

UNIVERSITATEA "POLITEHNICA" BUCUREȘTI  
Facultatea de Electronică și Telecomunicații

Angel Cațaron

**REȚELE NEURALE PENTRU  
RECUNOAȘTEREA VORBIRII**

Teză de doctorat

Conducător științific:  
Prof. dr. ing. Victor-Emil Neagoe

BUCUREȘTI  
2003



Angel Cațaron

**REȚELE NEURALE PENTRU  
RECUNOAȘTEREA VORBIRII**

Teză de doctorat

Conducător științific:  
Prof. dr. ing. Victor-Emil Neagoe



## Rezumat

Această lucrare investighează, pe de o parte, *cum poate fi implementat conceptul de concurență la nivelul unei colecții de rețele neurale* și, pe de altă parte, *modul în care importanța diferită a intrărilor influențează performanțele în recunoașterea ale unor tipuri de rețele neurale*. Recunoașterea formelor, domeniu în care se plasează și subiectul acestei teze, oferă un câmp foarte larg cercetării. Recunoașterea vorbirii grupează o gamă diversă de probleme precum recunoașterea cuvintelor izolate, identificarea vorbitorilor pe baza amprentei lor vocale, recunoașterea fonemelor etc. Rețelele neurale sunt un instrument care și-a demonstrat eficiența în rezolvarea unei palete largi de aplicații între care și cele de recunoaștere automată a vorbirii.

Majoritatea modelelor neurale abordează recunoașterea formelor ca o problemă unitară, globală, fără a face distincția între aportul individual diferit al intrărilor.

Este cunoscut că modularitatea și principiul *divide et impera* aplicate rețelelor neurale le pot îmbunătăți performanțele. Introducem modelul *Concurrent Neural Networks (CNN)* care îmbină paradigmele învățării supervizate și ne-supervizate și oferă o soluție primei probleme. Ideea concurenței este folosită la nivelul unei colecții de rețele neurale care sunt instruite independent pentru a rezolva subprobleme diferite. Recunoașterea se face prin identificarea rețelei neurale care furnizează cel mai bun răspuns. Așa cum o demonstrează și rezultatele experimentale, acuratețea recunoașterii este mai mare atunci când folosim modelul propus de noi față de cazurile în care nu folosim concurența.

A doua problemă are conotații mai largi deoarece se regăsește și în aplicații de *data mining*. Prin asocierea unei ponderi fiecărei intrări, se stabilește intensitatea cu care intrările influențează răspunsul sistemului și implicit capacitatea sa de recunoaștere corectă. Dezvoltăm algoritmul *Ordered Weighted Aggregation - Relevance LVQ (OWA-RLVQ)* care calculează ponderile sub forma unui set de coeficienți *Ordered Weighted Aggregation (OWA)*. Algoritmii *Mutual Information Relevance LVQ (MIRLVQ)* și *Energy Relevance LVQ (ERLVQ)* propuși în continuare stabilesc acești coeficienți prin maximizarea informației mutuale pe de o parte, și a unei mărimi definite cu ajutorul energiei informaționale, pe de altă parte. Avantajele algoritmilor noștri sunt confirmate experimental prin rate de recunoaștere în mod constant superioare față de alți algoritmi care rezolvă același tip de probleme.

Mulțumesc d-lui Prof. dr. ing. Victor-Emil Neagoe care în calitate de conducător științific m-a sprijinit constant prin colaborarea fructuoasă de-a lungul acestor ani, prin discuțiile pline de substanță și prin sfaturile fără de care nu aș fi putut finaliza această teză.

De asemenea, mulțumesc tuturor colegilor care au influențat într-un fel sau altul conținutul acestei teze. Mulțumesc profesorilor care mi-au transmis pasiunea și dorința de a progresa. Multe mulțumiri prietenilor care m-au încurajat să merg mai departe în timpul perioadei de pregătire a tezei.

Mulțumesc părinților pentru că prin educația pe care mi-au dat-o și prin modul în care m-au crescut m-au ajutat să ajung până aici.

Mulțumesc familiei care m-a susținut în toată această perioadă, mi-a oferit suportul moral, motivația de care am avut nevoie și care m-a învățat ca nimic nu este greu pentru că nu sunt singur.

Brașov, septembrie 2003

# Notății

Am folosit caractere **aldine** pentru a nota vectorii și matricile și caractere *cursive* pentru componentele matricilor și pentru scalari în general.

Vectorii care apar în text sunt vectori coloană:

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

iar

$$\mathbf{x} = \begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix}^t$$

este un vector linie. Prin  $t$  am notat matricea transpusă.

Indicii sunt plasați în partea dreaptă jos a literei cu care se notează vectorul sau scalarul. În anumite contexte, pentru claritatea expunerii am folosit și indici superiori însoțiți de explicațiile de rigoare pentru a evita confuzia cu operația de ridicare la putere. În această relație,  $p$  este un indice:

$$E^p = \frac{1}{2} \sum_{j \in O} (y_j^p - d_j^p)^2.$$

În expresiile ale căror valori se calculează iterativ, am folosit notația  $(t)$  pentru a desemna pasul iterației:

$$\Delta w_{ji}^{(t)} = \alpha \Delta w_{ji}^{(t-1)} + \eta \delta_j^{(t)} y_i^{(t)}.$$

Produsul scalar a doi vectori  $\mathbf{x}_1$  și  $\mathbf{x}_2$  este:

$$\mathbf{x}_1^t \mathbf{x}_2 = \sum_{i=1}^n x_{i1} x_{i2}.$$

Norma euclidiană a unui vector este notată prin:

$$\|\mathbf{x}\| = \sqrt{\mathbf{x}^t \mathbf{x}}.$$

Pentru matricea unitate  $n$ -dimensională folosim următoarele notații:

$$\mathbf{I}_n = \mathbf{I} = \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{pmatrix}.$$

Preferăm varianta fără indice pentru simplitate, iar dimensiunea matricii se va deduce din context.

# Cuprins

|          |   |           |
|----------|---|-----------|
| <b>1</b> | <b>Introducere</b>  | <b>11</b> |
| 1.1      | Recunoașterea formelor . . . . .                                  | 11        |
| 1.2      | Recunoașterea vorbirii . . . . .                                  | 14        |
| 1.3      | Structura lucrării . . . . .                                      | 18        |
| <b>2</b> | <b>Rețele neurale. Concepte fundamentale</b>                      | <b>21</b> |
| 2.1      | Introducere . . . . .   | 21        |
| 2.2      | Modelul McCulloch-Pitts . . . . .                                 | 22        |
| 2.3      | Perceptronul . . . . .  | 24        |
| 2.4      | Reguli de învățare . . . . .                                      | 26        |
| 2.5      | Modele de rețele neurale artificiale . . . . .                    | 29        |
| 2.6      | Calcul neural . . . . .   | 32        |
| 2.7      | Perceptronul multistrat . . . . .                                 | 33        |
| 2.7.1    | Algoritmul <i>backpropagation or error</i> . . . . .              | 36        |
| 2.7.2    | Parametrii învățării . . . . .                                    | 38        |
| 2.8      | Rețele neurale cu întârziere în timp . . . . .                    | 40        |
| 2.9      | Rețele neurale cu autoorganizare . . . . .                        | 43        |
| 2.9.1    | Învățarea hebbiană . . . . .                                      | 43        |
| 2.9.2    | Învățarea competitivă . . . . .                                   | 44        |
| 2.10     | Concluzii . . . . .   | 52        |
| <b>3</b> | <b>Recunoașterea vorbirii cu modele clasice de rețele neurale</b> | <b>53</b> |
| 3.1      | Introducere . . . . .   | 53        |
| 3.2      | Bazele de date Speechdata și Telephdata . . . . .                 | 55        |
| 3.3      | Experimente . . . . .   | 58        |
| 3.3.1    | Recunoașterea cuvintelor . . . . .                                | 60        |
| 3.3.2    | Recunoașterea vorbitorilor . . . . .                              | 61        |
| 3.4      | Concluzii . . . . .   | 61        |

|          |  |            |
|----------|--|------------|
| <b>4</b> | <b>Rețele neurale concurente</b>   | <b>63</b>  |
| 4.1      | Introducere . . . . .  | 63         |
| 4.2      | Rețele neurale concurente . . . . .  | 65         |
| 4.3      | Recunoașterea vorbitorilor cu rețele neurale<br>concurente. Experimente . . . . .                    | 68         |
| 4.4      | Recunoașterea vocalelor cu rețele neurale<br>concurente . . . . .                                    | 72         |
| 4.4.1    | Fonetică și fonologie în limba română . . . . .  | 72         |
| 4.4.2    | Experimente de recunoaștere a vocalelor . . . . .  | 76         |
| 4.5      | Concluzii . . . . .  | 86         |
| <b>5</b> | <b>Recunoaștere prin calculul relevanțelor intrărilor<br/>folosind operatori OWA</b>                 | <b>87</b>  |
| 5.1      | Introducere . . . . .  | 87         |
| 5.2      | Algoritmul Relevance LVQ (RLVQ) . . . . .  | 89         |
| 5.3      | Operatorii OWA . . . . .   | 91         |
| 5.3.1    | Problema generală a agregării . . . . .  | 91         |
| 5.3.2    | Conceptul de operator OWA . . . . .  | 92         |
| 5.3.3    | Instruirea operatorilor OWA . . . . .  | 93         |
| 5.4      | Algoritmul OWA - Relevance LVQ (OWA-RLVQ) . . . . .  | 95         |
| 5.5      | Experimente . . . . .  | 98         |
| 5.6      | Concluzii . . . . .  | 100        |
| <b>6</b> | <b>Recunoaștere prin calculul relevanțelor intrărilor<br/>folosind noțiuni de teoria informației</b> | <b>101</b> |
| 6.1      | Introducere . . . . .  | 101        |
| 6.2      | Entropia Shannon . . . . .   | 103        |
| 6.3      | Estimarea densităților de probabilitate cu<br>ferestre Parzen . . . . .                              | 105        |
| 6.4      | Informația mutuală pătratică . . . . .   | 112        |
| 6.5      | Algoritmul Mutual Information Relevance LVQ (MIRLVQ) . . . . .                                       | 115        |
| 6.6      | Energia informațională . . . . .   | 121        |
| 6.7      | Algoritmul Energy Relevance LVQ (ERLVQ) . . . . .  | 123        |
| 6.8      | Experimente . . . . .  | 128        |
| 6.9      | Concluzii . . . . .  | 130        |
| <b>7</b> | <b>Concluzii</b>   | <b>131</b> |
| 7.1      | Introducere . . . . .  | 131        |
| 7.2      | Contribuții . . . . .  | 131        |
| 7.3      | direcții viitoare de cercetare . . . . .   | 137        |

|                     |            |
|---------------------|------------|
| <b>Anexa A</b>      | <b>139</b> |
| <b>Anexa B</b>      | <b>141</b> |
| <b>Bibliografie</b> | <b>145</b> |



# Capitolul 1

## Introducere

Această lucrare abordează o temă recunoaștere a formelor folosind rețele neurale, cu aplicații în recunoașterea vorbirii. Sunt introduse mai multe modele originale de recunoaștere, iar dezvoltările teoretice și rezultatele experimentale demonstrează viabilitatea acestora. Spectrul aplicațiilor de recunoaștere a vorbirii abordate în teză este larg, începând cu recunoașterea fonemelor, continuând cu recunoașterea cuvintelor izolate și terminând cu recunoașterea vorbitorilor. Bazele de date sunt diverse, fiind vorba atât de colecții standard pentru limba engleză cât și de seturi alcătuite special pentru experimentele din această lucrare și conțin cuvinte din limba română. Am testat unele modele propuse de noi și pe seturi de date standard în recunoașterea formelor, cu scopul de a ilustra faptul că ele pot fi folosite și în alte contexte decât cel al recunoașterii vorbirii.

Acest capitol descrie principalele domenii care sunt abordate în lucrare, iar în final conține o prezentare a structurii tezei.

### 1.1 Recunoașterea formelor

Înțelegerea cuvintelor vorbite, recunoașterea unei fețe umane, citirea caracterelor scrise de mână sunt acțiuni firești, specific umane, în spatele cărora stau procese complexe care nu sunt încă deslușite, dar pe care dorim să le modelăm prin sisteme de recunoaștere automată. *Recunoașterea formelor* înseamnă analiza datelor care provin de la o intrare fizică în sistem și încadrarea acestora în anumite categorii. De la recunoașterea automată a vorbirii, identificarea amprentelor digitale până la recunoașterea caracterelor scrise de mână și identificarea secvențelor ADN, gama de aplicații a recunoașterii formelor este imensă, iar crearea unor mașini care să realizeze acest lucru automat și cu acuratețe este extrem de utilă. Pentru multe probleme de recunoaștere eforturile de gășire a

soluțiilor sunt, de fapt, influențate de cunoștințele pe care le avem despre modul în care sunt rezolvate acestea de creierul uman. Rețelele neurale artificiale reprezintă un suport care, în ultimii ani, s-a dovedit foarte performant.

Intrarea unui sistem de recunoaștere a formelor (fig. 1.1) este în majoritatea cazurilor un dispozitiv care preia informația primară și o convertește într-o colecție de date. Un astfel de dispozitiv poate fi alcătuit dintr-unul sau mai multe microfoane, una sau mai multe camere de luat vederi etc.

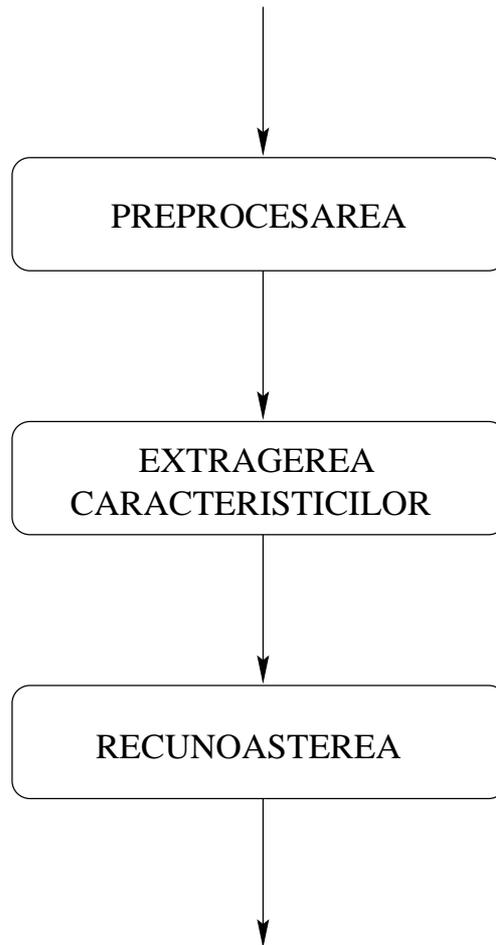


Figura 1.1: Arhitectura unui sistem de recunoaștere a formelor.

Colecția de date este sursa de informație a întregului sistem, iar modulul de preprocesare trebuie să segmenteze datele, adică să stabilească unde începe și unde se termină fiecare formă. Vom folosi în continuare termenul *formă* care are aceeași semnificație cu cea a cuvântului *pattern* din limba engleză. Segmentarea este una dintre cele mai dificile probleme în recunoașterea formelor.

În problemele de recunoaștere automată a vorbirii putem să ne propunem ca segmentarea să însemne găsirea unor entități de bază cum ar fi fonemele. Acestea ar putea fi apoi grupate astfel încât să alcătuiască silabe sau cuvinte. Fenomenul coarticulației face, însă, ca două foneme succesive să se influențeze reciproc, iar rezultatul este o variație de la pronunția standard prin modificarea spectrului de frecvențe specific fiecărui fonem, prin modificarea amplitudinii sau a duratei. Coarticularea induce astfel un element de confuzie și limita dintre foneme nu mai este foarte clară. Problema este delicată și poate influența negativ performanțele întregului sistem de recunoaștere a formelor.

Frontiera între extragerea caracteristicilor și recunoaștere nu este întotdeauna foarte clară [17]. Un modul ideal de extragere a caracteristicilor ar trebui să însemne o recunoaștere foarte ușoară, în timp ce un sistem foarte puternic de recunoaștere nu ar trebui să aibă nevoie de un modul de extragere a caracteristicilor foarte complex. Din punct de vedere practic, însă, este foarte dificil de construit module omnipotente.

Rolul modulului de extragere a caracteristicilor este de a genera un set de însușiri ale căror valori sunt similare obiectelor care fac parte din aceeași categorie și foarte diferite pentru obiectele din categorii diferite. Această afirmație ne conduce la ideea căutării trăsăturilor distinctive care sunt invariante la transformări irelevante asupra intrărilor. În recunoașterea vorbirii, dorim, de exemplu, ca însușirile să fie invariante la viteza de pronunție a cuvintelor de aceeași persoană sau de persoane diferite.

Extragerea trăsăturilor, ca și segmentarea, este o problemă dependentă de tipul aplicației și implică diverse cunoștințe specifice domeniului său. În urma acestei etape, se introduce un nivel de abstractizare a intrărilor și clasificatorul va deveni mai puțin dependent de aplicație.

Ca parte a întregului sistem, modulul de recunoaștere primește ca intrări vectori pe care îi asignează unei categorii sau unei clase. Implementarea perfectă a recunoașterii este, în majoritatea cazurilor, imposibilă, de aceea se preferă stabilirea unei marje de eroare în limitele căreia aceasta este considerată corectă. Dificultatea implementării poate să provină din complexitatea intrinsecă a problemei și atunci decizia este una netrivială.

Procesul prin care clasificatorul înmagazinează informația de care are nevoie pentru a fi capabil să asocieze o formă unei categorii se numește instruire. Instruirea presupune un model, un set de proceduri și date cu ajutorul cărora se stabilesc parametrii clasificatorului. Deși un sistem complex poate oferi o clasificare perfectă a datelor de instruire, se poate constata uneori că performanțele nu mai sunt aceleași pentru forme noi. De aceea, un clasificator bun va realiza întotdeauna echilibrul între un model nu prea simplu, care să poată face

distincția între clase și unul nu foarte complex, pentru a evita clasificarea slabă a formelor noi.

