importantă. Eroarea de recunoaștere a cuvintelor este măsura cea mai folosită. Datele disponibile sunt împărtjite în prealabil în date de antrenare și date de test. Sistemul se antrenează folosind doar datele de antrenare, iar la testare sistemului i se prezintă datele de test, pe baza cărora se calculează eroarea.

Pentru măsurarea erorii de recunoaștere a cuvintelor se consideră trei tipuri de erori:

- *substituire (SUBST)*: sistemul recunoaște un cuvânt incorrect în locul unuia corect
- *stergere (DEL)*: un cuvânt corect a fost omis din propoziția recunoscută
- *inserare (INS)*: un cuvânt care nu există în semnalul acustic a fost adăugat în propoziția recunoscută

Rata de eroare pe cuvinte se definește ca fiind

$$WER = \frac{SUBST + DEL + INS}{N} \quad (2.105)$$

unde  $N$  este numărul de cuvinte din propoziția prezentată sistemului pentru recunoaștere. Această problemă de aliniere este una clasică și este cunoscută sub numele de cel mai lung subșir comun [23]. Problema se rezolvă ușor cu metoda programării dinamice, algoritmul se bazează pe completarea unui tabel folosind următoarea recurență:

**Algoritm 10** Algoritmul WER

---

Pas 1. Inițializare  $d[0, 0] = 0$     $b[0, 0] = 0$

Pas 2. Recurență

Pentru fiecare  $i = 1 \dots n$

Pentru fiecare  $j = 1 \dots m$

Se calculează  $d[i, j]$  și  $b[i, j]$  conform ecuațiilor (2.106) și (2.107)

---

Pas 3.  $WER = \frac{d[n, m]}{n} \cdot 100\%$ , Reconstituire drum optim conform cu 2.108

---

$$d[i, j] = \min \begin{cases} d[i - 1, j] + 1 & (\text{stergere}) \\ d[i - 1, j - 1] & (\text{potrivire}) \\ d[i - 1, j - 1] + 1 & (\text{substituire}) \\ d[i, j - 1] + 1 & (\text{inserare}) \end{cases} \quad (2.106)$$

unde  $d[i, j]$  reprezintă eroarea minimă în alinierea subșirurilor  $A = (a_1, a_2, \dots, a_i)$  și  $B = (b_1, b_2, \dots, b_j)$ . Eroarea totală acumulată pentru şirurile  $A = (a_1, a_2, \dots, a_n)$  și  $B = (b_1, b_2, \dots, b_m)$  va fi  $d[n, m]$ . Dacă dorim să calculăm separat numărul de inserări, stergeri și substituiri, va trebui să utilizăm încă un tabel ajutător pentru a reține tipul erorii la fiecare  $d[i, j]$ . Astfel, dacă notăm

$$b[i, j] = \begin{cases} 1 & \text{stergere} \\ 2 & \text{inserare} \\ 3 & \text{potrivire} \\ 4 & \text{substituire} \end{cases} \quad (2.107)$$

vom putea determina exact frecvențele celor trei tipuri de erori. Pentru reconstruirea drumului optim format din stările  $(s_1, s_2, \dots, 0)$  se utilizează formulele

$$s_1 = B[n, m]$$

$$s_t = \begin{cases} B[i - 1, j] & \text{daca } s_{t-1} = 1 \\ B[i, j - 1] & \text{daca } s_{t-1} = 2 \quad \text{pentru } t = 2, \dots, s_t = 0 \\ B[i - 1, j - 1] & \text{daca } s_{t-1} = 3 \text{ sau } 4 \end{cases} \quad (2.108)$$

Algoritmul este definit în (alg. 10).

#### 2.6.4 Tehnici de căutare

Există două strategii de bază pentru căutarea propoziției celei mai probabile care satisfacă atât constrângerile acustice cât și cele lingvistice: strategia modulară și

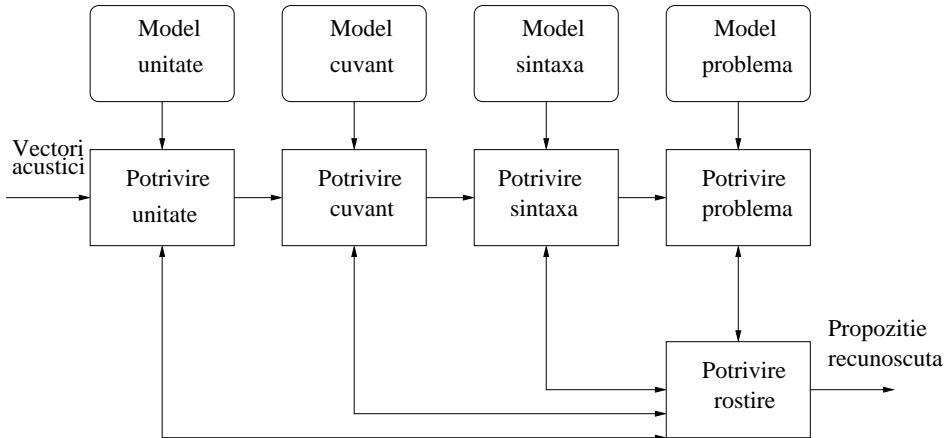


Figura 2.13: Un sistem ASR cu structură integrată

strategia integrată. În varianta integrată decizia asupra recunoașterii se ia luând în considerare toate sursele de cunoaștere (vezi figura 2.13). Utilizând toate sursele acustice, lexicale, sintactice și semantice se poate construi un automat finit, care va fi practic mașina de recunoaștere. Această strategie este extrem de des folosită în sistemele de recunoaștere a vorbirii.

Pe de altă parte, în abordarea secvențială, prezentată în figura 2.14, recunoașterea se obține prin potrivire acustică, potrivire lexicală, potrivire sintactică și potrivire semantică într-un mod secvențial. Avantajul acestei abordări este în primul rând simplitatea și posibilitatea de a implementa și testa separat modulele componente. Din punct de vedere computațional acest tip de sistem este mai abordabil. Există și o serie de neajunsuri la această abordare, dintre care cel mai grav ar fi faptul că, în fiecare modul în parte se iau o serie de decizii fără a lua în considerare restul modulelor, dar care influențează acestea, modulele fiind conectate în mod secvențial. Un prim exemplu pentru o asemenea sursă de eroare ar fi ca, în primul modul, care realizează potrivirea unităților (de exemplu foneme), dacă se determină secvența cea mai probabilă de foneme corespunzătoare secvenței acustice, fără să se ia în considerare lexicul limbajului, am putea obține secvențe de foneme inexistente în limbajul în cauză.

Căutarea poate fi făcută utilizând o singură trecere (single pass) în spațiul acustic sau utilizând mai multe treceri (multi pass). Spațiul de căutare poate fi expandat în mod static sau în mod dinamic. Expandarea statică a spațiului de căutare a devenit posibilă doar în ultimii ani, primii care au realizat un astfel de sistem au fost cei de la AT&T [62].

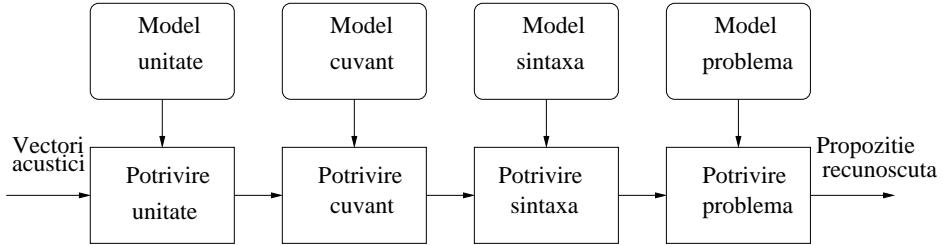


Figura 2.14: Un sistem ASR cu structură modulară

Tehnicile de căutare cele mai cunoscute sunt: căutare cadru-sincron cu o singură trecere, decodarea cu stivă.

#### 2.6.4.1 Căutare cadru-sincron cu o singură trecere

Acest tip de căutare își propune găsirea celui mai probabil drum într-o rețea finită. O căutare totală în lățime este de obicei impracticabilă datorită timpului uriaș de procesare. Pentru a reduce timpul de căutare se utilizează strategii de căutare suboptimale. O strategie este utilizarea aproximăției Viterbi, iar o altă strategie este căutarea cu fascicul. În căutarea cu fascicul cadru-sincron doar o parte din ipotezele aparținând unui fascicul sunt expandate pentru fiecare moment de timp. Algoritmul astfel obținut este unul care aproximează o căutare în lățime necesitând întregul spațiu de căutare pe tot parcursul căutării. Integrarea surselor de cunoștințe lexicale pot reduce semnificativ numărul ipotezelor viabile. Lățimea fasciculului se determină experimental și este dependentă atât de model cât și de problemă.

#### Definirea problemei-Aproximația Viterbi

La început algoritmul cu o singură trecere timp-sincron a fost proiectat pentru sisteme de recunoaștere a vorbirii cu un vocabular limitat, ca recunoașterea secvențelor de cifre. Acest algoritm însă s-a dovedit eficient și în sisteme având vocabulare cu peste 20.000 de cuvinte [64]. Spațiul de căutare poate fi considerat ca și o mașină finită de stări, definită în termenii modelelor acustici și ai modelului limbajului. Pentru o secvență de vectori acustici  $X = \{x_1 x_2 \dots x_T\}$  definim

secvența optimă de cuvinte prin

$$[w_1^N]_{opt} = \arg \max_{[w_1^N]} \left\{ P(w_1^N) \cdot \sum_{[s_1^T]} P(x_1^T; s_1^T | w_1^N) \right\} \quad (2.109)$$

unde suma se referă la toate secvențele de stări posibile  $[s_1^T]$  consistente cu secvența de cuvinte  $[w_1^N]$ . Aproximând suma cu maximul, obținem aproximarea Viterbi

$$[w_1^N]_{opt} = \arg \max_{[w_1^N]} \left\{ P(w_1^N) \cdot \max_{[s_1^T]} P(x_1^T; s_1^T | w_1^N) \right\}. \quad (2.110)$$

Altfel spus, secvența de cuvinte cea mai probabilă se aproximează prin secvența de stări cea mai probabilă. În aproximarea Viterbi, maximul nu va fi o estimare a sumei, doar cele două expresii conduc la aproximativ aceeași secvență de cuvinte.

În cazul aproximării Viterbi problema căutării se definește în modul următor:

Dorim să asociem cu fiecare vector acustic corespunzător momentului de timp  $t$  o pereche (stare, cuvânt). Această asociere poate fi privită ca o aliniere a timpului la secvența de perechi (stare, cuvânt):

$$(s_1, w_1), \dots, (s_t, w_t), \dots, (s_T, w_T) \quad (2.111)$$

Un exemplu de aliniere pentru problema recunoașterii unei secvențe de cuvinte este arătat în figura 2.15. La capetele cuvintelor avem nevoie de condiții de continuitate: să permitem tranzitii între stările  $S(v)$ , ultima stare a cuvântului  $v$  și prima stare a oricărui cuvânt  $w$ . Strategia programării dinamice ne va permite să evaluăm probabilitățile următoare

$$P(w_1 \dots w_N) \cdot P(x_1 x_2 \dots x_T; s_1 s_2 \dots s_T | w_1 w_2 \dots w_N) \quad (2.112)$$

într-o manieră stânga-dreapta și să determinăm secvența optimă de stări. Remarcăm faptul că, în această formulare a problemei secvența de stări și secvența de cuvinte sunt estimate simultan.

### Căutare cu fascicul cadru-sincron

Utilizarea metodei programării dinamice presupune că, soluția optimă a problemei se poate construi din soluții optime ale subproblemelor [23]. De aceea se

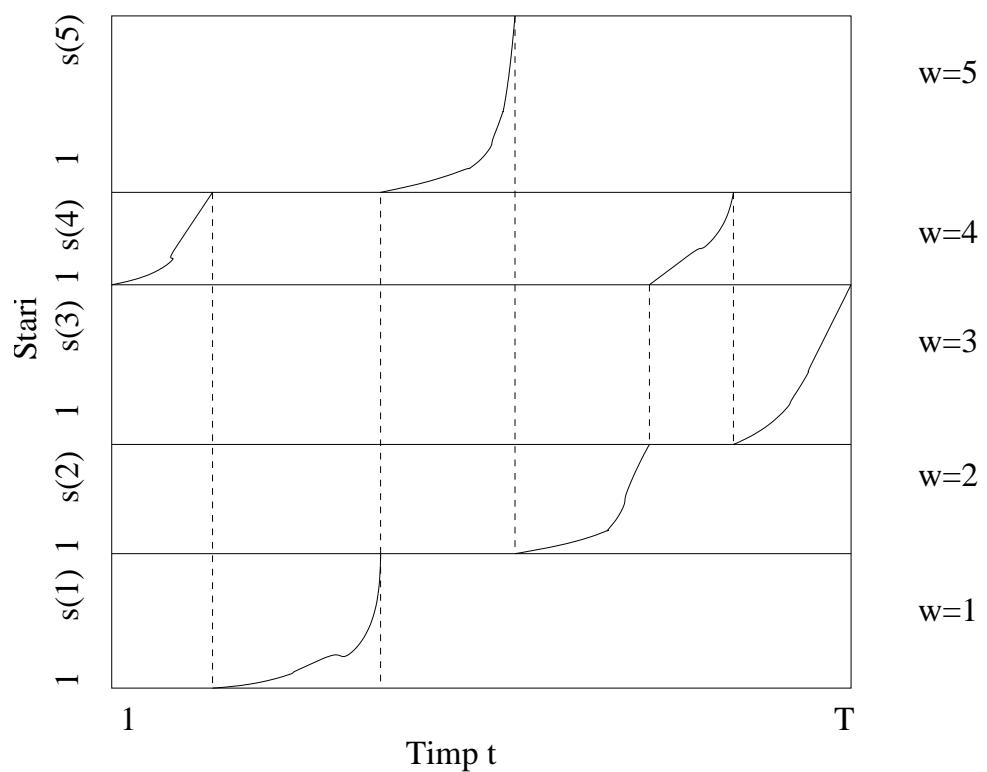


Figura 2.15: Aliniere timp-perechi stări cuvinte

utilizează un tablou pentru stocarea soluțiilor optimale ale subproblemelor. Pentru reconstituirea drumului se utilizează un tablou suplimentar. Notăm aceste tablouri în modul următor:

$Q(t, s; w)$  : scorul maxim obținut de drumul care se termină în starea  $s$  al cuvântului  $w$  la momentul de timp  $t$ .

$B(t, s; w)$  : timpul de început al drumului optim care se termină în starea  $s$  al cuvântului  $w$  la momentul de timp  $t$ .

Figura 2.16 ne arată că avem două tipuri de tranziții: între stările același cuvânt și între stări ale diferitelor cuvinte. Din această cauză vom defini două seturi de reguli pentru completarea tabelelor. Pentru stările același cuvânt definim următoarele relații recurente:

$$Q(t, s; w) = \max_{\sigma} \{q(x_t, s|\sigma; w) \cdot Q(t - 1, \sigma; w)\} \quad (2.113)$$

$$B(t, s; w) = B(t - 1, \sigma_{max}(t, s; w); w) \quad (2.114)$$

unde

$$q(x_t, s|\sigma; w) = a(s|\sigma; w) \cdot b(x_t|\sigma; w) \quad (2.115)$$

este probabilitatea condițională ca în starea  $\sigma$  a cuvântului  $w$  să fie observat vectorul acustic  $x_t$  și să se treacă în starea  $s$ .  $a(s|\sigma; w)$  este probabilitatea de tranziție din starea  $s$  în starea  $\sigma$ , iar  $b(x_t|\sigma; w)$  probabilitatea de emitere a vectorului acustic  $x_t$  în starea  $\sigma$  al cuvântului  $w$ , iar  $\sigma_{max}(t, s; w)$  este starea optimă precedentă a ipotezei  $(t, s; w)$ . Pentru a include capetele cuvintelor, introducem o stare specială  $s = 0$  utilizată pentru a începe un cuvânt nou. În momentul atingerii unui capăt de cuvânt, vom combina acest cuvânt cu cuvintele precedente. Pentru acest scop definim

$$H(w; t) := \max_v \{p(w|v) \cdot Q(t, S(v); v)\} \quad (2.116)$$

unde  $p(w|v)$  este probabilitatea condiționată a modelului limbaj de tip bigram, adică probabilitatea ca cuvântul  $v$  să fie urmat de cuvântul  $w$ . Cu simbolul  $S(v)$  se notează starea terminală a cuvântului  $v$ . În acest scop trebuie să transferăm atât scorul, cât și indicele de timp:

$$Q(t - 1, 0; w) = H(w; t - 1) \quad (2.117)$$

$$B(t - 1, 0; w) = t - 1. \quad (2.118)$$

Ecuațiile 2.117 și 2.118 presupun că prima dată sunt evaluate stările normale

---

**Algoritm 11** Algoritmul cu o singură trecere Single-best

---

Pentru  $t = 1 \dots T$  execută

Procesare la nivelul stărilor acustice

P1. Inițializare cu ecuațiile 2.117 și 2.118

P2. Aliniere, completarea tabelelor conform recurențelor 2.113 și 2.114

P3. Eliminarea ipotezelor necorespunzătoare

Procesarea la nivelul cuvintelor (stare specială)

Pentru fiecare pereche  $(w, t)$  execută

P1.  $H(w; t) = \max_v [p(w|v)Q_v(t, S(w); w)]$

P2.  $v_0(w; t) = \arg \max_v [p(w|v)Q_v(t, S(w); w)]$

se stochează predecesorul cel mai bun  $v_0 = v_0(w, t)$

se stochează marginea cea mai bună  $\tau_0 = B(t, S(v_0); v_0)$

---

$s = 1, \dots, S(w)$  pentru fiecare cuvânt  $w$  și după aceea se va evalua starea  $s = 0$ . Astfel  $Q(t, s; w)$  va conține atât probabilitățile furnizate de către modelele acustice cât și probabilitățile modelului de limbaj.

Algoritmul este dat în continuare (alg. 11). Secvența de vectori acustici este procesată de la stânga la dreapta. Algoritmul de căutare este o căutare în lățime, în care toate ipotezele sunt extinse în mod paralel pentru fiecare vector acustic.

Deoarece toate ipotezele se referă la aceeași secvență de intrare, scorurile acestora sunt comparabile. De aceea putem evita căutarea exhaustivă, eliminând ipotezele nepromițătoare. Pentru fiecare vector acustic se determină ipoteza cea mai bună, iar toate acele ipoteze care au scoruri mai mici cu un factor fixat, decât aceasta, sunt eliminate. Experimentele arată că, pentru acest tip de căutare cu fascicul, care utilizează la fiecare moment vectorul acustic și constrângerile impuse de modelul limbaj, doar un mic procent de ipoteze rămân viabile, toate celelalte fiind eliminate.

În algoritm 11. lexiconul a fost organizat într-un mod liniar. Fiecare cuvânt a fost reprezentat separat ca o secvență liniară de foneme, într-un mod independent de celealte foneme. Într-un vocabular mare (cu mai multe mii de cuvinte), există foarte multe cuvinte cu același prefix, ceea ce sugerează organizarea lexiconului sub formă de arbore.

Experimentele arată că, majoritatea ipotezelor sunt concentrate pe stări aparținând fonemelor de la începutul cuvintelor. De aceea este ideală organizarea vocabularului de pronunții sub formă de arbore.

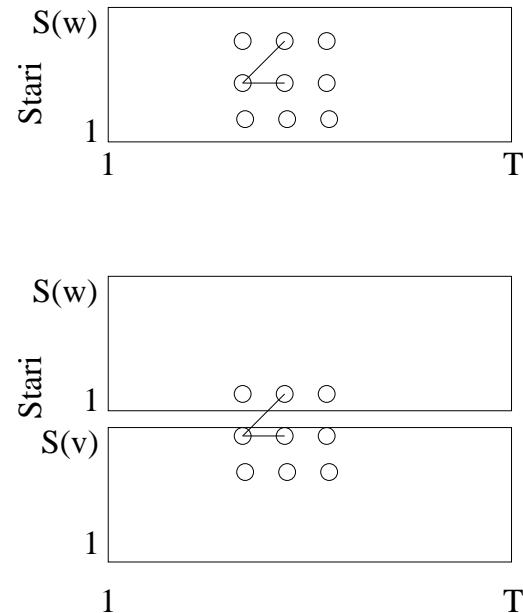


Figura 2.16: Tipuri de tranziții

#### 2.6.4.2 Decodarea cu stivă

Această tehnică este bazată pe strategia „primul cel mai bun” (Best-First). De obicei această strategie este implementată cu ajutorul unei stive pentru fiecare moment de timp, care menține o listă ordonată a ipotezelor. În practică s-a dovedit că, găsirea unei funcții euristice care să permită selectarea drumului optim este aproape imposibilă. Din această cauză căutarea cu stivă este în general mai puțin optimală decât căutarea cadru-sincron cu fascicul.

Decodarea cu stivă este o variantă a căutării  $A^*$  bazată pe algoritmul forward, unde funcția de evaluare este bazată pe probabilitatea forward. Această căutare este de fapt o căutare într-un arbore. Căutarea cu stivă, fiind un algoritm de căutare în arbore, caută acel drum în arbore, ale cărui ramuri corespund cuvintelor vocabularului  $V$ , nodurile neterminale corespund propozițiilor incomplete, iar cele terminale propozițiilor complete. Figura 2.17 arată un arbore de căutare pentru un vocabular format din 3 cuvinte  $V = \{w_1, w_2, w_3\}$ .

Un avantaj al decodorului cu stivă ar fi potrivirea cu algoritmul de antrenare forward-backward. Deoarece căutarea Viterbi este de fapt o căutare într-un graf, probabilitățile atașate drumurilor nu pot fi însumate, deoarece aceste drumuri pot corespunde unor secvențe diferite de cuvinte. În general căutarea Viterbi returnează secvența de stări optime în loc de secvența optimă de cuvinte. Pe

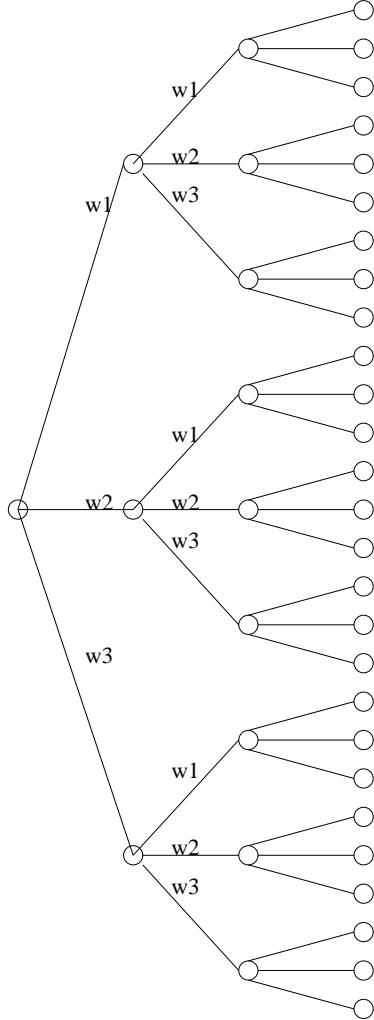


Figura 2.17: Arbore de căutare pentru decodorul cu stivă

de altă parte, decodorul cu stivă, fiind o căutare într-un arbore, pentru fiecare nod există un singur drum de la rădăcină la nodul respectiv, permitând utilizarea algoritmului forward pentru evaluarea modelelor de cuvinte.

Comparând cu algoritmul cadru-sincron Viterbi putem afirma despre acest algoritm următoarele:

- nu este cadru-sincron
- lucrează cu drumuri de lungimi diferite

Căutarea începe prin adăugarea tuturor ipotezelor de un singur cuvânt la o listă OPEN. Se alege ipoteza cea mai bună, se scoate din lista OPEN și se adaugă

toate ipotezele de drumuri la care care se poate extinde ipoteza cea mai bună. Se continuă căutarea până când se găsește drumul cel mai bun din toate drumurile din lista OPEN.

Deoarece căutarea cu stivă se realizează în mod asincron, este necesar un mecanism efectiv pentru a determina sfârșitul fonemelor, respectiv al cuvintelor. Astfel pentru a putea utiliza căutarea cu stivă sunt necesare următoarele:

- o funcție euristică pentru a estima secvența de vectori acustici încă neprelucrați
- o metodă pentru a determina sfârșitul fonemelor/cuvintelor

## 2.7 Baze de date

Pentru experimentele din această lucrare am folosit două baze de date de sunet. Baza de date TIMIT este o bază de date cu vorbitori în engleză, înregistrată în condiții ideale, de birou. Conține material vocal de la 630 de vorbitori, fiecare rostind 10 propoziții. Rata de eşantionare este 16kHz cu eşantioane de 16 biți. Întregul material fonetic a fost segmentat manual și împărțit în parte de antrenare și parte de testare. Partea de antrenare conține întregul material de la 462 vorbitori iar partea de testare întregul material audio de la 168 vorbitori. Fiecare înregistrare audio îi este atașat un fișier de date, care conține segmentarea în foneme indicând capetele în milisecunde ale acestora cât și etichetele fonemelor. În total sunt 61 de simboluri fonetice.

Cea de-a doua bază de date utilizată a fost OASIS Numbers, proiectată pentru recunoașterea numerelor din limba maghiară. Detalii despre această bază de date găsim în [47]. Frecvența de eşantionare este 22.05kHz cu eşantioane de 16 biți, condiții de birou. Partea segmentată a bazei de date este formată dintr-un material audio rostit de către 26 de vorbitori, dintre care un copil, 9 femei și 16 bărbați. Fiecare vorbitor citește de două ori cifrele 0 – 9, numerele 10 – 100, 1000, 1000000 și patru variante care se utilizează în formarea unor numere compuse cum ar fi 28. Orice număr din limba maghiară se poate forma din concatenarea acestor 26 de cuvinte. Segmentarea și anotarea bazei de date a fost făcută manual și s-au utilizat 31 de simboluri fonetice, ceea ce reprezintă doar o parte din fonemele limbii maghiare.

## Capitolul 3

# Recunoașterea fonetică a vorbitorului

### 3.1 Introducere

Recunoașterea vorbitorului reprezintă procesul automat al recunoașterii sau al verificării aceluia care vorbește, pe baza caracteristicilor individuale ce derivă din semnalul vocal al vorbitorului. Multe din aceste caracteristici sunt independente de conținutul lingvistic și sunt considerate în general surse de degradare în recunoașterea vorbirii.

Recunoașterea vorbitorului în general prezintă două probleme diferite: identificarea vorbitorului și verificarea vorbitorului. Identificarea vorbitorului este procesul în care se determină care din vorbitorii înregistrati a pronunțat fraza în cauză. Cu alte cuvinte, scopul este de a determina care din vocile cunoscute se apropie cel mai mult de vocea înregistrată. Verificarea vorbitorului este în esență o problemă de verificare de ipoteze, care va trebui să valideze sau să eliminate o identitate presupusă. Scopul este apoi de a determina dintr-o voce înregistrată, dacă persoana este într-adevăr cea care se susține a fi. În identificarea vorbitorului numărul alternativelor este egal cu mărimea populației (posibil plus unu în cazul problemei deschise, când vorbitorul de test nu face parte din vorbitorii înregistrați). În cazul verificării vorbitorului există două cazuri: acceptare și respingere, indiferent de mărimea populației. Structura simplificată a sistemelor de recunoaștere a vorbitorului este redată în figura 3.1. În blocul Extragerea parametrilor se extrag parametrii acustici (care în mod ideal vor fi specifici vorbitorului) din semnalul vocal. Cu comutatorul în poziția 1 ne aflăm în modul de antrenare, adică crearea modelelor pentru fiecare client potențial. Modelele statistice cele mai populare

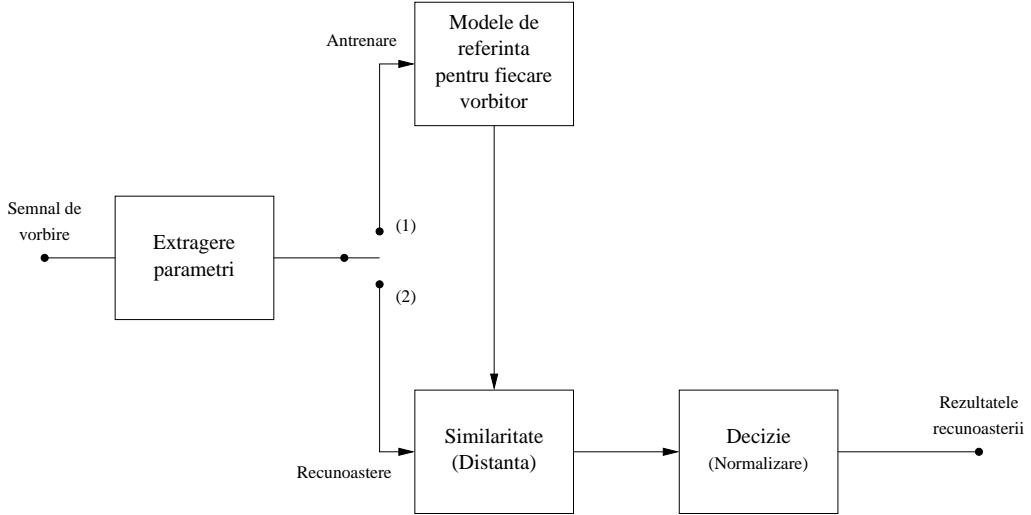


Figura 3.1: Structura unui sistem de recunoaștere a vorbitorului

sunt GMM [16, 71, 70], HMM [36], VQ [21, 45, 46] sau rețele neuronale ANN [15]. Cu comutatorul în poziția 2 vom recunoaște secvențele de test prin compararea parametrilor extrași din semnalul vocal cu modelele de referință ale vorbitorilor înregistrați. În cazul identificării vorbitorilor vom selecta pur și simplu vorbitorul cu gradul maxim de similaritate. În cazul verificării vorbitorului decizia este luată după calculul similarității și compararea acesteia cu o valoare de prag.