## 1.2 Recunoașterea vorbirii

Vorbirea este un proces care, pentru noi, este foarte natural și ușor. Încă de la vârste foarte mici căpătăm abilitățile necesare înțelegerii și transmiterii mesajelor vorbite. Învățarea se desfășoară în mai multe etape. Mai întâi recepționăm mesajul pe care apoi încercăm să îl reproducem. Înțelegerea sensului mesajului nu ne interesează în acest studiu. Pentru o învățare performantă și facilă este necesară existența unui set de obișnuințe preliminare, cum sunt structurile elementare de limbă și modul lor de pronunțare. Se presupune așadar că persoana este vorbitoare a limbii respective. În caz contrar, învățarea este ceva mai dificilă pentru că lipsesc cunoștințele preliminare. În ambele situații, însă, o ființă umană va fi capabilă să repete cu o mare acuratețe noțiunile pe care și le-a însușit în urma unui proces de învățare.

Tractul vocal și articulator sunt organe cu proprietăți neliniare, iar comportamentul lor este afectat de diferiți factori, de la genul persoanei până la starea ei emoțională. De aceea, felul în care sunt pronunțate cuvintele poate varia din punct de vedere al accentului, al articulării cuvintelor, al volumului sau al vitezei. În fine, pot interveni factorii perturbatori externi cum sunt zgomotul de fond sau ecoul, dar și caracteristicile electrice ale echipamentelor electronice care sunt folosite drept suport de stocare sau ale liniilor telefonice utilizate în transmisiile la distanță.

Înțelegerea vorbirii care pentru noi este o obișnuință, devine o problemă foarte complexă dacă dorim să o repetăm prin tehnici care țin de învățarea și recunoașterea automată. Aparent simplă, crearea unui sistem artificial care să recunoască limbajul vorbit s-a dovedit a fi, totuși, destul de dificilă.

Sistemele actuale de recunoaștere a vorbirii aparțin următoarelor categorii [16]:

- a) Sistemele care acceptă un vocabular de dimensiuni reduse, cu 10-100 de cuvinte;
- b) Sistemele care pot recunoaște cuvintele pronunțate izolat unul de celălalt și pentru care dimensiunea vocabularului poate atinge 10000 de cuvinte;
- c) Sisteme care acceptă la intrare fraze, recunoscând cuvintele chiar dacă acestea sunt pronunțate cursiv, fără pauze. Majoritatea acestor sisteme

au, în mod obișnuit, un vocabular restricționat la un anumit domeniu și recunosc 1000-5000 de cuvinte.

Majoritatea sistemelor actuale utilizate în aplicațiile practice fie folosesc un vocabular de dimensiuni reduse, fie recunosc cuvinte pronunțate izolat. Sistemele folosite pentru o comunicare naturală între om și calculator rămân un deziderat care nu a fost încă atins. Nici o mașină nu a reușit deocamdată să atingă performanțe maxime în aplicațiile practice.

Un sistem de recunoaștere a vorbirii va conduce întotdeauna la rezultate cu atât mai bune cu cât va fi folosit de mai puține persoane. Dacă persoanele care au antrenat sistemul sunt și cele care îl folosesc, atunci rezultatele sunt și mai bune. Pe de altă parte, utilizatorii care pronunță cuvintele cu claritate și mai rar se fac mai ușor înțeleși de un astfel de sistem. În contrast cu aceste considerații, aplicațiile de maximă utilitate și cu caracter universal care ar fi capabile să răspundă unor cerințe reale și complexe nu trebuie să folosească vocabular limitat, să impună o anumită modalitate de pronunție a cuvintelor sau o anumită topică a frazelor, să limiteze numărul persoanelor care pot folosi sistemul ori să fie influențat de perturbațiile mediului în care este folosit. Mai mult, un sistem ideal de recunoaștere a vorbirii ar trebui să fie capabil să se adapteze noilor tendințe ale limbii, să poată înmagazina și folosi informația legată de modificările lexicale, sintactice sau semantice. Toate acestea demonstrează complexitatea domeniului și dificultatea problemelor care se pun.

Parametrii vorbirii se obțin prin diverse metode de analiză [16], [72] care concentrează caracteristicile semnalului vocal. În felul acesta, se produc secvențe de vectori de trăsături care, în mod uzual, se numesc *pattern-uri ale vorbirii*. O abordare în recunoașterea formelor este cea conexionistă, care încorporează cunoștințele legate de trăsături într-o rețea densă de noduri și conexiuni, sub forma modelului de calcul cunoscut prin numele de *rețea neurală*. El combină elemente acustico-fonetice, elemente de recunoașterea formelor și de inteligență artificială.

Abordarea *acustico-fonetice* se bazează pe teoria conform căreia limba vorbită este alcătuită din unități fonetice distincte. Acestea sunt caracterizate de un set de proprietăți care se manifestă în semnalul vorbirii sau în spectru. Deși aceste proprietăți acustice ale unităților fonetice variază destul de mult în funcție de contextul în care se află și de vorbitor, se presupune că există un set de trăsături de bază care pot fi învățate și aplicate în situații practice. Primul pas în recunoașterea vorbirii prin prisma abordării acustico-fonetice este faza de segmentare și etichetare. Semnalul vocal este separat în regiuni distincte în care proprietățile acustice ale semnalului sunt reprezentative pentru una dintre

unitățile fonetice sau una dintre clasele de unități fonetice. Fiecarui astfel de segment i se atașează apoi o etichetă care marchează apartenența sa la o unitate fonetică sau clasă de unități fonetice. Cea de-a doua etapă constă în determinarea în interiorul unui cuvânt sau a unei secvențe de cuvinte a unităților fonetice etichetate la pasul anterior.

Abordarea recunoașterii vorbirii care se bazează pe *recunoașterea formelor* este una prin care vectorii de test, a căror apartenență la o anumită clasă nu este cunoscută, sunt comparați cu un set de vectori de referință. Decizia se ia pe baza unei măsuri de similaritate dintre vectorul necunoscut și fiecare vector de referință. Această metodă constă din doi pași: instruirea vectorilor vorbirii și recunoașterea unor vectori noi prin metoda comparației. Cu cât sunt folosite în setul de instruire mai multe versiuni ale formei care trebuie recunoscută, cu atât mai performantă va fi procedura de instruire în extragerea proprietăților acustice ale sale.

Abordarea prin prisma *inteligenței artificiale* a recunoașterii vorbirii este o metodă hibridă și exploatează idei și concepte din cele două abordări descrise mai sus, cea acustico-fonetice și cea bazată pe recunoașterea formelor. Ideea fundamentală este aceea de a încorpora în sistemul de recunoaștere cunoștințe care provin de la o varietate largă de surse. O tehnică ce poate fi încadrată în această clasă este folosirea unui sistem expert pentru segmentare și etichetare care realizează această operație nu numai pe baza informațiilor acustice, așa cum se întâmplă prin metodele acustico-fonetice. Un sistem expert poate integra informații fonetice, lexicale, sintactice sau semantice dar, în plus, este capabil să învețe și să se adapteze.

Folosirea rețelelor neurale reprezintă o abordare distinctă a recunoașterii vorbirii, dar poate fi privită și ca o structură care combină caracteristici ale tehnicilor descrise mai sus.

## Recunoașterea vocalelor

Producția vocală este o însiruire de unități fonetice care au caracteristici distincte și care stau la baza comunicării mesajelor între oameni. Unitățile fonetice (fonemele) au caracteristici fonice și articulatorii proprii care ajută la înțelegerea lor și care fac posibilă transmiterea mesajelor verbale. Modul lor de pronunție este, însă, influențat de caracteristicile fizice ale aparatului fonator propriu fiecărui individ. Aceste variații de pronunție nu influențează calitatea informației transmise între doi parteneri de dialog, însă creează probleme la recunoașterea automată a fonemelor. Unitățile fonetice alcătuiesc sistemul fonetic al unei limbi și sunt împărțite în două categorii: vocale și

consoane. Vocalele reprezintă sunete produse la trecerea liberă a coloanei de aer prin traiectul vocal, iar consoanele apar atunci când fluxul de aer întâlnește diferite obstacole. Vocalele sunt produse de vibrația coardelor vocale, dar și de buze, poziția limbii, a dinților și a altor organe vocale care influențează sunetul rezultat. Ele au în general o durată mai mare în comparație cu consoanele și sunt mai bine definite în domeniul spectral.

Sistemele de recunoaștere a vorbirii folosesc recunoașterea vocalelor ca o modalitate de creștere semnificativă a abilității lor de recunoaștere.

### Recunoașterea cuvintelor izolate

Despre sistemele care recunosc entități discrete, de regulă cuvinte, se spune că folosesc la recunoașterea cuvintelor izolate. În faza de recunoaștere, se presupune că vorbitorii utilizează în mod deliberat cuvinte rostite rar, cu pauze între ele. În mod uzual, acestea au minim 200 ms pentru ca sistemul să nu le confunde cu fonemele care în general sunt mai scurte. Faptul că limita dintre cuvinte este ușor de stabilit simplifică semnificativ recunoașterea. Limitele pot fi localizate prin diverse tehnici cum ar fi folosirea unor algoritmi pentru detecția punctelor de start și de sfârșit ale cuvintelor. Acest tip de sistem de recunoaștere este util aplicațiilor în care se utilizează comenzi sub forma unor cuvinte simple, dintr-un vocabular redus.

### Recunoașterea vorbitorilor

Expresia vocală este o caracteristică proprie fiecărui individ. Din acest motiv, este posibilă recunoașterea în condiții normale a unui vorbitor, de exemplu a partenerului unei conversații telefonice. Variațiile individuale între vorbitori sunt influențate de doi factori. În primul rând, este vorba de proprietățile fizice ale aparatului fonator care determină caracteristicile de frecvență ale semnalului vocal. Pe de altă parte, o frază nu este pronunțată la fel de doi vorbitori. Întotdeauna vor apărea diferențe în debitul verbal, diferențe de ton, de accent etc. Așa cum sugerează și termenul, *recunoașterea unui vorbitor* se referă la problema identificării unui individ dintr-o populație, fără a interesa mesajul transmis.

Atunci când recunoașterea vorbitorilor se face pe baza unui vocabular limitat, bine determinat, sau pe baza unui text impus, este vorba de tehnica numită *fixed-text speaker recognition*. Dacă recunoașterea este independentă de text, este vorba de *free-text speaker recognition*.

În strânsă legătură cu recunoașterea vorbitorului este *verificarea vorbitorului*. De această dată problema constă în determinarea identității declarate de

un individ. De regulă, se compară distanța dintre expresia sa vocală și referința sa personală.

### 1.3 Structura lucrării

Teza este organizată în șapte capitole care conțin părți teoretice și descrieri ale experimentelor de recunoaștere a formelor, iar contribuțiile originale sunt argumentate în detaliu. Experimentele de recunoaștere vocală sunt diverse, pornind de la recunoașterea cuvintelor izolate, continuând cu recunoașterea vorbitorilor și terminând cu recunoașterea vocalelor. Bazele de date folosite acoperă o gamă largă, unele fiind publice, altele alcătuite special pentru experimentele din această lucrare. În unele experimente am utilizat și baze de date care nu sunt destinate strict recunoașterii vorbirii, demonstrând astfel că algoritmi propuși pot fi folosiți și pentru probleme generale de recunoaștere a formelor.

Capitolul 2 face o prezentare a tipurilor fundamentale de rețele neurale și a tehnicilor de instruire care stau la baza modelelor introduse în celelalte capitole. Neuronul McCulloch-Pitts reprezintă punctul de pornire al domeniului conexiionist. Deși este foarte simplu ca structură, acest neuron artificial inspirat ca structură și funcționare de procesele neurale care au loc în creierul uman s-a dovedit a avea un mare potențial de calcul. Pot fi implementate cu ajutorul său chiar și funcții logice oricât de complicate. Perceptronul introdus de Rosenblatt este un model mai evoluat, fiind capabil să învețe să clasifice vectori liniar separabili. Perceptronul multistrat constă dintr-o colecție de perceptroni elementari, organizați ierarhic. Este un model adaptiv avansat care poate fi instruit și care poate implementa funcții oricât de complicate. Datorită acestei proprietăți, poate fi utilizat ca și clasificator pentru vectori neseparabili liniar. Algoritmul de instruire *backpropagation of error* se bazează pe tehnica gradientului negativ și permite ajustarea iterativă, supervizată parametrilor perceptronului multistrat pornind de la un set de date de instruire. Rețeaua neurală cu întârziere în timp este un perceptron multistrat extins pentru probleme de recunoaștere a vorbirii prin introducerea unor elemente de întârziere a semnalelor de intrare. Învățarea competitivă constă dintr-un proces de selecție în urma căruia vectorii de instruire sunt grupați după caracteristicile comune. Algoritmul *self-organizing feature map* este nesupervizat și realizează o transformare a unui hiperspațiu cu un număr mare de dimensiuni într-unul unidimensional sau bidimensional, iar reprezentarea în plan a acestei transformări se numește hartă cu autoorganizare. Algoritmii din clasa *learning vector quantization* sunt supervizați și aproximează setul de date de instruire printr-o mulțime finită de

vectors numiți prototipuri.

În capitolul 3 este prezentat un set de experimente de recunoaștere a cuvintelor izolate și a vorbitorilor folosind trei modele clasice de rețele neurale: perceptronul multistrat, rețeaua neurală cu întârziere în timp și harta cu autoorganizare. Cele două baze de date, Speechdata și Telephdata, pe care le-am folosit în aceste teste constau din cuvinte izolate pronunțate de diverși vorbitori mai întâi în condiții de studio, fără perturbații, și apoi prin telefon. Pentru preprocesarea semnalelor vocale se folosește o cunoscută metodă de calcul al coeficienților cepstrali. Problema instruirii și cea a recunoașterii informației din aceste două baze de date sunt deosebit de complicate și rezultatele experimentale demonstrează acest lucru. Recunoașterea cuvintelor izolate conduce la rezultate satisfăcătoare în majoritatea cazurilor, însă recunoașterea vorbitorilor ridică probleme deosebite.

Rețelele neurale concurente introduse în capitolul 4 reprezintă un model conexiunist de recunoaștere care pornește de la observația că o rețea neurală este mai eficientă atunci când implementează o problemă mai puțin complexă sau când este specializată pe o subproblemă a problemei globale de recunoaștere. Modelul propus în acest capitol este o colecție de rețele neurale care lucrează în paralel, iar decizia în recunoaștere se ia conform regulii *câștigătorul ia tot*. Fiecare rețea neurală componentă este instruită individual cu un set propriu de date. Fiind vorba de module care rezolvă probleme cu o complexitate intrinsecă mai scăzută, numărul de neuroni pentru fiecare modul poate fi redus semnificativ. Am instruit rețele neurale concurente având drept unități elementare perceptroni multistrat, rețele neurale cu întârziere în timp și hărți cu autoorganizare pentru aceeași problemă de recunoaștere a vorbitorilor din capitolul 3, iar performanțele pentru cele două baze de date au crescut semnificativ. Cu același model am implementat și recunoașterea unui subset de șase vocale din limba română, baza de date fiind alcătuită prin segmentarea manuală a semnalelor vocale, în scopul extragerii secvențelor vocalizate. Preprocesarea constă din calculul coeficienților cepstrali care sunt vectorii de instruire. Rețeaua neurală concurentă care folosește hărți cu autoorganizare specializate pe recunoașterea câte unei vocale are performanțe mai bune decât o hartă unică, instruită cu toate vocalele.

În capitolul 5, problema recunoașterii este abordată dintr-o perspectivă diferită pentru că ne punem problema creșterii performanțelor prin selectarea acelor caracteristici ale vectorilor de instruire și de test care contribuie cel mai mult la o recunoaștere corectă. Fiecărei componente a vectorilor de intrare  $i$  se asociază câte un coeficient care reprezintă importanța acestuia, iar coeficienții sunt instruiți în așa fel încât să minimizeze eroarea de cuantizare a

unei rețele neurale de tip LVQ. Algoritmul RLVQ calculează printr-o tehnică euristică acești coeficienți. Operatorii OWA sunt o familie de operatori de agregare care se bazează pe reordonarea criteriilor care trebuie satisfăcute într-un proces de clasificare. Algoritmul OWA-RLVQ propus de noi calculează dinamic relevanțele intrărilor ca ponderi ale unui operator OWA, făcând totodată o legătură importantă între algoritmul RLVQ și modelul matematic al operatorilor OWA. Performanțele algoritmului nostru sunt comparate cu ratele de recunoaștere obținute cu LVQ și RLVQ și rezultatele înregistrate de OWA-RLVQ sunt superioare pentru toate cele trei baze de date utilizate în teste. Dovedim astfel că metoda noastră este competitivă atât pentru probleme de recunoaștere vocală cât și pentru probleme generale de recunoaștere a formelor.

Capitolul 6 introduce două noi metode de calcul al relevanțelor pentru componentele vectorilor de intrare cu ajutorul informației mutuale pe de o parte și a energiei informaționale pe de altă parte. Estimarea pe baza unor date de instruire a informației mutuale definită cu ajutorul entropiei Shannon ridică probleme computaționale serioase. Se pot folosi metode alternative cum ar fi cea care calculează entropia Renyi pornind de la estimarea densităților de probabilitate cu ferestre Parzen. Informația mutuală pătratică dintre două variabile aleatoare obținută fie prin similaritatea cu entropia Renyi pătratică, fie pe baza inegalității lui Pinsker devine un criteriu care poate fi maximizat iterativ într-o problemă de gradient pozitiv. Algoritmul MIRLVQ propus de noi calculează iterativ relevanțele vectorilor de intrare maximizând informația mutuală pătratică dintre o variabilă aleatoare continuă ale cărei realizări sunt calculate pe baza datelor de instruire și o variabilă aleatoare discretă ale cărei realizări sunt clasele din care fac parte vectorii de instruire. Energia informațională este o mărime pe baza căreia se calculează măsura  $o(X, Y)$ , unde  $X$  și  $Y$  sunt două variabile aleatoare, care este similară informației mutuale. Al doilea algoritm propus de noi, numit ERLVQ, calculează relevanțele componentelor din vectorii de intrare prin maximizarea măsurii  $o(X, Y)$  și obținem o formulă de actualizare similară celei din MIRLVQ. Experimentele realizate cu aceleași baze de date folosite în capitolul 5 arată că noua metodă este superioară din punct de vedere al ratei de recunoaștere algoritmului OWA-RLVQ.

Capitolul 7 care este o sinteză prezintă contribuțiile pe care teza le aduce la dezvoltarea domeniului recunoașterii formelor și posibilele direcții viitoare de dezvoltare a rezultatelor obținute în lucrare.

## Capitolul 2

# Rețele neurale. Concepte fundamentale

*Rețelele neurale constau din elemente de procesare numite neuroni care pot fi organizate în structuri regulate, pe nivele ierarhice. Adaptarea parametrilor rețelelor neurale la datele de intrare se poate face printr-un proces de instruire. Acest capitol prezintă arhitecturile neurale fundamentale pe care le-am folosit la dezvoltarea modelelor originale propuse în teză.*

### 2.1 Introducere

Rețelele neurale se pot defini, pe de o parte, ca o clasă de algoritmi matematici deoarece o rețea poate fi privită ca o reprezentare grafică pentru o largă categorie de algoritmi. Pe de altă parte, rețelele neurale imită rețelele neurale biologice proprii organismelor vii. Privite prin prisma cunoștințelor limitate pe care le avem despre rețelele neurale biologice, cea mai potrivită definiție se apropie de cea algoritmică.

Deși inspirate din biologie, există mari diferențe între rețelele neurale artificiale și cele naturale, neexistând încă modele care să concureze cu succes performanțele creierului uman. Calculul neural se realizează pe o rețea densă de noduri și conexiuni. Nodurile lucrează în mod colectiv, în paralel, și se numesc *neuroni artificiali* sau, simplu, *neuroni*. Aceștia sunt configurați în arhitecturi regulate, pot fi organizați pe nivele ierarhice și sunt permise conexiuni în cadrul aceluiași nivel sau către nivele adiacente. Puterea unei conexiuni se exprimă printr-o valoare numerică numită *pondere* care poate fi modificată. În acest capitol vom vedea care sunt principalele modele neurale și care sunt cele mai importante metode folosite la instruirea neuronilor și a rețelelor de neuroni.

## 2.2 Modelul McCulloch-Pitts

Rezultatele primelor cercetări referitoare la calculul inspirat de procesele naturale din creierul uman au fost publicate în anul 1943 de către Warren S. McCulloch și Walter Pitts [55]. Ei au pornit de la ideea că activitatea nervoasă are un caracter de tipul *totul sau nimic*, ceea ce înseamnă că procesele cerebrale ar putea fi descrise pe baza calculului logic. Aceste considerații au condus la crearea primului model neural care acceptă la intrările  $x_i, i = 1, \dots, n$  semnale 0 sau 1, iar ieșirea  $y$  se calculează pe baza regulii

$$y = \begin{cases} 1, & \text{dacă } (\sum_{i=1}^n w_i x_i) \geq T \\ 0, & \text{dacă } (\sum_{i=1}^n w_i x_i) < T \end{cases} \quad (2.1)$$

Mărimile  $w_i, i = 1, \dots, n$  se numesc *ponderi sinaptice* și au valori 1 sau -1, iar  $T$  se numește *prag al neuronului*.

Modelul general al neuronului artificial este cel din figura 2.1.

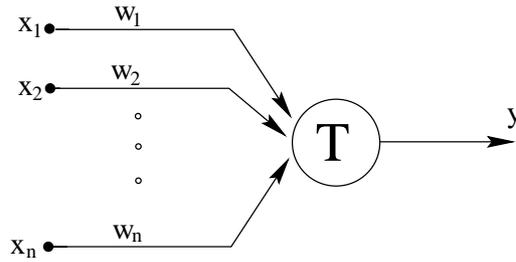
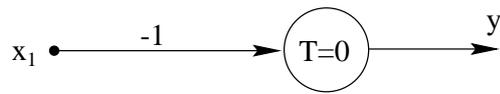


Figura 2.1: Modelul neuronului artificial.

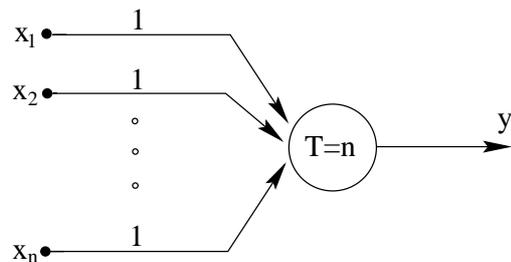
Deși este foarte simplu ca structură, potențialul de calcul al neuronului artificial este mare. Folosind scheme cu mai mulți neuroni și alegând valorile potrivite pentru ponderi și praguri, se pot implementa de la operațiile logice fundamentale până la funcțiile logice oricât de complicate. Deoarece orice funcție booleană poate fi redusă la una dintre formele canonice conjunctivă sau disjunctivă, ea poate fi implementată de operațiile NOT, AND și OR. Aceste operații elementare sunt implementate de modelele din figura 2.2 [97].

Folosind neuronul McCulloch-Pitts și proprietatea sa de întârziere a semnalului cu o unitate de timp, se pot construi circuite digitale secvențiale. În aceste condiții, o buclă dinspre ieșirea neuronului înspre intrarea sa va crea imediat modelul unei celule de memorie, ca în figura 2.3 [97].

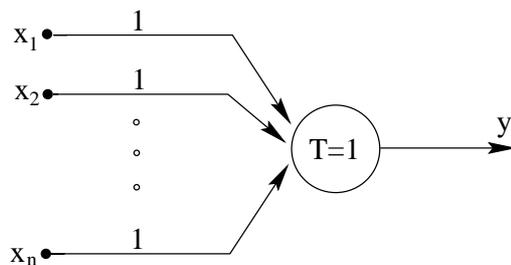
Am arătat cum se pot implementa porțile NOT, AND, OR și celula de memorie, dar scopul neuronilor artificiali nu este cel de a reconstrui elementele



(a)



(b)



(c)

Figura 2.2: Funcții logice implementate folosind neuronul McCulloch-Pitts: (a)NOT, (b)AND, (c)OR.

esențiale ale unui sistem digital. Vom vedea, însă, cum se poate exploata puterea de calcul a modelului neural și cum se poate extinde acesta pentru a construi entități ce pot fi supuse unui proces de învățare.

Deși oferă o metodă elegantă de implementare a funcțiilor logice și este caracterizat de definiția sa foarte precisă, neuronul propus de McCulloch și Pitts face o serie de simplificări care îl îndepărtează de modelul biologic. El operează într-un spațiu discret deoarece permite doar intrări și ieșiri din domeniul  $\{0, 1\}$ , dar și în timp discret când neuronii dintr-o rețea sunt sincronizați. Ponderile și pragurile sale sunt, de asemenea, fixate și sunt stabilite în faza proiectării.

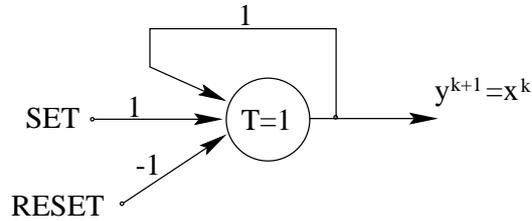


Figura 2.3: Celula de memorie implementată cu neuronul McCulloch-Pitts.

## 2.3 Perceptronul

*Perceptronul* [97] introdus de Frank Rosenblatt în 1958 a fost prima mașină capabilă să învețe. El reprezintă cea mai simplă formă de rețea neurală care poate clasifica vectori liniar separabili, adică vectori care, atunci când sunt reprezentați într-un hiperspațiu, pot fi separați cu un hiperplan. Adaptarea perceptronului constă în ajustarea ponderilor sinaptice și a pragurilor în așa fel încât, în final, să clasifice corect vectorii din mulțimea de instruire.

În principiu, operațiile pe care le realizează perceptronul sunt cele pe care le implementează neuronul McCulloch-Pitts. Așadar, un perceptron este un element care realizează o combinație liniară a intrărilor, urmat de un element de procesare care calculează valoarea ieșirii.

Operația de bază implementată de un perceptron este:

$$y = f \left( \sum_{i=1}^n w_i x_i - \theta \right).$$

În relația de mai sus,  $f$  se numește *funcție de activare a perceptronului*, iar  $\theta$  este *pragul perceptronului*. Domeniul de definiție al funcției  $f$  este mulțimea valorilor activării funcției  $f$ ,  $\sum_{i=1}^n w_i x_i$ . Această sumă se poate rescrie prin produsul scalar  $\mathbf{w}^t \mathbf{x}$ . Am notat prin

$$\mathbf{w} = [w_1 \ w_2 \ \dots \ w_n]^t$$

și

$$\mathbf{x} = [x_1 \ x_2 \ \dots \ x_n]^t$$

valorile ponderilor, respectiv ale intrărilor perceptronului.

În practică una dintre intrări este fixă, având de regulă valoarea -1, iar ponderea corespunzătoare ei este pragul perceptronului.

Dacă funcția de activare este *signum*:

$$f(\mathbf{w}^t \mathbf{x}) = \text{sgn}(\mathbf{w}^t \mathbf{x}) = \begin{cases} -1, & \text{dacă } \mathbf{w}^t \mathbf{x} < 0 \\ 1, & \text{dacă } \mathbf{w}^t \mathbf{x} \geq 0 \end{cases}, \quad (2.2)$$

atunci acest model va implementa un perceptron discret. Dacă funcția de activare este continuă:

$$f(\mathbf{w}^t \mathbf{x}) = \frac{2}{1 + e^{-\lambda \mathbf{w}^t \mathbf{x}}} - 1, \quad (2.3)$$

atunci se implementează un perceptron continuu. În relația (2.3),  $\lambda$  se numește *câștigul perceptronului*.

Pentru că funcțiile de activare descrise de relațiile (2.2) și (2.3) iau valori în intervalul  $[-1, 1]$ , adică atât valori pozitive cât și valori negative, cele două modele se vor numi *perceptron bipolar discret* și *perceptron bipolar continuu*. La implementarea unui *perceptron unipolar discret* și a unui *perceptron unipolar continuu*, care au valori de ieșire în intervalul  $[0, 1]$ , putem folosi funcțiile:

$$f(\mathbf{w}^t \mathbf{x}) = \begin{cases} 0, & \text{dacă } \mathbf{w}^t \mathbf{x} < 0 \\ 1, & \text{dacă } \mathbf{w}^t \mathbf{x} \geq 0 \end{cases} \quad (2.4)$$

și

$$f(\mathbf{w}^t \mathbf{x}) = \frac{1}{1 + e^{-\lambda \mathbf{w}^t \mathbf{x}}}. \quad (2.5)$$

Cu scopul de a descrie procesul de învățare a unui perceptron discret, să presupunem că dispunem de două mulțimi de vectori liniar separabili care aparțin claselor  $X_1$  și  $X_2$ . Prin reuniunea celor două mulțimi, alcătuim întregul set de vectori. Pentru a implementa corect această separare liniară, perceptronul va trebui să aibă un set de ponderi  $\mathbf{w}$  pentru care

$$\begin{cases} \mathbf{w}^t \mathbf{x} < 0, & \text{când } \mathbf{x} \in X_1 \\ \mathbf{w}^t \mathbf{x} \geq 0, & \text{când } \mathbf{x} \in X_2. \end{cases}$$

Prin instruirea perceptronului putem găsi, în condițiile teoremei de convergență a perceptronului discret enunțată mai jos, un  $\mathbf{w}$  care îndeplinește relația anterioară pentru orice vector  $\mathbf{x}$ . Pentru aceasta, la instruirea unui perceptron ca și clasificator liniar trebuie să aplicăm iterativ următorii pași fiecărui vector  $i = 1, \dots, n$ :

**Pasul 1.** Dacă vectorul  $\mathbf{x}_i$  este clasificat corect, atunci nu se fac corecții asupra lui  $\mathbf{w}$ .

**Pasul 2.** Dacă vectorul  $\mathbf{x}_i$  nu este clasificat corect, atunci

$$\begin{cases} \mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \eta \mathbf{x}_i, & \text{dacă } \mathbf{w}^{(t)t} \mathbf{x}_i < 0 \text{ și } \mathbf{x}_i \in X_1 \\ \mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \mathbf{x}_i, & \text{dacă } \mathbf{w}^{(t)t} \mathbf{x}_i \geq 0 \text{ și } \mathbf{x}_i \in X_2. \end{cases}$$

În aceste relații,  $t$  este pasul de instruire,  $\eta$  este un parametru ajustabil care se numește *rată de învățare*, iar  $\mathbf{w}^{(t+1)}$  este vectorul de ponderi care rezultă în urma ajustării la pasul  $t + 1$ .

Teorema de convergență a perceptronului discret cu un increment fixat arată că instruirea cu o mulțime liniar separabilă de vectori se realizează într-un număr finit de pași.

**Teorema de convergență a perceptronului discret.** *Fie submulțimile de vectori  $X_1$  și  $X_2$  liniar separabile. Instruirea unui perceptron discret cu acești vectori are loc într-un număr  $t_0$  de iterații, adică*

$$\mathbf{w}^{(t_0)} = \mathbf{w}^{(t_0+1)} = \dots$$

*este vectorul soluție.*

## 2.4 Reguli de învățare

Datele sunt stocate într-o rețea neurală printr-un *proces de învățare*. Acesta constă dintr-o recalculare a parametrilor rețelei, forțându-se un anumit răspuns pentru o intrare specifică. Un răspuns particular poate să genereze sau nu o corecție.

### Învățare supervizată și învățare nesupervizată

Procesul de învățare se bazează pe conceptul de *feedback*. Acesta conduce la introducerea unei relații între intrarea și ieșirea sistemului. Există două tipuri de învățare: cea cu supervizare și cea fără supervizare.

În învățarea supervizată presupunem că fiecare ieșire a rețelei este comparată cu o ieșire dorită  $d$  asociată vectorului introdus. Diferențele dintre răspunsul dorit și cel generat de rețea servesc drept măsură a erorii care este folosită pentru a corecta parametrii rețelei.

În învățarea fără supervizare, răspunsul dorit nu se cunoaște și din acest motiv nu se poate stabili o măsură a erorii în funcție de care să se corecteze comportarea rețelei. De aceea, mecanismul de învățare trebuie să prevadă o modalitate de ajustare a parametrilor pe baza similitudinilor între vectori. Rețeaua trebuie să descopere singură proprietăți, regularități ale vectorilor care să conducă la separare. Această modalitate de adaptare se numește *autoorganizare*.

### Regula generală de învățare

Un neuron este un element adaptiv care își modifică ponderile în funcție de intrare, de ieșire și de răspunsul dorit pentru fiecare vector pe baza regulii de instruire utilizată. Atunci când nu există răspuns dorit, ponderile se modifică în funcție doar de intrare și de ieșire, deci nesupervizat.

În studiul rețelelor neurale a fost adoptată o regulă generală de instruire [97] care, prin diferite particularizări, produce o serie de noi reguli. Pentru un neuron ca cel din figura 2.1, vectorul ponderilor  $\mathbf{w} = [w_1 w_2 \dots w_n]^t$  se modifică proporțional cu produsul dintre intrările  $\mathbf{x}$  și *semnalul de învățare*  $s$ . Semnalul de învățare se calculează ca o funcție de vectorul  $\mathbf{w}$ , de intrarea  $\mathbf{x}$  și, atunci când este cazul, de ieșirea dorită  $d$ . Odată făcute aceste notații, regula de actualizare a ponderilor neuronului se scrie astfel:

$$\mathbf{w}(t+1) = \eta s(\mathbf{w}(t), \mathbf{x}(t), d(t)) \mathbf{x}(t) \quad (2.6)$$

pentru cazul continuu la momentul  $t+1$  și sub forma:

$$\mathbf{w}^{(t+1)} = \eta s(\mathbf{w}^{(t)}, \mathbf{x}^{(t)}, d^{(t)}) \mathbf{x}^{(t)} \quad (2.7)$$

în cazul discret la pasul  $t+1$ .

Principalele reguli de instruire a neuronilor se obțin prin particularizarea regulii generale.

## Regula lui Hebb

Această regulă este nesupervizată și implementează *postulatul învățării* formulat de Hebb în cartea sa apărută în 1949 [31] care propune ideea învățării asociative la nivel celular:

*Dacă celula A excită în mod repetat celula B făcând-o să genereze un impuls, atunci are loc un proces de creștere (schimbare metabolică) într-una sau ambele celule, astfel încât eficiența excitării lui B de către A crește.*

Cu alte cuvinte, vectorii de intrare cu frecvența cea mai mare vor avea influența cea mai puternică asupra ponderilor conexiunilor.

Această regulă particularizează semnalul de învățare astfel:

$$s = f(\mathbf{w}^t \mathbf{x}),$$

adică semnalul de învățare este chiar ieșirea neuronului. Atunci:

$$\Delta \mathbf{w} = \eta f(\mathbf{w}^t \mathbf{x}) \mathbf{x},$$

unde  $\eta$  este o constantă de învățare.

## Regula perceptronului

Așa cum a fost definită de Rosenblatt, regula de instruire a perceptronului consideră că semnalul de învățare este diferența dintre răspunsul dorit și cel actual al neuronului:

$$s = d - y. \quad (2.8)$$

Din (2.6) și (2.8) obținem formula de ajustare a ponderilor perceptronului:

$$\Delta \mathbf{w} = \eta [d - \text{sgn}(\mathbf{w}^t \mathbf{x})] \mathbf{x} \quad (2.9)$$

unde ieșirea se calculează pe baza funcției bipolare discrete

$$y = \text{sgn}(\mathbf{w}^t \mathbf{x}).$$

Se observă că această regulă se aplică doar pentru neuroni cu răspuns binar, iar ponderile sunt ajustate doar dacă ieșirea actuală  $y$  este incorectă. Această observație ne permite să rescriem relația (2.9) ținând cont și de faptul că ieșirile dorite pot lua doar valorile 1 și -1 astfel:

$$\Delta \mathbf{w} = \begin{cases} -2\eta \mathbf{x}, & \text{când } d = -1 \text{ și } \text{sgn}(\mathbf{w}^t \mathbf{x}) = 1 \\ 2\eta \mathbf{x}, & \text{când } d = 1 \text{ și } \text{sgn}(\mathbf{w}^t \mathbf{x}) = -1 \end{cases}.$$

## Regula delta

Regula delta este aplicabilă doar pentru funcții continue de activare și învățare supervizată. Poate fi dedusă căutând varianta de scădere cea mai rapidă a erorii pătratice medii dintre ieșirea actuală  $y$  a neuronului și ieșirea dorită  $d$ . Eroarea pătratică medie

$$E = \frac{1}{2}(d - y)^2$$

este echivalentă cu:

$$E = \frac{1}{2} [d - f(\mathbf{w}^t \mathbf{x})]^2,$$

iar gradientul ei va fi:

$$\nabla E = -(d - y) f'(\mathbf{w}^t \mathbf{x}) \mathbf{x}. \quad (2.10)$$

Pentru ca algoritmul să aibă eficiența maximă, ponderile vor fi modificate în sensul gradientului negativ al erorii astfel:

$$\Delta \mathbf{w} = -\eta \nabla E, \quad (2.11)$$

unde  $\eta$  este o constantă pozitivă. Din relațiile (2.10) și (2.11) obținem formula generală de modificare a ponderilor pentru regula de instruire delta:

$$\Delta \mathbf{w} = \eta (d - y) f'(\mathbf{w}^t \mathbf{x}) \mathbf{x}. \quad (2.12)$$

Comparând regula generală de instruire în cazul continuu dată de relația (2.6) și regula delta exprimată prin relația (2.12), observăm că în acest caz semnalul de învățare este:

$$s = [d - f(\mathbf{w}^t \mathbf{x})] f'(\mathbf{w}^t \mathbf{x}).$$

Regula delta a fost introdusă de McClelland și Rumelhart în 1986 și se mai numește și regula de instruire a perceptronului continuu, putând fi extinsă și pentru instruirea rețelelor neurale multistrat.