Tehnologia de verificare a vorbitorului are multe aplicații: face posibilă verificarea identității persoanelor prin control acces vocal, fiind considerată mult mai convenabilă decât folosirea codurilor PIN, cartelelor, cheilor sau altor metode artificiale. Alte aplicații potențiale ar fi: tranzacții bancare prin telefon, acces la baze de date, comerț electronic prin telefon, control de securitate pentru zone cu informații confidențiale. În sfârșit verificarea și identificarea vorbitorului pot fi folosite în aplicații militare.

Verificarea vorbitorului poate fi bazată pe metode dependente de text sau independente de text. În cazul metodelor dependente de text vorbitorul folosește aceeași frază sau parolă atât în antrenare, cât și în recunoaștere, iar în cazul metodelor independente de text nu există această constrângere. Verificarea vorbitorului prin metoda dependență de text se va folosi de informații lexicale în afară de caracteristicile vocale ale vorbitorului, pe când metodele independente de text vor folosi doar acestea din urmă. În general, din cauza variațiilor acustice-fonetice mari ale înregistrărilor independente de text, este nevoie de un material de antrenare mai vast pentru a caracteriza corect vorbitorul, comparativ cu metodele

dependente de text. Metodele dependente de text se bazează în general pe tehnici de potrivire a formelor, în care axa timpului secvenței de intrare este aliniată cu fiecare formă de referință, iar similaritatea calculată se acumulează de la începutul până la sfârșitul rostirii. Deoarece această metodă folosește direct identitățile vocale asociate fiecărei foneme sau silabe, în general atinge rate de recunoaștere mai mari decât sistemele independente de text. Există câteva aplicații în care nu pot fi folosite cuvinte predeterminate. De exemplu rostiri ale aceluiași cuvânt nu pot fi comparate în investigații criminale. Foarte des apare cazul în care vorbitorii nu sunt cooperanți și își schimbă intenționat rata și felul vorbirii. De altfel oamenii pot recunoaște vorbitorii indiferent de conținutul rostirii.

Ambele metode de verificare au o problemă serioasă: vocea vorbitorului poate fi înregistrată și apoi redată. Pentru a preveni această problemă, apare soluția de verificare a vorbitorului cu text promptat constând într-un text sintetic și vorbitorul este rugat să repete o secvență aleatoare de cuvinte cheie, care sunt luate dintr-un dicționar restrâns (de exemplu cifre). În acest caz recunoașterea vorbirii independentă de vorbitor, se efectuează înaintea verificării vorbitorului, pentru a se asigura că vorbitorul într-adevăr a repetat ceea ce i s-a cerut. Astfel rostirile preînregistrate nu vor fi de nici un folos impostorului, deoarece acesta nu va mai putea prevedea secvența cerută. Până și acest sistem poate fi păcălit cu un sistem avansat care reproduce cuvintele cheie în ordinea cerută. Astfel s-a propus o nouă soluție în care secvențele aleatoare promptate sunt luate dintr-un vocabular foarte mare.

În servicii reale este de dorit ca un utilizator să-și poată alege singur parola la care se efectuează verificarea: verificarea vorbitorului prin parolă definită de utilizator. Această metodă permite securitate mărită, deoarece accesul fraudulos necesită cunoașterea parolei clientului, iar utilizatorul își poate schimba parola oricând dorește.

Cercetările recente în domeniul recunoașterii vorbitorului încearcă să utilizeze informații de nivel înalt, caracteristic vorbitorului, cum ar fi stilul de vorbire, duratele specifice ale fonemelor vorbitorului, informații lexicale și informații legate de prozodie. Desigur, aceste informații de nivel înalt nu vor înlocui informațiile acustice, ci le vor extinde. O primă încercare în acest sens a fost făcută de către [28], cu modelarea vorbitorilor prin metode cunoscute din recunoașterea vorbirii și anume cu utilizarea n-gramelor. Metoda bazată pe duratele specifice ale diferitelor unități lingvistice, specifice vorbitorilor, a fost utilizată cu succes în lucrarea [34]. Modele structurate fonetic au fost utilizate în [33], iar lucrarea [60] prezintă investigații în scopul fuzionării a mai multor tipuri de caracteristici

de nivel înalt pentru reprezentarea vorbitorului.

Scopul nostru este de a studia posibilitățile de integrare a conținutului fonetic în modelul vorbitorului. Pentru acesta vom începe cu studierea conținutului din punct de vedere fonetic al unui model clasic, format din componente gaussiene. După acesta vom construi modele de vorbitori utilizând doar câte o anumită clasă fonetică, cum ar fi vocalele sau fricativele. Pe baza acestor modele vom putea determina puterea de discriminare a acestor clase fonetice pentru aplicații ca recunoașterea vorbitorilor. În final vom construi modele de vorbitori structurate fonetic și vom compara performanța acestor modele cu cele clasice. O publicație proprie dedicată studiului conținutului fonetic ale modelelor GMM de vorbitori este [8], iar publicații legate de puterea de discriminare a claselor fonetice sunt [8, 9, 11].

### 3.2 Modelarea vorbitorului

În această secțiune vom prezenta diferite metode de modelare a vorbitorilor. De fapt, metodele au fost prezentate în capitolul 2, aici vom prezenta doar utilizarea lor pentru rezolvarea acestei probleme specifice, identificarea vorbitorului.

În cazul utilizării metodei VQ prezentată în 2.4.4, fiecare model de vorbitor s-a obținut separat utilizând partea de antrenare a materialului înregistrat de la fiecare vorbitor. Astfel s-a obținut câte un dicționar de coduri, care reprezintă clasele acustice ale vorbitorului. Presupunând că un sistem de recunoaștere are  $N$  vorbitori înregistrați, vom nota modelul pentru vorbitorul  $i$  cu

$$\lambda_i = \left( m_1^{(i)}, m_2^{(i)}, \dots, m_M^{(i)} \right),$$

unde cu  $M$  am notat dimensiunea modelului. Fiecare vorbitor a fost reprezentat cu un model de aceeași dimensiune.

În faza de identificare a vorbitorului notăm cu  $X = \{x_1, x_2, \dots, x_T\}$ ,  $x_i \in R^d$  vectorii acustici extrași din semnalul vocal al vorbitorului care a cerut identificarea. Identificarea necesită calculul a  $N$  distanțe, adică se va calcula distanța mulțimii  $X$  față de fiecare model de vorbitor. Pentru modelul  $\lambda_i$  distanța este definită de către următoarea formulă

$$d(X, \lambda_i) = \frac{1}{T} \sum_{k=1}^T \min_{j=1..M} d_E(x_k, m_j^{(i)}) \quad (3.1)$$

unde  $d_E$  este metrica Euclidiană definită în  $R^d$ . Vorbitorul se va identifica, ca

fiind acel vorbitor, pentru care distorsiunea definită de 3.1 este minimă, adică

$$Id = \arg \min_{i=1,N} d(X, \lambda_i) \quad (3.2)$$

Cel de-al doilea model matematic utilizat pentru reprezentarea unui vorbitor a fost modelul mixturilor gaussiene prezentată în secțiunea 2.4.6. și în acest caz modelele formate din ponderi, valori medii și matrice de covarianțe sunt notate cu  $\{\lambda_1, \lambda_2, \dots, \lambda_N\}$ . Obiectivul este de găsi acel vorbitor, pentru care probabilitatea a posteriori pentru secvența de vectori acustici  $X = \{x_1, x_2, \dots, x_T\}$  este maximă. Prin formulă se exprimă în modul următor

$$Id = \arg \max_{1 \leq i \leq N} P(\lambda_i | X) = \arg \max_{1 \leq i \leq N} \frac{p(X | \lambda_i)}{p(X)} P(\lambda_i). \quad (3.3)$$

Considerând că fiecare vorbitor are aceeași probabilitate, iar  $p(X)$  constant, ecuația 3.3 devine

$$Id = \arg \max_{1 \leq i \leq N} p(X | \lambda_i). \quad (3.4)$$

Pentru a calcula probabilitatea unei secvențe de vectori acustici într-un model de vorbitor, pentru simplificarea calculelor, vom presupune că acești vectori acustici sunt independenți. Astfel calculul probabilității se realizează cu formula

$$p(X | \lambda_i) = \prod_{t=1}^T p(x_t | \lambda_i). \quad (3.5)$$

Pentru a evita produsul unor valori subunitare foarte mici, în loc de probabilitate se calculează logaritmul probabilității și se împarte cu  $T$ , numărul vectorilor acustici din secvența  $X$ . Formula exactă utilizată în aplicație:

$$\frac{1}{T} \log p(X | \lambda_i) = \frac{1}{T} \sum_{t=1}^T \log p(x_t | \lambda_i), \quad (3.6)$$

iar pentru calculul  $p(x_t | \lambda_i)$  se utilizează formula 2.25.

Au fost făcute o serie de experimente pentru a testa o parte din algoritmii de grupare prezentați în secțiunile precedente, în scopul recunoașterii vorbitorilor. Publicațiile relevante în acest domeniu sunt [1, 2, 3, 4]. În [3] experimentele au fost făcute pe o bază de date proprie cu 66 de vorbitori, 37 vorbitori nativi români și 29 vorbitori nativi maghiari. Din cauza condițiilor de înregistrare (zgomot și utilizarea diferitelor tipuri de microfoane) rezultatele au fost sub cele obișnuite (în jur de 80%). Mai mult, deoarece nu au fost date suficiente pentru estimarea

parametrilor modelelor GMM, metoda VQ s-a dovedit mai robustă pentru acest corpus.

Primele rezultate obținute pe o bază de date publică TIMIT au fost publicate în [2]. În aceste experimente semnalul a fost împărțit în cadre de 32ms cu o suprapunere de 20 de ms și din fiecare cadru au fost extrași primii 12 de coeficienți Mel-cepstrali [40]. Pentru experimente s-au folosit doar 168 de vorbitori dintr-un total de 630 de vorbitori în limba engleză. Antrenarea modelelor vorbitorilor a fost făcută aproximativ pe 24s de vorbire de la fiecare vorbitor, iar pentru testare s-au folosit aproximativ 6s de vorbire de la fiecare vorbitor în parte. Au fost antrenate atât modele VQ cât și modele GMM, cele GMM în acest caz atingând rata de recunoaștere 100%.

Articolul [4] este de fapt un studiu privind eficiența diferenților algoritmilor de grupare pentru identificarea vorbitorilor. Au fost comparați algoritmul generalizat K-means, algoritmul de cuantizare vectorială binară, algoritmul LVQ și GMM. Toate experimentele au fost făcute pe toți 630 de vorbitori din corpusul TIMIT, obținându-se pentru fiecare tehnică de grupare o eroare de identificare sub 1%. Parametrii acustici au fost cei obișnuiți, 12 coeficienți mel-cepstrali.

Toate studiile efectuate în vederea recunoașterii vorbitorului au arătat că în cazul când semnalul vocal este înregistrat cu microfoane de calitate, în condiții de birou, se pot obține rezultate excepționale utilizând coeficienți cepstrali, aproape cu orice algoritm de grupare. Algoritmul de grupare va determina clasele acustice specifice vorbitorului și aceste clase vor fi reprezentate în mod specific algoritmului, centroid în cazul algoritmului k-means și valoare medie, respectiv matrice de covarianță în cazul modelelor GMM. Modelele GMM mai conțin un set de ponderi, care ajută la ponderarea acestor clase acustice.

### 3.3 Distribuția conținutului fonetic

#### 3.3.1 Modele VQ

În acest capitol vom prezenta o analiză statistică a conținutului fonetic al unui model de vorbitor. Pentru simplitate am ales modelul VQ de vorbitor prezentat în 2.4.4. Acest model este format din centroizi, care reprezintă grupe de vectori acustici similari. Obiectivul nostru este de a verifica, dacă acest model separă vectorii acustici în categorii fonetice, sau nu. Pașii analizei pentru fiecare vorbitor au fost următorii:

1. Am extras vectorii acustici din datele vorbitorului, excluzând porțiunile de

- liniște, obținând mulțimea  $X = \{x_1, x_2, \dots, x_T\}$ ,  $x_i \in R^d$ .
2. Am etichetat vectorii acustici cu categoria din care fac parte aceștia și am obținut mulțimea  $Y = \{(x_1, f_1), (x_2, f_2), \dots, (x_T, f_T)\}$ , unde  $f_i \in \{A, F, N, S, V, W\}$ . Etichetele din această mulțime reprezintă categoriile fonetice prezentate în tabelul 5.
  3. Utilizând mulțimea  $Y$  am aplicat un algoritm VQ pentru un număr  $M$  de grupe, fixat în prealabil și am analizat conținutul acestor grupe din punctul de vedere al categoriilor fonetice.

Figura 3.2 conține distribuția categoriilor fonetice în modele de vorbitori formate din 2, 3, 4, 5, respectiv 6 clustere. Prima diagramă arată distribuția în două clustere. În acest caz primul cluster este populat de nazale, vocale și semivocale și cel de-al doilea de africate, fricative și plozive. Această tendință de separare este explicabilă cunoscând proprietățile acustice ale acestor categorii fonetice. Categoria africatelor este cea mai mică și vectorii aparținând acestei categorii fac parte întotdeauna din același cluster.

Ultima figură conține exact 6 clustere, cât este și numărul de categorii de foneme. Analizând această figură se poate vedea că, marea majoritate a fricativelor este concentrată în același cluster, doar o mică parte a acestora este distribuită în restul clusterelor. Același lucru este valabil și pentru categoria nazalelor. Vocalele, semivocalele și plozivele sunt distribuite în mai multe clustere.

Acest experiment ne arată că, grupând vectorii acustici pe baza similarităților dintre acestea, obținem modele de vorbitori ale căror clustere nu separă vectorii acustici în categorii de foneme.

Varietatea acustică prezentă într-o categorie fonetică se poate caracteriza prin varianță. Am calculat această varianță pentru categoriile fonetice ale fiecărui vorbitor și am determinat media aritmetică pe categorii fonetice. Pentru calcule am utilizat cei 462 de vorbitori din partea de test a bazei de date TIMIT și vectorii acustici cu 12 coeficienți Mel-cepstrali. Valorile obținute se pot vedea în tabelul 3.

### 3.3.2 Modele GMM

În acest experiment am urmărit tendința de separare în componente gaussiene ale fonemelor din datele înregistrate de la un vorbitor. Din cauză că, cadrele de la început și de la sfârșit sunt afectate de coarticulare, am utilizat din fiecare fonem doar un singur cadru, cadrul din mijloc. Fie  $X = \{x_1, x_2, \dots, x_T\}$  vectorii

### 3. Recunoașterea fonetică a vorbitorului

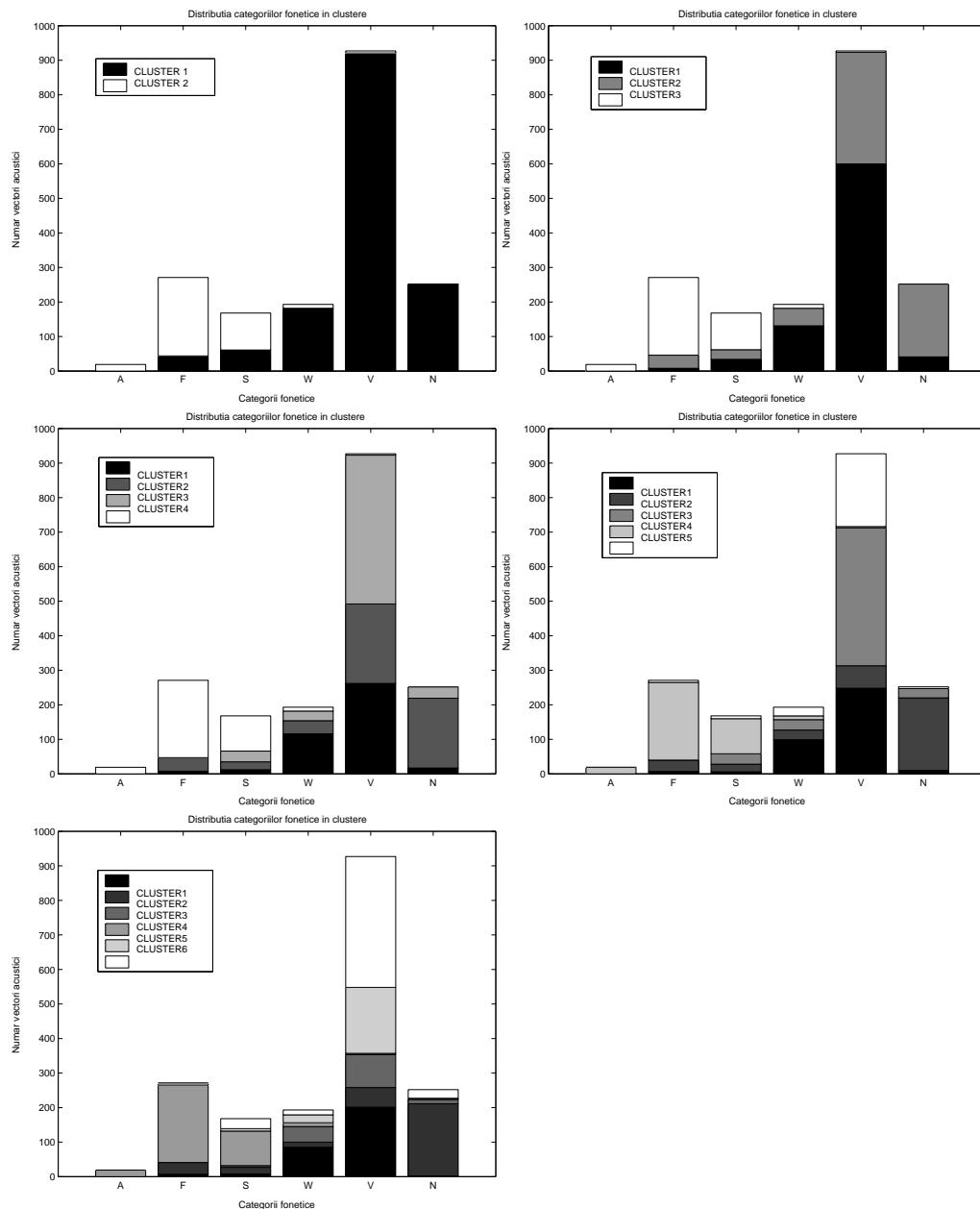


Figura 3.2: Distribuția categoriilor fonetice în clustere

Categorie fonetică	Varianță
Africate	23.22
Fricative	47.35
Plozive	37.85
Semivocale	46.76
Vocale	41.33
Nazale	20.53

Tabelul 3: Varianțe pe categorii fonetice

acustici obținuți de la un vorbitor, selectând din fiecare fonem doar cadrul din mijloc. Selectia cadrelor din mijloc este prezentată în figura 3.3.

Am etichetat fiecare vector acustic obținut dintr-un cadru extras din mijlocul fonemului cu eticheta acestuia. Ca etichete am utilizat codurile celor 39 de grupe de foneme prezentate în tabelul 4 și am obținut

$$Y = \{(x_1, f_1), (x_2, f_2), \dots, (x_T, f_T)\}, \quad (3.7)$$

unde  $f_i \in \{0, 1, \dots, 38\}$ . Am creat modele GMM cu câte 32 de componente gaussiene. După crearea modelelor, fiecare vector acustic a fost etichetat cu indicele acelei componente gaussiene care ne-a furnizat probabilitatea maximă pentru acel vector acustic. În final am obținut mulțimea

$$Z = \{(x_1, f_1, g_1), (x_2, f_2, g_2), \dots, (x_T, f_T, g_T)\}, \quad (3.8)$$

unde  $g_i \in \{1, 2, \dots, 32\}$ . Mulțimea  $\{(f_1, g_1), (f_2, g_2), \dots, (f_T, g_T)\}$  a fost reprezentată sub formă de diagramă. Figurile 3.4 și 3.5 prezintă diagramele astfel obținute pentru vorbitorul feminin *felcθ*, respectiv unul masculin *mdabθ*. Diagramele ne arată că, componentele gaussiene nu sunt pure, ele conținând vectori acustici din mai multe foneme. Plozivele și fricativele de obicei vor fi în aceeași componentă, iar vocalele sunt împărățiate în mai multe componente.

### 3.4 Modele GMM pure

Pentru a putea crea modele de vorbitori utilizând doar vectori acustici aparținând unei categorii fonetice specifice, vectorii acustici extrași din datele vorbitorilor trebuie etichetați cu categoria lor fonetică. Această etichetă se poate obține în două moduri: prin utilizarea unei baze de date segmentate, cum ar fi TIMIT, care conține și segmentarea fonetică a materialului audio sau prin crearea unui

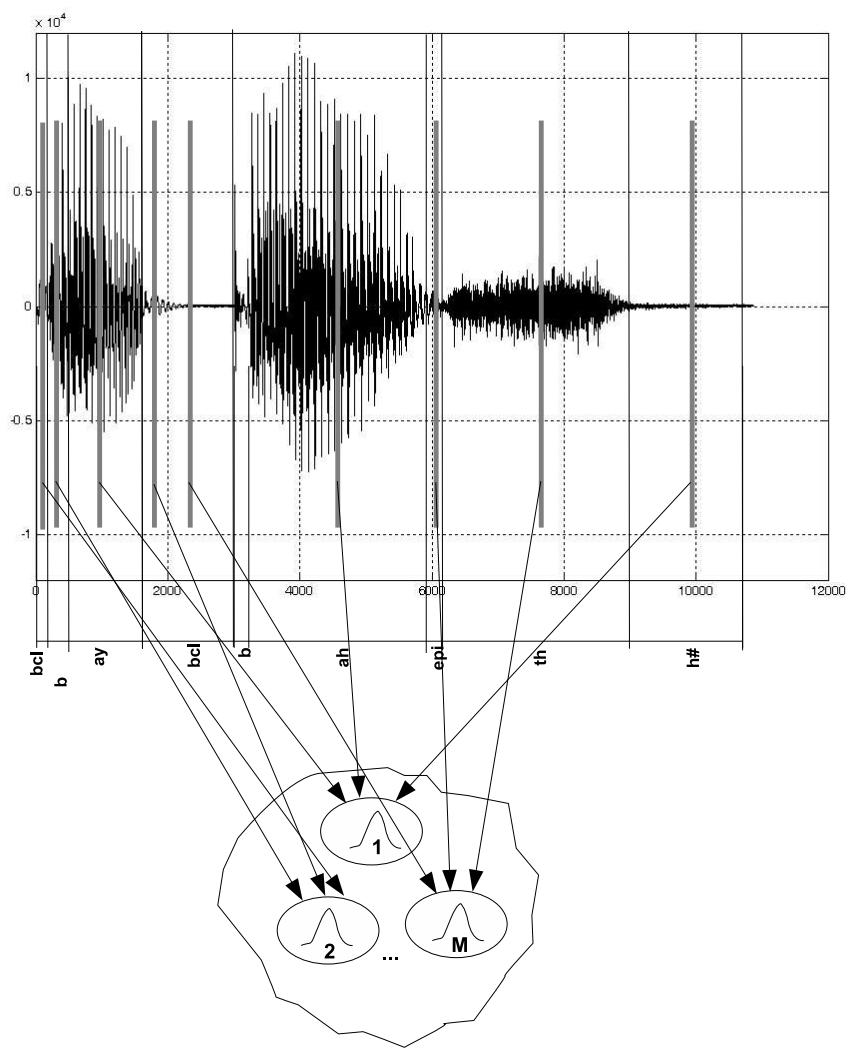


Figura 3.3: Selectia cadrelor din mijlocul fonemelor

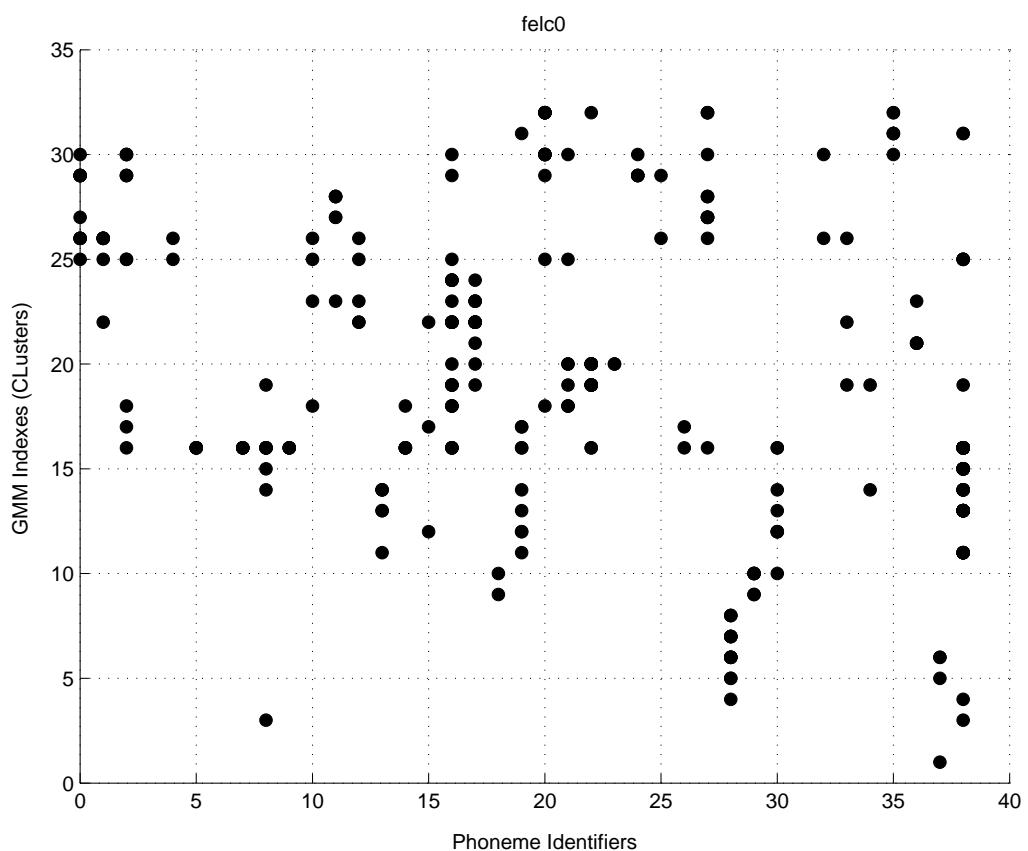


Figura 3.4: felc0-Distribuția fonemelor în componentele gaussiene ale unui model de vorbitor

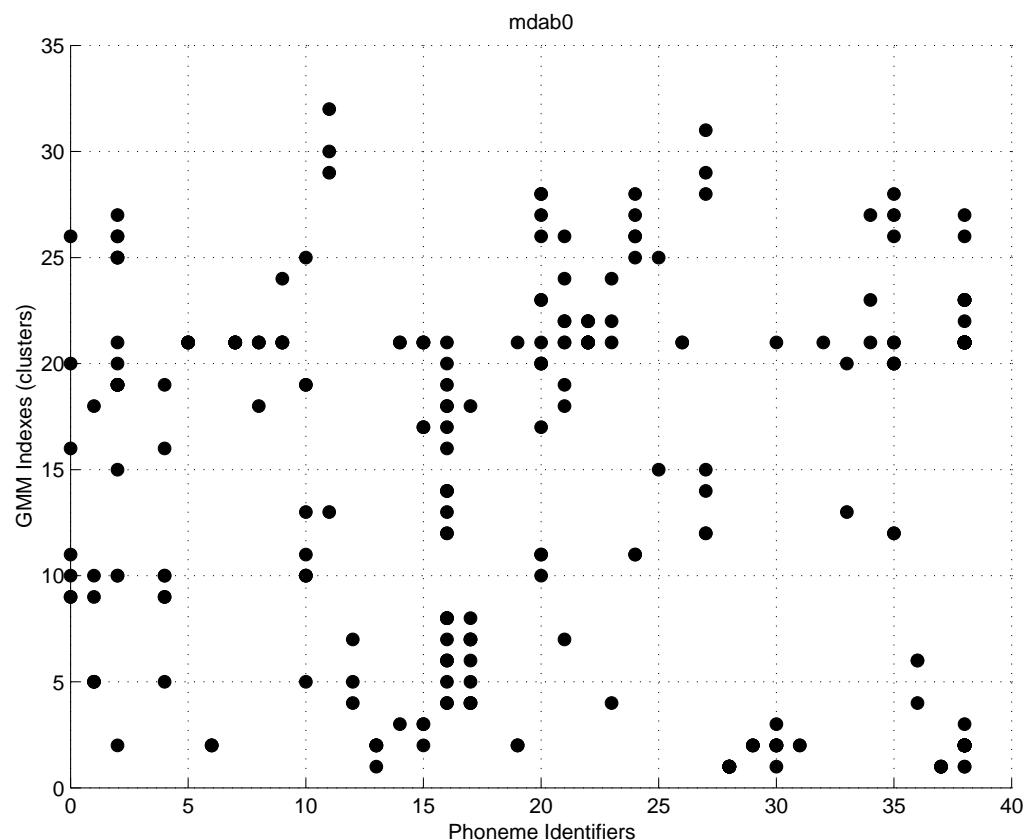


Figura 3.5: mdab0-Distribuția fonemelor în componentele gaussiene ale unui model de vorbitor

clasificator, care face clasificarea vectorilor acustici în categorii fonetice. Noi am abordat ambele metode, dar în acest capitol, pentru a nu fi nevoiți să considerăm erorile datorate acestui clasificator, vom prezenta doar rezultate obținute pentru prima variantă. Acuratețea sistemului nostru de clasificare în categorii fonetice a fost publicată în [7] și vom reveni la prezentarea acestor rezultate în capitolul 4.