### Regula *câștigătorul ia tot*

Această regulă se folosește pentru instruirea nesupervizată a unor rețele de neuroni, fiind un exemplu de învățare competitivă. Spre deosebire de regulile prezentate anterior, regula *câștigătorul ia tot* se aplică asupra unei colecții de  $p$  neuroni organizați pe un singur strat.

Învățarea se bazează pe faptul că unul dintre neuroni, de exemplu neuronul  $m$ , are răspuns maxim pentru vectorul  $\mathbf{x}$ , fiind declarat *câștigător*. Câștigătorul este selectat pe baza următorului criteriu:

$$\mathbf{w}_m^t \mathbf{x} = \max_{i=1,2,\dots,p} (\mathbf{w}_i^t \mathbf{x}).$$

Drept urmare, doar ponderile

$$\mathbf{w}_m = [w_{m1} \ w_{m2} \ \dots \ w_{mn}]^t$$

care converg către el vor fi actualizate în următorul pas al învățării. Formula de ajustare a ponderilor este:

$$\Delta \mathbf{w}_m = \alpha (\mathbf{x} - \mathbf{w}_m),$$

unde  $\alpha$  este o constantă pozitivă cu valori mici.

De regulă, actualizarea ponderilor se face pentru neuronii dintr-o vecinătate a neuronului câștigător. Ponderile se inițializează cu valori aleatoare și sunt normalizate pe parcursul învățării.

## 2.5 Modele de rețele neurale artificiale

O rețea neurală poate fi definită ca o interconectare de neuroni cu proprietatea că ieșirile unor neuroni sunt conectate prin intermediul conexiunilor sinaptice cu intrările altor neuroni sau ale lor însele.

Pornind de la această definiție, se pot imagina diverse arhitecturi de rețele neurale. Cele mai populare sunt rețelele *feedforward* și rețelele *feedback* [30], [97].

## Rețele neurale feedforward

O arhitectură elementară feedforward constă dintr-o colecție de neuroni care permit transmiterea semnalelor într-un singur sens, dinspre intrări înspre ieșiri.

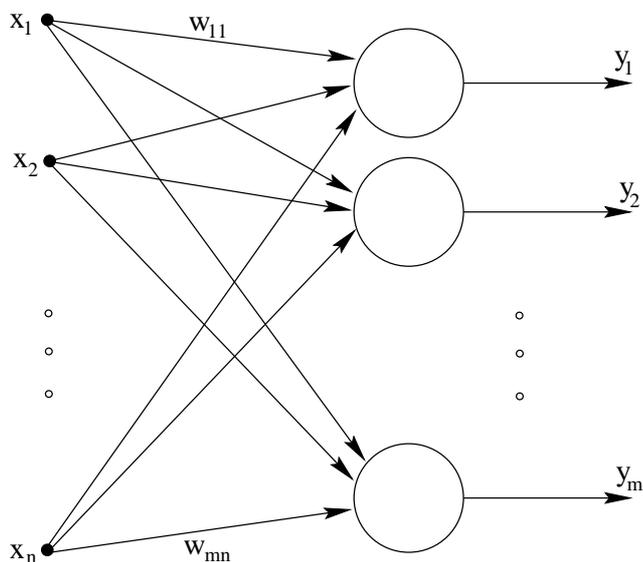


Figura 2.4: Rețea neurală feedforward cu un singur strat de neuroni.

În figura 2.4 am reprezentat o rețea neurală feedforward cu un singur strat de neuroni. Notăm cu  $\mathbf{x}$  și  $\mathbf{y}$  vectorii intrărilor și ieșirilor:

$$\mathbf{x} = [x_1 \ x_2 \ \dots \ x_n]^t$$

$$\mathbf{y} = [y_1 \ y_2 \ \dots \ y_m]^t.$$

Dacă notăm cu  $w_{ji}$  ponderea conexiunii care leagă neuronul  $j$  cu intrarea  $i$ , atunci matricea  $\mathbf{W}$  a ponderilor se va scrie astfel:

$$\mathbf{W} = \begin{bmatrix} \mathbf{w}_1 \\ \mathbf{w}_2 \\ \vdots \\ \mathbf{w}_m \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1n} \\ w_{21} & w_{22} & \dots & w_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{m1} & w_{m2} & \dots & w_{mn} \end{bmatrix}.$$

Folosind aceste notații, ieșirile neuronilor se vor calcula conform relației:

$$y_i = f(\mathbf{w}_i^t \mathbf{x}), \quad i = 1, \dots, m.$$

Rețelele neurale feedforward pot fi extinse grupând neuronii pe straturi. Arhitecturile multistrat se obțin prin cascada mai multor straturi de neuroni.

Stratul final se numește *strat de ieșire*, iar celelalte straturi de neuroni se numesc *straturi ascunse*. Intrările  $\mathbf{x}$  și ieșirile  $\mathbf{y}$  se numesc *vectori de intrare* și *vectori de ieșire*. În timpul instruirii, ieșirile rețelei sunt comparate cu un set de ieșiri dorite și, în funcție de diferențele care apar, sunt modificate valorile ponderilor  $\mathbf{W}$ .

## Rețele neurale feedback

Rețelele neurale feedback se obțin din rețele neurale feedforward prin conectarea ieșirilor rețelei la intrările sale, ca în figura 2.5. Prin buclele care se închid peste rețea, fiecare ieșire va fi controlată de valorile celorlalte ieșiri. Ieșirea de la momentul  $t$  controlează ieșirea de la momentul  $t + \Delta$ :

$$y_j(t + \Delta) = \sum_{i=1}^n w_{ji} y_i(t), \quad j = 1, \dots, n \quad (2.13)$$

unde diferența de timp  $\Delta$  este introdusă de elementele de întârziere plasate pe buclele de feedback. Intrările  $x(t)$  sunt folosite pentru inițializarea rețelei, după care dispar și sistemul rămâne autonom. Din acest moment, el se autoîntreține datorită legăturilor dintre ieșiri și intrări.

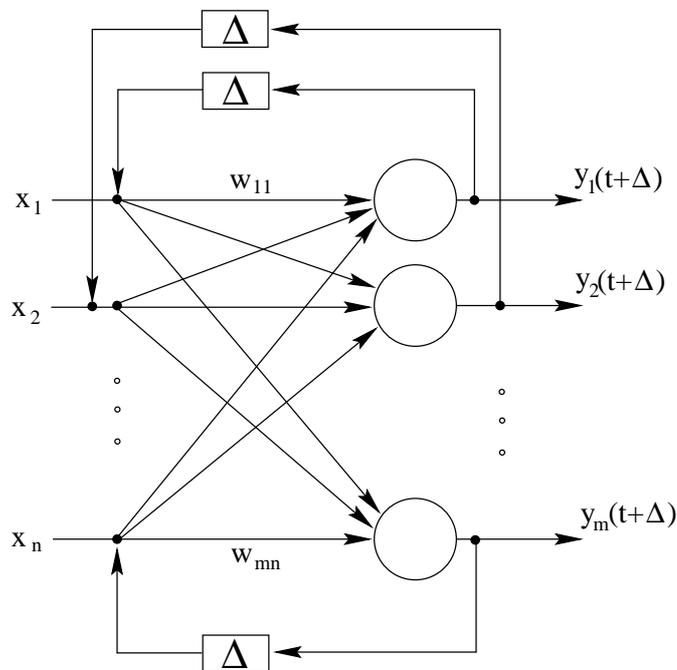


Figura 2.5: Rețea neurală feedback.

În situația în care ieșirile sunt înregistrate la intervale determinate de timp, rețeaua se numește cu *timp discret*. Din acest motiv, putem rescrie ecuația (2.13) asimilând întârzierea  $\Delta$  cu unitatea și momentele de timp cu întregi pozitivi astfel:

$$y_j^{(t+1)} = \sum_{i=1}^n w_{ji} y_i^{(t)}, \quad j = 1, \dots, n. \quad (2.14)$$

Rețeaua neurală din figura 2.5 se numește *recurentă* pentru că răspunsul său la momentul de timp  $t + 1$  depinde de întregul istoric, începând cu momentul  $t = 0$ . Relația (2.14) arată că pentru a calcula  $y_j^{(t+1)}$  trebuie să calculăm mai întâi  $y_j^{(t)}$ . Pentru aceasta, este nevoie de  $y_j^{(t-1)}$  etc.

Rețelele recurente lucrează de obicei cu o reprezentare discretă a datelor, folosind neuroni cu funcție de activare discretă. Un sistem cu intrări discrete și o reprezentare discretă a datelor se numește *automat*. Rețelele neurale recurente din această clasă pot fi considerate niște automate. Ecuația (2.14) calculează starea rețelei la momentele  $t = 1, 2, \dots$ , adică descrie secvența de *tranziții ale stărilor*. Rețeaua își începe tranziția la momentul 0 cu  $\mathbf{x}^0$  și evoluează până la o eventuală stare de echilibru care se numește *attractor*. Un attractor constă dintr-o singură stare sau dintr-un număr limitat de stări.

## 2.6 Calcul neural

Așa cum am văzut în paragrafele anterioare, o rețea neurală calculează ieșirea  $\mathbf{y}$  corespunzătoare intrării  $\mathbf{x}$  după ce în prealabil au fost stabilite structura rețelei și valorile ponderilor ei. Acest proces constă de fapt în *regăsirea* informației care a fost înmagazinată anterior.

Să presupunem că o rețea a stocat o mulțime de vectori. Ulterior, dacă rețelei  $i$  se prezintă un vector din mulțimea de instruire aceasta îl va asocia cu cel memorat. Un vector incomplet dar asemănător unuia folosit la învățare va reprezenta un bun punct de pornire pentru regăsire. Acest proces se numește *autoasociere*.

O rețea neurală poate să stocheze, însă, asocieri între perechi de vectori astfel încât la regăsire, la prezentarea unui vector sau a unei versiuni ușor distorsionate a lui, rețeaua va genera perechea sa. Acest mecanism se numește *heteroasociere*.

Calculul neural poate fi privit și prin prisma diferenței dintre *clasificare* și *recunoaștere*.

Clasificarea este un caz special de heteroasociere. Să presupunem că avem un număr de vectori împărțiți în mai multe clase. La aplicarea unui vector dintre cei folosiți la instruire, clasificatorul va genera o informație care arată

clasa din care face parte vectorul. Heteroasocierea se face în această situație între vector și codul asociat clasei.

Dacă, însă, vectorul folosit în procesul de regăsire nu se potrivește exact cu nici unul din mulțimea de instruire, atunci avem de a face cu o recunoaștere. Pattern-ul va fi introdus în clasa pentru care similaritatea este maximă.

Una dintre caracteristicile cele mai importante ale rețelelor neurale este capacitatea de generalizare. O rețea face o bună generalizare atunci când reușește să recunoască în mod corect vectori noi.

## 2.7 Perceptronul multistrat

Acesta constă dintr-o mulțime de noduri care alcătuiesc stratul de intrare, unul sau mai multe straturi ascunse și un strat de ieșire, după cum se poate vedea în figura 2.7 [97]. Semnalul de intrare se propagă prin rețea dinspre intrări înspre ieșiri, parcurgând fiecare strat în parte. Perceptronul multistrat așa cum îl descriem aici reprezintă o generalizare a perceptronului introdus de Rosenblatt.

Perceptronul multistrat a fost folosit cu succes în mai multe aplicații cu diverse grade de dificultate prin instruire supervizată folosind algoritmul *back-propagation of error* (*propagarea în urmă a erorii*). Acest algoritm ajustează ponderile unei rețele neurale multistrat pe baza unei metode de minimizare a erorii dintre răspunsul actual al rețelei și răspunsul dorit.

Un perceptron multistrat are trei caracteristici importante.

- a) Modelul fiecărui neuron din rețea include o funcție neliniară care calculează ieșirea acestuia. Pentru algoritmul de propagare în urmă a erorii este important ca această funcție să fie diferentiabilă în orice punct. Câteva exemple de funcții care satisfac această cerință sunt cele din figura 2.6 a), d) și e).
- b) O rețea neurală conține unul sau mai multe straturi ascunse. Acestea îmbunătățesc performanțele în învățare ale rețelei pentru că măresc capacitatea de stocare a caracteristicilor vectorilor de instruire. În același timp, straturile ascunse sporesc capacitatea de generalizare a rețelelor.
- c) Ponderile sinaptice conferă rețelei un grad ridicat de conectivitate, de aceea are un potențial de procesare paralelă distribuită suficient de mare pentru a rezolva probleme dificile.

Caracteristicile de mai sus împreună cu posibilitatea de a învăța din experiență dau forță acestui model. Primele rezultate obținute în studierea rețe-

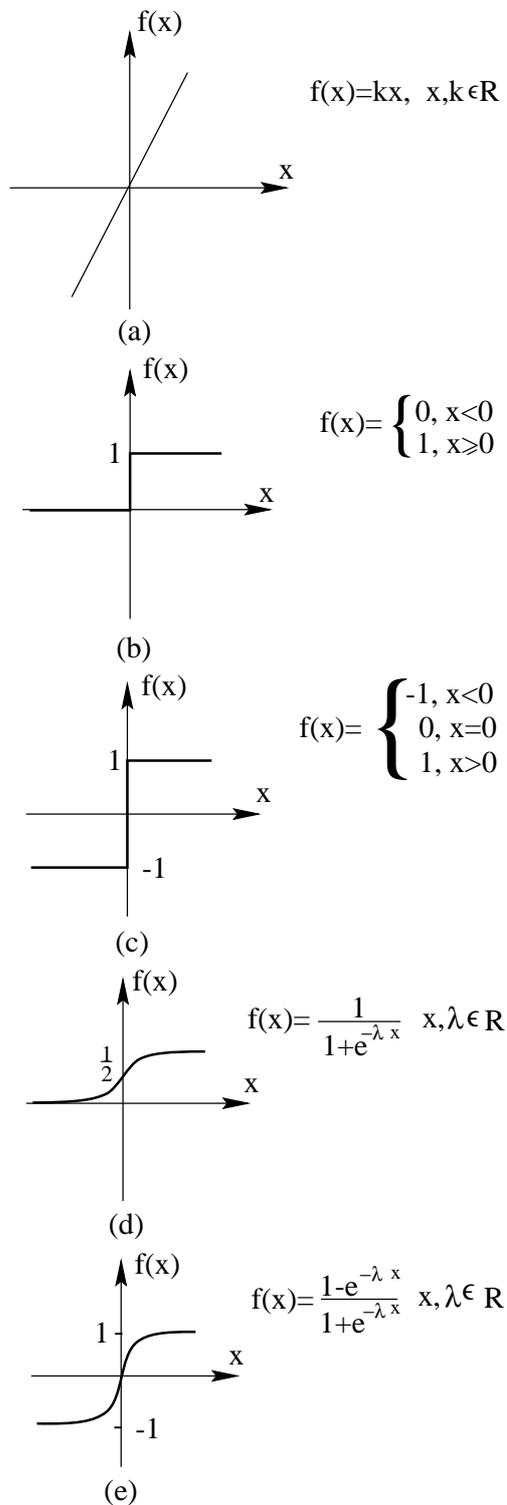


Figura 2.6: Funcții de transfer: (a)funcția liniară; (b)funcția treaptă; (c)funcția signum; (d)sigmoida unipolară; (e)sigmoida bipolară.

lelor neurale feedforward multistrat au apărut în 1958, atunci când Rosenblatt a prezentat pentru prima dată perceptronul. Evoluția acestui domeniu s-a dovedit a fi, însă, greoaie datorită lipsei unui algoritm performant de instruire a rețelelor de neuroni artificiali.

Publicarea în anul 1986 a lucrării *Parallel Distributed Processing* [75] de către Rumelhardt și McClelland a reprezentat o relansare a cercetărilor în domeniul conexiunilor. Algoritmii *backpropagation* oferă o metodă eficientă de instruire a rețelelor feedforward multistrat, cu toate că nu se poate afirma că ea conduce la o soluție pentru orice problemă solvabilă.

Pentru a prezenta algoritmul de instruire a perceptronului multistrat prin propagarea în urmă a erorii vom folosi modelul din figura 2.7.

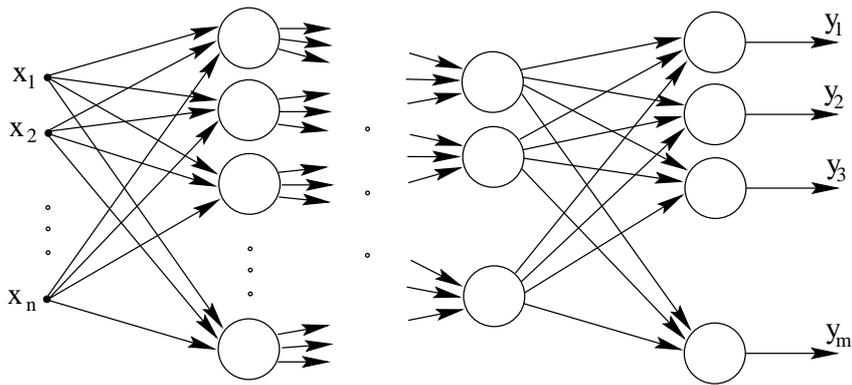


Figura 2.7: Perceptronul multistrat.