Baza de date TIMIT conține 61 foneme distincte. Grupând fonemele similară se obțin 39 de grupe, care pot fi grupate în categorii de foneme ca vocale (V), nazale (N), fricative (F), africate (A), plozive (S), semivocale (W) și liniște, împreună cu părțile de închidere ale plozivelor (O). Pentru notația fonemelor se utilizează alfabetul fonetic ARPAbet descris în capitol doi și sunt prezentate în tabelul 4. Tabelul 5 prezintă numărul fonemelor, respectiv a grupelor de foneme pe categorii.

Am realizat o statistică privind cantitatea medie de date de antrenare, respectiv privind cantitatea medie de date de test medie existent pentru un vorbitor. Aceste date sunt prezentate în tabelul 6. Deoarece sunt vorbitori pentru care africatele lipsesc ori din partea de antrenare, ori din partea de testare, această categorie nu a fost luată în considerare. Pentru fiecare vorbitor am creat câte un model GMM utilizând vectorii acustici etichetați cu această categorie. După crearea modelelor utilizând datele de test, am măsurat rata de identificare a vorbitorilor, acesta fiind raportul dintre numărul de vorbitori identificați corect și numărul total de vorbitori. Pentru a determina numărul de componente optim pentru fiecare categorie în parte, la modelarea fiecărei categorii am pornit cu un singur gaussian și am crescut numărul acestora până când rata de identificare a atins maximul. Tabelul 6 prezintă pentru fiecare categorie în parte numărul de componente optim utilizat, rata de identificare, precum și un coeficient normalizat de pondere calculat din rata de identificare, care va fi folosit în crearea modelelor GMM structurate fonetic.

### 3.5 Modele GMM structurate fonetic

Testarea unui model structurat fonetic se desfășoară conform figurii 3.6. Testarea unui asemenea model se face similar cu testarea modelului standard, adică se calculează probabilitatea fiecărui vector acustic în fiecare componentă gaussiană, se înmulțește această probabilitate cu ponderea atașată fiecărei componente și se însumează, obținându-se probabilitatea vectorului acustic în model. Formula

Fonem	Cod	Cod grupă	Categorie	Fonem	Cod	Cod grupă	Categorie
aa	0	0	V	em	18	21	N
ae	1	1	V	en	19	22	N
ah	2	2	V	eng	20	23	N
ao	3	0	V	m	37	21	N
aw	4	3	V	n	38	22	N
ax	5	2	V	ng	39	23	N
ax-h	6	2	V	nx	40	22	N
axr	7	11	V	el	17	20	W
ay	8	4	V	hh	28	15	W
eh	16	10	V	hv	29	15	W
er	22	11	V	l	36	20	W
ey	23	12	V	r	47	27	W
ih	30	16	V	w	57	35	W
ix	31	16	V	y	58	36	W
iy	32	17	V	bcl	21	38	O
ow	41	24	V	dcl	27	38	O
oy	42	25	V	epi	44	38	O
uh	53	32	V	gcl	26	38	O
uw	54	33	V	h#	27	38	O
ux	55	33	V	kcl	35	38	O
ch	11	6	A	pau	44	38	O
jh	33	18	A	pcl	45	38	O
dh	14	8	F	q	46	38	O
f	24	13	F	tcl	51	38	O
s	48	28	F	b	9	5	S
sh	49	29	F	d	12	7	S
th	52	31	F	dx	15	9	S
v	56	34	F	g	25	14	S
z	59	37	F	k	34	19	S
zh	60	29	F	p	43	26	S
				t	50	30	S

Tabelul 4: Foneme TIMIT

Categorie	Număr grupe de foneme	Număr foneme
Vocale (V)	14	20
Nazale (N)	3	7
Fricative (F)	7	8
Semivocale (W)	5	7
Plozive (S)	7	7
Africate (A)	2	2
Altele (O)	1	10

Tabelul 5: Categorii de foneme TIMIT

Categorie	Antrenare	Testare	Nr. mixt.	Rata id.	Pondere
Vocale	9.65s	2.27s	<b>8</b>	<b>95.39%</b>	0.37
Nazale	1.34s	0.38s	<b>1</b>	<b>70.31%</b>	0.27
Fricative	3.27s	0.95s	4	44.60%	0.16
Semivocale	2.38s	0.47s	4	41.74%	0.16
Plozive	1.43s	0.35s	4	10.47%	0.04

Tabelul 6: Rata de identificare pentru 630 de vorbitori utilizând GMM-uri fonetic pure

exactă conform figurii 3.6

$$p(x|\lambda) = \sum_{i=1}^5 w_i p_i(x), \quad \sum_{i=1}^5 w_i = 1 \quad (3.9)$$

unde  $p_i(x)$  este probabilitatea calculată în submodelul GMM corespunzător categoriilor fonetice. Pentru un submodel atașat unei categorii fonetice probabilitatea se calculează conform formulei următoare

$$p_i(x) = \sum_{k=1}^{n_i} c_{ik} b_{ik}(x), \quad \sum_{k=1}^{n_i} c_{ik} = 1, \quad (3.10)$$

unde  $b_{ik}(x)$  este densitatea de probabilitate  $d$ -dimensională dată de formula 2.23.

În experimentele efectuate cu aceste modele structurate fonetic, pentru ponderile atașate submodelelor am folosit două abordări: ponderi egale, ponderi determinate pe baza puterii de discriminare a categoriilor fonetice. Fiecare categorie fonetică a fost modelată cu câte o componentă gaussiană, cu excepția vocalelor, pentru care am folosit două componente gaussiene, în total șase. În primul experiment, notat cu modele PSGMM1, am utilizat ponderi egale pentru fiecare com-

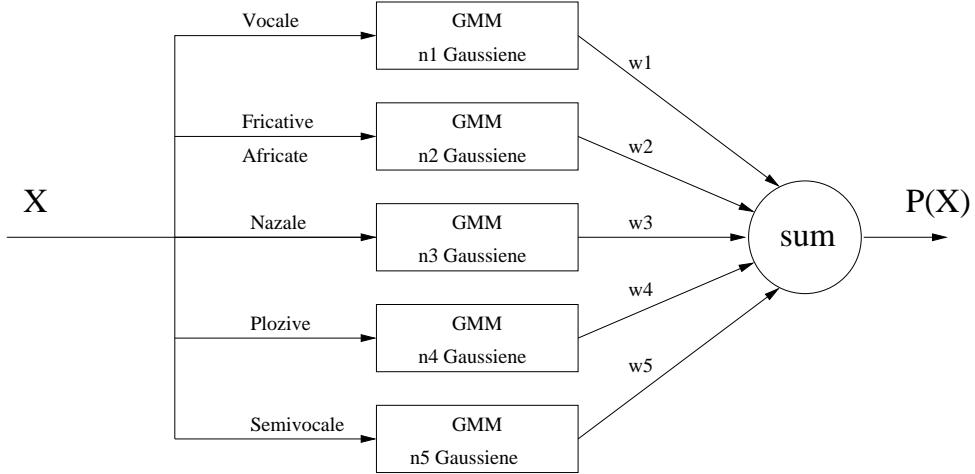


Figura 3.6: Calculul similarității într-un model de vorbitor structurat fonetic

ponentă gaussiană, iar în cel de-al doilea experiment, notat cu modele PSGMM2, ponderile utilizate au fost cele din ultima coloană a tabelului 6. Deoarece vocalele au fost modelate cu două componente gaussiene, ponderile obținute pentru aceste două componente,  $c_{11}$  și  $c_{12}$ , au fost înmulțite cu  $w_1 = 0.37$  satisfăcându-se astfel toate constrângerile necesare pentru model. Deoarece în partea de test a materialului vocal am avut două propoziții rostite, am repetat experimentul câte o dată pentru fiecare propoziție și am calculat media aritmetică a ratei de recunoaștere. Propozițiile rostite au o durată medie de 3 secunde. Experimentele au fost repete și pentru cele două propoziții, obținându-se în acest caz o îmbunătățire semnificativă a ratei de recunoaștere. Pentru comparație am inclus și rata de recunoaștere obținută pentru un model standard GMM având același număr de componente, adică 6.

Tabelul 7 prezintă rezultatele comparative.

Tip model	Nr.mixt.	Rata id. 3s test	Rata id. 6s test
GMM	6	93.02%	98.58%
PSGMM1	6	92.14%	98.74%
PSGMM2	6	92.06%	98.16%

Tabelul 7: Rate de recunoaștere pentru modele de vorbitori structurate fonetic - TIMIT 630 de vorbitori

Așa cum arată rezultatele, utilizând doar informații relative la categoria fonetică a vectorilor acustici, modelele create au performanțe similare cu cele standard.

Am realizat un studiu cu privire la vorbitorii identificați în mod incorect utilizând modelele standard GMM și cele două modele structurate fonetic. Tabelul 8 prezintă lista vorbitorilor identificați în mod incorect. Așa cum se poate vedea, ponderile nu influențează în mod semnificativ rata recunoașterii, mai mult, listele vorbitorilor identificați în mod greșit de către metodele PSGMM1 și PSGMM2 sunt similare. Nu același lucru se poate spune, dacă comparăm modelul standard cu cel structurat fonetic. În acest caz listele conțin vorbitori diferenți. De aici rezultă că, există posibilitatea de a combina cele două modele.

GMM	PSGMM1	PSGMM2
mgwt0	fram1	fram1
mmdm2	mcrc0	mjvw0
mthc0	mrdd0	mcrc0
mwar0	fdxw0	mrfk0
fmmh0	mdps0	mdlh0
mlbc0	mrfk0	fkkh0
fkkh0	fjxp0	flood0
mjpg0	fkkh0	mrmh0
mcre0	flood0	
	mcre0	
	mrmh0	

Tabelul 8: Lista identificatorilor vorbitorilor identificați în mod incorect

Pentru a putea caracteriza cele două tipuri de modele de vorbitori, am reprezentat probabilitățile cu care au fost identificați vorbitorii (vezi figura 3.7). Am ales o listă de vorbitori identificați în mod corect de ambele metode. Așa cum rezultă și din figură, modelul standard (GMM) indică aproape întotdeauna o probabilitate mai mare, decât modelele structurate fonetic (PSGMM1, PSGMM2).

## 3.6 Concluzii

În aceast capitol am prezentat diferite tehnici de modelare a vorbitorului. Toate modelele prezentate în acest capitol au fost implementate de către autor și testate pe baze de date publice. Publicații relevante pentru această problemă sunt:

[1, 3]: Compară metodele GMM și VQ pentru identificarea vorbitorului pe o bază de date proprie cu 66 de vorbitori. Deoarece nu au fost date suficiente pentru estimarea robustă a parametrilor modelului GMM, în acest caz modelele VQ au avut o performanță ușor ridicată față de performanțele modelelor GMM.

[2]: Compară metodele GMM și VQ pentru 168 de vorbitori din baza de date

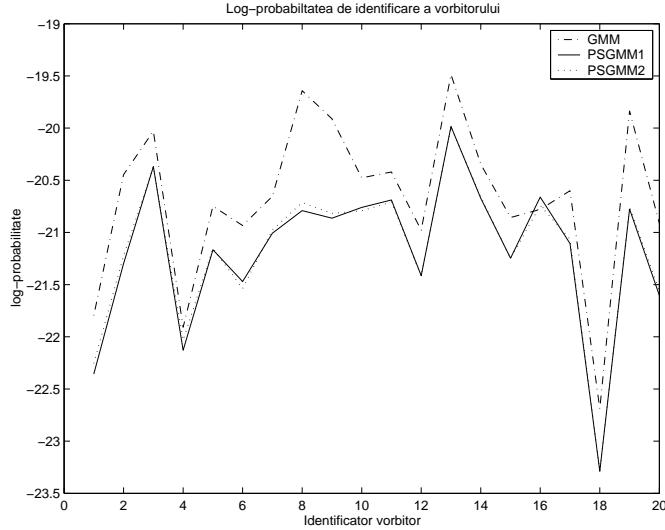


Figura 3.7: Probabilități de identificare a vorbitorilor pentru primii 20 de vorbitori

TIMIT. Modelele GMM s-au dovedit superioare în acest caz.

[4]: Lucrarea prezintă comparația diferenților algoritmilor de clustering pentru recunoașterea vorbitorului. Au fost comparate algoritmii: Randomized Generalized Lloyd Algorithm (RGLA), Modified Generalized Lloyd Algorithm (MGLA), Learning Vector Quantization (LVQ) și Gaussian Mixture Models (GMM). Concluzia acestei lucrări ar fi că, orice algoritm de clustering implementat în mod adecvat este corespunzător pentru modelarea vorbitorilor. Contează mai mult initializarea clusterelor, decât algoritmul în sine.

[6]: Lucrarea studiază puterea de discriminare a categoriilor fonetice pe un număr de 100 de vorbitori din baza de date TIMIT.

[8]: În această lucrare am comparat coeficienții mel cepstrali (MFCC) cu coeficienții cepstrali obținuți din parametrii LPC (LPCC). Toate măsurătorile au fost făcute pentru toți cei 630 de vorbitori din baza de date TIMIT. Am studiat și efectul utilizării doar a cadrelor din mijloc din fiecare fonem asupra recunoașterii vorbitorului. Am verificat, dacă clusterele unui model de vorbitor sunt sau nu fonetic pure.

Parametrii LPCC s-au dovedit ușor mai bune decât parametrii MFCC. Utilizând doar un singur cadru din fiecare fonem rezultatele au scăzut doar cu un singur procent față de utilizarea tuturor cadrelor. În sfârșit am arătat că un model de vorbitor nu separă fonemele vorbitorului în clustere diferite, cu toate că fiecare fonem are clusterele specifice. De obicei vocalele sunt împărățiate în mai multe

clustere, iar plozivele și fricativele în maxim două clustere. În acest experiment numărul clusterelor a fost 32 și cel al fonemelor 39.

[9, 11]: În aceste două lucrări am studiat puterea de discriminare a categoriilor fonetice. În acest scop am creat modele pure de vorbitori, adică conținând vectori acustici dintr-o singură categorie fonetică. Am ordonat categoriile fonetice din punct de vedere al puterii de discriminare al vorbitorului. Nazalele și vocalele s-au dovedit a fi cele mai bune din acest punct de vedere.

Rezultatele obținute pentru modele de vorbitori structurate fonetic urmează să fie publicate.



## Capitolul 4

# Clasificarea fonemelor cu metode statistice

### 4.1 Introducere

Pentru a construi un sistem de recunoaștere automată a vorbirii este necesară modelarea unei părți a vorbirii, din care se poate compune orice propoziție iar numărul acestor părți de vorbire să fie relativ restrâns. S-au propus mai multe asemenea unități, cum ar fi silaba, demisilaba, fenonul. Sistemul TANGORA de la IBM este bazat pe modele de tip fenon [25]. Totuși, majoritatea sistemelor comerciale pentru recunoașterea vorbirii sunt bazate pe foneme. Pentru a modela coarticularea prezentă în vorbire se utilizează trifoneme, adică triplete de foneme. Aceste modele se numesc și modele dependente de context, deoarece modelează nu numai fonemul ci și contextul acestuia. Într-un sistem de recunoașterea vorbirii o componentă utilă ar fi recunoașterea fonemelor. Recunoașterea însă, pe lângă modelarea fonemelor, necesită și tehnici de căutare prezentate în capitolul 2.6.4. Clasificarea fonemelor este utilă pentru evaluarea modelelor. Dacă la acestea se adaugă tehnici corespunzătoare de căutare, se poate ajunge și la recunoașterea fonemelor.

Au fost publicate câteva studii prezentând rezultate privind clasificarea fonemelor pe baza de date TIMIT utilizând modele neuronale recurente [73] sau modele Markov ascunse [22]. Majoritatea articolelor prezintă și rezultate experimentale pentru recunoașterea fonemelor [51, 73, 78]. Clasificarea presupune cunoașterea capetelor fonemelor, iar recunoașterea nu are informații legate de poziția fonemelor în semnalul vocal. Bineînțeles, recunoașterea este o problemă mai complexă.

Un studiu propriu, care efectuează clasificarea fonemelor din baza de date TIMIT este [5], iar articolul [7] propune o nouă schemă de clasificare pe două nivele. În continuare descriem modelele folosite în articolele menționate cu rezultate experimentale ușor îmbunătățite față de cele publicate. Aceste îmbunătățiri se datorează faptului că a fost îmbunătățită modalitatea de calcul a diferențelor de ordin unu și doi ale coeficienților cepstrali.

## 4.2 Analiza acustică

Forma de undă a fost împărțită pe cadre scurte de 20 de milisecunde cu o suprapunere de 50% între cadrele consecutive. Din fiecare cadru am extras 13 coeficienți Mel-cepstrali, conform procedurii descrise în [40]. Am utilizat 28 de filtre triunghiulare, trece-bandă, așezate conform scalei Mel. Am calculat coeficienții cepstrali cu formula

$$C_p = \sum_{r=1}^{28} S_r \cos \left( p \cdot \left( r - \frac{1}{2} \right) \cdot \frac{\pi}{28} \right), \quad 0 \leq p < 28 \quad (4.1)$$

unde  $S_r$  reprezintă logaritmul energiei la ieșirea filtrului cu numărul de ordine  $r$ . Am utilizat doar primii 13 coeficienți,  $C_0, C_1, \dots, C_{12}$ . Coeficienții dinamici au fost obținuți cu următoarea formulă

$$\Delta C_m(t) = \frac{\sum_{k=-K}^K k C_m(t+k)}{\sum_{k=-K}^K k^2} \quad (4.2)$$

În formula 4.2 pentru constanta  $K$  am utilizat valoarea 2. Coeficienții cepstrali împreună cu derivatele de ordin unu și doi, în total 39 de coeficienți, au format un vector acustic.

## 4.3 Evaluarea clasificării

Clasificarea a fost efectuată pe cele două baze de date segmentate TIMIT și OASIS. Rata de recunoaștere pentru un fonem este dată de formula

$$R_k = \frac{\text{numar\_foneme\_k\_corect\_clasificate}}{\text{numar\_total\_foneme\_k}} \quad (4.3)$$

iar rata globală de recunoaștere de formula

$$R_{tot} = \frac{\text{numar\_foneme\_corect\_clasificate}}{\text{numar\_total\_foneme}} \quad (4.4)$$

## 4.4 Rezultate experimentale

### 4.4.1 TIMIT

Baza de date TIMIT conține 61 foneme distințe. Grupând fonemele similare se obțin 39 de grupe, care pot fi grupate în categorii de foneme ca vocale (V), nazale (N), fricative (F), africate (A), plozive (S), semivocale (W) și liniște împreună cu părțile de închidere ale plozivelor (O). Pentru notația fonemelor se utilizează alfabetul fonetic ARPANET descris în capitolul 2 și este prezentat în tabelul 4. Ca și modele au fost utilizate modele GMM, valorile medii fiind inițializate cu valorile medii determinate de algoritmul de cuantizare vectorială binară (alg. 5). Măsurătorile au fost făcute odată pe setul de 61 foneme, după care pe setul de 39 de foneme. Modelele diferitelor foneme au avut același număr de componenete, acestea variind între 1 și 32. Rezultatele obținute pentru cele două seturi de foneme sunt prezentate în tabelele 10 și 11. Toate experimentele au fost făcute pe întreaga bază de date, adică antrenarea făcându-se cu fonemele a 462 vorbitori, iar evaluarea cu datele celor 168 de vorbitori rămași. Numărul exact de foneme din setul de test și cel de antrenare este prezentat în tabelul 9.

Figurile 4.1, 4.2 și 4.3 prezintă ratele de clasificare obținute pentru cele 61 de foneme din baza de date TIMIT cu 8, 16 respectiv 32 de componente.

Din matricea de confuzie am determinat erorile cele mai frecvente. Tabelul 12 conține primele 10, cele mai frecvente erori. Aceste erori sunt similare cu cele publicate în [73], obținute cu rețele neurale recurente. După eticheta fonemului, simbolul din paranteză reprezintă categoria fonemului.

Pe baza tabelului 12 și a matricei de confuzie (vezi figura 4.4), din care reiese că majoritatea fonemelor clasificate incorrect au fost confundate cu alte foneme din propria categorie, se propune o clasificare pe două niveluri prezentată în figura 4.5.

Pentru a realiza clasificatorul pe două niveluri, în primul pas am creat un clasificator pentru cele șapte categorii de foneme. Categoriile de foneme au fost modelate cu mixturi gaussiene, odată cu 32 de componente și încă odată cu 64 de componente. Performanțele acestui clasificator sunt prezentate în tabelul 13 și în figura 4.6. Ratele cele mai bune de clasificare au furnizat categoriile fonetice: africate, nazale și vocale.

Rezultatele obținute cu clasificatorul pe două niveluri sunt prezentate în tabelul 14. În paranteză sunt ratele de clasificare obținute cu metoda obișnuită, prezentate și în tabelul 10. Așa cum arată aceste rezultate, utilizând un clasificator cu două niveluri, performanța sistemului nu scade în mod semnificativ, dar

Fonem	Test	Antr.	Fonem	Test	Antr.	Fonem	Test	Antr.
aa	1133	3064	eng	5	38	nx	360	971
ae	1407	3997	epi	536	1464	ow	777	2136
ah	879	2306	er	800	2046	oy	263	684
ao	1156	2940	ey	806	2282	p	956	2588
aw	216	729	f	912	2216	pau	391	952
ax	1346	3610	g	750	2017	pcl	965	2644
ax-h	118	375	gcl	808	2223	q	1243	3590
axr	1383	3407	h#	3360	9240	r	2525	6539
ay	852	2390	hh	356	957	s	2639	7475
b	878	2181	hv	369	1154	sh	796	2238
bcl	776	1909	ih	1709	5051	t	1534	4364
ch	259	822	ix	2945	8642	tcl	2334	6644
d	1242	3548	iy	2710	6953	th	267	751
dcl	1643	4942	jh	372	1209	uh	221	535
dh	1053	2826	k	1614	4874	uw	170	555
dx	940	2709	kcl	1964	5859	ux	580	1908
eh	1440	3853	l	2356	5801	v	710	1994
el	343	951	m	1524	3903	w	1239	3140
em	47	124	n	2501	7068	y	634	1715
en	251	723	ng	414	1330	z	867	3773
						zh	44	151

Tabelul 9: Frecvența fonemelor în mulțimea de test și cea de antrenare

Număr componente	Rata de recunoaștere
1	45.55%
2	48.95%
4	52.98%
8	57.61%
16	60.84%
<b>32</b>	<b>63.34%</b>

Tabelul 10: Rata de recunoaștere pentru setul TIMIT-61

Număr de componente	Rata de recunoaștere
1	52.68%
2	55.12%
4	58.81%
8	63.09%
16	65.48%
<b>32</b>	<b>68.13%</b>

Tabelul 11: Rata de recunoaștere pentru setul TIMIT-39

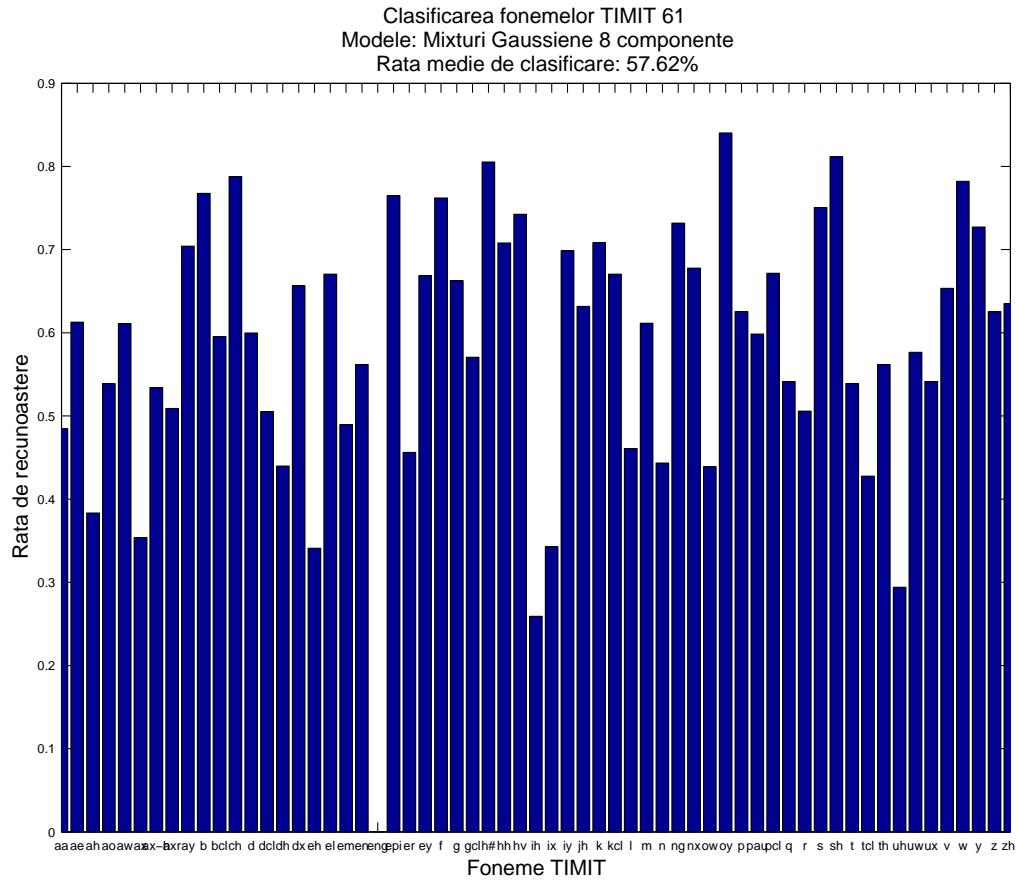


Figura 4.1: Clasificare foneme - TIMIT 61 - GMM 8

Nr.	Fonem inițial	Fonem recunoscut	Procent de eroare
1	/ix/ (V)	/ih/ (V)	1.59%
2	/r/ (W)	/er/ (V)	1.39%
3	/ix/ (V)	/ax/ (V)	1.35%
4	/s/ (F)	/z/ (F)	1.29%
5	/z/ (F)	/s/ (F)	1.29%
6	/ih/ (V)	/ix/ (V)	1.25%
7	/eh/ (V)	/ae/ (V)	0.95%
8	/eh/(V)	/ah/ (V)	0.95%
9	/n/ (N)	/m/ (N)	0.95%
10	/r/ (W)	/axr/ (V)	0.95%

Tabelul 12: Cele mai frecvente erori - TIMIT 61

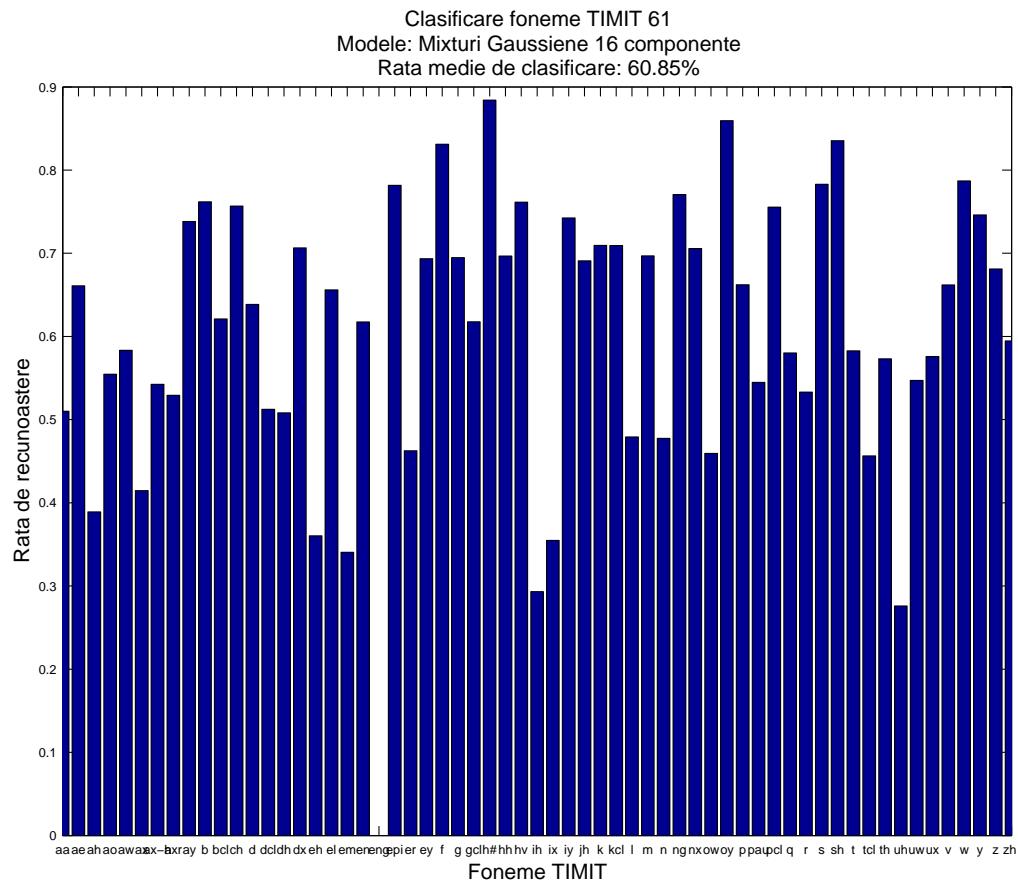


Figura 4.2: Clasificare foneme - TIMIT 61 - GMM 16

Categorie fonetică	GMM 32	GMM 64
Vocale (V)	83.99%	85.10%
Nazale (N)	88.85%	89.95%
Fricative (F)	78.82%	81.20%
Semivocale (W)	77.26%	80.02%
Plozive (S)	82.46%	84.80%
Africate (A)	92.55%	92.55%
Liniște (O)	87.90%	88.30%
<b>Medie</b>	<b>83.68%</b>	<b>85.13%</b>

Tabelul 13: Clasificare categorii fonetice TIMIT

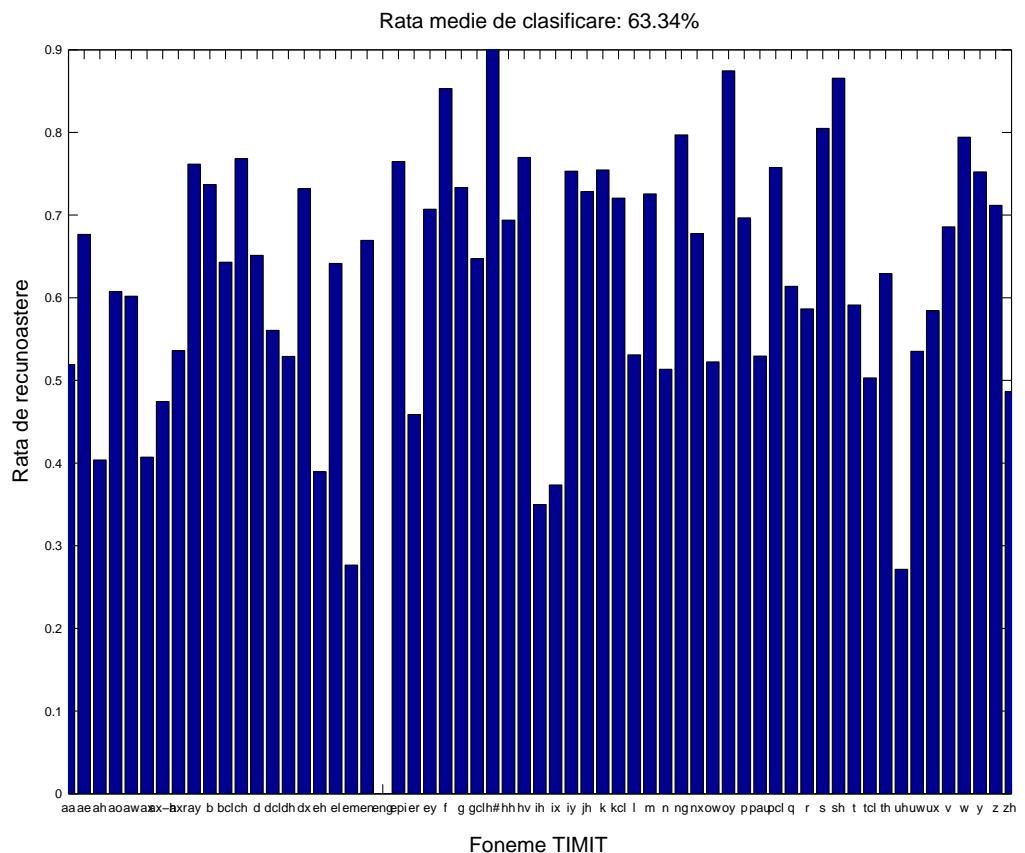


Figura 4.3: Clasificare foneme - TIMIT 61 - GMM 32

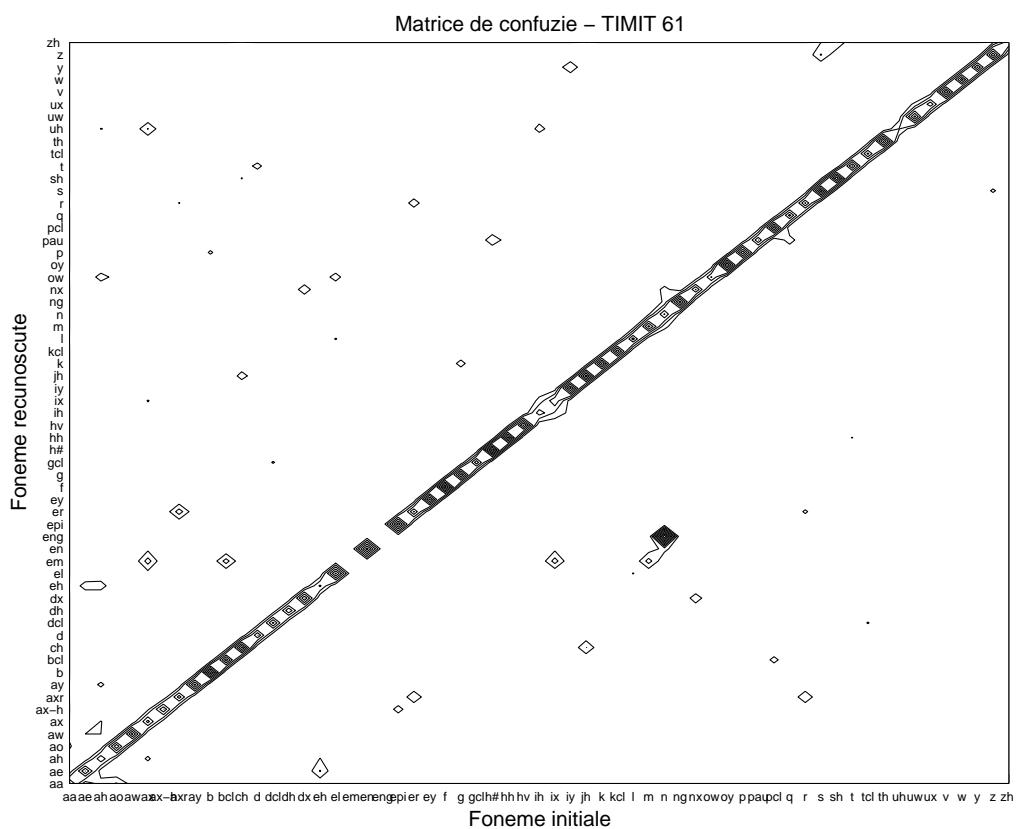


Figura 4.4: Matricea de confuzie TIMIT 61

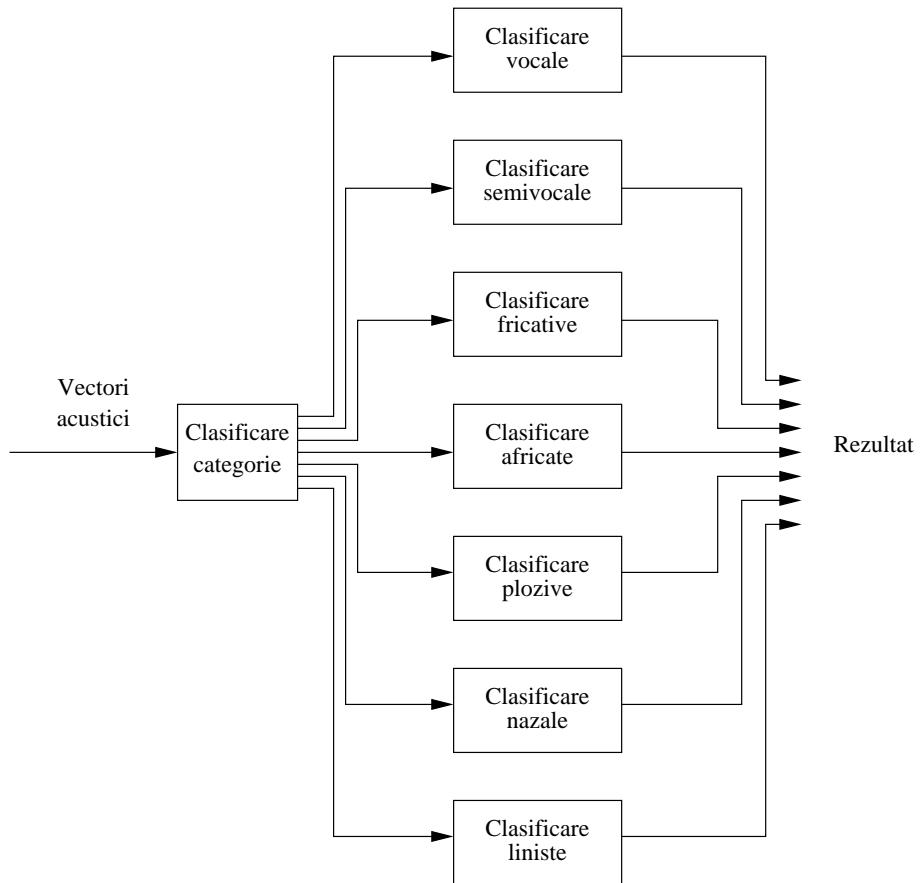


Figura 4.5: Clasificator cu două niveluri

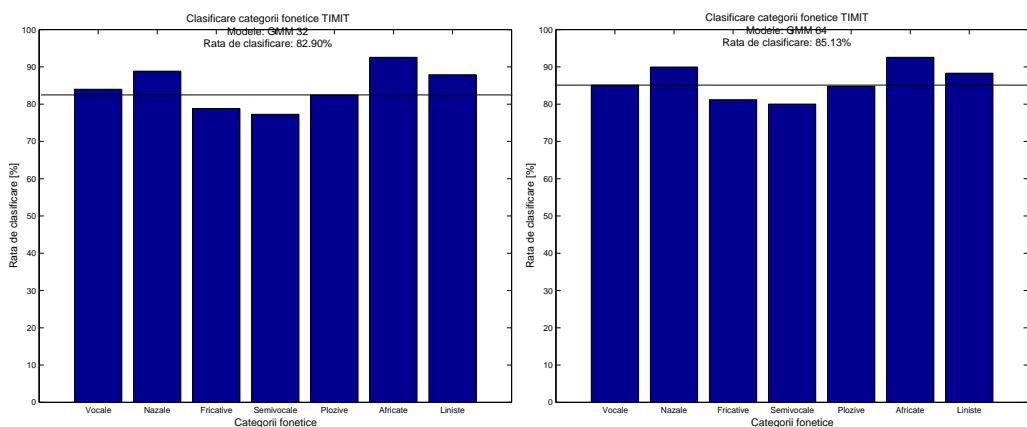


Figura 4.6: Clasificare categorii fonetice TIMIT

timpul de procesare se reduce considerabil. În cazul vocalelor se vor face 7+20 de comparații în loc de 61. În primul nivel întotdeauna se fac 7 comparații, deoarece în total avem șapte categorii, iar pe al doilea nivel numărul comparațiilor depinde de numărul fonemelor din categoria respectivă. În total sunt 20 de vocale. În cazul nazalelor se vor face 7+7 comparații, iar în cazul africatelor 7+2 comparații.

	Model categ. GMM 32	Model categ. GMM 64
Model fonem GMM 8	55.51% (57.61%)	56.64% (57.61%)
Model fonem GMM 16	58.00% (60.84%)	59.21% (60.84%)
Model fonem GMM 32	59.61% (63.34%)	60.84% (63.34%)

Tabelul 14: Clasificare foneme TIMIT 61 cu clasificatorul pe două niveluri

Ultimele experimente au fost făcute pentru clasificarea cadrelor (vectorilor acustici). Acest tip de clasificare se mai numește și clasificare cadru cu cadru (frame by frame classification). și acest tip de clasificare urmărește evaluarea modelelor. Experimentul a fost repetat de trei ori. Prima dată s-a folosit câte un cadru, după aceea am crescut cu doi numărul cadrelor din secvența de clasificat. Tabelul 15 prezintă rezultatele obținute în cazul modelelor gaussiene cu 32 componente.

Nr. cadre	Rată de clasificare cadru cu cadru
1	44.29%
3	46.88%
5	49.21%
7	49.30%
9	48.27%

Tabelul 15: Clasificare cadru cu cadru-TIMIT 61

#### 4.4.2 OASIS Numbers

Rezultatele prezentate în această lucrare se referă la toate cele 31 foneme, chiar dacă câteva pot fi fuzionate, reprezentând alofoanele aceluiași fonem. În articolele [48, 76] rezultatele sunt prezentate referitor la cele 28 de foneme. Rezultatele noastre în cazul aplicării metodei GMM sunt mai bune decât cele prezentate în [48] (79.20%) cu aceeași metodă, deși ei au lucrat doar cu 28 de foneme, grupând alofoanele. Cele mai bune rezultate au fost obținute de ei prin aplicarea transformării LDA pentru vectorii acustici (86.76%). În [76] modelarea fonemelor a fost realizată cu pachetul HTK, utilizând modele Markov continue cu 3 stări ascunse

și obținând o rată de recunoaștere de 90.48%.

Tabelul 16 prezintă frecvențele de apariție a diferitelor foneme în partea de antrenare respectiv testare.

Fonem	Antrenare	Test	Fonem	Antrenare	Test
-	519	156	i:	40	12
-:	40	12	j	80	24
'd'	118	36	J	80	24
'k'	200	60	l	200	60
't'	399	120	l:	40	12
'ts	200	60	m	120	36
+	120	36	n	559	168
~	2080	624	O	240	72
2	80	24	o	160	48
:2	40	12	o:	40	12
A:	120	36	r	160	48
E	600	180	s	160	48
e:	160	48	u	80	24
h	306	96	u:	40	12
i	240	72	v	240	72
			z	240	72

Tabelul 16: Frecvența fonemelor în mulțimea de antrenare și test

Tabelul 17 prezintă ratele de clasificare obținute pe această bază de date folosind ca modele de foneme mixturi gaussiene cu 1, 2, 4, 8, 16 și 32 de componente.

Figurile 4.7 și 4.8 prezintă matricea de confuzie respectiv ratele de clasificare pentru cele 31 de foneme OASIS obținute cu modele având 16 componente gaussiene.

Din matricea de confuzie am determinat erorile cele mai frecvente. Tabelul 18 conține primele 20, cele mai frecvente erori. Putem constata că erorile cele mai frecvente sunt între foneme, care au câte două forme, una scurtă și una lungă. În cazul limbilor finugorice fiecare vocală are atât formă lungă cât și formă scurtă. Pentru recunoașterea acestor limbaje, modelarea duratei fonemelor este o problemă critică.

Ultimele experimente au fost făcute pentru clasificarea cadrelor (vectorilor acustici). Acest tip de clasificare se mai numește și clasificare cadru cu cadru (frame by frame classification). Și acest tip de clasificare urmărește evaluarea modelelor. Experimentul a fost repetat de cinci ori. Prima dată s-a folosit câte un cadru, după aceea am crescut cu doi numărul cadrelor din secvența de clasificat. Cadrul din mijloc determină întotdeauna apartenența la o fonemă a secvenței de

Număr componente	Rata de recunoaștere
1	80.40%
2	87.78%
4	90.84%
8	92.27%
<b>16</b>	<b>92.74%</b>
32	90.84%

Tabelul 17: Clasificare foneme OASIS 31

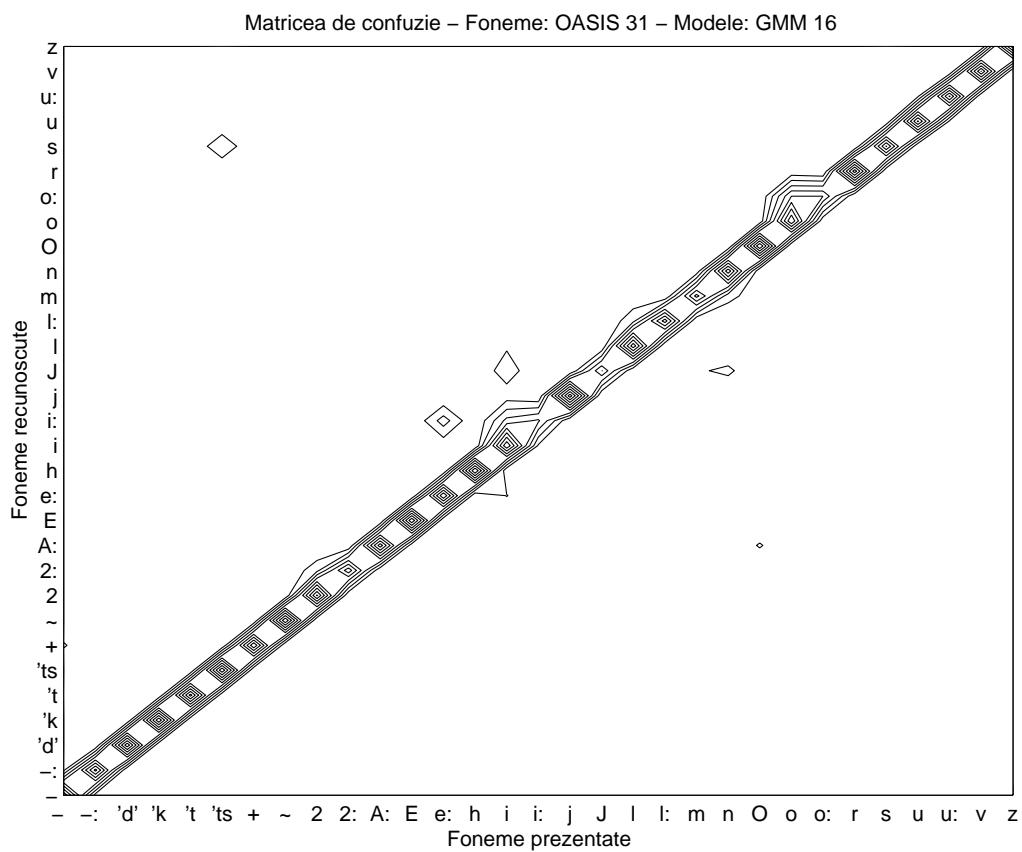


Figura 4.7: Matricea de confuzie OASIS 31

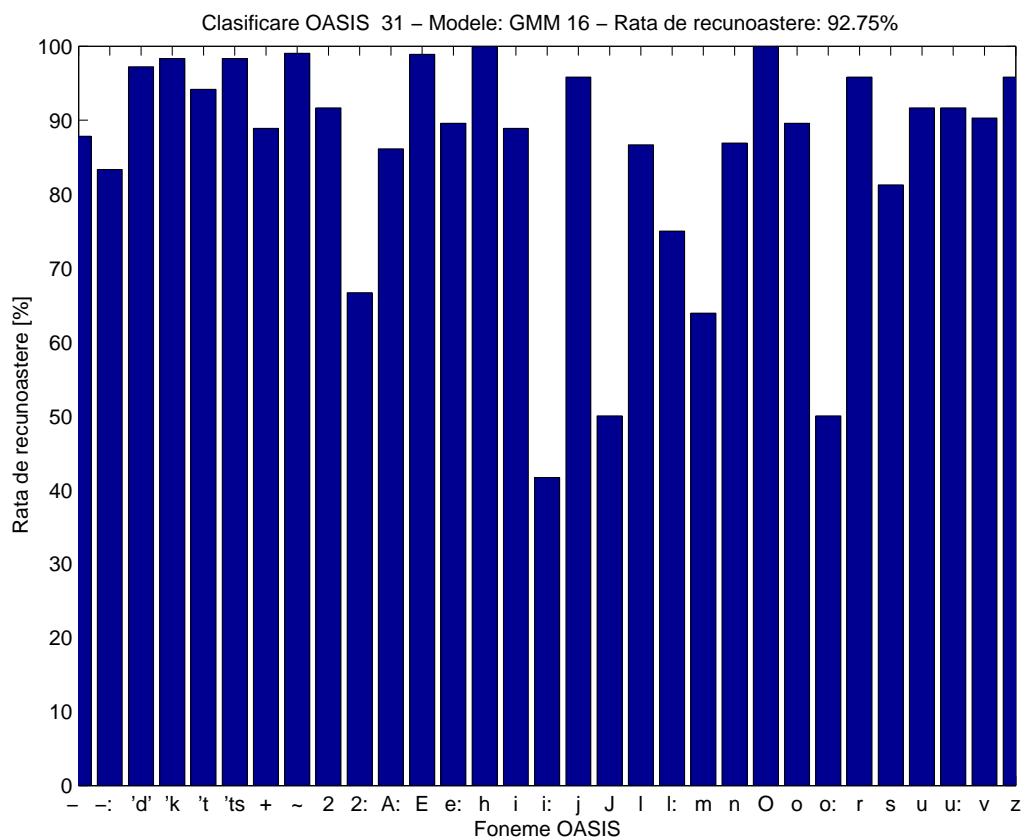


Figura 4.8: Ratele de clasificare OASIS 31, modele cu 16 componente gaussiene

Nr.	Fonem initial	Fonem recunoscut	Procent de eroare
1.	/-/	/-:/	5.35%
2.	/s/	/'ts/	5.35%
3.	/i/	/i:/	3.57%
4.	/o:/	/o/	3.57%
5.	/e:/	/i/	2.97%
6.	/l/	/i/	2.97%
7.	/m/	/n/	2.97%
8.	/-/	/'ts/	2.38%
9.	-	+	2.38%
10.	+	-	2.38%
11.	/A:/	/O/	2.38%
12.	/i:/	/i/	2.38%
13.	/J/	/i/	2.38%
14.	/n/	/~/	2.38%
15.	/n/	/i/	2.38%
16.	/n/	/m/	2.38%
17.	/'t/	/-/-	1.78%
18.	/~/	/h/	1.78%
19.	/i:/	/e:/	1.78%
20.	/J/	/n/	1.78%

Tabelul 18: Cele mai frecvente erori - TIMIT 61

cadre. Tabelul 19 prezintă rezultatele obținute în cazul modelelor gaussiene cu 16 componente.

Număr cadre	Rată de clasificare
1	63.36%
3	70.44%
5	74.03%
7	76.10%
9	77.36%

Tabelul 19: Clasificare cadru cu cadru - OASIS 31

## 4.5 Concluzii

Clasificarea fonemelor se utilizează în recunoașterea vorbirii pentru evaluarea modelelor. Cu modele de foneme performante avem șansa de a construi sisteme de recunoașterea vorbirii cu o acuratețe ridicată. În acest capitol am prezentat sistemul realizat de autor pentru clasificarea fonemelor. Ca și modele de foneme am utilizat modele GMM și am arătat că aceste modele ating performanțele modelelor HMM cu 3 stări. Cu aceasta am demonstrat faptul că probabilitățile de tranziție între stările unui model Markov ascuns sunt nesemnificative comparativ cu probabilitățile de emisie ale stărilor.

Pentru clasificarea fonemelor din baza de date TIMIT, toate articolele prezintă rezultate obținute pe setul de 39 de foneme, de aceea am putut compara doar rezultatele obținute pe această mulțime. Tabelul 20 este un tabel comparativ care se referă la clasificarea celor 39 de grupuri de foneme.

Articol	Vector acustic	Model	#mixt.	Clasif.
[22]	MFCC26	HMM-3	15	58.97%
[56]	LPCC24	HMM-3	24	57.10%
[61]	MFCC36	HMM-3	nespec.	63.00%
<b>Această lucrare</b>	<b>MFCC39</b>	<b>GMM</b>	<b>32</b>	<b>68.13%</b>

Tabelul 20: Rezultate comparative pentru clasificarea fonemelor TIMIT 39

În afară de clasificarea fonemelor și a grupurilor de foneme am clasificat și categoriile fonetice. Toate rezultatele au fost publicate în articolele [5, 7]. Articolul [7] introduce un nou tip de clasificator, clasificatorul cu două niveluri. Rezultatele prezentate în acest capitol sunt îmbunătățite față de cele publicate, deoarece am îmbunătățit extragerea coeficienților cepstrali, derivele de ordin 1 și 2.

Pe lângă baza de date TIMIT, am testat performanțele sistemului și pe baza de date OASIS Numbers. Am comparat și pe această bază de date rezultatele noastre cu cele obținute de către alte grupuri de cercetători. Rezultatele obținute de autor cu modele GMM sunt la fel de bune ca cele obținute cu modele HMM. Tabelul 21 prezintă rezultatele comparative. Cele două lucrări, cu care am comparat rezultatele noastre, au utilizat un set de 29 de foneme, grupând alofoanele. Noi am făcut experimentele fără această grupare, utilizând toate cele 31 de foneme.

Articol	Vector acustic	Model	#mixt.	Clasif.
[77]	<b>MFCC39</b>	HMM	24	90.48%
[48]	<b>MFCC39+LDA</b>	GMM	3	86.76%
<b>Această lucrare</b>	<b>MFCC39</b>	<b>GMM</b>	<b>16</b>	<b>92.74%</b>

Tabelul 21: Rezultate comparative pentru clasificarea fonemelor OASIS 31

## Capitolul 5

# Recunoașterea fonemelor cu modele GMM

### 5.1 Introducere

Analiza acustică a fost identică cu cea de la clasificarea fonemelor (vezi 4.2). Ca modele am utilizat modele Markov ascunse cu observații continue descrise în secțiunea 2.5.4.1. Fie  $O = \{o_1, o_2, \dots, o_T\}$  secvența de observații, pentru care trebuie determinată secvența de foneme, pe care aceasta o reprezintă. Dacă notăm cu  $F$  mulțimea tuturor secvențelor de foneme posibile, problema recunoașterii se poate formula, ca determinarea acelei secvențe de foneme  $\hat{f} \in F$ , care verifică următoarea egalitate

$$\hat{f} = \arg \max_{f \in F} P(f|O) = \arg \max_{f \in F} \frac{P(O|f) \cdot P(f)}{P(O)} \quad (5.1)$$

unde  $P(f)$  este modelul limbaj (la nivel de foneme). Presupunând că toate secvențele de observații sunt echiprobabile, ecuația 5.1 devine

$$\hat{f} = \arg \max_{f \in F} P(O|f) \cdot P(f) \quad (5.2)$$

Deoarece fiecare fonem a fost modelat cu un model Markov cu o singură stare, problema găsirii celei mai probabile secvențe de stări este identică cu găsirea celei mai probabile secvențe de foneme, problemă care este rezolvată de algoritmul Viterbi și care se bazează pe calcularea probabilităților înainte pentru fiecare secvență de timp conform ecuației

$$\alpha_t(j) = \max_{1 \leq i \leq n} \alpha_{t-1}(i) \cdot a_{ij} \cdot b_j(o_t). \quad (5.3)$$

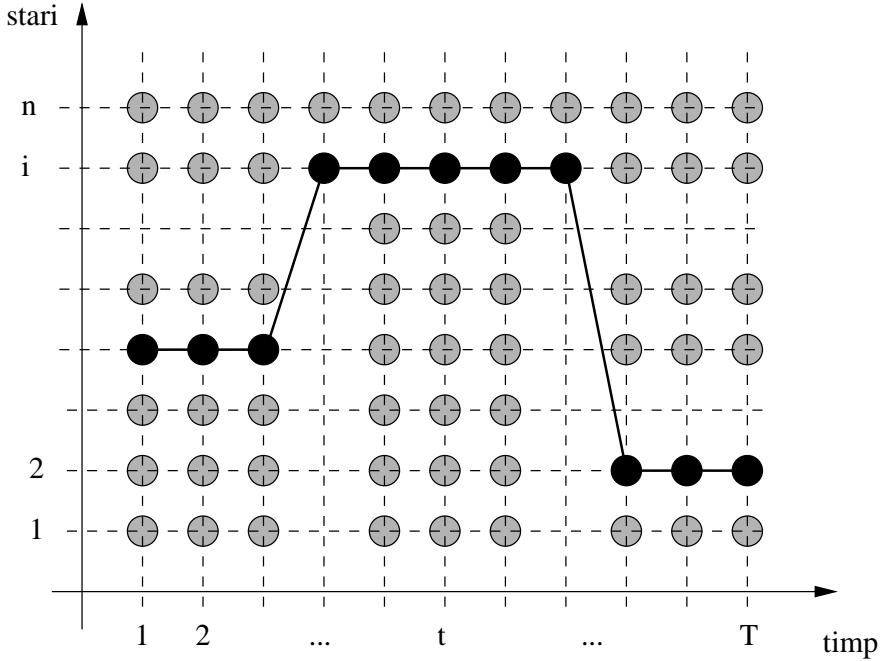


Figura 5.1: Algoritmul Viterbi - Căutare în spațiul stărilor

Algoritmul Viterbi realizează o căutare în spațiul stărilor (fonemelor), determinând secvența cea mai probabilă de stări (foneme) pentru un sir de vectori acustici. Practic realizează segmentarea sirului de vectori acustici în stări (foneme). Rezultatul căutării este reprezentat în figura 5.1.