În forma generală, rețeaua este complet conectată, ceea ce înseamnă că fiecare neuron din rețea are câte o legătură cu fiecare neuron de pe stratul anterior celui pe care se află el.

În timpul instruirii, rețeaua neurală este parcursă de două tipuri de semnale [30]. *Semnalul funcției implementate de rețea* este generat de stimulii aplicați la intrare și se propagă prin fiecare strat în parte, apărând la ieșire sub forma semnalului de ieșire. *Semnalul de eroare* este folosit în algoritmul de instruire și este calculat pe baza diferențelor dintre ieșirile actuale ale rețelei și ieșirile dorite. Nivelul acestui semnal se calculează dinspre ieșiri înspre intrări, de aceea convenim să spunem că semnalul de eroare se propagă în sens invers semnalului funcției implementate de rețea.

### 2.7.1 Algoritmul *backpropagation of error*

Vom prezenta detaliat algoritmul *backpropagation of error* [30] datorită importanței sale deosebite.

Considerăm că avem o rețea neurală cu un strat de intrare, un strat ascuns și unul de ieșire, fără a afecta gradul de generalitate a algoritmului. Vom vedea că relațiile de calcul deduse în acest fel se aplică în mod similar pentru rețele cu oricâte straturi ascunse.

Dorim, așadar, să aflăm ponderile unei rețele neurale cu un număr fixat de straturi și de neuroni pe fiecare strat astfel încât aceasta să implementeze funcția dorită de noi care aplică vectorii de intrare în vectorii de ieșire. Datorită caracterului neliniar al elementelor de procesare pe care le-am numit neuroni, vom deduce o metodă iterativă de calcul a valorii ponderilor. Metoda se bazează pe minimizarea unei funcții globale de eroare  $E$ .

Pentru fiecare neuron calculăm valoarea activării după care, pe baza acesteia, calculăm ieșirea:

$$a_j^p = \sum_{i \in I} w_{ji} y_i^p \quad (2.15)$$

$$y_j^p = f_{act}(a_j^p), \quad j \in J. \quad (2.16)$$

În aceste relații,  $p$  este indicele vectorului curent din mulțimea de instruire pentru care se calculează activarea și ieșirea. Am notat cu  $I$  mulțimea intrărilor neuronului, iar cu  $J$  mulțimea ieșirilor tuturor neuronilor din rețea. După cum se vede, ieșirea unui neuron poate deveni intrare pentru alți neuroni. Pentru neuronii din primul strat ascuns, intrările  $y_i^p$  sunt chiar intrările rețelei. Matricea

$$\mathbf{W} = \{w_{ji}\}_{j \in O, i \in I}$$

păstrează valorile ponderilor sinaptice pentru conexiunile dintre neuroni, iar  $a_j^p$  reprezintă activarea neuronului  $j$  atunci când neuronul primește la intrare vectorul  $p$ . Mulțimea  $O$  este o submulțime a lui  $J$  și conține ieșirile neuronilor de pe stratul de ieșire. Funcția de activare  $f_{act}$  trebuie să fie derivabilă deoarece algoritmul de propagare în urmă a erorii calculează ponderile sinaptice ținând cont și de gradientul lui  $f_{act}$ . Cele mai folosite funcții de activare sunt sigmoida unipolară și biploară (fig. 2.6). De regulă, fiecare neuron are una dintre intrări fixată la valoarea -1, motiv pentru care ponderea corespunzătoare acesteia poate fi interpretată ca prag al neuronului.

Măsura adaptării rețelei neurale la setul de vectori de instruire este dată de funcția de eroare. Algoritmul *backpropagation* utilizează funcția de eroare pătratică medie care are expresia:

$$E^p = \frac{1}{2} \sum_{j \in O} (d_j^p - y_j^p)^2. \quad (2.17)$$

$E^p$  este eroarea individuală pentru un vector,  $d_j^p$  este ieșirea dorită pentru neuronul  $j$  și vectorul  $p$ , iar  $y_j^p$  este ieșirea neuronului  $j$  atunci când a fost folosit vectorul  $p$ .

Instruirea rețelei constă din modificarea fiecărei ponderi proporțional cu influența sa asupra erorii în sensul reducerii lui  $E$ :

$$\Delta w_{ji}^p = -\eta \frac{\partial E^p}{\partial w_{ji}}. \quad (2.18)$$

În relația de mai sus,  $\eta$  este o constantă numită *rata de instruire*. Semnul minus arată că modificarea ponderii se face în sensul gradientului negativ al erorii.

Relația (2.18) poate fi rescrisă ținând cont de (2.16) astfel:

$$\Delta w_{ji}^p = \eta \frac{\partial E^p}{\partial y_j^p} \frac{\partial y_j^p}{\partial a_j^p} \frac{\partial a_j^p}{\partial w_{ji}}. \quad (2.19)$$

Observăm că:

$$f'_{act}(a_j^p) = \frac{\partial y_j^p}{\partial a_j^p}$$

și

$$\frac{\partial a_j^p}{\partial w_{ji}} = y_i^p,$$

astfel că relația (2.19) devine:

$$\Delta w_{ji}^p = \eta \frac{\partial E^p}{\partial y_j^p} f'_{act}(a_j^p) y_i^p. \quad (2.20)$$

În acest moment, se tratează separat cazul în care  $y_j^p$  aparține unui neuron de pe stratul de ieșire și cel în care aparține unui neuron din straturile ascunse.

Primul caz, cel în care  $j$  este un neuron de ieșire, este mai ușor de rezolvat. Dacă înlocuim (2.17) în (2.20) putem scrie:

$$\Delta w_{ji}^p = \eta (d^p - y^p) f'_{act}(a_j^p) y_i^p, \quad j \in O. \quad (2.21)$$

A doua situație este cea în care neuronul  $j \notin O$ , adică se găsește pe un strat ascuns și nu se poate defini un răspuns dorit pentru el:

$$\Delta w_{ji}^p = \eta \left( \sum_{k \in K} \frac{\partial E^p}{\partial y_k^p} \frac{\partial y_k^p}{\partial a_k^p} \frac{\partial a_k^p}{\partial y_j^p} \right) f'_{act}(a_j^p) y_i^p.$$

Relația se poate rescrie astfel:

$$\Delta w_{ji}^p = \eta \left( \sum_{k \in K} \frac{\partial E^p}{\partial y_k^p} f'_{act}(a_k^p) w_{kj} \right) f'_{act}(a_j^p) y_i^p. \quad (2.22)$$

Recursia care apare în această ecuație arată că  $\frac{\partial E^p}{\partial y_m^p}$  pentru un strat  $m$  se calculează în funcție de valorile  $\frac{\partial E^p}{\partial y_n^p}$  obținute pe stratul  $n$  care îl succede pe  $m$ . Aceasta înseamnă că ajustările ponderilor se vor calcula din relațiile de mai sus, parcurgând rețeaua de la ieșiri înspre intrări. În mod obișnuit, se face notația:

$$\delta_j^p = \frac{\partial E^p}{\partial a_j^p} = \begin{cases} (d_j^p - y_j^p) f'_{act}(a_j^p), & j \in O \\ \sum_{k \in K} (\delta_k^p w_{kj}) f'_{act}(a_j^p), & j \notin O \end{cases},$$

iar  $\delta_j^k$  se numește *gradient local*.

Cu această notație, relațiile (2.21) și (2.22) se pot scrie într-o formă compactă astfel:

$$\Delta w_{ji}^p = \eta \delta_j^p y_i^p \quad (2.23)$$

care este regula de calcul ce stă la baza instruirii rețelelor neurale prin propagarea în urmă a erorii.

## 2.7.2 Parametrii învățării

Convergența algoritmului de instruire prin propagarea în urmă a erorii este asigurată de o serie de factori. Vom prezenta modul în care învățarea este influențată de ponderile inițiale, de constanta de instruire, de factorul momentum și de modul de ajustare a ponderilor [97].

### Ponderile inițiale

Ponderile rețelei care urmează a fi instruită sunt inițializate cu valori aleatoare și mici. Acestea afectează puternic soluția finală.

Inițializarea cu valori egale a ponderilor rețelei conduce la o comportare nesatisfăcătoare, cu excepția cazului când acestea sunt perturbate în mod aleator în timpul instruirii.

Instruirea unei rețele neurale poate ajunge într-un punct în care eroarea să se stabilizeze la o valoare insuficient de mică sau chiar să înceapă să crească. Experimentele [30] arată că în astfel de situații continuarea instruirii este inutilă. Valorile ponderilor pot începe să crească, în timp ce calitatea rețelei scade. De aceea, se preferă reluarea învățării prin reinițializarea ponderilor.

## Constanta de instruire

Convergența algoritmului *backpropagation* depinde în mod semnificativ de valoarea constantei de instruire  $\eta$ . De obicei, valoarea optimă a acestei constante depinde de problema care trebuie rezolvată.

Deși instruirea rețelelor neurale prin algoritmul *backpropagation* este o metodă eficientă, suprafața de eroare are, adeseori, proprietăți care încetinesc convergența. În astfel de situații, o valoare mai mare a lui  $\eta$  accelerează instruirea. În cazurile în care ajustarea ponderilor se face cu valori prea mari, existând pericolul depășirii soluției și a unui comportament instabil al rețelei, se preferă micșorarea valorii constantei de învățare. Aceste considerații conduc la concluzia că  $\eta$  trebuie ales experimental pentru fiecare problemă și că rata de instruire poate fi modificată dinamic în timpul procesului. Cu toate acestea, experimentele [30] au arătat că o valoare cuprinsă între  $10^{-3}$  și 10 este potrivită pentru majoritatea cazurilor.

## Factorul momentum

O metodă de accelerare a convergenței algoritmului *backpropagation* și de prevenire a comportamentului instabil este modificarea regulii delta prin introducerea *factorului momentum* conform formulei:

$$\Delta w_{ji}^{(t)} = \alpha \Delta w_{ji}^{(t-1)} + \eta \delta_j^{(t)} y_i^{(t)}, \quad (2.24)$$

unde  $t$  și  $t - 1$  indică pasul curent de instruire și cel anterior, iar  $\alpha$  este factorul momentum care este o constantă pozitivă. Primul termen al relației (2.24) se numește *termenul momentum* și scalează cea mai recentă ajustare a ponderilor. După o serie de  $t$  pași de instruire folosind termenul momentum, modificarea actuală se mai poate calcula și conform formulei:

$$\Delta w_{ji}^{(t)} = -\eta \sum_{k=0}^t \alpha^{t-k} \frac{\partial E^{(k)}}{\partial w_{ji}^{(k)}}, \quad (2.25)$$

unde am ținut cont de (2.18) și (2.23).

Ajustarea  $\Delta w_{ji}^{(t)}$  reprezintă suma unei serii exponențiale. Pentru ca aceasta să fie convergentă, trebuie impusă condiția  $0 \leq |\alpha| < 1$ .

Atunci când  $\frac{\partial E^{(k)}}{\partial w_{ji}^{(k)}}$  are același semn de-a lungul mai multor iterații consecutive,  $\Delta w_{ji}^{(k)}$  crește iar ponderea  $w_{ji}$  este ajustată cu o valoare mai mare. În acest fel, factorul momentum tinde să accelereze convergența algoritmului. Atunci când  $\frac{\partial E^{(k)}}{\partial w_{ji}^{(k)}}$  are semne opuse în iterații consecutive,  $\Delta w_{ji}^{(k)}$  scade iar ponderea  $w_{ji}$  este ajustată cu o valoare mai mică. În acest caz, factorul momentum are un efect de stabilizare a oscilațiilor.

## Ajustare cumulativă și ajustare incrementală a ponderilor

Algoritmul de instruire prin propagarea în urmă a erorii poate realiza o mică ajustare a ponderilor după prezentarea fiecărui vector. Această tehnică se numește *actualizare incrementală*. Algoritmul backpropagation poate fi conceput, însă, să ajusteze ponderile după o prezentare a întregului set de vectori, caz în care avem de a face cu o *actualizare cumulativă*.

Avantajul folosirii primei metode este acela că algoritmul implementează, într-adevăr, o metodă a gradientului descendent al erorii. În plus, o astfel de modalitate de instruire nu necesită memorarea componentelor calculate pentru fiecare vector care, apoi, să participe la calcularea ajustării globale.

Oricare ar fi modalitatea de instruire folosită, un element foarte important este acela că ordinea în care sunt aleși vectorii este bine să fie aleatoare.

## 2.8 Rețele neurale cu întârziere în timp

Rețelele neurale cu întârziere în timp (*time delay neural networks*) [90] sunt o extensie a perceptronului multistrat și au fost create pentru rezolvarea problemelor de recunoaștere a vorbirii.

Semnalul asociat vorbirii se caracterizează prin modificări rapide în timp, de aceea trebuie prelucrat cu ajutorul unei rețele neurale adaptate acestui comportament. O astfel de rețea neurală trebuie să dispună de o memorie [18] pentru a fi capabilă să surprindă evoluția în timp a semnalului. Tehnica adoptată pentru a rezolva această problemă constă în introducerea unor elemente de întârziere în structura legăturilor sinaptice ale rețelei.

### Arhitectura rețelelor neurale cu întârziere în timp

Neuronul, elementul de bază folosit de o rețea neurală, calculează suma ponderată a intrărilor pe care o transmite apoi unui bloc ce implementează o funcție neliniară care furnizează răspunsul. Neuronul este modificat într-o rețea neurală cu întârziere în timp prin introducerea celulelor de întârziere  $\Delta_1, \Delta_2 \dots \Delta_N$  (figura 2.8). În felul acesta, anumite componente ale intrării vor influența în mod repetat răspunsul neuronului. Se poate spune, astfel, că rețeaua neurală cu întârziere în timp memorează cele mai recente evenimente, fiind adaptată evoluției în succesiune temporală a intrărilor, așa cum este necesar în recunoașterea vorbirii.

Un vector folosit la instruirea unei rețele neurale pentru recunoașterea vorbirii are, de regulă, un număr foarte mare de componente. Modelul *time delay neural network* permite ca și atunci când se folosesc vectori de dimensiuni mari

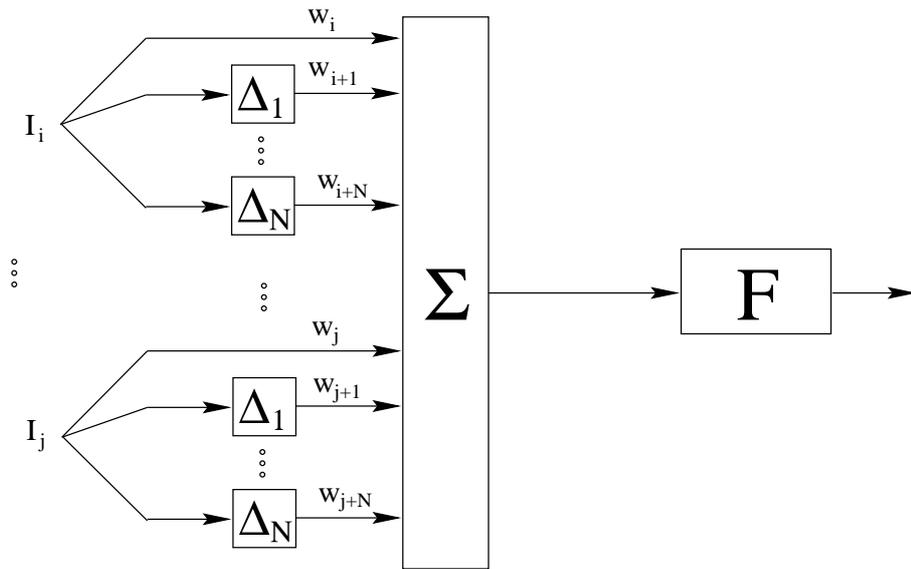


Figura 2.8: Modelul neuronului folosit de rețelele neurale cu întârziere în timp.

numărul de neuroni și de conexiuni sinaptice să rămână, totuși, mic. La un moment dat, doar o parte a componentelor vectorului influențează ieșirea rețelei. Acesta este parcurs de o fereastră care selectează intrarea curentă, așa cum ilustrează figura 2.9.

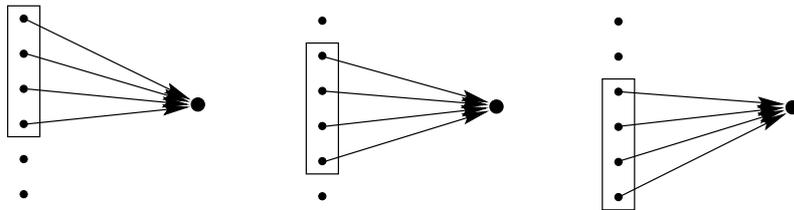


Figura 2.9: Secvența în care componentele unui vector devin intrări ale neuronului într-o rețea neurală cu întârziere în timp.

Maniera descrisă mai sus de parcurgere a unui vector poate fi extinsă pentru a construi o rețea neurală cu întârziere în timp ca cea din figura 2.10

Numai componentele vectorului de intrare care se găsesc în fereastră sunt transmise prin intermediul conexiunilor sinaptice spre stratul următor. Fiecărei ferestre dintr-un strat sursă îi corespunde o fereastră în stratul următor. Dimensiunea ferestrei care parcurge un strat este constantă, iar numărul de ferestre din stratul destinație este determinat de dimensiunea vectorului care se aplică

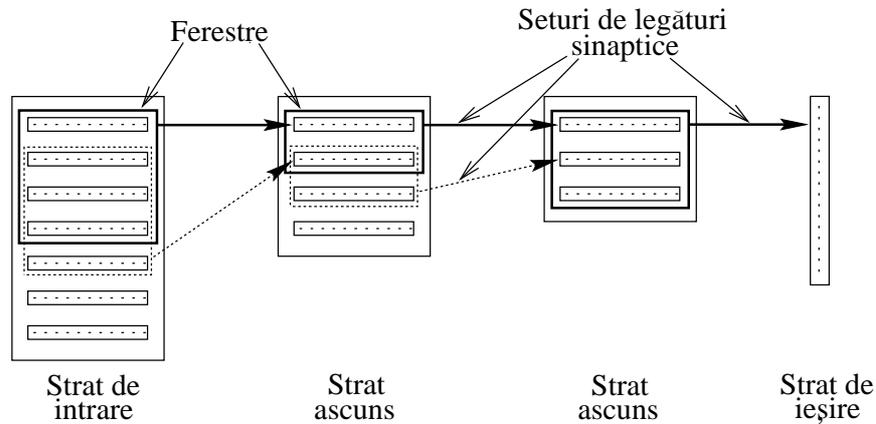


Figura 2.10: Arhitectura unei rețele neurale cu întârziere în timp.

stratului sursă și de dimensiunea ferestrei care baleiază vectorul.