În această lucrare a fost utilizată forma logaritmată a formulei 5.3, adică

$$\log \alpha_t(j) = \max_{1 \leq i \leq n} \{ \log \alpha_{t-1}(i) + \log a_{ij} \} + \log b_j(o_t) \quad (5.4)$$

Deoarece am utilizat modele HMM cu o singură stare, nu există probabilitate de trecere de la o stare la alta (lipsește  $a_{ij}$ ). Aplicând algoritmul Viterbi clasic am obținut rezultate foarte slabe, numărul inserărilor fiind foarte mare. De aceea, am introdus în formula 5.4 în loc de  $a_{ij}$  o penalizare de inserare notată cu

$$I_{ij} = \begin{cases} \beta, & \text{daca } i = j \\ 0, & \text{caz contrar} \end{cases} \quad (5.5)$$

cu care am reușit să controlăm numărul inserărilor. De fapt această valoare empirică a fost utilizată și pentru balansarea erorilor de ștergere și inserare. Astfel

formula 5.4 se modifică în

$$\log \alpha_t(j) = \max_{1 \leq i \leq n} \{\log \alpha_{t-1}(i) + I_{ij}\} + \log b_j(o_t) \quad (5.6)$$

Constanta  $\beta$  a fost determinată în mod empiric. În secțiunile următoare prezentăm determinarea acesteia pentru cele două baze de date TIMIT și OASIS Numbers.

## 5.2 Rezultate experimentale

Pentru evaluarea rezultatelor am utilizat metrica standard, rata de eroare la nivel de foneme, notată cu PER. Această măsură calculează diferența dintre secvența de foneme de intrare și secvența de foneme recunoscută. Distanța se calculează cu algoritmul prezentat în secțiunea 2.6.3, doar că în acest caz se notează cu PER în loc de WER, fiind vorba de recunoașterea fonemelor și nu a cuvintelor.

$$PER = 100 \cdot \frac{INS + SUBST + DEL}{N} \quad (5.7)$$

unde  $N$  reprezintă numărul de foneme în propoziția de recunoscut. Acuratețea sistemului se va calcula cu formula

$$ACCURACY = 100 - PER \quad (5.8)$$

O altă măsură utilizată este numărul corect de foneme returnate de către sistemul de recunoaștere. Această valoare se calculează tot cu algoritmul 10. Noi o vom nota în continuare cu *CORRECT*.

### 5.2.1 TIMIT

Măsurătorile în acest caz au fost făcute pentru două seturi de date de test, numite și **timit\_test\_core** respectiv **timit\_test**. Setul de test **timit\_test\_core** conține în total 192 de propoziții, iar **timit\_test** conține câte 10 propoziții de la 168 de vorbitori, deci în total 1680 de propoziții. Rezultatele pentru setul **timit\_test** sunt aproximativ cu 4% mai bune decât cele pentru setul **timit\_test\_core**. Tabelele 22 și 23 prezintă rezultatele obținute pentru recunoașterea celor 61 de foneme din baza de date TIMIT obținute pe seturile de test **timit\_test\_core** respectiv **timit\_test**. În ambele cazuri am folosit ca și modele de foneme mixturi gaussiene cu câte 32 de componente. Setul **timit\_test\_core** conține în total 7525 de foneme, iar setul **timit\_test** în total 65825 de foneme.

Mentionăm că rezultatele au fost obținute fără a utiliza model de limbaj la nivel de foneme. Acest lucru înseamnă că nu a fost făcută nici un fel de constrângere, orice fonem poate să urmeze după oricare altul. Conținutul tabelului 23 este reprezentat grafic în figura 5.2.

$\beta$	DEL	INS	SUBST	ACCURACY	CORRECT
10	6.18%	9.49%	38.77%	45.55%	55.04%
11	7.41%	7.68%	37.76%	47.14%	54.81%
13	9.75%	4.81%	36.37%	49.06%	53.87%
<b>20</b>	<b>18%</b>	<b>1.26%</b>	<b>30.88%</b>	<b>49.86%</b>	<b>51.12%</b>

Tabelul 22: Recunoașterea fonemelor - **timit\_test\_core**

$\beta$	DEL	INS	SUBST	ACCURACY	CORRECT
11	6.79%	7.39%	34.54%	51.27%	58.66%
13	9.13%	4.76%	33.12%	52.98%	57.74%
<b>15</b>	<b>11.42%</b>	<b>3.24%</b>	<b>31.66%</b>	<b>53.67%</b>	<b>56.91%</b>
20	17.16%	1.31%	28.12%	53.40%	54.71%

Tabelul 23: Recunoașterea fonemelor - **timit\_test**

### 5.2.2 OASIS Numbers

În această secțiune prezentăm rezultatele sistemului pentru recunoașterea fonemelor din baza de date OASIS Numbers. Fonemele au fost modelate cu mixturi gaussiene cu 16 componente și setul de test conține în total 2316 foneme. Tabelul 24 prezintă rezultatele obținute pentru diferite valori ale parametrului  $\beta$ . Figura 5.3 prezintă sub formă grafică rezultatele din tabelul 24.

$\beta$	DEL	INS	SUBST	ACCURACY	CORRECT
11	0.52%	27.93%	5.18%	66.36%	94.30%
13	0.82%	21.80%	5.53%	71.85%	93.65%
20	2.37%	11.27%	6.00%	80.35%	91.62%
<b>25</b>	<b>4.06%</b>	<b>7.77%</b>	<b>5.78%</b>	<b>82.38%</b>	<b>90.15%</b>
27	4.83%	7.08%	5.83%	82.25%	89.33%
30	6.43%	5.82%	5.70%	82.03%	87.86%

Tabelul 24: Recunoașterea fonemelor - **oasis\_test**

Rezultatele noastre se compară în mod favorabil cu cele publicate în [80] pe aceeași bază de date. Fără utilizarea cunoștințelor rezultate din vocabular, adică

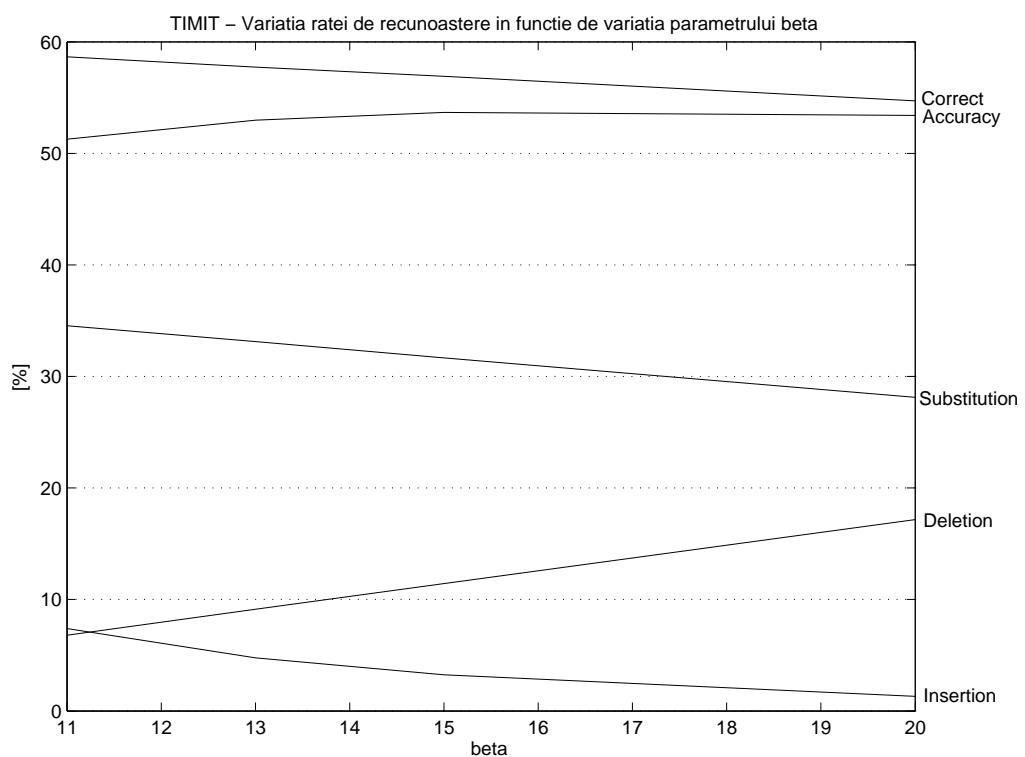


Figura 5.2: TIMIT-Variatia ratei de recunoastere în funcție de variația parametrului  $\beta$

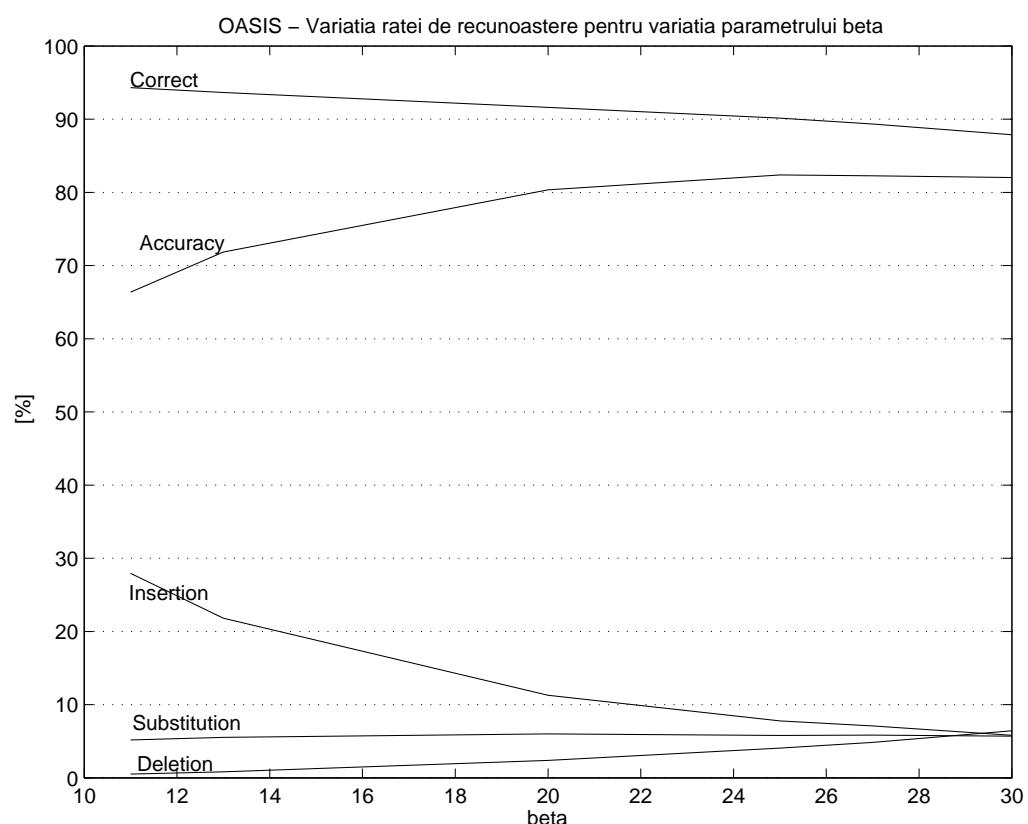


Figura 5.3: OASIS-Variatia ratei de recunoastere în funcție de variația parametrului  $\beta$

orice fonem poate să urmeze după oricare altul, creatorii acestei baze de date au obținut rata de recunoaștere maximă de 82.05%. Noi am obținut rata maximă de 82.38%, iar cu balansarea erorilor de inserare respectiv de ștergere, 82.03% (ultimul rând din tabelul 24). Utilizând un model bigram la nivel de foneme această rată de recunoaștere a fost ridicată la 96.87%.

### 5.3 Concluzii

În acest capitol am prezentat sistemul realizat de autor pentru recunoașterea fonemelor. De fapt acest capitol se bazează pe cel precedent, clasificarea fonemelor, deoarece modelele de foneme utilizate au fost evaluate în capitolul precedent. Pentru a putea detecta capetele fonemelor în vorbirea continuă este nevoie de un algoritm eficient de căutare. În literatura de specialitate sunt recomandați algoritmii: Viterbi clasic, Viterbi fascicular, respectiv algoritmul cu stivă. Noi am aplicat algoritmul Viterbi clasic. Din cauză că, acesta producea un număr mare de erori de inserare pentru modelele de foneme independente de context, am introdus o modificare în acest algoritm. Această modificare se referă la introducerea unei constante care controlează trcerea de la o stare la alta. Cu cât este mai mare valoarea acestei constante, cu atât este mai defavorizată trcerea la o stare (fonem) nouă. Cu această constantă se poate regla și raportul dintre numărul de erori de tip inserare, respectiv ștergere. Cu cât este mai mare valoarea acestei constante, cu atât este mai mare numărul erorilor de tip ștergere, deci este direct proporțională cu acest tip de eroare. Cu cât este mai mare valoarea acestei constante, cu atât este mai mic numărul erorilor de tip inserare, deci este invers proporțională cu acesta. Am prezentat o serie de experimente cu diferite valori ale acestei constante atât pentru baza de date TIMIT, cât și pentru OASIS Numbers. Așa cum arată experimentele, valoarea acestui parametru depinde de baza de date, adică de durata medie a fonemelor din limba respectivă.

Rezultatele au fost comparate doar în cazul bazei de date OASIS, deoarece pentru TIMIT nu există rezultate publicate pentru setul de 61 de foneme utilizând modele statistice. Toate rezultatele privind recunoașterea fonemelor pentru baza de date TIMIT, utilizând modele statistice, sunt publicate pentru setul restrâns de 39 de foneme. Doar articolele care prezintă experimentele realizate cu rețele neuronale prezintă rezultate pentru toate cele 61 de foneme. Însă aceste modele nu sunt comparabile cu modelele de foneme GMM independente de context, eventual cu cele dependente de context, deoarece rețelele neuronale descrise au capacitatea de a învăța contextul. Primul articol care publică rezultate pentru

setul **timit\_test** cu 39 de foneme este [51], obținând 58.77% cu modele independente de context și fără model de limbaj. În articolul [78] rezultatele sunt prezentate tot pentru setul de 39 de foneme cu modele hibride (HMM/ANN) obținând 58.40% rată de recunoaștere cu vectori acustici similari. Rezulatul nostru cel mai bun pentru setul de 61 de foneme 53.67% se compară în mod favorabil cu aceste rezultate.

Pentru baza de date OASIS, rezultatele noastre cele mai bune (82.38%) sunt aproape identice cu cele obținute de autorii bazei de date (82.05%) cu același tip de vectori acustici și același tip de modele GMM. Singura diferență în ambele cazuri este algoritmul de decodare, al nostru fiind o variantă Viterbi foarte eficient.

## Capitolul 6

# Modelarea duratei fonemelor

### 6.1 Introducere

Unul dintre punctele slabe ale modelelor Markov ascunse este că, nu oferă o reprezentare potrivită pentru structura temporală a vorbirii. Acest lucru se datorează faptului că, probabilitatea de a rămâne în aceeași stare descrește exponențial cu timpul, adică aceste modele incorporează în mod implicit un model de durată exponențial. Probabilitatea de a rămâne în starea  $i$  exact de  $t$  ori se poate exprima cu formula

$$d_i(t) = a_{ii}^t(1 - a_{ii}) \quad (6.1)$$

în cazul modelelor Markov clasice, care permit două tipuri de tranziții dintr-o stare.  $a_{ii}$  este probabilitatea de a rămâne în starea  $i$ , iar  $(1 - a_{ii})$  este probabilitatea de a trece în starea următoare. Figura 6.1 ne arată structura unui astfel de model Markov cu trei stări.

Durata poate fi inclusă în modelele Markov. O primă încercare în acest sens a fost făcută în [54], care modifică algoritmul Viterbi pentru a include și duratele. Fie  $O = \{o_1, o_2, \dots, o_T\}$  secvența de observații care trebuie decodată într-o secvență de foneme și cu  $\alpha_t(j)$  notăm probabilitatea maximă de a fi în starea  $j$  și observând secvența de observații  $o_1, o_2, \dots, o_t$ . Atunci valorile  $\alpha_t(j)$  se pot determina cu ajutorul formulei următoare

$$\alpha_t(j) = \max_{1 \leq i \leq n} \left\{ \max_{1 \leq \tau \leq D} \left\{ \alpha_{t-\tau}(i) a_{ij} d_j(\tau) \prod_{\Theta=0}^{\tau-1} b_j(o_{t-\Theta}) \right\} \right\} \quad (6.2)$$

unde  $D$  este durata maximă admisă pentru o fonemă. Acest algoritm crește complexitatea algoritmului Viterbi cu un factor  $D$ , care poate fi inadmisibil pentru

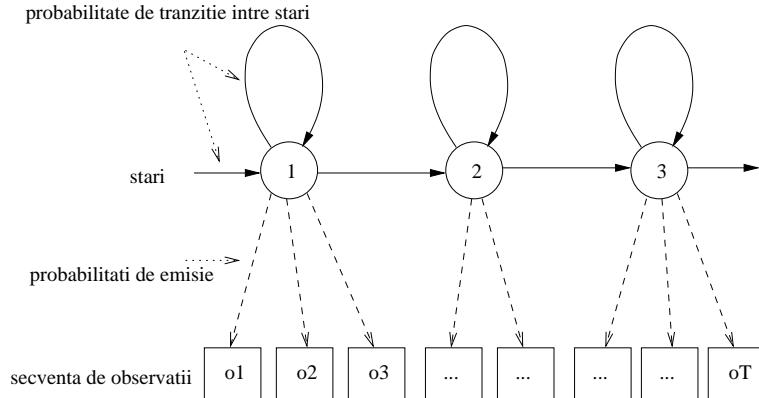


Figura 6.1: Model Markov ascuns cu trei stări

recunoașterea vorbirii continue, însă acceptabil pentru recunoașterea cuvintelor izolate.

Inițial au fost utilizate distribuții neparametrice. Acest lucru înseamnă că, fiecare  $d_i(\tau)$  a fost estimat pentru fiecare  $1 \leq \tau \leq D$ . Acest tip de estimare însă necesită o cantitate foarte mare de date de antrenare. Din această cauză au fost propuse câteva distribuții parametrice împreună cu formulele de estimare ale acestora. Distribuțiile propuse sunt: Poisson, Gamma și cea gaussiană.

## 6.2 Rezultate experimentale

### 6.2.1 OASIS Numbers

În această secțiune prezentăm rezultatele sistemului pentru recunoașterea fonemelor din baza de date OASIS Numbers utilizând modelarea duratei. Ca modele acustice pentru foneme au fost folosite modele Markov ascunse cu o singură stare.

Figura 6.2 prezintă duratele minime, respectiv maxime a celor 31 de foneme OASIS. Durata medie este marcată cu steluță. Duratele minime, respectiv cele maxime vor fi folosite în ecuația 6.4.

O primă încercare este să verificăm comportarea sistemului fără a utiliza vre-un model pentru durata fonemelor, însă impunând niște durate minime și maxime pentru acestea. Pentru acesta formula 6.2 a fost simplificată rezultând în

$$\alpha_t(j) = \max_{1 \leq i \leq n} \left\{ \max_{\tau_{min} \leq \tau \leq \tau_{max}} \left\{ \alpha_{t-\tau}(i) \cdot a_{ij} \cdot \prod_{\Theta=0}^{\tau-1} b_j(o_{t-\Theta}) \right\} \right\} \quad (6.3)$$

pentru  $1 \leq j \leq n$  și  $1 \leq t \leq T$ , unde  $\tau_{min}$  și  $\tau_{max}$  sunt duratele permise

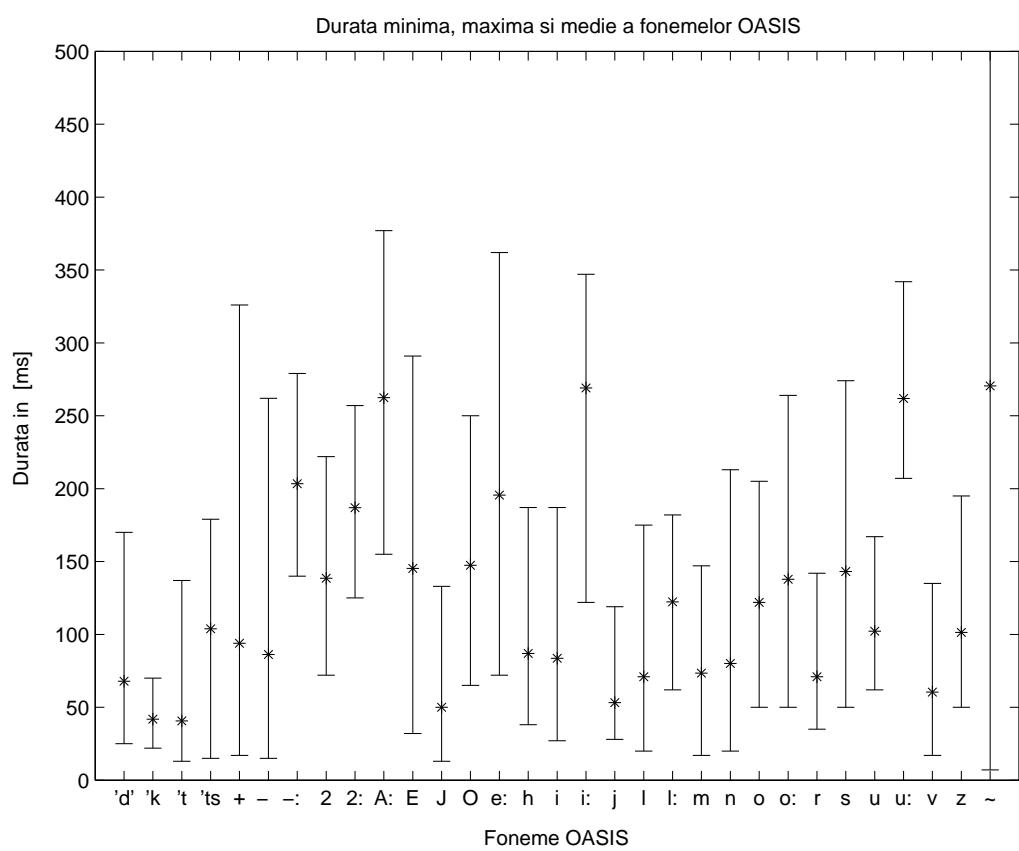


Figura 6.2: Durata minimă, maximă și medie a fonemelor OASIS

pentru unitățile fonetice. Deoarece fonemele sunt modelate cu o singură stare, semnificația  $a_{ij}$  se modifică, acesta reprezentând probabilitatea de trecere de la fonemul  $f_i$  la fonemul  $f_j$ . Noi am utilizat doar valori 0 și 1, 0 reprezentând faptul că  $f_i$  nu poate fi urmat de  $f_j$ , iar valoarea de 1 contrariul. Elementele matricii  $a_{ij}$  se determină pe baza vocabularului și reprezintă un model bigram al limbajului, la nivel fonetic.

Pentru a reduce cantitatea de calcule am calculat, tot pe baza înregistrărilor segmentate la nivel fonetic, durata maximă, respectiv durata minimă pentru fiecare fonem în parte (vezi figura 6.2). Aceste valori se notează cu  $\tau_{min}(j)$  respectiv  $\tau_{max}(j)$  pentru  $j = 1 \dots n$ . Utilizând aceste valori, formula 6.3 se modifică în

$$\alpha_t(j) = \max_{1 \leq i \leq n} \left\{ \max_{\tau_{min}(j) \leq \tau \leq \tau_{max}(j)} \left\{ \alpha_{t-\tau}(i) \cdot a_{ij} \cdot \prod_{\Theta=0}^{\tau-1} b_j(o_{t-\Theta}) \right\} \right\} \quad (6.4)$$

Tabelul 25 prezintă rezultatele obținute pentru recunoaștere indicând duratele minime și maxime utilizate. Aceste rezultate au fost publicate în [10]. Toate măsurările au fost făcute pe două seturi de modele, unul folosind 8 componente gaussiene, iar al doilea folosind 16 componente gaussiene. Duratele  $\tau_{min} \dots \tau_{max}$  reprezintă faptul că, în acest caz am folosit duratele specifice pentru fiecare fonem în parte.

Nr. mixt.	Durate	CORR.	ACC.	INS	SUBST	DEL
8	1..50	97.23%	11.70%	85.53%	2.63%	0.12%
8	2..50	96.93%	60.70%	36.22%	2.84%	0.21%
8	3..50	94.68%	76.29%	18.39%	4.23%	1.07%
8	$\tau_{min} \dots \tau_{max}$	92.96%	76.64%	16.32%	5.95%	1.07%
16	1..50	97.40%	24.30%	73.10%	2.54%	0.04%
16	2..50	96.37%	68.65%	27.72%	3.45%	0.17%
16	3..50	96.64%	77.54%	17.09%	4.66%	0.69%
16	$\tau_{min} \dots \tau_{max}$	92.65%	76.94%	15.71%	6.21%	1.12%

Tabelul 25: Recunoașterea fonemelor - OASIS - 31

### 6.2.2 TIMIT

Figura 6.3 prezintă duratele minime respectiv maxime a celor 61 de foneme TIMIT. Durata medie este marcată cu steluță.

Distribuția duratelor pentru fiecare fonem în parte sunt reprezentate cu câte o histogramă (vezi anexa 1).

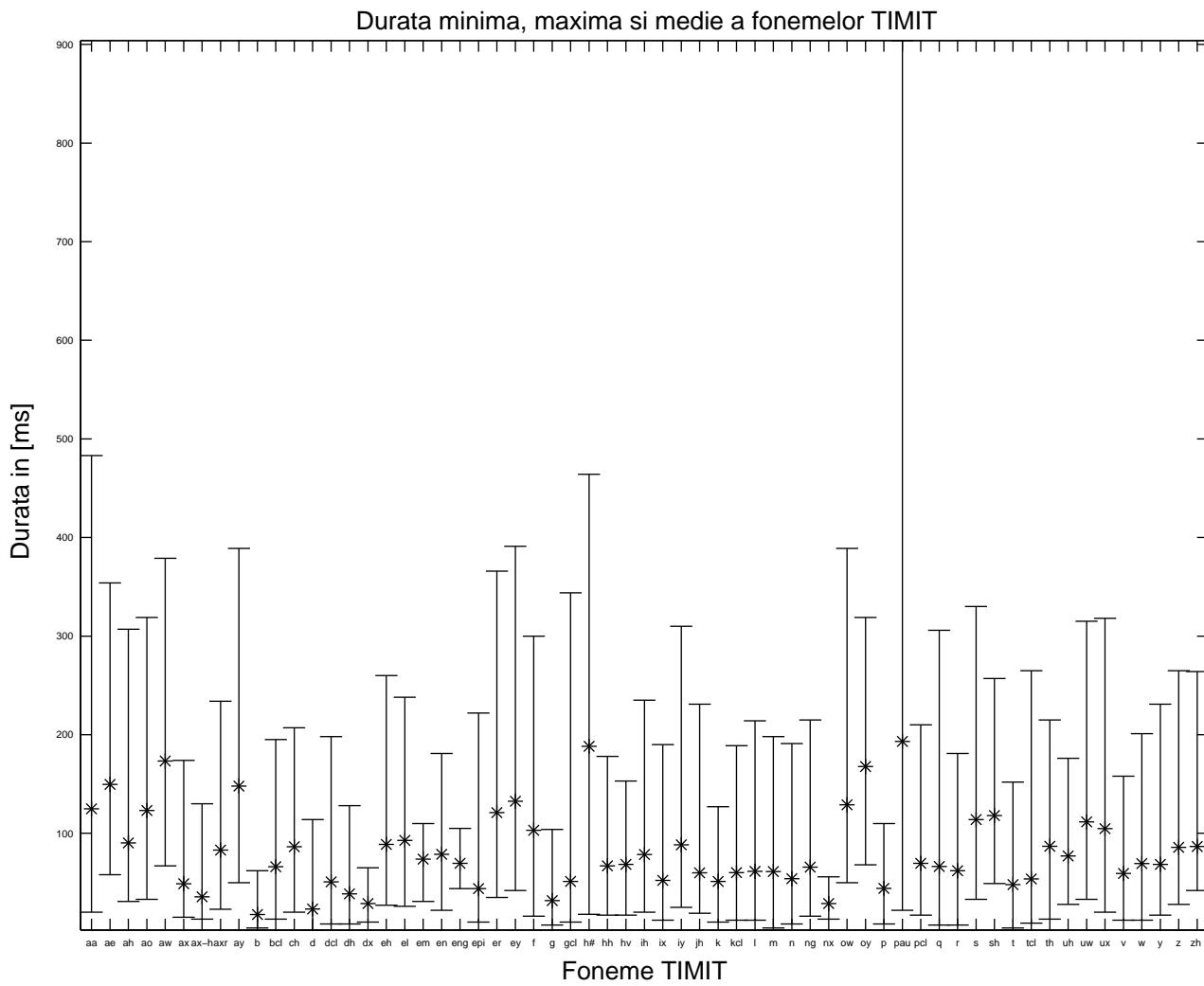


Figura 6.3: Durata minimă, maximă și medie a fonemelor TIMIT

### 6.3 Concluzii

Pe baza duratelor specifice ale fonemelor am stabilit o relație între parametrul empiric și durata medie de fonem. Tabelul 26 prezintă duratele medii, respectiv valorile optime ale parametrului  $\beta$ . Pentru parametrul empiric am considerat acea valoare, pentru care acuratețea de recunoaștere era maximă (vezi tabelele 23 și 24). Relația cu care se poate calcula acest parametru empiric ar fi

$$\beta = \frac{\text{durata\_medie\_fonem}}{5} \quad (6.5)$$

	Durată medie	$\beta$
TIMIT	80 ms	15
OASIS	125 ms	25

Tabelul 26: Durata medie de fonem și valoare  $\beta$

O concluzie referitoare la duratele experimentate pentru baza de date OASIS ar fi că, impunerea unei dure minime este foarte importantă, acesta poate să furnizeze rezultate la fel de bune ca și includerea în formulă a duratelor specifice fiecarei foneme ( $\tau_{min}$  și  $\tau_{max}$ ). De remarcat este și faptul că, rezultatele obținute în capitolul precedent pentru recunoașterea fonemelor din baza de date OASIS au fost mai bune (82.38%), decât cele obținute cu modelarea explicită a duratei (77.54%). Ambele rezultate au fost obținute cu modele de foneme formate din 16 componente gaussiene. O parte din aceste rezultate a fost publicată în [10].

## Capitolul 7

# Recunoașterea cuvintelor prin modelare fonetică

### 7.1 Introducere

În acest capitol prezentăm rezultatele experimentale obținute pentru recunoașterea cuvintelor pronunțate în mod izolat. Au fost făcute două rânduri de experimente, în primul modelând cuvinte întregi și în al doilea modelând fonemele.

În cel de-al doilea caz cuvintele au fost reprezentate cu ajutorul unui vocabular care conține pronunții fonetice ale cuvintelor. Există diferite abordări ale creării unui asemenea vocabular și anume [72]:

- fiecare cuvânt are o singură reprezentare fonetică
- fiecare cuvânt poate avea mai multe reprezentări fonetice
- fiecare cuvânt poate avea mai multe reprezentări fonetice, fiecărei reprezentări îi corespunde câte o probabilitate

Pentru exemplificare considerăm cuvântul din limba engleză *bottle*. Acest cuvânt admite mai multe pronunții corecte cum ar fi: [ b aa dx el ], [ b aa t el ], [ b aa dx ax 1 ], [ b aa t ax 1 ]. Introducând multe variante de pronunție crește dimensiunea vocabularului, însă ridică acuratețea sistemului. Probabilitățile atașate pronunților se determină în funcție de frecvența acestor pronunții.

O altă problemă, care se ridică, se referă la reprezentarea eficientă a vocabularului. Există două abordări:

- liniar

- organizat sub formă de arbore (trie)

Fiecare metodă are avantaje și dezavantaje. Avantajul modelului liniar ar fi aplicarea simplă a modelului de limbaj, iar dezavantajul este că, ocupă multă memorie. Modelul organizat sub formă de arbore reprezintă o singură dată prefizurile comune ale cuvintelor, deci este mult mai avantajos din punctul de vedere al memoriei, însă apar dificultăți în cazul aplicării modelului de limbaj.

## 7.2 Modele de cuvinte

Modelarea cuvintelor poate fi făcută în mai multe moduri, chiar dacă ne referim doar la metodele statistice. Sistemele de recunoașterea cuvintelor izolate, care conțin un număr redus de cuvinte, cum ar fi cifrele, pot utiliza și metoda alinierii dinamice în timp pentru recunoașterea cuvintelor. Descrierea acestor tipuri de modele se găsesc în [26, 67, 79]. Modelele Markov ascunse cu observații discrete furnizează rezultate foarte bune pentru această problemă [68]. O serie de rezultate au fost publicate și pentru limba română pentru modelarea cuvintelor folosind atât modelele statistice, cât și diferite tipuri de rețele neuronale [20, 31, 79].

Noi am utilizat modele Markov ascunse cu o singură stare (GMM) pentru a demonstra că, și aceste modele simple, care reprezintă doar conținutul acustic al cuvintelor, fără a modela distribuția acestora în timp, pot discrimina un număr de zece cuvinte. Antrenarea cuvintelor (cifrelor) s-a făcut cu câte 40 de exemplare, 20 de vorbitori rostesc de câte două ori fiecare cifră. Testarea fiecărui cuvânt s-a făcut pe 12 exemplare, câte două de la fiecare dintre cei 6 vorbitori. Figura 7.1 prezintă procesul de evaluare, iar rezultatele sunt prezentate în tabelul 27. Pentru a determina numărul de mixturi potrivite am început cu o singură mixtură pentru fiecare cuvânt, dublând numărul acestora la fiecare experiment. Cu 16 mixturi am obținut o rată de recunoaștere 100%.

Număr mixturi	Rata de recunoaștere
1	83%
2	82%
4	94%
8	98%
16	100%

Tabelul 27: Recunoașterea cuvintelor cu modele GMM - OASIS-10 cifre

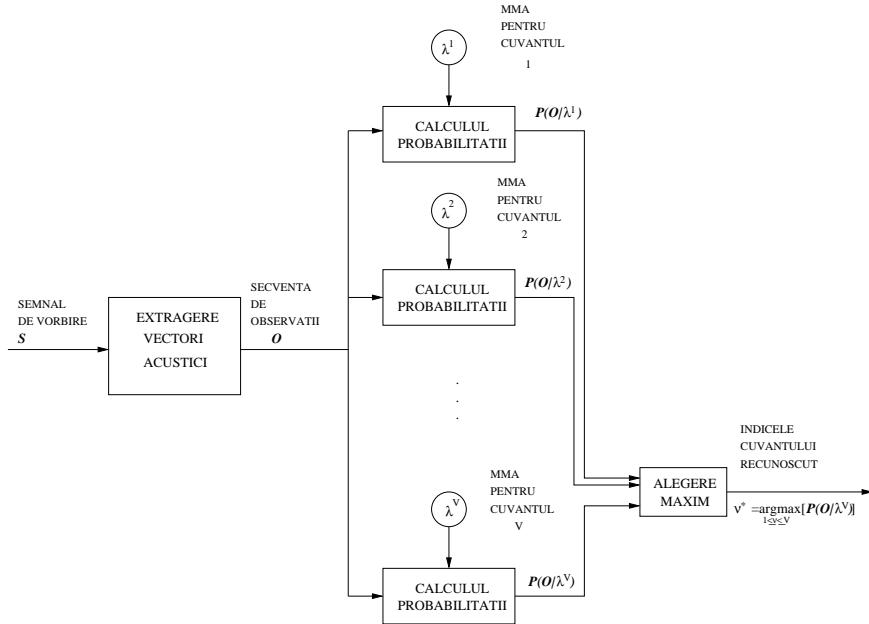


Figura 7.1: Recunoașterea cuvintelor izolate (după [67])

### 7.3 Modele de foneme

În vocabular am inclus mai multe pronunții pentru același cuvânt, dar numărul acestor pronunții a fost limitat la două. Pentru reprezentarea vocabularului am ales cea de-a doua variantă, adică reprezentarea sub formă de arbore. Acest tip de arbore se numește și *trie*, denumirea provenind de la information retrieval. Acest arbore permite organizarea eficientă a vocabularilor, mai mult, este structura de date, care permite o căutare eficientă de timp constant  $\Theta(1)$ , timpul de căutare al unui cuvânt depinzând doar de lungimea acestuia, nu și de dimensiunea vocabularului. Acest tip de căutare este numit și căutare digitală [50]. Pentru decodare am utilizat algoritmul Viterbi cu parametrul empiric, introdus în secțiunea 5.1. Figura 7.2 prezintă o parte din vocabular organizat sub formă de trie.

Recunoașterea de cuvinte a fost făcută conform figurii 7.3. În primul pas am obținut vectorii acustici din semnalul vocal (vezi secțiunea 4.2), după care am aplicat o variantă a algoritmului Viterbi. Acest algoritm ne-a furnizat o secvență de foneme, care eventual a fost postprocesată aplicând niște corecții pe baza regulilor fonetice și apoi căutată în vocabularul reprezentat sub formă de arbore. Căutarea se începe întotdeauna de la rădăcina arborelui și se caută un drum de la rădăcină la un nod, care reprezintă sfârșitul unui cuvânt. Dacă secvența de foneme nu a fost găsită, atunci acesta nu va fi identificată cu nici un cuvânt.

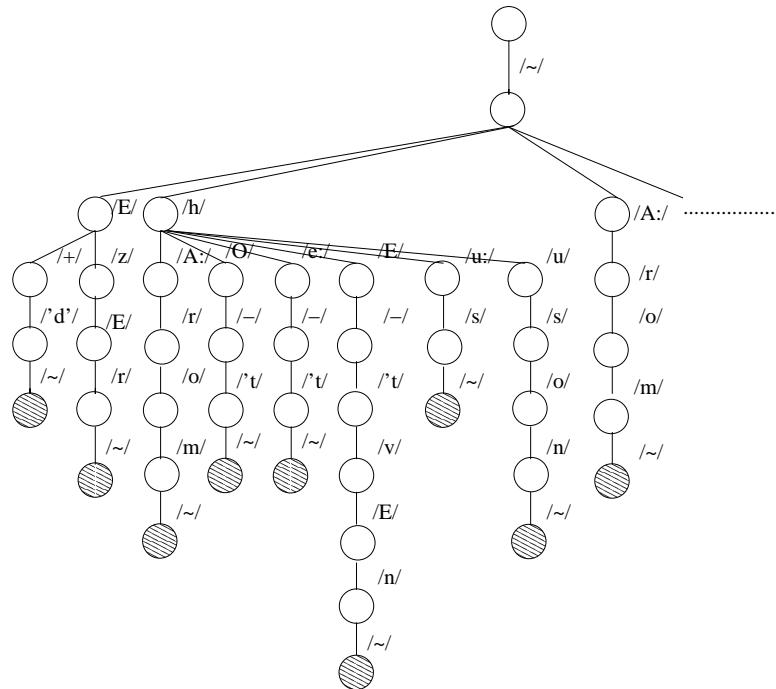


Figura 7.2: OASIS-Vocabular organizat sub formă de trie

În capitolul Recunoașterea fonemelor am prezentat diferitele tipuri de erori de recunoaștere. Am văzut că, erorile de inserare sunt cele mai ridicate. În cazul acestor erori cea mai frecventă este inserarea diferitelor foneme, care reprezintă liniștea la începutul, dar mai ales la sfârșitul rostirilor. În faza de postprocesare deocamdată am comasat acele foneme, care reprezintă linștea. Exemplu:

Secvență de foneme la intrare: ~ E z E r ~

Secvență de foneme decodată : ~ E z E r - : ~

După comasare : ~ E z E r ~

Tabelul 28 prezintă eroarea la nivel de cuvinte obținută pentru setul de test, care conține 312 cuvinte rostită de către 6 vorbitori diferenți. Tabelul prezintă rezultatele cu și fără comasarea fonemelor de tip liniște. Vocabularul conține 34 de cuvinte, care reprezintă de fapt 26 de cuvinte, câteva cuvinte admit mai multe pronunții fonetice.

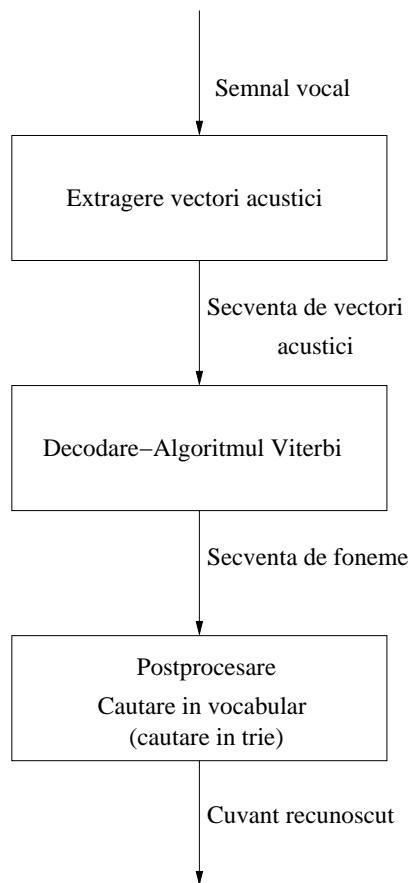


Figura 7.3: Recunoașterea cuvintelor cu modele de foneme

Tip recunoaștere	DEL	INS	SUBST	REC. FO- NEME	CORR.	REC. CU- VINTE
fără postprocesare	2.37%	11.27%	6.00%	80.35%	91.62%	84.62%
cu postprocesare	2.46%	7.64%	5.52%	84.36%	92.01%	86.54%

Tabelul 28: Recunoaștere cuvinte - OASIS 26 de cuvinte -  $\beta = 20$

## 7.4 Concluzii

În acest capitol am prezentat rezultate experimentale pentru recunoașterea cuvintelor rostite izolat. Sunt prezentate două rânduri de experimente, primul introducând modelele GMM pentru modelarea cifrelor. Pentru un vocabular conținând doar 10 cuvinte modelele GMM au fost adecvate. Folosind 16 componente gausiene, am obținut rata de recunoaștere 100%. Pentru un vocabular cu mai multe cuvinte și cuvinte mai similare, acest model nu ar fi potrivit.

Partea a doua a capitolului utilizează cunoștințele dobândite la recunoașterea fonemelor și prezintă sistemul de recunoașterea cuvintelor izolate bazată pe foneme. Noutatea față de recunoașterea fonemelor este, că este necesară introducerea modelării vocabularului. Noi am optat pentru modelul mai eficient din punct de vedere al spațiului de memorie, și am utilizat un vocabular sub formă de arbore. Rata de recunoaștere a cuvintelor, fără introducerea limbajului de model este de 84.62%. Cu introducerea cunoștințelor fonetice acest rezulat poate fi îmbunătățit. Noi am aplicat o simplă postprocesare pentru comasarea erorilor de tip inserare, ceea ce a crescut deja cu 2% rata de recunoaștere, reducând-se mai ales erorile de tip inserare.

## Capitolul 8

# Concluzii finale

În această lucrare am prezentat structura generală a sistemelor de recunoaștere vorbitorului și recunoașterea vorbirii continue. În acest scop am realizat o aplicație cu următoarele module, care se poate utiliza atât pentru recunoașterea vorbitorului cât și la recunoașterea vorbirii:

- modul pentru extragerea vectorilor acustici (MFCC, LPCC)
- modul pentru tehnici de grupare conținând tehnicile: k-means, cuantizare vectorială binară ierarhică, LVQ, GMM
- modul pentru modele Markov ascunse discrete (HMM discret)
- modul pentru reprezentarea lexiconului
- modul pentru tehnici de decodare

Menționăm că, toate componentele acestui sistem sunt independente de limba, care trebuie recunoscută. Experimentele au fost făcute pentru două baze de date, una pentru limba engleză TIMIT și una pentru limba maghiară OASIS Numbers. Aceste baze de date au fost puse la dispoziția noastră de către *Research Group on Artificial Intelligence, University of Szeged*.

### 8.1 Contribuții la recunoașterea vorbitorului

În cazul recunoașterii vorbitorului a fost studiată recunoașterea independentă de text. Am studiat performanțele modelelor de vorbitori create cu ajutorul diferitelor tehnici de grupare. Tipul de model ales depinde și de numărul de vorbitori, precum și de cantitatea de material audio disponibil. Modelele GMM au mai

multi parametri, deci necesita o cantitate mai mare de date pentru estimarea robustă a parametrilor modelelor. Au fost studiate modelele de vorbitori din punct de vedere al conținutului fonetic. Deoarece mulți autori stabilesc dimensiunea modelului (numărul componentelor gaussiene, respectiv numărul clusterelor) în funcție de numărul fonemelor unei limbi, am verificat dacă fonemele, respectiv categoriile fonetice sunt separate în componente ale acestor modele. Pentru a elibera efectul de coarticulare, din foneme au fost luate în considerare doar mijlocul acestora, care nu sunt afectate de coarticulare. Surprinzător, modelele nu separă toate categoriile de foneme, mai ales vocalele și semivocalele sunt împreștiinate în mai multe componente. Pornind de la acest experiment, am verificat performanțele modelelor structurate pe categorii fonetice. Concluzia este că, utilizând doar categoria fonetică și construind modele cu 6 componente gaussiene, aceste modele vor avea performanțe similare cu modelele obișnuite. Am constituit lista vorbitorilor identificați în mod incorrect atât la modelele obișnuite GMM cât și la modelele structurate pe categorii fonetice. Deoarece cele două sisteme produc liste diferite, aceste două tipuri de modelări ar putea fi combinate în scopul realizării unui sistem mai performant. Acest aspect rămâne o direcție de cercetare.

Un alt aspect studiat în această lucrare a fost puterea de discriminare a diferențelor categorii fonetice. Am reușit să ordonăm grupele fonetice pe baza puterii lor de discriminare, nazalele și vocalele demonstrându-se foarte bune în acest scop.

## 8.2 Contribuții la recunoașterea vorbirii

Deși există mai multe abordări ale acestei probleme, componente, care trebuie realizate în cazul acestor sisteme reprezintă deja un standard. Prima dată trebuie selectată unitatea acustică, care va fi modelată și apoi se antrenează modelele acustice. Problema următoare, care trebuie rezolvată, este modelarea limbajului și reprezentarea vocabularului limbii. După ce aceste două probleme sunt rezolvate, se poate implementa o tehnică de decodare, care este de fapt o problemă de căutare cu o serie de constrângeri.

Am realizat o serie de experimente pe cele două baze de date disponibile, clasificând foneme, grupe de foneme și categorii de foneme. Am reușit să obținem cu modele GMM performanțe similare cu cele ale modelelor Markov ascunse cu trei stări, astfel demonstrând efectul neglijabil al probabilităților de tranziție între stările unui model de fonem. Pentru cele 61 de foneme TIMIT am obținut o rată de clasificare de 63.34%, pentru cel 39 de grupe de foneme TIMIT o rată de 68.13%, ambele cu modele de foneme mixturi gaussiene cu 32 de componente. În cazul

bazei de date OASIS Numbers, rata de clasificare este de 92.74% tot cu modele GMM, în acest caz modelele având 16 de componente. Numărul componentelor este puternic influențat de cantitatea de date de antrenare. Pentru reducerea calculelor aici s-a propus un clasificator cu două niveluri ale cărui performanțe au fost testate pe baza de date TIMIT și care, deși este mai eficient din punct de vedere computațional, nu reduce în mod semnificativ rata de clasificare.

În cazul recunoașterii fonemelor am modificat algoritmul Viterbi pentru modele de foneme cu o singură stare. Introducând o constantă, a cărei valori se determină în mod empiric, am reușit să reducem erorile de tip inserare. Recunoașterea a fost testată doar pentru varianta de 61 de foneme, rezultate mult mai bune s-ar fi putut obține cu gruparea alofoanelor (varianta cu 39 de grupe). În acest caz pe TIMIT am obținut 53.67, ca rată maximă de recunoaștere. Pentru OASIS Numbers rata maximă a fost 82.38%.

Studiul duratei fonemelor este o problemă importantă, mai ales în cazul recunoașterii limbilor, care au pentru același sunet atât variantă scurtă, cât și lungă, durata fiind singura diferență între acestea. Limbile finugorice fac parte din categoria acestor limbi. Această problemă este însă una mult studiată și pentru limba engleză. Am prezentat o analiză statistică a celor două baze de date disponibile din punctul de vedere al duratei fonemelor și pe baza cunoștințelor dobândite am propus o variantă pentru modelarea duratei fonemelor. Modelarea duratei se face doar la nivel de decodare, însă această variantă ridică foarte mult timpul de decodare. Din această cauză rezultate experimentale sunt prezentate doar pentru OASIS Numbers, care conține doar cuvinte. Pentru baza de date TIMIT acest tip de decodare este nepractic, ridicând insuportabil timpul necesar pentru decodare. Rezultatele experimentale ne arată că, impunerea unei dure minime este chiar mai importantă decât impunerea duratelor specifice fiecărei foneme.

În final sunt prezentate rezultate pentru recunoașterea cuvintelor. În cazul utilizării unui număr redus de cuvinte, cum ar fi cifrele, se recomandă folosirea modelelor de cuvinte. Cu modelele având 16 componente gaussiene am obținut o rată de recunoaștere 100% în cazul unui sistem independent de vorbitor. Dacă avem de recunoscut un număr mai mare de cuvinte, atunci se recomandă crearea modelelor de cuvinte din modelele de foneme ale limbii respective. Rezultatul este de 84.36% pe 26 de cuvinte, câteva dintre ele fiind foarte similare, deosebindu-se doar printr-un fonem. În acest caz nu au fost utilizate cunoștințe lingvistice, orice fonem poate să urmeze după oricare altul.

Majoritatea rezultatelor prezentate în această lucrare au fost publicate în reviste de specialitate și în volume ale unor conferințe naționale și internaționale,

acestea fiind totodată citate în capitolele corespunzătoare din lucrare.

### 8.3 Direcții de cercetare

Recunoașterea vorbitorului:

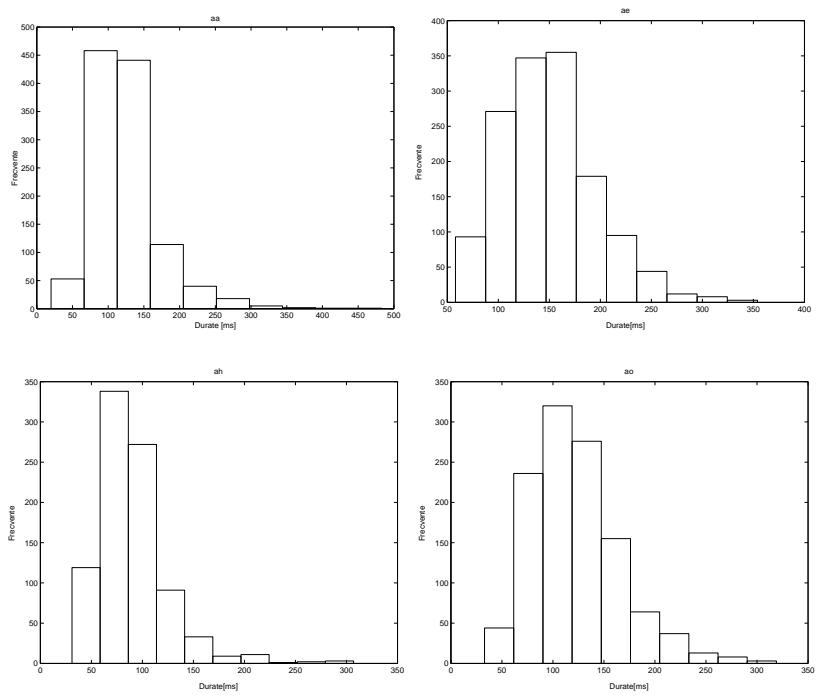
- O continuare firească ar fi construirea modelelor de vorbitori prin adaptarea modelelor generale de foneme. Pentru această investigare ar fi nevoie de o altă bază de date, deoarece în baza de date TIMIT există foneme care lipsesc sau sunt foarte slab reprezentate (cazul africatelor).
- O altă direcție de cercetare ar fi verificarea identității pe baza unor informații lingvistice, cum ar fi durata unor foneme sau cuvinte specifice.
- În final, ar fi de dorit combinarea mai multor informații atât de nivel inferior cât și de nivel superior. Aceasta iarăși necesită crearea unei baze de date care permite aceste studii.

Recunoașterea vorbirii:

- În cazul recunoașterii fonemelor, deoarece modelele Markov ascunse reprezintă o modelare total neadecvată a duratelor fonemelor, ar fi necesar un studiu complet al tuturor posibilităților de integrare al modelării acestor durează fără a mări complexitatea algoritmului de decodare. Actualmente toate încercările de modelare a duratei fonemelor măresc complexitatea calculelor.

## Anexa A

# Duratele fonemelor din baza de date TIMIT



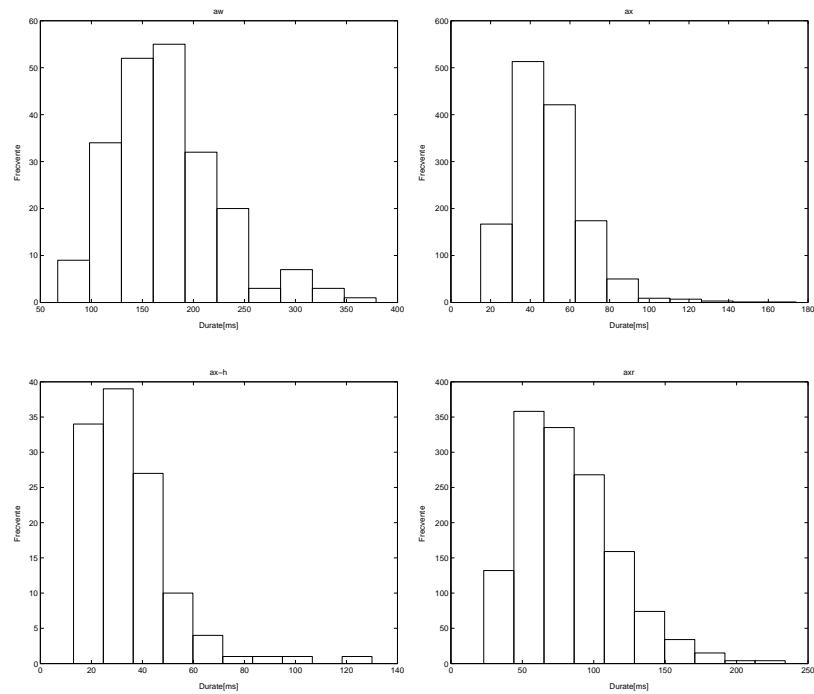


Figura A.1: Histograme ale duratelor fonemelor: /aa/, /ae/, /ah/, /ao/, /aw/, /ax/, /ax-h/, /axr/

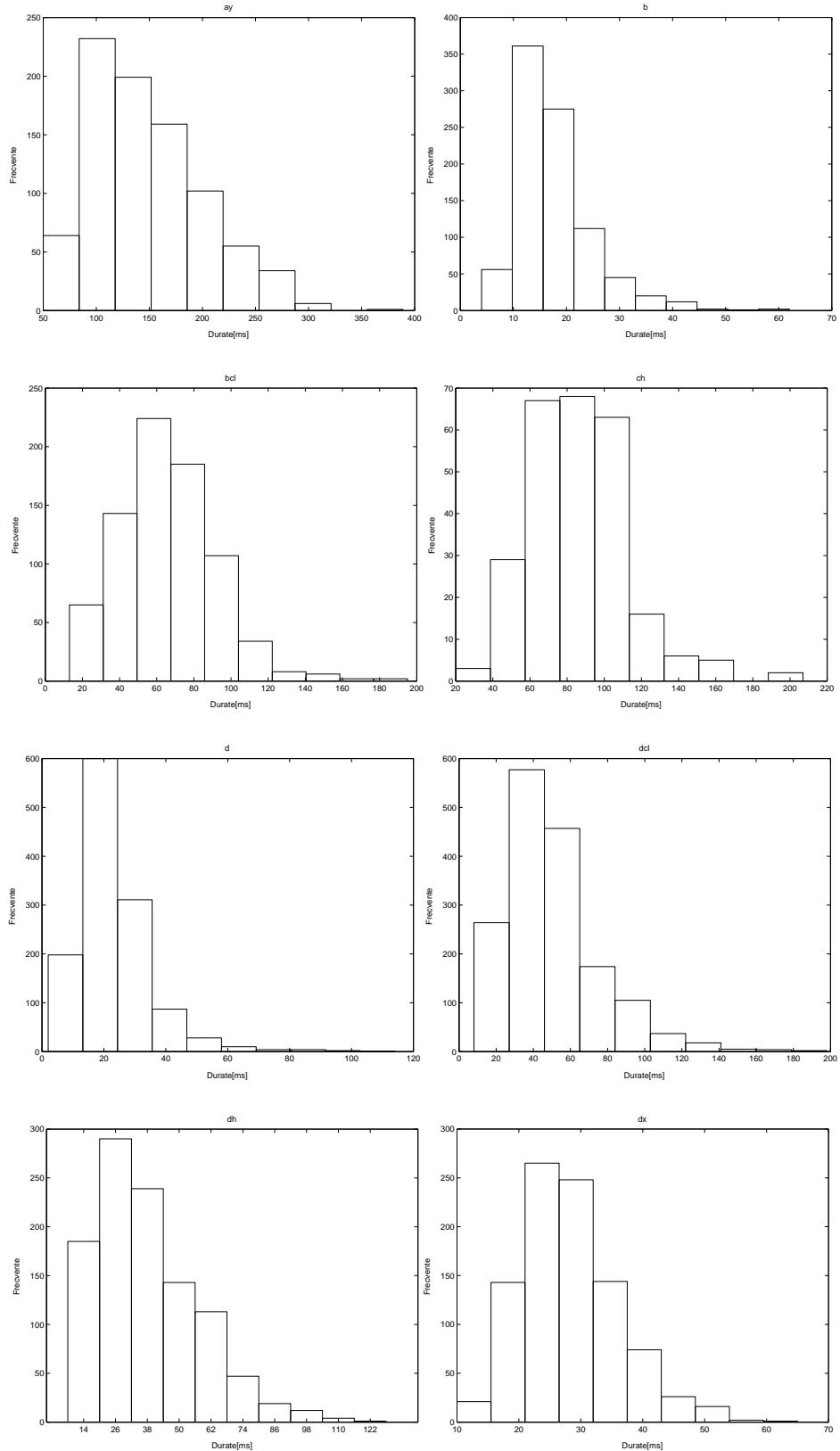


Figura A.2: Histogramme ale duratelor fonemelor: /ay/, /b/, /bcl/, /ch/, /d/, /dcl/, /dh/, /dx/