În exemplul nostru, vectorul este structurat sub forma unei matrici. Fereastra care îl baleiază selectează o secvență a lui, aceasta devenind intrarea actuală. Pe baza intrării selectate se calculează ieșirile neuronilor din stratul următor. Fiecare deplasare a ferestrei din stratul de intrare înseamnă adăugarea noilor componente la vectorul care se formează în stratul de ieșire. Fereastra baleiază vectorul sursă până la epuizarea acestuia, deci până când se calculează toate componentele ieșirii din stratul destinație.

### Instruirea unei rețele neurale cu întârziere în timp

Pentru acest model se folosește același procedeu de învățare ca și pentru perceptronul multistrat. Prin instruirea care folosește propagarea în urmă a erorii se realizează două treceri prin rețea. Mai întâi, se aplică rețelei neurale un vector de intrare care constituie informația propagată înainte. Ponderile sunt ajustate în timpul propagării înapoi a semnalului de eroare care este calculat în funcție de ieșirile rețelei și de ieșirile dorite. Această procedură se repetă până când rețeaua converge și produce ieșirile dorite.

Cele două etape ale algoritmului backpropagation, propagarea înainte și propagarea înapoi, se desfășoară ca evenimente separate. Modificarea ponderilor se face prin medierea valorilor calculate pentru fiecare fereastră în parte. Așadar, un pas de instruire se consideră încheiat doar atunci când ferestrele au baleiat complet intrările tuturor straturilor de neuroni.

Prin această procedură, rețeaua înmagazinează caracteristicile acustico-fonetice ale vectorilor de intrare.

## 2.9 Rețele neurale cu autoorganizare

### 2.9.1 Învățarea hebbiană

O caracteristică importantă a rețelelor neurale este că au capacitatea de a învăța și, prin învățare, își îmbunătățesc performanțele. Atunci când învățarea este supervizată, este dirijată de perechile intrare-ieșire pe care algoritmul trebuie să le recunoască. Scopul algoritmilor nesupervizați sau cu autoorganizare este de a descoperi caracteristicile semnificative ale datelor de intrare fără ca ele să fie transmise din exterior. Acești algoritmi înglobează reguli cu caracter local, care influențează modificările legăturilor sinaptice dintre un neuron și cei aflați în imediata sa vecinătate, permițând sintetizarea unor corespondențe între intrări și ieșiri. Tehnica de adaptare a caracteristicilor rețelei neurale prin autoorganizare se inspiră din neurobiologie într-o măsură mult mai mare decât învățarea supervizată, deoarece acest proces de organizare este fundamental pentru creier.

Învățarea hebbiană constă din modificări ale ponderilor sinaptice pe baza interacțiunilor dintre neuronii vecini, o rețea putând avea mai mulți neuroni activi la un moment dat. Această idee este implementată prin regula Hebb de instruire.

Problema cheie este cum se poate ajunge la o structură viabilă prin autoorganizare, iar răspunsul este dat de Turing [30] care arată că *ordinea globală poate proveni din interacțiuni locale*. Această observație ilustrează faptul că interacțiunile locale aparent aleatoare dintre neuronii învecinați ai unei rețele pot conduce, în final, la un comportament global coerent, proces care constituie esența autoorganizării.

Pentru a ajunge în faza de autoorganizare a rețelei, care este opusă stabilizării, trebuie ca feedback-ul dintre modificările sinaptice și răspunsul rețelei la stimulii de intrare să fie pozitiv. Primul principiu al autoorganizării formulat de von der Malsburg [30] exprimă acest lucru: *modificările asupra ponderilor sinaptice tind să se autoamplifice*. Autoamplificarea se face pe baza semnalelor locale, adică prin semnalul presinaptic și cel postsinaptic, acest mecanism fiind, de fapt, o reformulare a postulatului învățării al lui Hebb.

Pe de altă parte, pentru ca sistemul să se stabilizeze este nevoie de o competiție pentru un set limitat de resurse care pot fi, de exemplu, resurse energetice sau număr de intrări. Cu alte cuvinte, creșterile unor ponderi sinaptice trebuie compensate de scăderi ale altora care pot chiar să dispară. Această observație se poate sintetiza în cel de-al doilea principiu al autoorganizării formulat de von der Malsburg [30]: *limitarea resurselor conduce la competiție între*

*sinapse, rezultând o selecție a celor mai puternice în detrimentul celor mai puțin adaptate.*

Sinapsele nu pot genera singure comportamente favorabile, fiind nevoie de cooperare între sinapsele care converg către un neuron pentru ca acesta să fie activat. Se poate formula, astfel, cel de-al treilea principiu al autoorganizării (von der Malsburg [30]) care arată că *modificările în ponderile sinaptice tind să fie rezultatul cooperării.*

Pentru ca învățarea prin autoorganizare să fie eficientă, trebuie ca datele de intrare să înglobeze un volum semnificativ de informație redundantă. Aceasta este, de fapt, cea care alimentează procesul global de instruire deoarece redundanța din vectorii de instruire este transferată parametrilor rețelei.

## 2.9.2 Învățarea competitivă

Această secțiune prezintă o clasă specială de rețele neurale cunoscute sub numele de *hărți cu autoorganizare (self-organizing maps)*[42]. Ele se bazează pe *învățarea competitivă* astfel încât, în urma unui proces de selecție, doar un neuron sau un grup de neuroni vor fi activați la un moment dat. Instruirea acestor rețele se bazează pe regula *câștigătorul ia tot*.

Într-o hartă cu autoorganizare, neuronii sunt plasați de regulă sub forma unei matrici bidimensionale. În timpul procesului de învățare ei devin selectivi la anumiți vectori de intrare sau la anumite clase de vectori. Poziția neuronilor care răspund fiecărei clase tinde să se stabilizeze pe parcursul instruirii, formându-se o hartă a vectorilor de intrare prin proiecție bidimensională. Acești vectori se vor grupa în urma instruirii conform unor caracteristici interne ale lor, de aceea rezultatul se mai numește și *hartă cu autoorganizare a caracteristicilor*.

### Cuantizarea vectorială

Spațiul intrărilor unei rețele neurale care acceptă ca intrări vectori cu componente reale poate fi considerat infinit. Cuantizarea vectorială este o metodă care produce aproximarea acestui spațiu printr-un număr  $K$  finit de vectori  $\mathbf{w}_i$ ,  $i = 1 \dots K$  numiți și *prototipuri*. Mulțimea de prototipuri se mai numește și *codebook*. Odată ales acest set de vectori, aproximarea intrării  $\mathbf{x}$  înseamnă alegerea vectorului  $\mathbf{w}_a$  care este cel mai apropiat de  $\mathbf{x}$  în sensul distanței euclidiene:

$$\|\mathbf{x} - \mathbf{w}_a\| = \min_i \|\mathbf{x} - \mathbf{w}_i\|, \quad i = 1 \dots K.$$

Principalul avantaj al cuantizării vectoriale constă în reducerea volumului de informație referitoare la mulțimea intrărilor posibile în rețea.

Fiecărui vector  $\mathbf{w}_i$ ,  $i = 1 \dots K$  i se poate asocia o etichetă. În acest fel, alegerea lui  $\mathbf{w}_a$  ca cel mai apropiat dintre ei de vectorul de intrare  $\mathbf{x}$  înseamnă asignarea etichetei acestui vector lui  $\mathbf{x}$ , adică recunoașterea automată a lui  $\mathbf{x}$ .

Pe de altă parte, pentru că numărul vectorilor  $\mathbf{w}_i$ ,  $i = 1 \dots K$  este finit, alegerea lui  $\mathbf{w}_a$  în locul lui  $\mathbf{x}$  pentru prelucrările ulterioare implică introducerea unei distorsiuni, adică a unei erori de cuantizare. În mod evident, cu cât  $K$  este mai mare, cu atât distorsiunea este mai mică.

### Self-organizing feature map

Algoritmul *self-organizing feature map* (SOFM) [42] transformă spațiul multidimensional al vectorilor de intrare într-unul unidimensional sau bidimensional printr-o grupare conform trăsăturilor comune. În urma instruirii, un vector de intrare va provoca activarea unui grup de neuroni bine definit și localizat în matricea ieșirilor. Neuronii care se activează sunt organizați în jurul unui neuron central care generează o valoare de activare maximă. Cu cât un neuron este mai depărtat de neuronul central, cu atât valoarea lui de activare va fi mai mică.

La construirea unei rețele neurale care implementează o hartă cu autoorganizare trebuie implementat un mecanism care compară rezultatele funcțiilor discriminant ale neuronilor și îl selectează pe cel cu răspunsul cel mai mare pentru fiecare vector de instruire. Procesul de adaptare trebuie să decurgă în așa fel încât neuronii selectați să devină și mai bine adaptați la vectorii de intrare.

Vom folosi o rețea neurală ca cea din figura 2.11 în scopul de a prezenta algoritmul de instruire a hărților cu autoorganizare.

Intrările sunt conectate cu toți neuronii de pe stratul de ieșire. Notăm vectorii intrărilor

$$\mathbf{x} = [x_1 \ x_2 \ \dots \ x_M]^t,$$

iar ieșirile vor fi

$$\mathbf{y} = [y_1 \ y_2 \ \dots \ y_N]^t,$$

unde  $M$  este numărul intrărilor în rețeaua neurală, iar  $N$  numărul elementelor care alcătuiesc harta. Matricea ponderilor va fi:

$$\mathbf{W} = \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1M} \\ \vdots & \vdots & \ddots & \vdots \\ w_{N1} & w_{N2} & \dots & w_{NM} \end{bmatrix}.$$

Pentru găsierea celei mai bune potriviri a vectorului de intrare  $\mathbf{x}$  cu matricea ponderilor  $\mathbf{w}_j$  corespunzătoare neuronului cu indicele  $j$ , se compară produsele

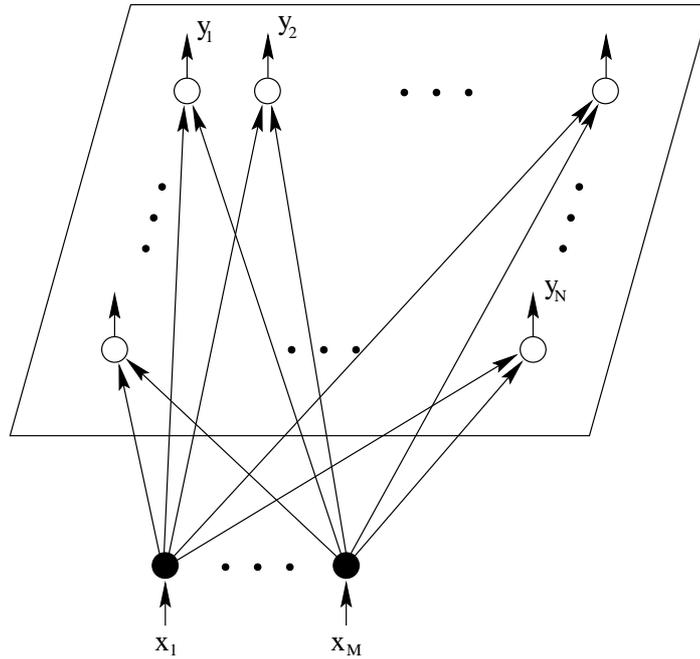


Figura 2.11: Modelul rețelei neurale folosite de SOFM.

scalare  $\mathbf{w}_j \mathbf{x}$ ,  $j = 1 \dots N$  și se alege cel mai mare. Prin această selecție se determină centrul zonei cu activare maximă pentru vectorul  $\mathbf{x}$ .

Algoritmul adaptiv normalizează vectorii de ponderi  $\mathbf{w}_j$ . În această situație, găsirea celei mai bune potriviri poate fi înlocuită cu găsirea distanței euclidiene minime dintre vectori. Astfel se stabilește indicele neuronului care îndeplinește această condiție, adică:

$$i = \text{acel } j \text{ pentru care } \|\mathbf{x} - \mathbf{w}_j\| \text{ este minim, } j = 1 \dots N.$$

Neuronul  $i$  care satisface condiția de mai sus se numește *neuron câștigător* pentru patternul  $\mathbf{x}$ . Prin această relație, spațiul continuu al intrărilor este transformat în spațiul discret al mulțimii neuronilor. În funcție de aplicația implementată, răspunsul rețelei poate fi ori indicele neuronului câștigător, ori vectorul ponderilor sinaptice corespunzătoare acestuia. În același context, răspunsul fiecărui neuron poate fi distanța euclidiană dintre vectorul de intrare și vectorul său de ponderi.

Algoritmul SOFM folosește noțiunea de *vecinătate* a neuronului câștigător (fig. 2.12) care este mulțimea de neuroni din jurul lui asupra cărora se aplică tehnicile de instruire. Dimensiunea vecinătății nu este fixată dinainte, iar forma ei poate fi pătrată, circulară etc.

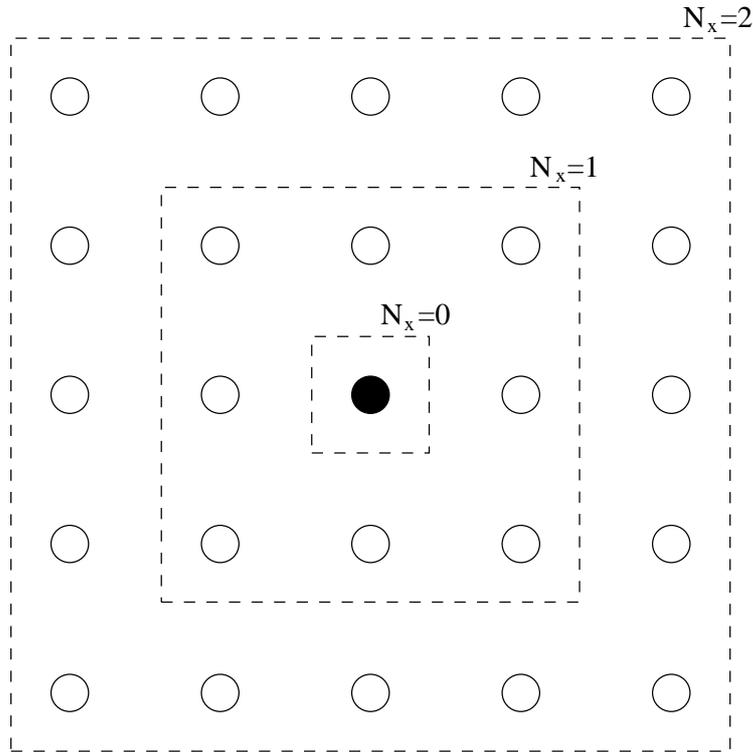


Figura 2.12: Vecinătăți pătrate ale neuronului câștigător.

Experimental [42] s-a dovedit că alegerea unei vecinătăți largi care se micșorează în timp conduce la cele mai bune rezultate în instruire. Această modalitate de selecție a neuronilor asupra cărora se aplică formula de ajustare a ponderilor este echivalentă mai întâi cu o organizare a zonelor de activare pentru fiecare clasă de vectori. Apoi, pe măsura micșorării vecinătăților, aceste zone și, implicit, neuronii câștigători devin din ce în ce mai specializați și mai selectivi. În final, vecinătatea are dimensiune 0, conținând doar neuronul câștigător.

Pentru ca rețeaua să se autoorganizeze, ponderile sinaptice trebuie să se modifice și în funcție de vectorii de intrare. Postulatul lui Hebb al învățării arată că valoarea unei ponderi sinaptice crește în funcție de activitățile presinaptică și postsinaptică [30]. Această regulă este doar parțial satisfăcătoare. Pentru că modificările au loc doar în sensul creșterii, ponderile finale pot ajunge la saturație. Acest lucru este evitat modificând regula lui Hebb prin introducerea unui termen de uitare [30]:

$$-g(y_j)\mathbf{w}_j.$$

În expresia de mai sus,  $\mathbf{w}_j$  este vectorul ponderilor sinaptice ale neuronului  $j$ , iar  $g(y_j)$  este o funcție cu valori pozitive, parametrul său fiind răspunsul

$y_j$  al neuronului. Folosind aceste date, ecuația diferențială care stă la baza formalismului discret este:

$$\frac{\partial \mathbf{w}_j}{\partial t} = \eta y_j \mathbf{x} - g(y_j) \mathbf{w}_j, \quad j = 1, 2 \dots N. \quad (2.26)$$

În relația (2.26),  $\eta$  este rata de instruire. Condiția care se pune asupra funcției  $g$  este

$$g(y_j) = 0 \quad \text{pentru } y_j = 0, \forall j = 1, 2 \dots N,$$

astfel încât pentru un răspuns nul al neuronului să nu se producă nici o modificare a ponderilor.

Experimentele [30] arată că dacă neuronii folosesc o funcție de activare neliniară care le limitează inferior și superior ieșirile, după un număr suficient de mare de pași răspunsul neuronilor din interiorul vecinătății va tinde spre limita maximă, iar răspunsul celor din exterior va tinde către cea minimă. În aceste condiții, se poate considera că:

$$y_j = \begin{cases} 1, & j \text{ activ și se află în interiorul vecinătății} \\ 0, & j \text{ inactiv și se află în exteriorul vecinătății.} \end{cases}$$

Se poate folosi, astfel, o variantă simplificată a funcției  $g(y_j)$ :

$$g(y_j) = \begin{cases} \alpha, & \text{dacă neuronul } j \text{ este activ} \\ 0, & \text{dacă neuronul } j \text{ este inactiv,} \end{cases}$$

unde  $\alpha$  este o constantă pozitivă.