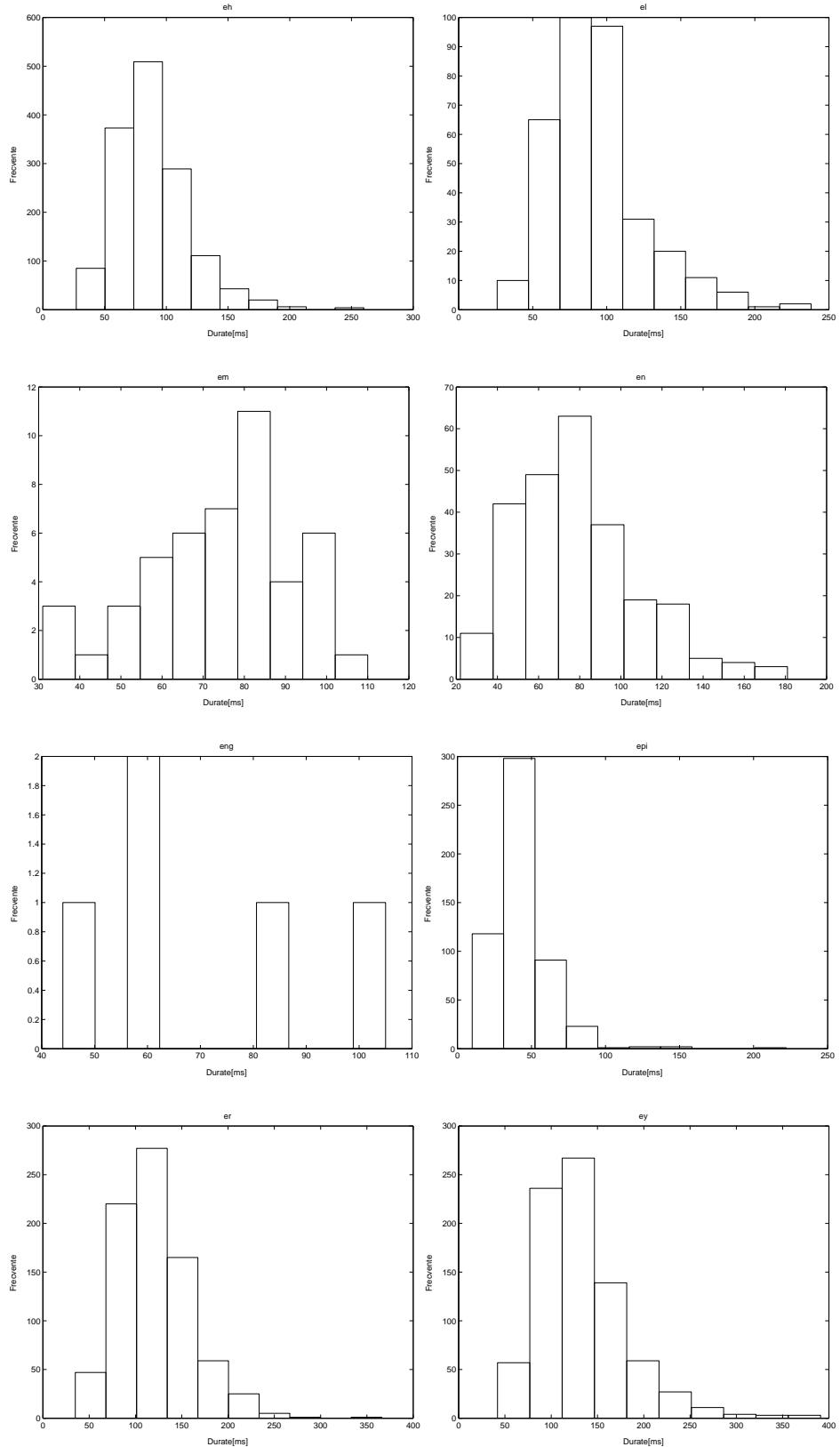


Figura A.3: Histogramme ale duratelor fonemelor: /eh/, /el/, /em/, /en/, /eng/, /epi/, /er/, /ey/

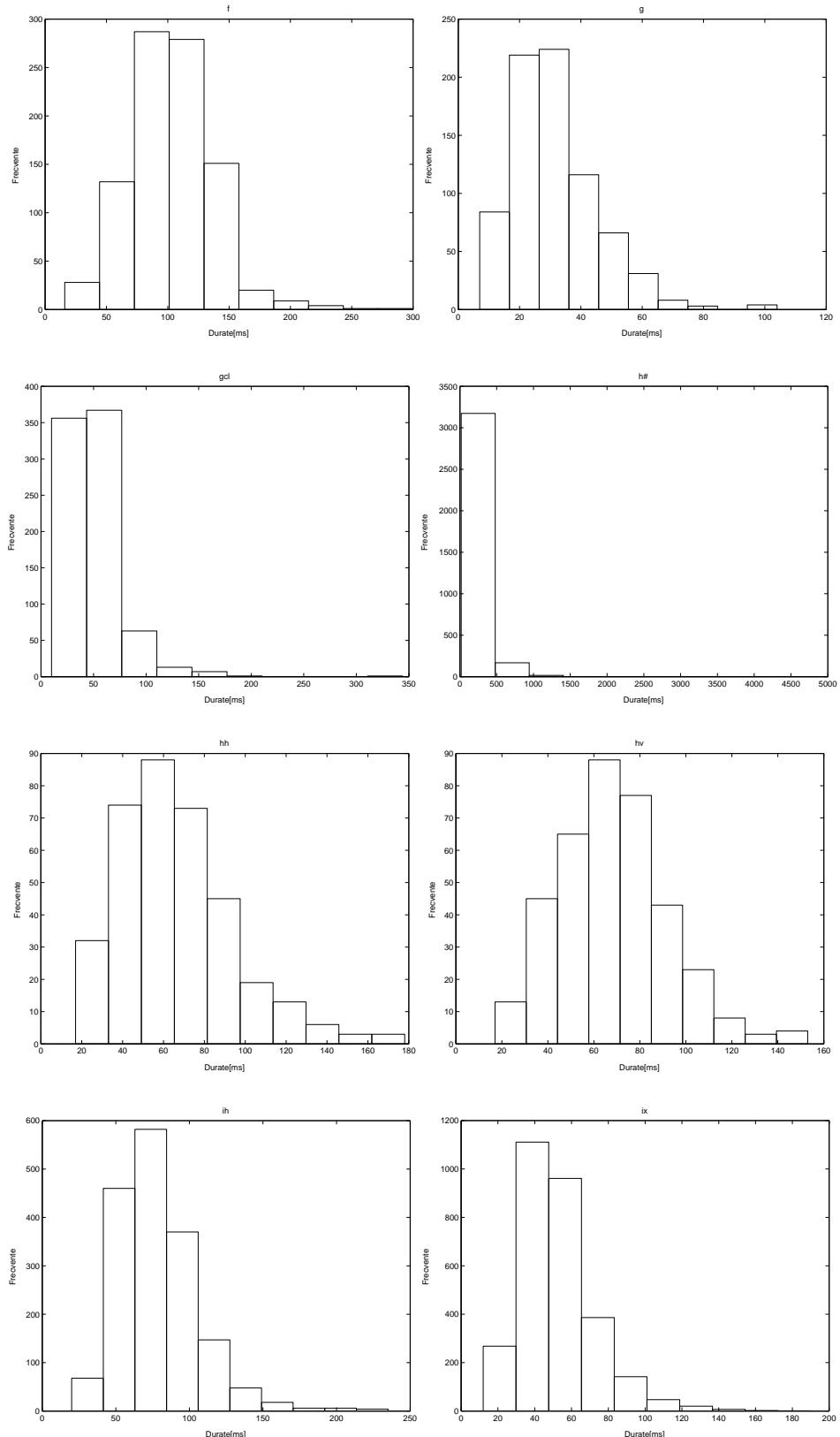


Figura A.4: Histogramme ale duratelor fonemelor: /f/, /g/, /gcl/, /h#//, /hh/, /hv/, /ih/, /ix/

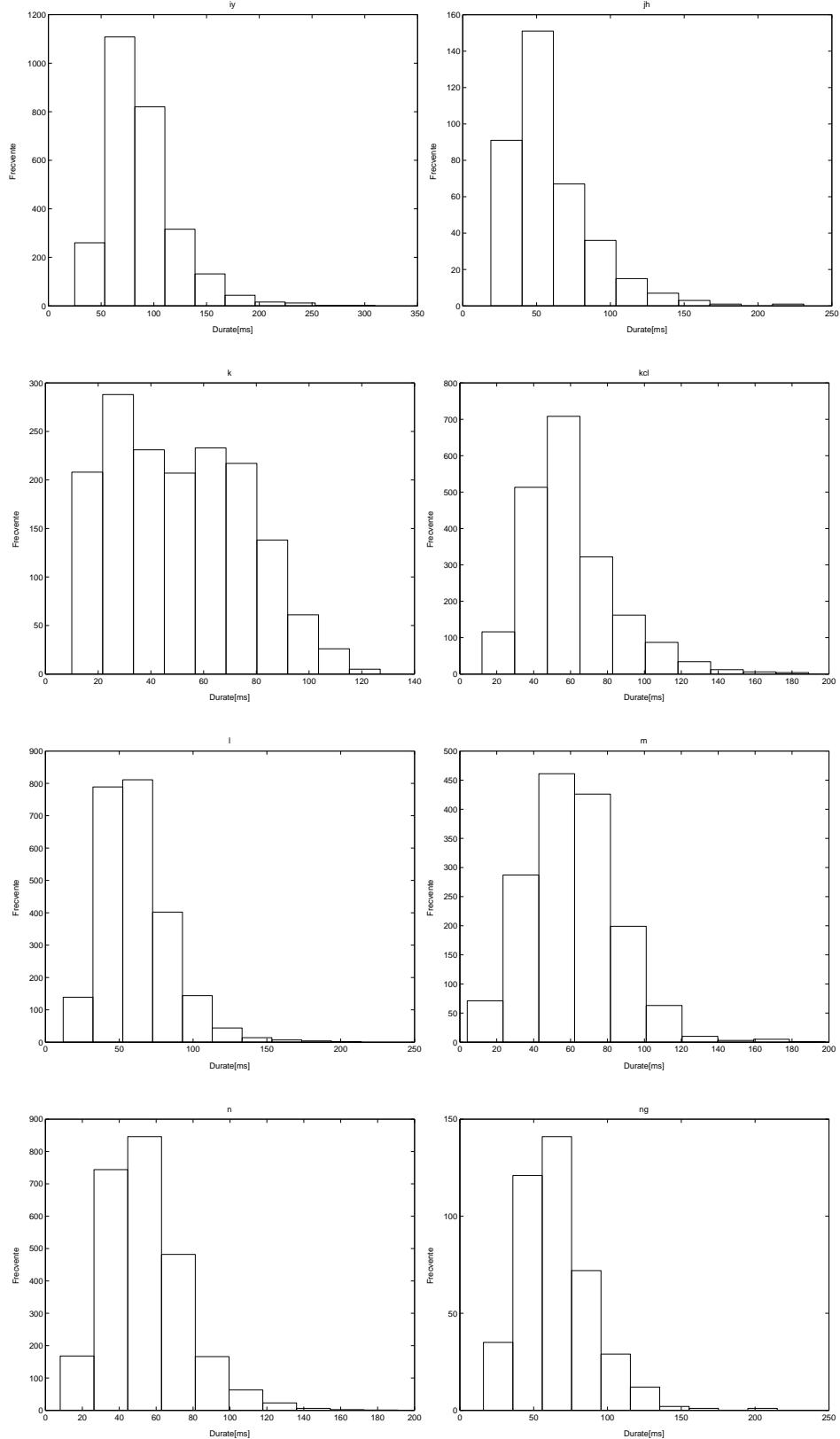


Figura A.5: Histogramme ale duratelor fonemelor: /iy/, /jh/, /k/, /kcl/, /l/, /m/, /n/, /ng/

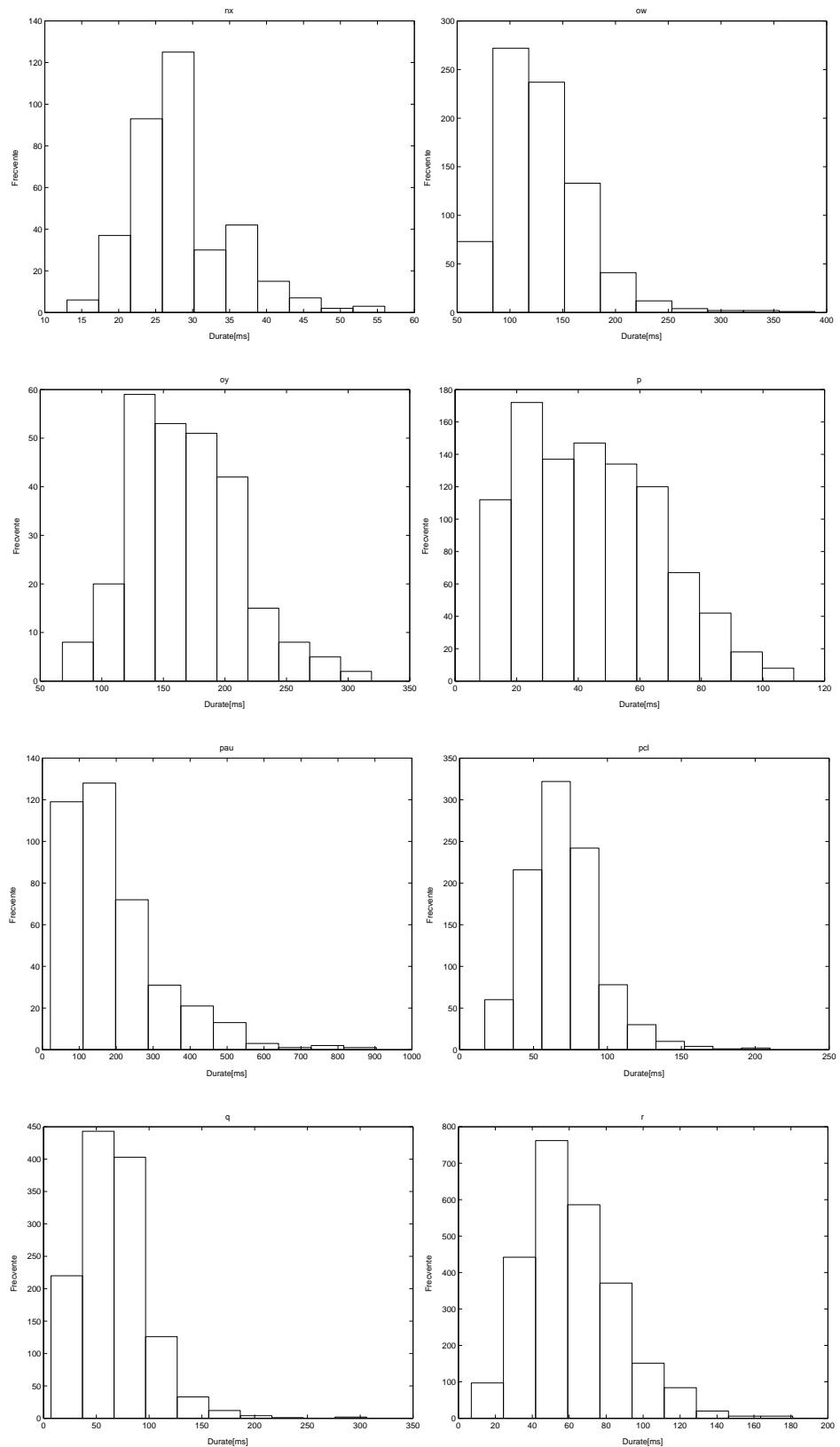


Figura A.6: Histogramme ale duratelor fonemelor: /nx/, /ow/, /oy/, /p/, /pau/, /pcl/, /q/, /r/

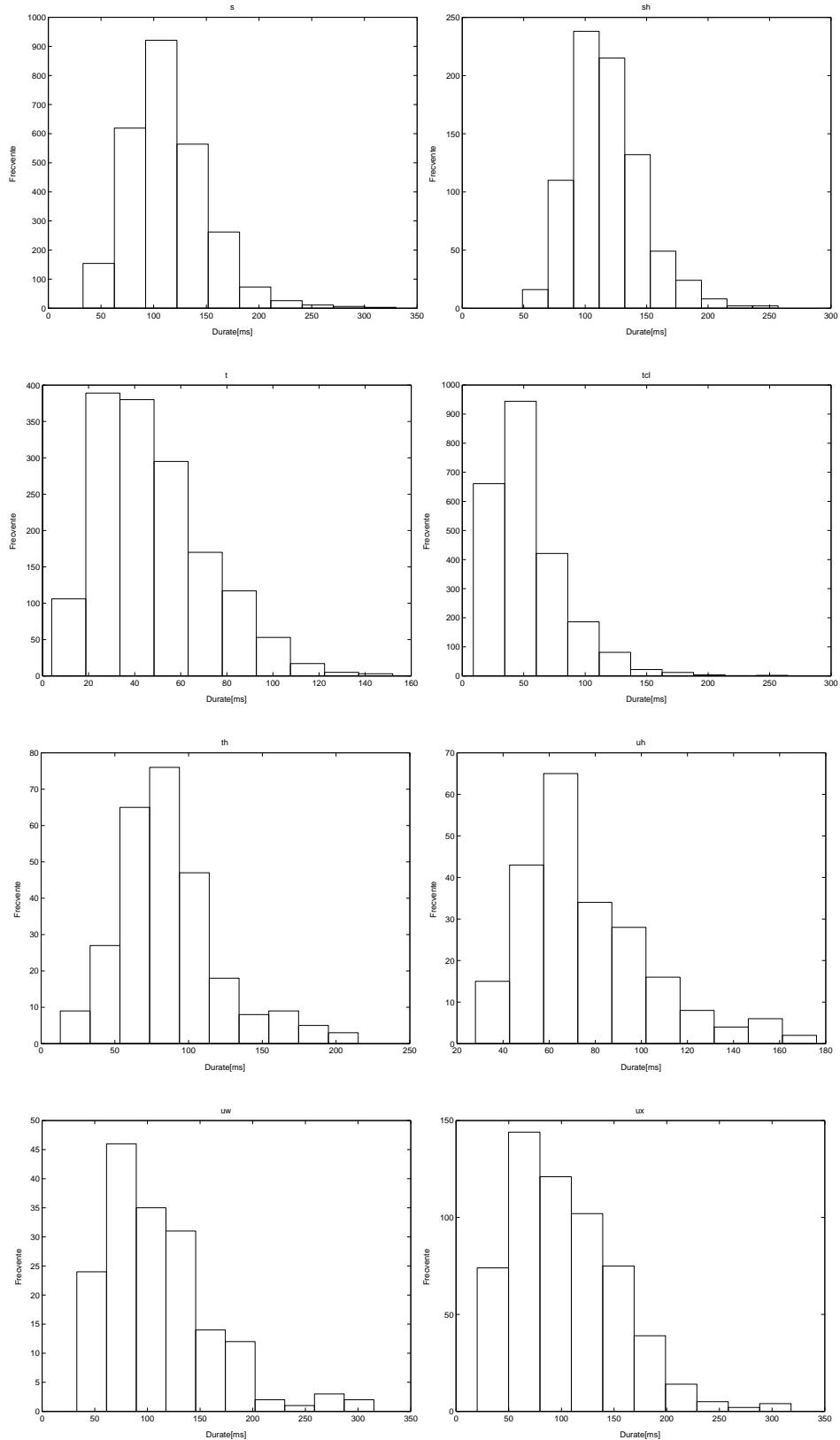


Figura A.7: Histogramme ale duratelor fonemelor: /s/, /sh/, /t/, /tcl/, /th/, /uh/, /uw/, /ux/

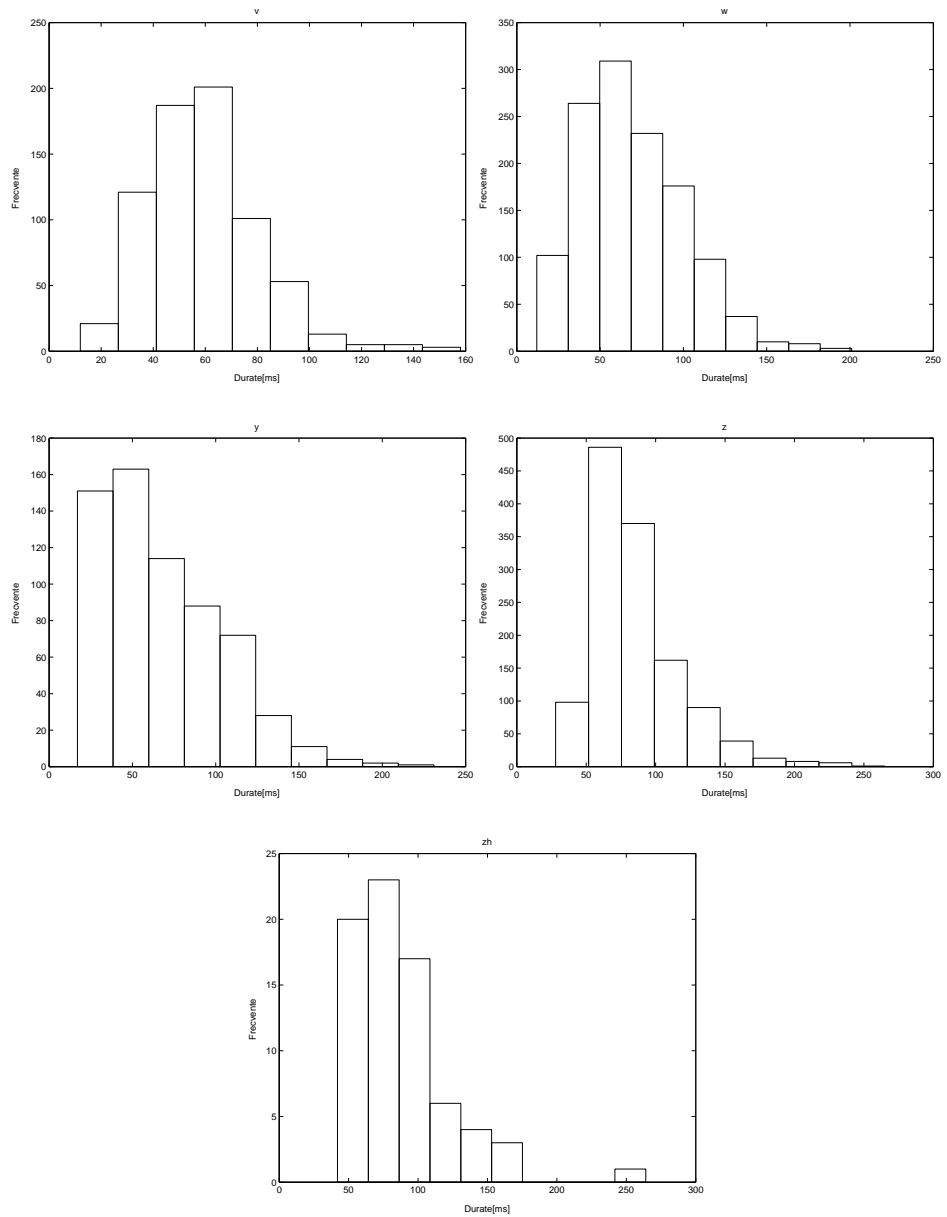
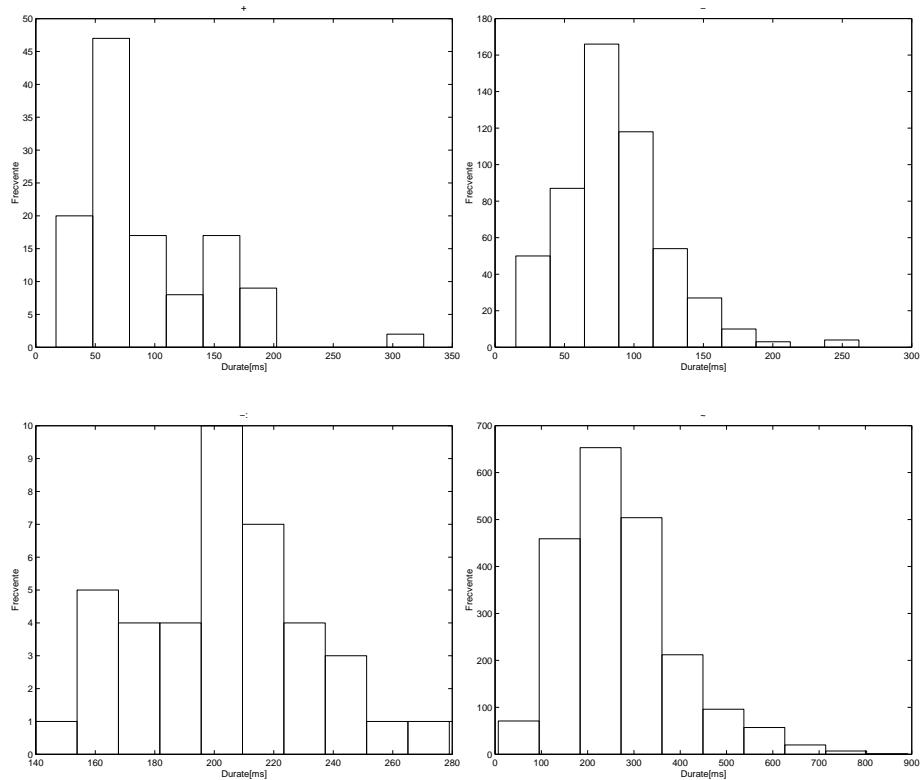


Figura A.8: Histograme ale duratelor fonemelor: /v/, /w/, /y/, /z/, /hz/



## Anexa B

# Duratele fonemelor din baza de date OASIS



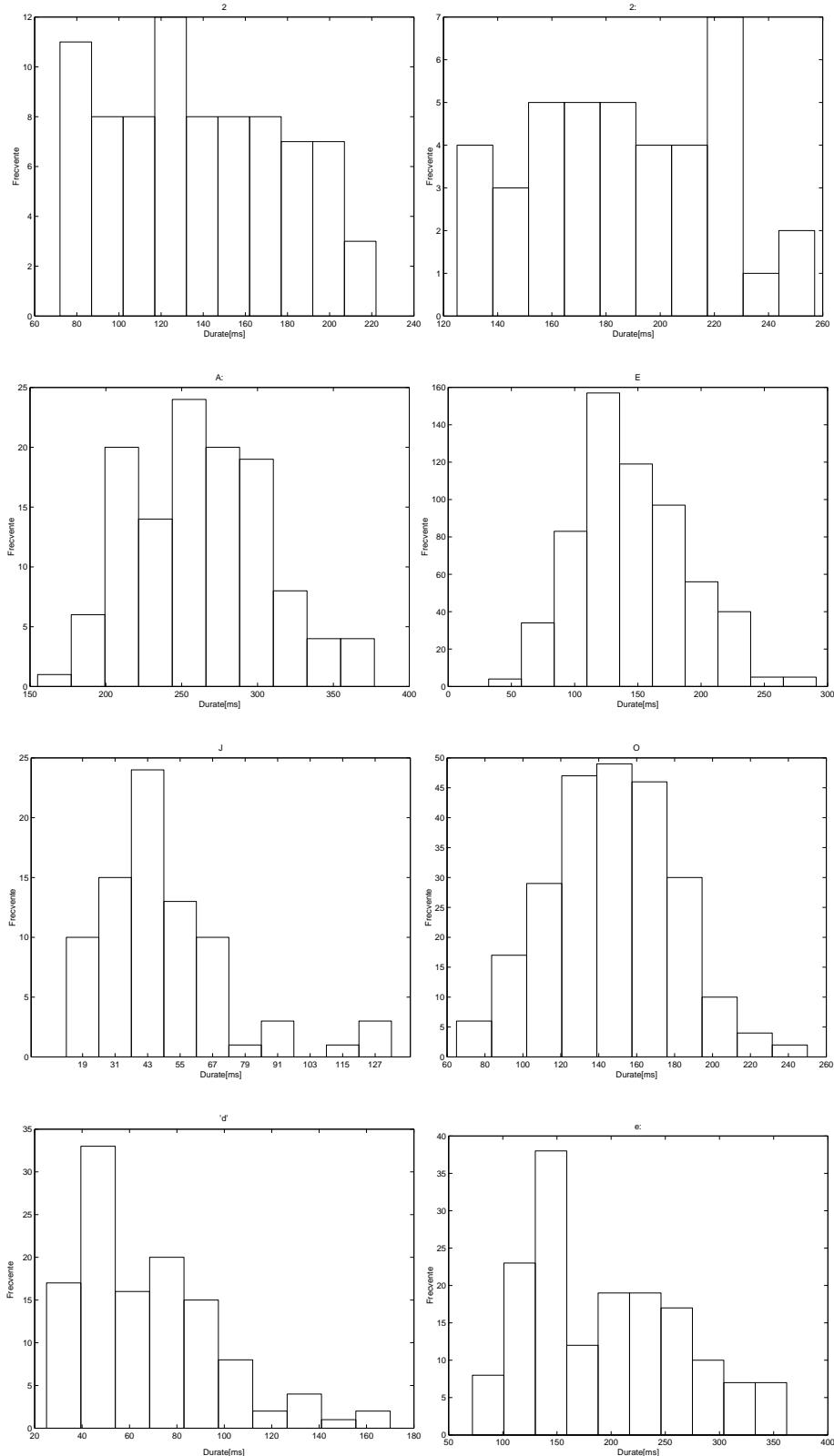


Figura B.1: Histogramme ale duratelor fonemelor: /ə/, /ə:/, /A/, /E/, /J/, /O/, /d/, /e:/

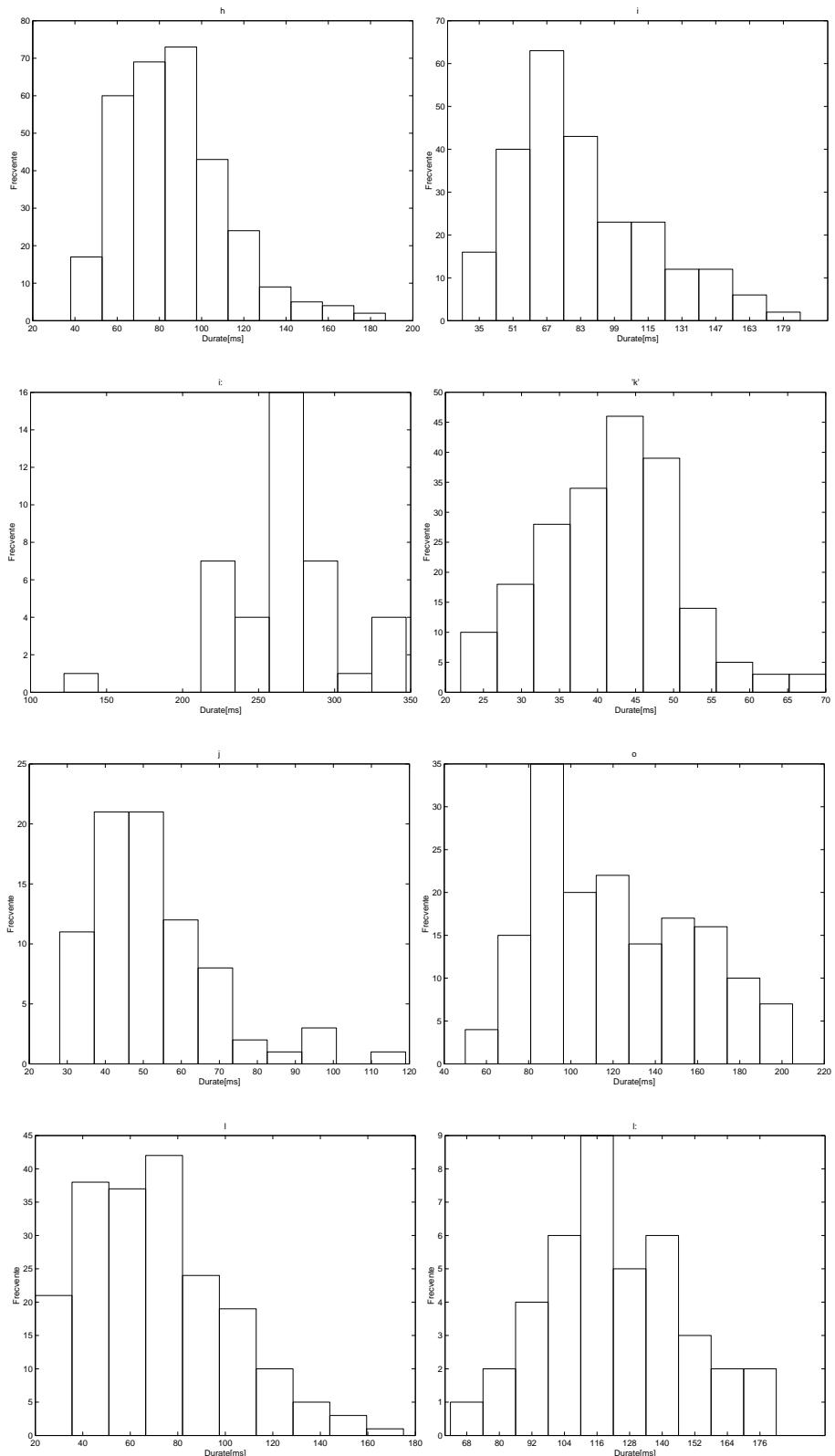


Figura B.2: Histogramme ale duratelor fonemelor: /h/, /i/, /i:/, /'k/, /j/, /o/, /l/, /l:/

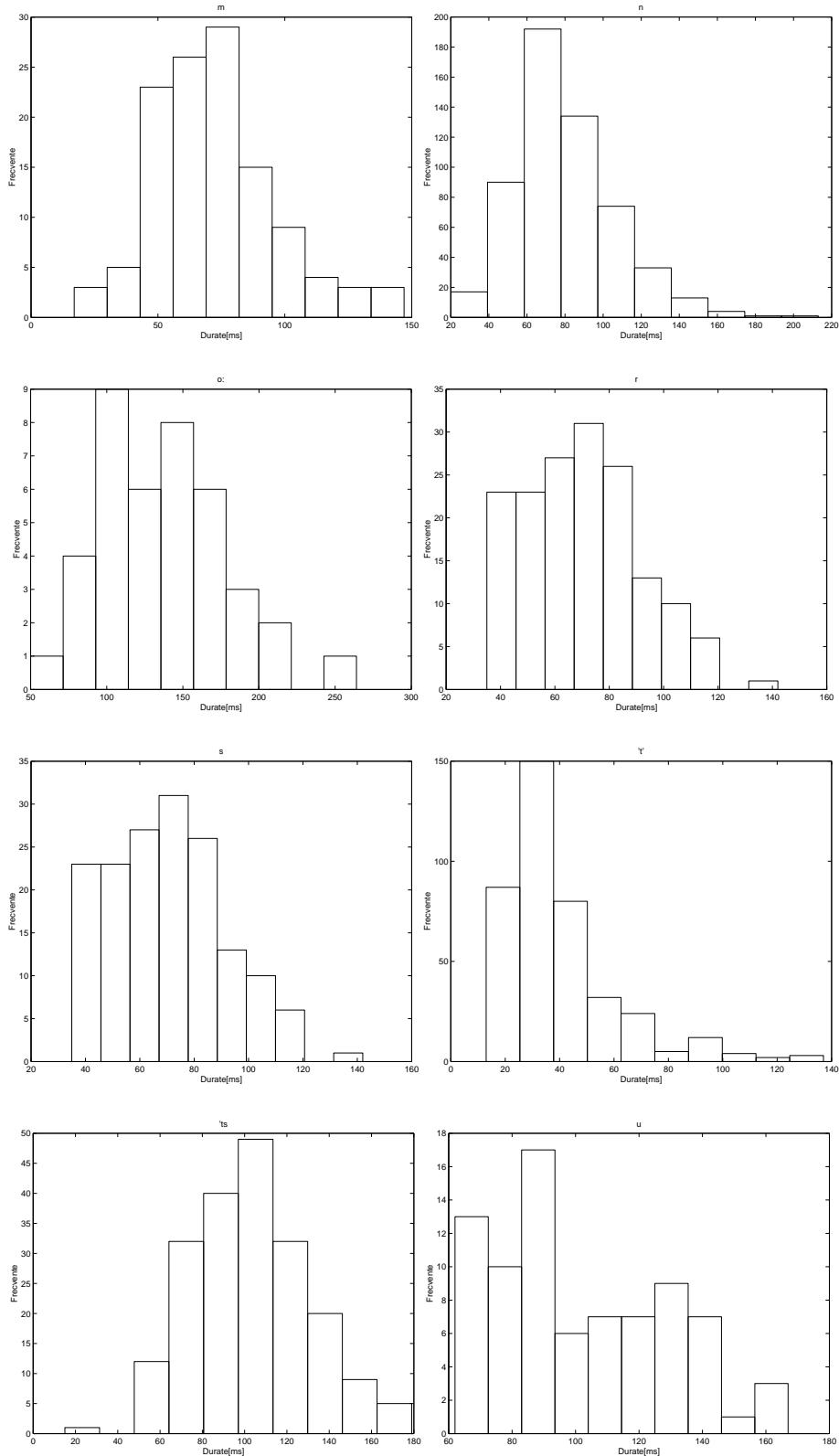


Figura B.3: Histogramme ale duratelor fonemelor: /m/, /n/, /o:/, /r/, /s/, /t/, /'ts/, /u/

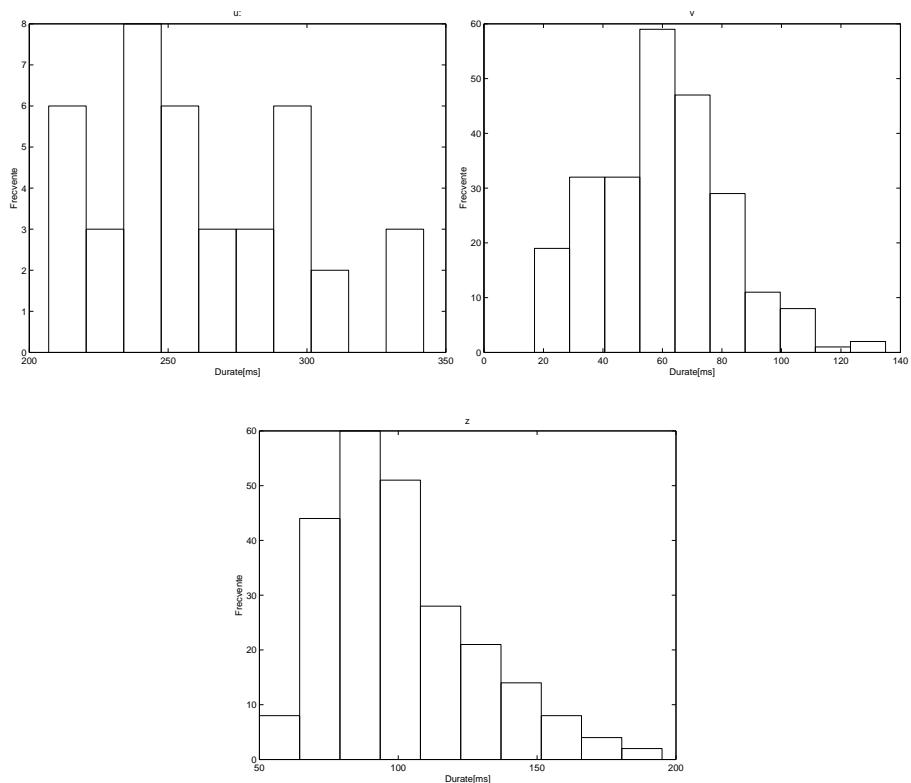


Figura B.4: Histogramme ale duratelor fonemelor: /u:/, /v/, /z/



## **Anexa C**

### **Foneme TIMIT**

Pentru segmentarea bazei de date TIMIT s-au utilizat 61 de simboluri. Părțile de închidere ale plozivelor s-au notat separat, astfel simbolurile **bcl**, **dcl**, **gcl**, **pcl**, **tck**, **kcl** reprezintă părțile de închidere ale fonemelor **b**, **d**, **g**, **p**, **t**, **k**.

Simbolul **h#** s-a utilizat pentru marcarea părților de liniște de la începutul, respectiv sfârșitul rostirilor. Există încă două simboluri pentru liniște: **pau** a fost utilizat pentru marcarea liniștii dintre cuvinte, iar **epi** este un fel de liniște, care apare între pronunțarea unei fricative și a unei semivocale.

Categorie	Simbol	Cuvânt	Transcriere fonetică
Plozive	b	bee	BCL B iy
	d	day	DCL D ey
	g	gay	GCL G ey
	p	pea	PCL P iy
	t	tea	TCL T iy
	k	key	KCL K iy
	dx	muddy, dirty	m ah DX iy, dcl d er DX iy
	q	bat	bcl b ae Q
Africate	jh	joke	DCL JH ow kcl k
	ch	choke	TCL CH ow kcl k
Fricative	s	sea	S iy
	sh	she	SH iy
	z	zone	Z ow n
	zh	azure	ae ZH er
	f	fin	F ih n
	th	thin	TH ih n
	v	van	V ae n
	dh	then	DH e n
Nazale	m	mom	M aa M
	n	noon	N uw N
	ng	sing	s ih NG
	em	bottom	b aa tcl t EM
	en	button	b ah q EN
	eng	washington	w aa sh ENG tcl t ax n
	nx	winner	w ih NX axr

Tabelul 29: Foneme TIMIT-1

Categorie	Simbol	Cuvânt	Transcriere fonetică
Semivocale	l	lay	L ey
	r	ray	R ey
	w	way	W ey
	y	yacht	Y aa tcl t
	hh	hay	HH ey
	hv	ahead	ax HV eh dcl d
	el	bottle	bcl b aa tcl t EL
Vocală	iy	beet	bcl b IY tcl t
	ih	bit	bcl b IH tcl t
	eh	bet	bcl b EH tcl t
	ey	bait	bcl b EY tcl t
	ae	bat	bcl b AE tcl t
	aa	bott	bcl b AA tcl t
	aw	bout	bcl b AW tcl t
	ay	bite	bcl b AY tcl t
	ah	but	bcl b AH tcl t
	ao	bought	bcl b AO tcl t
	oy	boy	bcl b OY
	ow	boat	bcl b OW tcl t
	uh	book	bcl b UH kcl k
	uw	boot	bcl b UW tcl t
	ux	toot	tcl t UX tcl t
	er	bird	bcl b ER dcl d
	ax	about	AX bcl b aw tcl t
	ix	debit	dcl d eh bcl b IX tcl t
	axr	butter	bcl b ah dx AXR
	ax-h	suspect	s AX-H s pcl p eh kcl k tcl t

Tabelul 30: Foneme TIMIT-2



## **Anexa D**

### **Vocabulary OASIS Numbers**

Transcriere fonetică	Număr
E + 'd'	1
E z E r	1000
h A: r o m	3
A: r o m	3
h O r m i n - 'ts	30
O r m i n - 'ts	30
h O - 't	6
O - 't	6
h O - 't v O n	60
O - 't v O n	60
h e: - 't	7
e: - 't	7
h E - 't v E n	70
E - 't v E n	70
h u: s	20
u: s	20
h u s o n	20-
u s o n	20-
'k e: - 't E z E r	2000
'k e: - s A: z	200
'k E - : 't 2:	2
'k i l E n - 'ts	9
'k i l E n - 'ts v E n	90
m i l i o:	1000000
n e: + 'd'	4
n E + 'd' v E n	40
n u l: O	0
J j o l - 'ts	8
J j o l - 'ts v O n	80
2 - 't	5
2 - 't v E n	50
s A: z	100
't i: z	10
't i z E n	10-

Tabelul 31: Vocabular alofonic OASIS Numbers

# Bibliografie

- [1] Antal, M., Statistical methods in speaker recognition, *Lucrările sesiunii naționale de comunicări științifice*, Univ. Petru Maior, Facultatea de Știinte Economice și Administrative, vol. 2, pp. 446-451, 2003.
- [2] Antal, M., Statistical Methods for Speaker Identification, *Proceedings of the 6th International Conference on Applied Informatics*, Jan. 27-31, Eger, Hungary, pp. 65-74, 2004.
- [3] Antal, M., Soós, A., Unsupervised Classification for Designing Speaker Identification Systems, *Studia Universitatis Babes-Bolyai Mathematica*, Vol. XLIX, No.1, pp. 3-14, 2004.
- [4] Antal, M., A Comparison of Parametric Clustering Techniques used in Speaker Identification, *IJSIT Lecture Notes of International Conference on Intelligent Knowledge Systems*, Aug. 16-20, Assos, Turkey, Vol.1, No.1, pp. 19-25, 2004.
- [5] Antal, M., Speaker Independent Phoneme Classification in Continuous Speech, *Studia Univ. Babeş-Bolyai Informatica*, Vol. XLIX, No. 2, pp. 55-64, 2004.
- [6] Antal, M., Phonetic Aspects of Speaker Identification, *15th International Conference in Computer Science and Education SzamOkt-2005*, Cluj-Napoca, Romania, Mar. 17-20, pp. 231-238, 2005.
- [7] Antal, M., Toderean, G., Phoneme Classification Using Two-Level Classification Scheme, *Proceedings of the 3rd Conference of Speech Technology and Human-Computer-Dialog*, Cluj-Napoca, Romania, May 13-14, pp. 231-238, 2005.

- [8] Antal, M., Toderean, G., Phonetic Analysis of GMM-based Speaker Models, Biometric Authentication Workshop, Cluj-Napoca, May 26-27, pp. 235-240, 2005.
- [9] Antal, M., Toderean, G., Speaker Recognition and Broad Phonetic Groups, Proceedings of the 24th IASTED International Multi-Conference on Signal Processing, Pattern Recognition and Applications, Febr. 15-17, Innsbruck, Austria, pp. 155-158, 2006.
- [10] Antal, M., Phoneme Recognition for ASR, Proc. of the International Conference Communications 2006, June 8-10, Bucharest, Romania, pp. 123-126, 2006.
- [11] Antal, M., Toderean, G., Broad Phonetic Classes Expressing Speaker Individuality, Studia Univ. Babes-Bolyai Informatica, Vol. LI, No. 1, pp. 49-58, 2006.
- [12] Aubert, X. L., A Brief Overview of Decoding Techniques for Large Vocabulary Continuous Speech Recognition, ASR-2000, Paris, France, , Sept. 18-20, pp. 91-97, 2000.
- [13] Baum, L., Petrie, T., Statistical Inference for Probabilistic Functions of Finite State Markov chains, Annals of Mathematical Statistics, vol. 37, pp. 1554-1563, 1966.
- [14] Bechetti, L. P. Ricotti, Speech Recognition, Theory and C++ Implementation, John Wiley & Sons, 1999.
- [15] BenZeghiba, M. F., Bourlard, H., On the Combination of Speech and Speaker Recognition, Proc. Eurospeech, Geneva, Switzerland, pp. 1361-1364, 2003.
- [16] Bimbot, F., Bonastre, J-F, Fredouille, C., Gravier, G., Margin-Chagnolleau, I., Meignier, S., Merlin, T., Ortega-Gracia, J., Petrovska-Delacretaz, D., Reynolds, D. A., A Tutorial on Text-Independent Speaker Verification, EURASIP Journal on Applied Signal Processing 4, 430-451, 2004.
- [17] Chagnolleau, I.M., Bonastre, J-F., Bimbot, F., Effect of utterance duration and phonetic content on speaker identification using second order statistical methods, Eurospeech, pp. 337-340, 1995.
- [18] Bourlard, H., Morgan, N., Connectionist Speech Recognition: A hybrid approach, Kluwer Academic Publishers, 1994.

- [19] Bourlard, H., Hermansky, H., Morgan., N., Towards Increasing Speech Recognition Error Rates, *Speech Communication*, Vol. 18., pp. 205-231, 1996.
- [20] Burileanu, D., Sima, M., Negrescu, C., Croitoru, V., Robust Recognition of Small-Vocabulary Telephone Quality Speech, *Speech Technology and Human-Computer Dialogue*, Ed, Academiei Române, Bucureşti, pp. 145-154, 2003.
- [21] Campbell, J.P., Speaker recognition: A tutorial, *Proc. IEEE*, vol. 85, no. 9., pp. 1437-1462, 1997.
- [22] Chengalvarayan, R., Deng, L., Speech Trajectory Discrimination Using the Minimum Classification Error learning, *IEEE Transaction on Speech and Audio Processing*, Vol. 6, No. 6, pp. 505-515, 1998.
- [23] Cormen, T. H., Leiserson, C. E., Rivest, R. R., *Introducere în algoritmi*, Agora, 2000.
- [24] Cover, T. M., Thomas, J. A., *Elements of information theory*, Wiley New York, 1991.
- [25] Das, S. K., Picheny, M. A., Issues in Practical Large Vocabulary Isolated Word Recognition: The IBM Tangora System, *Automatic Speech and Speaker Recognition, Advanced Topics*, ed. by Lee, C-H, Soong, F. K., Paliwal, K. K., Kluwer Academic Publisher, 1996.
- [26] Deller, J.R., Hansen, J. H. L., Proakis, J. G., *Discrete-Time Signal Processing of Speech Signals*, John Wiley & Sons, 2000.
- [27] Dempster, A. P., Laird, N. M., Rubin, D. B., Maximum Likelihood from Incomplete Data via the EM Algorithm, *Journal of the Royal Statistical Society*, vol. 39, pp. 1-88, 1977.
- [28] Doddington, G., Speaker Recognition based on Idiolectal Differences between Speakers, *Eurospeech*, pp. 2521-2524, 2001.
- [29] Duda, R. O. , Hart, P. E., *Pattern Classification*, John Wiley & Sons, 2000.
- [30] Dumitrescu, D., Costin, H., *Rețele neuronale*, Teora, 1996.
- [31] Dumitru, C. O., Gavat, I., Statistical, Neural and Hybrid Methods for Speech Recognition in Romanian Language, *Communications*, Proc. of the International Conference Communications, June 8-10, Bucharest, Romania, pp. 139-142, 2006.

- [32] Eatock, J.P., Mason, J. S., A Quantitative Assessment of the Relative Speaker Discriminating Properties of Phonemes, International Conference on Acoustics, Speech, and Signal Processing, pp. 333-336, 1994.
- [33] Falthauser, R., Ruske, G., Improving Speaker Recognition Performance Using Phonetically Structured Gaussian Mixture Models, Eurospeech, pp.751-754, 2001.
- [34] Ferrer, L., Bratt, H., Gadde, V. R. R., Kajarekar, S. S., Shriberg, E., Sonmez, K., Stolcke, A., Venkataraman, A., Modeling duration patterns for speaker recognition, Eurospeech, pp. 20017-2020, 2002.
- [35] Fukunaga, K., Introduction to Statistical Pattern Recognition, Academic Press, 1990.
- [36] Furui, S., An Overview of Speaker Recognition Technology, Automatic Speech and Speaker Recognition, Advanced Topics, ed. by Lee, C-H, Soong, F. K., Paliwal, K. K., Kluwer Academic Publisher, pp. 31-56, 1996.
- [37] Furui, S., Recent Advances in Speaker Identification, Pattern Recognition Letters, vol. 18, no. 9, pp. 859-872, 1997.
- [38] Furui, S., Digital Speech Processing, Synthesis and Recognition, Marcel Dekker, New York, 2001.
- [39] House, A. S., and K. N. Stevens, Analog studies of the nasalization of vowels, Journal of Speech and Hearing Disorders, vol. 21, pp. 218-232, 1956.
- [40] Huang, X., Acero, A., Hon, H-W., Spoken Language Processing, Prentice Hall, 2001.
- [41] Jain, A. K., Duin, R. P. W., Mao, J., Statistical Pattern Recognition, IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol.22, No. 1, pp. 4-37, 2000.
- [42] Jang, J-S. R., Sun, C.-T., Mizutani, E., Neuro-Fuzzy and Soft Computing, Prentice-Hall NJ., 1997
- [43] Jelinek, F. , Statistical Methods for Speech Recognition, MIT Press, 1997.
- [44] Jurafsky, D., Martin, J. H., Speech and Language Processing, Prentice Hall, 2000.

- [45] Kinnunen, T., Kilpelainen, T., Franti, O., Comparison of clustering algorithms in speaker identification, Proc. IASTED International Conference on Signal Processing and Communications, pp.222-227, 2000.
- [46] Kinnunen, T., Franti, P., Speaker Discriminative Weighting Method for VQ-based Speaker Identification, Proc. 3<sup>rd</sup> International Conference on audio and video-band biometric person authentication, pp.150-156, Halmstad, Sweden, 2001.
- [47] Kocsor, A., Kuba, A. Jr., Tóth, L., An Overview of the OASIS Speech Recognition Project, Proceedings of the 4th International Conference on Applied Informatics, August 30 - September 3, Eger-Noszvaj, Hungary, pp. 94-102, 1999.
- [48] Kocsor, A., Tóth, L., Kuba, A., Kovács, K., Jelasity, M., Gyimóthy, T., Csirik J., A Comparative Study of Several Feature Transformation and Learning Methods, International Journal of Speech Technology, pp. Vol. 3, Nr 3/4, pp. 253-262, 2000.
- [49] Kohonen, T., Self organizing Maps, Springer, 2001.
- [50] Knuth, D. E., Arta programării calculatoarelor, Vol. 3., Sortare și căutare, Teora, 2002.
- [51] Lee, K-F., Hon, H-W., Speaker-independent Phone Recognition Using Hidden Markov Models, IEEE Transactions on Acoustics, Speech, and Signal Processing, Vol. 37, No. 11, 1989.
- [52] Lee, C-H., Giachin, E., Rabiner, L. R., Pierrickini, R., Rosenberg, A. E., Improved Acoustic Modeling for Continuous Speech Recognition, Computer Speech and Language, vol. 6, no. 2, pp. 103-127, 1992.
- [53] Lee, C-H., Soong, F. K., Paliwal, K. K. (editors), Automatic Speech and Speaker Recognition, Kluwer Academic Publishers, 1996.
- [54] Levinson, S. E., Ljolje, A., Miller, L. G., Continuous Recognition from Phonetic Transcription, Proceedings of a workshop on Speech and Natural Language, Pennsylvania, U.S., pp. 190-199, 1990.
- [55] Lin, Q., Jan, E-E., Che, C., Yuk, D-S, Flanagan, J. L., Selective use of the speech spectrum and a VQGMM method for speaker identification, Fourth

- International Conference on Spoken Language Processing, pp. 1321-1324, 1996.
- [56] Logan, B., Moreno, P., Factorial HMMs for acoustic modeling, ICASSP, vol. 2, pp. 813-816, 1998.
  - [57] Lovell,B.,C., Caelli, T., Hidden Markov Models for spatio-temporal pattern recognition, Handbook of Pattern Recognition and Computer Vision, Chen, C.H., Wang, P.S.P. editors, World Scientific, pp. 25-40, 2005.
  - [58] Manning, C., Schütze, H., Foundations of Statistical Natural Language Processing, MIT Press Cambridge, 1999.
  - [59] Mari, J. F., Fohr, D., Junqua, J-C., A second order HMM for high performance word and phoneme-based speech recognition, IEEE Transactions on Speech and Audio Processing, vol. 23, no. 2, pp. 435-438, 1996.
  - [60] Mason, M., Vogt, R., Baker, B., Shridharan, S., The QUT NIST 2004 Speaker Verification System: A fused acoustic and high-level approach, Proc. of the 10th Australian International Conference on Speech Science & Technology, pp. 398-403, 2004.
  - [61] Merwe, R., Variations on Statistical Phoneme Recognition - a hybrid approach, master thesis, 1997.
  - [62] Mohri, M., Riley, M., Hindle, D., Ljolje, A., Pereira, F., Full Expansion of Context-Dependent Networks in Large Vocabulary Speech Recognition, Proc. ICASSP, pp. 665-668, 1998.
  - [63] Nedic, B., Bourlard, H., Recent Developments in Speaker Verification at IDIAP, IDIAP Research Report 00-26, 2000.
  - [64] Ney, H., Aubert, X., Dynamic Programming Search Strategies: From Digit Strings to Large Vocabulary Word Graphs, Automatic Speech and Speaker Recognition, Advanced Topics, ed. by Lee, C-H, Soong, F. K., Paliwal, K. K., Kluwer Academic Publisher, pp. 385-411, 1996.
  - [65] Patane, G., Russo, M., Fully Automatic Clustering System, IEEE Transactions on Neural Networks, Vol. 13, No. 6, November, pp. 1285-1298, 2002.
  - [66] Pylkkonen, J., Phone Duration Modeling Techniques in Continuous Speech Recognition, Master thesis Helsinki University of Technology, 2004.

- [67] Rabiner, L., Juang, B.H., Fundamentals of Speech Recognition, Prentice-Hall, Englewood Cliffs, NJ, 1993.
- [68] Rabiner, L. R., A tutorial on hidden markov models and selected applications in speech recognition, Readings in speech recognition, pp. 267-296, 1990.
- [69] Rabiner, L. R., Juang, B.-H., Lee, C. -H., An overview of Automatic Speech Recognition, Automatic Speech and Speaker Recognition, Advanced Topics, ed. by Lee, C-H, Soong, F. K., Paliwal, K. K., Kluwer Academic Publisher, pp. 1-30, 1996.
- [70] Reynolds, D. A., Speaker identification and verification using Gaussian mixture speaker models, Speech Communications 17, pp. 91-108, 1995.
- [71] Reynolds, D. A., Automatic Speaker Recognition using Gaussian Mixture Speaker Models, The Lincoln Laboratory Journal, vol. 8, nr. 2, pp. 173-192, 1995.
- [72] Riley, M. D., Ljolje, A., Automatic Generation of Detailed Pronunciation Lexicons, Automatic Speech and Speaker Recognition, Advanced Topics, ed. by Lee, C-H, Soong, F. K., Paliwal, K. K., Kluwer Academic Publisher, pp. 285-301, 1996.
- [73] Robinson, T., Fallside, F., A Recurrent Error Propagation Network Speech Recognition System, Computer Speech and Language, Vol. 5, No. 3, pp. 259-274, 1991.
- [74] Stapert, R., Mason, J.S., Speaker Recognition and the Acoustic Speech Space, Odyssey Speaker Recognition Workshop, Crete, pp. 195-199, 2001.
- [75] Stolojanu, G., Podaru, V., Cetină, F., Prelucrarea numerică a semnalului vocal, Editura Militară, 1984.
- [76] Szarvas, M., Matsunaga, S., Improving Phoneme Classification Performance Using Observation Context-Dependent Segment Models, International Journal of Speech Technology, pp. Vol. 3, Nr 3/4, pp. 263-276, 2000.
- [77] Szarvas, M., Fegyó, T., Mihajlik, P., Tatai, P., Automatic Recognition of Hungarian: Theory and Practice, International Journal of Speech Technology, pp. Vol. 3, Nr 3/4, pp. 237-251, 2000.
- [78] Schwarz, P., Matejka, P., Cernocky, J., Recognition of Phoneme Strings using TRAP Technique, Eurospeech, Geneve, pp. 825-828, 2003.

- [79] Toderean, G., Caruntu, A., Metode de recunoașterea vorbirii, Risoprint, Cluj-Napoca, 2005.
- [80] Tóth, L., Kocsor, A., Csirik, J., On naive Bayes in Speech Recognition, In. J. Appl. Math. Comput. Sci., Vol. 15, No. 2, pp. 287-294, 2005.