Ecuația (2.26) devine, în aceste condiții:

$$\frac{\partial \mathbf{w}_j}{\partial t} = \begin{cases} \eta \mathbf{x} - \alpha \mathbf{w}_j, & j \text{ în interiorul vecinătății} \\ 0, & j \text{ în afara vecinătății.} \end{cases}$$

Fără a se pierde generalitatea, se alege același factor de scalare astfel încât  $\eta = \alpha$ . Relația care se obține arată că vectorul de ponderi  $\mathbf{w}_j$  tinde să se apropie de vectorul intrărilor  $\mathbf{x}$ :

$$\frac{\partial \mathbf{w}_j}{\partial t} = \begin{cases} \eta (\mathbf{x} - \mathbf{w}_j), & j \text{ în interiorul vecinătății} \\ 0, & j \text{ în afara vecinătății.} \end{cases}$$

Dacă transformăm relația continuă într-una discretă, obținem formula finală care arată cum se calculează ponderile unei hărți cu autoorganizare:

$$\mathbf{w}_j^{(t+1)} = \begin{cases} \mathbf{w}_j^{(t)} + \eta (\mathbf{x} - \mathbf{w}_j^{(t)}), & j \text{ în interiorul vecinătății} \\ \mathbf{w}_j^{(t)}, & j \text{ în afara vecinătății,} \end{cases}$$

unde  $t$  reprezintă pasul instruirii.

O variantă a algoritmului este cea în care rata de instruire se modifică de-a lungul procesului de instruire. În prima fază, cea de *ordonare*, în care se stabilesc zonele de răspuns pentru fiecare clasă de vectori, rata de instruire este mai mare pentru a permite o evoluție rapidă. În faza a doua, cea de *convergență*, în care ponderile sunt rafinate, rata de instruire este redusă pentru o bună acuratețe a rezultatelor.

## Learning vector quantization

Dacă hărțile cu autoorganizare sunt folosite pentru clasificarea vectorilor prin gruparea neuronilor în clase discrete, problema devine una de decizie. Hărțile aproximează semnalele de intrare sau funcțiile de densitate de probabilitate care descriu mulțimile de semnale de intrare prin vectori codebook plasați în spațiul de intrare în așa fel încât să minimizeze o eroare de cuantizare globală bazată pe distanța euclidiană. Când semnalele de intrare ajung să fie clasificate, se folosesc mai mulți vectori codebook pentru a reprezenta fiecare clasă și identitatea vectorilor în clasă nu mai este atât de importantă. Devine importantă decizia la nivelul frontierei dintre clase. Această strategie introdusă de Kohonen [42] constă dintr-o serie de trei algoritmi intitulati generic *learning vector quantization* (LVQ).

Prima variantă numită LVQ1 este cea mai cunoscută și mai utilizată. Se consideră că fiecărei clase îi sunt asignați mai mulți vectori codebook  $\mathbf{w}_j$  și că fiecare dintre ei este etichetat cu simbolul asociat clasei corespunzătoare. Regiunea corespunzătoare fiecărei clase este definită prin compararea distanțelor euclidiene dintre vectorii de intrare  $\mathbf{x}_i$ ,  $i = 1, \dots, M$  și prototipurile  $\mathbf{w}_j$ ,  $j = 1, \dots, N$ . Eticheta celui mai apropiat prototip  $\mathbf{w}_j$  definește clasificarea intrării  $\mathbf{x}_i$ .

Pentru a calcula o poziționare optimă a prototipurilor  $\mathbf{w}_j$ ,  $j = 1, \dots, N$  în spațiul intrărilor, se aleg valorile inițiale ale vectorilor din codebook prin diverse metode, una dintre acestea fiind algoritmul SOFM. Prototipurile sunt apoi etichetate cu ajutorul unui număr de eșantioane pentru care se cunoaște cărei clase aparțin, eticheta asociată fiecărui prototip fiind dată de eticheta predominantă a intrărilor din vecinătatea sa. Acest procedeu se numește *calibrare*. Acuratețea clasificării este îmbunătățită prin algoritmul iterativ pe care îl prezentăm în continuare. Ideea fundamentală a acestuia este de a îndepărta prototipurile de suprafețele de decizie în scopul de a delimita clasele cât mai bine.

Fie vectorul de intrare  $\mathbf{x}_i$  și  $\mathbf{w}_a$  prototipul cel mai apropiat de el în metrică

euclidiană. Atunci algoritmul LVQ1 de actualizare a lui  $\mathbf{w}_a$  este următorul:

$$\begin{aligned} \mathbf{w}_a^{(t+1)} &= \begin{cases} \mathbf{w}_a^{(t)} + \eta(\mathbf{x}_i - \mathbf{w}_j^{(t)}), & \text{când } \mathbf{x}_i \text{ este clasificat corect} \\ \mathbf{w}_a^{(t)} - \eta(\mathbf{x}_i - \mathbf{w}_j^{(t)}), & \text{când } \mathbf{x}_i \text{ nu este clasificat corect} \end{cases} \\ \mathbf{w}_j^{(t+1)} &= \mathbf{w}_j^{(t)}, \text{ pentru } j \neq a. \end{aligned} \quad (2.27)$$

Rata de instruire  $\eta$  este un număr real, pozitiv, subunitar ( $0 < \eta < 1$ ) și poate să fie fixată sau poate să se modifice în timp. În cea de-a doua variantă, se poate porni de la o valoare apropiată de unitate care este micșorată către 0 de-a lungul unui ciclu relativ mare de pași, de exemplu 100000. Se obține în felul acesta o metodă de adaptare fină a prototipurilor, care se stabilizează ușor la valorile finale.

Kohonen a introdus o variantă optimizată a acestui algoritm numită *Optimized Learning Vector Quantization* (OLVQ1) [44] care asociază câte o rată de instruire fiecărui prototip, iar rezultatul este o convergență mai rapidă:

$$\begin{aligned} \mathbf{w}_a^{(t+1)} &= \begin{cases} \mathbf{w}_a^{(t)} + \eta_a^{(t)}(\mathbf{x}_i - \mathbf{w}_j^{(t)}), & \text{când } \mathbf{x}_i \text{ este clasificat corect} \\ \mathbf{w}_a^{(t)} - \eta_a^{(t)}(\mathbf{x}_i - \mathbf{w}_j^{(t)}), & \text{când } \mathbf{x}_i \text{ nu este clasificat corect} \end{cases} \\ \mathbf{w}_j^{(t+1)} &= \mathbf{w}_j^{(t)}, \text{ pentru } j \neq a. \end{aligned} \quad (2.28)$$

Relația (2.28) se poate scrie sub forma:

$$\mathbf{w}_a^{(t+1)} = (1 - s^{(t)}\eta_a^{(t)})\mathbf{w}_a^{(t)} + s^{(t)}\eta_a^{(t)}\mathbf{x}_i^{(t)},$$

unde  $s^{(t)} = +1$  dacă  $\mathbf{x}_i$  este clasificat corect și  $s^{(t)} = -1$  în caz contrar. Atunci valoarea optimă a ratei de instruire va fi [44]:

$$\eta_a^{(t+1)} = \frac{\eta_a^{(t)}}{1 + s^{(t)}\eta_a^{(t)}}.$$

O variantă modificată a acestui algoritm este LVQ2.1. De această dată, în procesul de instruire se folosesc două prototipuri,  $\mathbf{w}_j$  și  $\mathbf{w}_k$ , cele mai apropiate de vectorul de intrare  $\mathbf{x}_i$ , care se găsesc în clase diferite și care vor fi modificate simultan. Fără a afecta generalitatea problemei, se presupune că cele două clase cărora le aparțin prototipurile au distribuții diferite, iar planul de decizie care inițial este stabilit la jumătatea distanței dintre  $\mathbf{w}_j$  și  $\mathbf{w}_k$  trebuie deplasat către poziția sa reală. Așa cum se poate vedea și din figura 2.13, se definește o fereastră centrată pe mijlocul acestei distanțe, iar algoritmul se aplică doar vectorilor de intrare care se găsesc în această fereastră.

Dacă se lucrează în spații multidimensionale, atunci limitele acestei ferestre sunt hipersfere. Fereastra se definește în termenii distanțelor relative  $d_j$  și  $d_k$  dintre  $\mathbf{x}_i$  și  $\mathbf{w}_j$ , respectiv  $\mathbf{w}_k$ , fiind caracterizată de un parametru constant

$$s = \frac{1 - l}{1 + l},$$

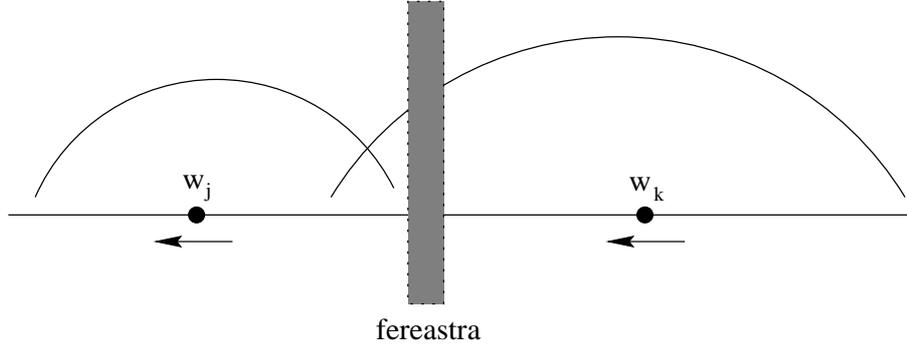


Figura 2.13: Fereastra folosită în algoritmi LVQ2.1 și LVQ3. Curbele reprezintă distribuția eșantioanelor din cele două clase.

unde  $l$  este lățimea relativă a ferestrei în punctul său cel mai îngust. Pentru aplicațiile practice, în [44] se recomandă ca  $l$  să ia valori între 0,1 și 0,2. Se consideră că vectorul  $\mathbf{x}_i$  se găsește în interiorul ferestrei dacă

$$\min\left(\frac{d_j}{d_k}, \frac{d_k}{d_j}\right) > s.$$

Formula de actualizare a prototipurilor conform algoritmului LVQ2.1 este [44]:

$$\begin{aligned} \mathbf{w}_j^{(t+1)} &= \mathbf{w}_j^{(t)} - \eta^{(t)}(\mathbf{x}^{(t)} - \mathbf{w}_j^{(t)}) \\ \mathbf{w}_k^{(t+1)} &= \mathbf{w}_k^{(t)} + \eta^{(t)}(\mathbf{x}^{(t)} - \mathbf{w}_k^{(t)}), \end{aligned} \quad (2.29)$$

unde  $\mathbf{w}_j$  și  $\mathbf{w}_k$  sunt prototipurile cele mai apropiate de  $\mathbf{x}_i$ , în condițiile în care  $\mathbf{x}_i$  și  $\mathbf{w}_k$  aparțin aceleiași clase, iar  $\mathbf{x}_i$  și  $\mathbf{w}_j$  sunt în clase diferite. În toate celelalte cazuri avem:

$$\mathbf{w}_m^{(t+1)} = \mathbf{w}_m^{(t)}, \quad m \neq j, k. \quad (2.30)$$

Acest algoritm corectează pozițiile în spațiu ale prototipurilor  $\mathbf{w}_j$  și  $\mathbf{w}_k$  în așa fel încât planul de decizie poziționat inițial la jumătatea distanței dintre cei doi vectori să se deplaseze către limita dintre cele două distribuții, făcând ca în final să coincidă cu limita de decizie Bayes.

O altă variantă modificată a algoritmului LVQ este LVQ3. Algoritmul LVQ2.1 se bazează pe ideea de modificare a limitei dintre două clase fără a lua în considerare efectele produse de continuarea rulării acestuia și după ce poziția corectă a limitei dintre clase a fost atinsă. În [42] se arată că, deoarece corecțiile se fac proporțional cu distanțele dintre  $\mathbf{x}_i$  și  $\mathbf{w}_j$ , respectiv  $\mathbf{x}_i$  și  $\mathbf{w}_k$ , modificarea lui  $\mathbf{w}_k$ , prototipul care face parte din clasa corectă, este mai amplă decât cea a lui  $\mathbf{w}_j$ . Rezultatul este o scădere monotonă a distanței  $\|\mathbf{w}_j - \mathbf{w}_k\|$ , iar vectorii codebook nu mai sunt plasați optim pentru a descrie cu acuratețe

forma granițelor dintre clase. Din acest motiv, este absolut necesar ca prototipurile să continue să aproximeze distribuția claselor indiferent de numărul de pași ai algoritmului. Algoritmul LVQ3 este următorul [44]:

$$\begin{aligned}\mathbf{w}_j^{(t+1)} &= \mathbf{w}_j^{(t)} - \eta^{(t)}(\mathbf{x}^{(t)} - \mathbf{w}_j^{(t)}) \\ \mathbf{w}_k^{(t+1)} &= \mathbf{w}_k^{(t)} + \eta^{(t)}(\mathbf{x}^{(t)} - \mathbf{w}_k^{(t)}),\end{aligned}\tag{2.31}$$

unde  $\mathbf{w}_j$  și  $\mathbf{w}_k$  sunt prototipurile cele mai apropiate de  $\mathbf{x}_i$ , în condițiile în care  $\mathbf{x}_i$  și  $\mathbf{w}_k$  aparțin aceleiași clase, iar  $\mathbf{x}_i$  și  $\mathbf{w}_j$  sunt în clase diferite. Dacă  $\mathbf{x}_i$ ,  $\mathbf{w}_j$  și  $\mathbf{w}_k$  aparțin aceleiași clase și  $l \in \{j, k\}$ , actualizarea se face conform formulei:

$$\mathbf{w}_l^{(t+1)} = \mathbf{w}_l^{(t)} + \epsilon\eta^{(t)}(\mathbf{x}_i^{(t)} - \mathbf{w}_l^{(t)}).\tag{2.32}$$

Modificarea constă în introducerea unei relații de actualizare suplimentare variabilă pentru cazul în care cele două prototipuri aparțin aceleiași clase. În [44] se arată că valoarea parametrului  $\epsilon$  poate fi între 0,1 și 0,5, iar valoarea optimă depinde de dimensiunea ferestrei, fiind mai mică pentru ferestre mai înguste. Algoritmul are tendința de a se stabili, iar continuarea instruirii nu afectează plasarea optimă a prototipurilor. Cu toate acestea, în [76] se arată că atunci când se asignează câte un singur prototip fiecărei clase, LVQ3 este similar cu LVQ2.1.

Se impune observația că în timp ce LVQ1 actualizează un singur prototip în timpul unui pas de instruire, algoritmi LVQ2.1 și LVQ3 modifică doi vectori codebook simultan.

## 2.10 Concluzii

În acest capitol am prezentat modelele neurale fundamentale și principalele metode care se folosesc pentru instruire. Am văzut care sunt cele mai importante caracteristici ale neuronului McCulloch-Pitts și ale perceptronului lui Rosenblatt și am studiat principalele reguli de instruire a unui neuron. Am arătat că rețelele neurale feedforward și feedback sunt cele mai utilizate arhitecturi. Am prezentat apoi perceptronul multistrat și algoritmul *backpropagation of error* pentru instruirea supervizată a sa. Rețelele neurale cu întârziere în timp sunt o extensie a perceptronului multistrat și pot fi folosite pentru instruirea vectorilor cu caracteristici temporale. În final, am văzut că hărțile cu autoorganizare și algoritmi tip LVQ sunt cele mai importante modele de învățare competitivă. În capitolul următor vom analiza printr-un studiu comparativ rezultatele obținute în recunoașterea cuvintelor izolate cu perceptroni multistrat, cu rețele neurale cu întârziere în timp și cu hărți cu autoorganizare.

## Capitolul 3

# Recunoașterea vorbirii cu modele clasice de rețele neurale

*Acest capitol prezintă un studiu comparativ asupra performanțelor în recunoașterea cuvintelor izolate și a vorbitorilor a perceptronilor multistrat, a rețelelor neurale cu întârziere în timp și a hărților cu autoorganizare. Au fost folosite două baze de date care conțin informații prelucrate prin analiză cepstrală după înregistrarea semnalelor vocale mai întâi într-un mediu fără perturbații, iar apoi după ce au fost transmise printr-o linie telefonică. Cu toate că pentru recunoașterea cuvintelor izolate rezultatele au fost bune, testele de recunoaștere a vorbitorilor au demonstrat că aceste modele nu oferă întotdeauna soluții satisfăcătoare.*

### 3.1 Introducere

Problemele computaționale asociate recunoașterii vorbirii și succesul limitat al diverselor tehnici convenționale au condus la dezvoltarea unor abordări bazate pe rețelele neurale. Metodele descrise în literatură diferă în special prin modul în care sunt convertite semnalele vocale în diverse formate cu scopul de a fi utilizate ca intrări ale rețelelor neurale, prin categoriile de aplicații implementate, de exemplu recunoașterea cuvintelor sau a fonemelor, recunoaștere dependentă sau independentă de vorbitor etc. și prin tipul de rețea neurală folosită.

Pentru ca o rețea să fie dinamică și astfel să fie capabilă să proceseze informație vocală trebuie să dispună de memorie [18]. Cea mai simplă strate-

gie care poate fi folosită este de a prezenta o secvență din fluxul de date la intrarea rețelei neurale. Aceasta este o strategie statică deoarece nu surprinde natura temporală a datelor. Stratul de intrare poate fi privit ca un buffer care păstrează secvența temporală curentă. Există două metode prin care se poate schimba conținutul buffer-ului: informația din registru se poate deplasa pentru a fi înlocuită treptat cu date care sosesc pe fluxul de intrare sau poate conține eşantioane obținute prin ferestruirea fluxului de date, ferestrele putându-se suprapune. În [7], [23] se folosește această abordare care, în mod tipic, are ca suport perceptronul multistrat. Abordarea statică are câteva dezavantaje [18]: impune vectorilor durată fixă, nu ia în considerare caracteristicile temporale ale intrărilor, iar perceptronul multistrat nu se comportă satisfăcător pentru date noi. De altfel, s-a arătat că rețelele neurale feedforward nu sunt capabile să învețe relații temporale care, din acest motiv, trebuie programate în avans [24].

Rețelele neurale *time-delay*, cu întârziere în timp, sunt caracterizate prin faptul că semnalul de intrare este întârziat, iar ieșirile sunt calculate în funcție de intrările curente și cele anterioare. În [90] se descrie acest model care este folosit și în câteva teste de recunoaștere a vorbirii. Modelul prezentat acolo are o proprietate importantă legată de structura dinamică a vorbirii: este invariant la translație, adică informația învățată de rețea nu depinde de deplasarea în timp.

Un model diferit de cele anterioare este harta cu autoorganizare [45] care a fost folosită de Kohonen pentru recunoașterea automată a vorbirii. Cu ajutorul *hărții fonotopice (phonotopic map)* bazate pe acest model, Kohonen a implementat un sistem de transcriere automată a vorbirii numit *phonetic typewriter* capabil să extragă un text după dictare [40].

Majoritatea acestor lucrări oferă rezultate, rate de recunoaștere mai bune decât cele obținute prin alte tehnici. De exemplu, în [23] este prezentat un sistem de recunoaștere a vorbirii care permite recunoașterea unui vocabular format din 45 de cuvinte din limba germană, independent de vorbitor. Experimentele descrise aici arată că rata de recunoaștere este de până la 91% pentru vorbitori de același sex și de până la 72% pentru un set de date care cuprinde cuvinte pronunțate atât de bărbați cât și de femei. Este, însă, dificil să comparăm în mod direct rezultatele obținute de diverși autori pentru că în majoritatea cazurilor acestea se bazează pe seturi particulare de date. Este de dorit ca pentru toate modelele testate să folosim aceleași date de instruire, obținând astfel o imagine mult mai fidelă a performanțelor lor. Scopul acestui capitol este de a prezenta o comparație a rezultatelor obținute în recunoașterea vorbirii prin trei tehnici diferite [66]: perceptronul multistrat (MLP), rețeaua

neurală cu întârziere în timp (TDNN) și harta cu autoorganizare (SOM), pentru aceeași bază de date. Pentru că aceste tipuri de rețele neurale sunt foarte cunoscute, le-am asociat atributul "clasice".

### 3.2 Bazele de date Speechdata și Telephdata

În testele de recunoaștere a vorbirii cu MLP, TDNN și SOM am folosit două baze de date care au fost realizate la Northern Ireland BioEngineering Center (NIBEC). Ele cuprind pronunțiile în limba engleză ale cifrelor 0 până la 9 la care au fost adăugate *nought*, *oh* și *zero* care sunt variante ale cifrei 0. Pentru crearea fiecărei baze de date au fost alese câte 25 de persoane care au repetat cele 12 cuvinte de câte 20 de ori, rezultând un număr total de 6000 de pronunții. La alcătuirea primei baze de date au fost aleși 19 bărbați și 6 femei, iar pentru cea de-a doua 15 bărbați și 10 femei. Toate persoanele au avut vârsta cuprinsă între 20 și 60 de ani și au folosit limba engleză cu accent nord-irlandez.

Pentru prima bază de date care se numește Speechdata, informațiile au fost înregistrate într-un mediu controlat, în care zgomotele și interferențele au fost reduse la minim. Fiecare vorbitor a fost înregistrat într-un studio izolat fonic, pe un casetofon profesional, cu un microfon de înaltă calitate. Semnalul vocal a fost apoi digitizat folosind un convertor analog/digital pe 12 biți cu frecvența de eșantionare de 7,5 KHz, după care informația a fost stocată pe suport magnetic.

Cea de-a doua bază de date numită Telephdata este destinată recunoașterii automate a vorbirii printr-o rețea telefonică, fiind mai apropiată de condițiile unei aplicații reale. Semnalul vocal a fost, de această dată, captat și stocat în timp real. Vorbitorii au fost așezați pe rând la un pupitru dotat cu un aparat telefonic și un calculator. Ei au pronunțat la telefon cuvintele pe care le-au citit de pe ecran. Semnalul telefonic recepționat a fost introdus în sistemul de digitizare prezentat mai sus.

Crearea mulțimii de vectori se bazează pe folosirea unui banc de filtre ale cărui benzi de trecere sunt calculate conform unei scale neliniare numită *scala mel*. Ieșirile generate de acest banc de filtre au fost proiectate în domeniul cepstral printr-o transformare cosinus discretă. Coeficienții cepstrali împreună cu o colecție de *coeficienți delta*, calculați pe baza unor diferențe ale primilor, alcătuiesc informația folosită la instruirea rețelelor neurale.

Oricare ar fi sistemul de recunoaștere a vorbirii folosit, etapa de preprocesare a semnalului vocal este foarte importantă. Scopul ei este de a extrage componentele de bază ale semnalului acustic și de a le coda într-un set compact și eficient de parametri.

Scala de frecvențe *mel* a fost realizată în urma studierii modului în care

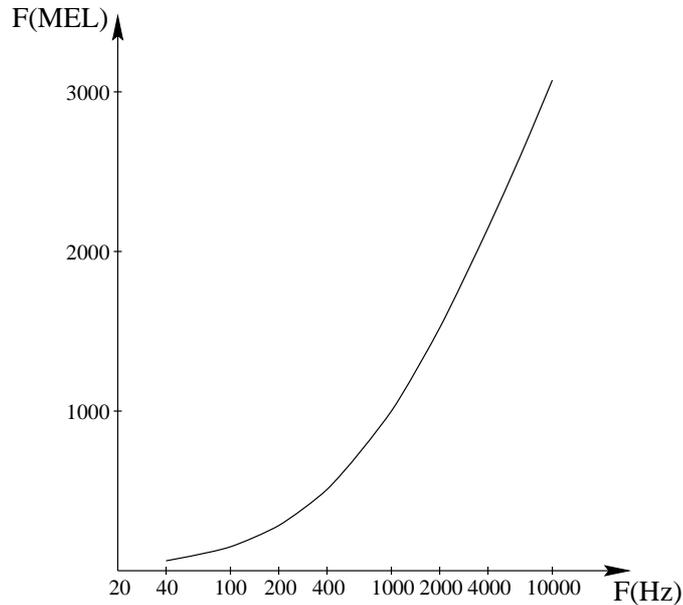


Figura 3.1: Scala de frecvențe *mel* realizată în urma studierii modului de percepție a sunetelor de sistemul auditiv uman.

sunetele sunt percepute de sistemul auditiv uman. *Mel* este o unitate de măsură a înălțimii percepute de ureche umană pentru un ton. Este cunoscut faptul că sunetele de frecvențe scăzute sunt distinse mai greu de către om. La fel se întâmplă și cu sunetele de frecvență înaltă. În funcție de nivelul de antrenare și de starea de sănătate, frecvențele maxime audibile pot să difere de la un individ la altul. Stevens și Volkman au stabilit în 1940 această scală pe baza unui experiment descris în continuare [16]. În mod arbitrar au stabilit frecvența de referință de 1000 Hz și i-au atribuit 1000 mel. Ascultătorilor li s-a cerut apoi să modifice frecvența fizică a sunetului până când intensitatea percepută a devenit de două ori mai mare decât referința. Această frecvență a fost etichetată cu 2000 mel. Repetând procedeul pentru frecvențe mai mari de 1000 Hz, dar și mai mici decât acest prag, a fost stabilită prin mediere scala din figura 3.1.

Aceasta este aproximativ liniară pentru frecvențe sub 1 KHz și logaritmică peste acest nivel. Pentru aproximarea ei, în faza de preprocesare a bazei de date s-a folosit relația:

$$F_{\text{MEL}} = 2595 \log \left( 1 + \frac{F_{\text{Hz}}}{700} \right).$$

Pornind de la rezultatul prezentat mai sus, a fost realizat un banc de filtre, construcție care se bazează pe capacitatea sistemului auditiv de a percepe diferitele frecvențe din spectrul audio. Pentru frecvențe sub 1 KHz lățimile

benzilor sunt de 100 Hz și cresc logaritmice peste 1KHz (figura 3.2).

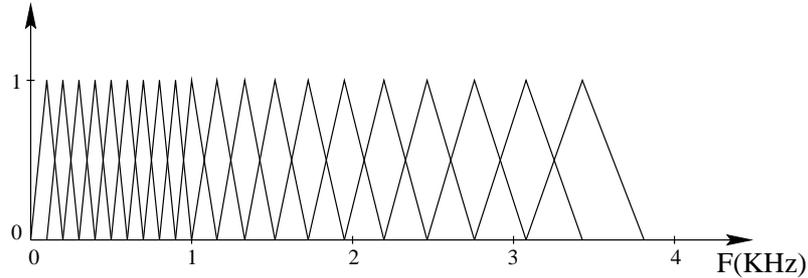


Figura 3.2: Banc de filtre a cărui construcție se bazează pe capacitatea sistemului auditiv uman de a percepe diferitele frecvențe din spectrul audio.

Semnalul vocal a fost împărțit în cadre de 20 ms folosind ferestre Hamming. Fiind suprapuse, durata lor a fost stabilită la 30 ms. Semnalul ferestruit a fost aplicat bancului de filtre, iar ieșirea fiecăruia dintre ele,  $X_j$ , a fost obținută prin agregarea componentelor spectrale din banda de frecvență  $j$ .

Conform modelului general acceptat, semnalul vocal este alcătuit dintr-o secvență de excitație compusă printr-o relație de convoluție cu răspunsul la impuls al tractului vocal. Noi avem acces doar la semnalul de ieșire, însă este deosebit de utilă eliminarea uneia dintre componente pentru studiul celeilalte. Deoarece combinația rezultată nu este liniară, separarea devine o problemă destul de dificilă.

Să presupunem că dorim să studiem un semnal  $x(n)$  de frecvență joasă care este perturbat prin suprapunere de un semnal de frecvență înaltă  $w(n)$ . Semnalul rezultat va fi  $y(n) = x(n) + w(n)$ . Putem folosi transformata Fourier care este un operator liniar pentru a transfera semnalul din domeniul timp în domeniul frecvență. Spectrul este o reprezentare a semnalului în care putem separa acele componente specifice semnalului util de cele ale perturbației. Dacă scopul este studiul semnalului  $x(n)$ , putem folosi un filtru trece-jos care extrage componentele dorite. Considerând că  $O$  este un operator care realizează filtrarea, atunci

$$O\{y(n)\} = O\{x(n) + w(n)\} = O\{x(n)\} + O\{w(n)\} \approx x(n).$$

Atunci când semnalele sunt combinate prin convoluție, nu se poate trage o concluzie despre influența filtrului asupra fiecărei componente. De aceea a fost adoptată utilizarea cepstrului în probleme de recunoașterea vorbirii [16], [72]. Folosind această metodă, componentele semnalului vor fi separate în cepstru și

vor fi combinate liniar. Pentru un semnal  $y(n)$ , cepstrul se definește astfel:

$$c_y(n) = F^{-1}\{\log |F\{y(n)\}|\} = \frac{1}{2\pi} \int_{-\pi}^{\pi} \log |Y(\omega)| e^{j\omega n} d\omega$$

unde  $F\{\cdot\}$  este o transformare Fourier discretă. Pentru că  $\log |Y(\omega)|$  este par,  $F\{\cdot\}$  se înlocuiește cu o transformare cosinus.

Pe baza acestor considerații, coeficienții cepstrali folosiți în bazele de date testate au fost calculați conform următoarei relații:

$$C_i = \sum_{j=1}^N \left( \log |X_j| \cdot \cos \left( \frac{\pi i}{N} (j - 0,5) \right) \right), \quad 1 \leq i \leq M$$

unde  $N$  este numărul de filtre din bancul de filtre, iar  $M$  este numărul dorit de coeficienți cepstrali. Fiecare coeficient delta,  $d_P$ , al cadrului  $P$  se calculează conform formulei [16]:

$$d_P = \frac{\sum_{i=1}^M i(C_{P+i} - C_{P-i})}{2 \sum_{i=1}^M i^2},$$

unde  $C_P$  este coeficientul cepstral pentru cadrul  $P$ . Pentru începutul și sfârșitul cuvântului se folosesc relațiile [16]:

$$d_P = C_{P+1} - C_P, \quad P < M$$

și

$$d_P = C_P - C_{P-1}, \quad P \geq N_F - M$$

unde  $N_F$  este numărul de cadre din cuvânt.

Pentru Speechdata și Telephdata, valorile lui  $N$  și  $M$  au fost 16 și 8, vectorul unui cadru având 18 componente care constau din 8 coeficienți cepstrali, 1 coeficient numit *log-energy* obținut prin logaritmare a energiei semnalului din fereastra Hamming curentă și 9 coeficienți delta. Semnalul vocal pentru fiecare cuvânt a fost procesat în 15 ferestre Hamming și s-au obținut câte 270 de coeficienți pentru un cuvânt.

### 3.3 Experimente

Vom descrie în acest subcapitol rezultatele obținute la testarea capacității de recunoaștere a perceptronului multistrat, a rețelelor neurale cu întârziere în timp și a hărților cu autoorganizare folosind bazele de date Speechdata și Telephdata.

De regulă, condițiile experimentelor nu sunt aceleași și compararea obiectivă a performanțelor diverselor tipuri de rețele neurale este dificil de realizat prin

sintetizarea rezultatelor disponibile în literatură. Atunci când bazele de date diferă, comparația nu se poate face corect. Acesta este motivul pentru care am considerat necesară o analiză pe baza unui singur set de date.

Experimentele prezentate în acest paragraf au urmărit două obiective: testarea performanțelor acestor rețele neurale în recunoașterea vorbirii și în recunoașterea vorbitorilor. Am separat datele de instruire disponibile în cele două baze de date într-un set folosit la instruire și altul folosit la recunoaștere. În setul de antrenare am inclus 40% din date, iar la testare le-am utilizat pe celelalte 60%.

Rețelele neurale pe care le-am testat în acest capitol sunt modele de referință în domeniu și a fost dezvoltată o paletă foarte largă de implementări ale algoritmilor lor de instruire. Multe dintre aceste implementări sunt publice și sunt disponibile în rețeaua Internet cu scopul declarat de a facilita activitatea de cercetare.

Stuttgart Neural Network Simulator (SNNS) [98] este un mediu eficient și flexibil pe care l-am folosit la implementarea perceptronilor multistrat și a rețelelor neurale cu întârziere în timp. Are patru componente importante: nucleul de simulare, interfața grafică (*GUI*), simulatorul programabil (*Batchman*) și compilatorul (*snns2c*). Nucleul de simulare realizează operațiile asupra rețelelor neurale. Interfața grafică este un modul care controlează simularea în timpul rulării. Oferă, de asemenea, o reprezentare grafică a rețelelor neurale. Interfața grafică poate fi folosită direct pentru a crea, manipula și vizualiza rețelele neurale în diverse moduri. Simulatorul folosește un limbaj de programare special, integrat în SNNS. Permite execuția operațiilor în background și implementarea rețelelor neurale fără interfața grafică. Acest limbaj de programare este similar unui limbaj de nivel înalt. Toate opțiunile care sunt accesibile prin interfața grafică sunt disponibile și în *Batchman*.

Pachetul de programe SOM\_PAK [43] a fost realizat de o echipă condusă de T. Kohonen și implementează toate operațiile necesare unei rulări corecte a instruirii și aplicării hărților cu autoorganizare. Am folosit acest set de programe în experimentele de recunoaștere a cuvintelor izolate și a vorbitorilor cu SOM. Modulele de inițializare (*randinit*, *lininit*), instruire (*vsom*) și testare (*qerror*) pot fi rulate atât automat, folosind programul *vfind* care solicită în mod interactiv toți parametrii necesari, sau separat, direct de la consolă. Acest pachet oferă, în plus, o serie de module pentru calibrare, vizualizare și testare a hărților cu autoorganizare.

### 3.3.1 Recunoașterea cuvintelor

Cele două baze de date conțin câte 12 cuvinte pronunțate de câte 25 de vorbitori. Fiecare vorbitor a repetat fiecare cuvânt de câte 20 de ori, astfel încât atât Speechdata cât și Telephdata conțin câte 6000 de repetiții. De fiecare dată primele 8 repetiții ale fiecărui cuvânt au fost folosite la instruire și ultimele 12 la test. Datorită metodei de creare a bazelor de date, vectorii au fost organizați sub forma unor matrici de 15x18 componente, deci rețelele neurale folosite au avut câte 270 de intrări.

Atât pentru perceptronul multistrat cât și pentru rețelele neurale cu întârziere în timp am folosit configurații cu un singur strat ascuns. Numărul de neuroni de pe stratul ascuns a fost de aproximativ 100 de neuroni, cu mici variații care nu au influențat semnificativ capacitatea de învățare.

Pentru că numărul de cuvinte este 12, am stabilit ca acesta să fie și numărul de neuroni de ieșire pentru MLP și TDNN. Fiecare ieșire se activează pentru câte o singură clasă de vectori, generând o ieșire apropiată de 1. Ieșirile care nu se activează generează valori apropiate de 0. În faza de instruire am folosit chiar valorile 1 și 0 pentru ieșirile dorite, iar la testare am acceptat un prag de 0,3. Aceasta înseamnă că o ieșire între 0,7 și 1 este asimilată lui 1, iar una între 0 și 0,3 este 0.

Hărțile cu autoorganizare folosite pentru recunoașterea cuvintelor au avut 17x15 noduri. Instruirea a avut loc în două etape. Pentru prima, cea de organizare a hărții, am rulat 10000 de pași ai algoritmului, iar pentru a doua, cea a accentuării zonelor, am parcurs 120000 de pași. Prin calibrare se stabilește pe baza distanței euclidiene clasa căreia îi aparține fiecare neuron. Aceasta se realizează statistic, prin prezentarea la intrările hărții instruite a unor vectori despre care se știe din ce clase fac parte. Nodurile rămase neetichetate nu sunt folosite în faza de test.

Rezultatele experimentelor de recunoaștere a cuvintelor pentru fiecare tip de rețea neurală și pentru ambele baze de date sunt prezentate în tabelul 3.1.

Tabelul 3.1: Rezultatele experimentelor de recunoaștere a cuvintelor.

|            | MLP   | TDNN  | SOM   |
|------------|-------|-------|-------|
| Speechdata | 93,3% | 95,1% | 97,1% |
| Telephdata | 76,3% | 90,5% | 96,8% |

### 3.3.2 Recunoașterea vorbitorilor

O problemă mai dificilă pe care am implementat-o cu modelele clasice de rețele neurale a fost recunoașterea vorbitorilor. Ne-am propus să instruiim MLP și TDNN cu vectori grupați după criteriul identității persoanei. Astfel, cuvintele din cele două baze de date au fost separate în 25 de clase care corespund celor 25 de vorbitori. Acum nu ne mai interesează care sunt cuvintele pe care le-au pronunțat subiecții pentru că această informație devine nerelevantă. Repetarea, însă, a acelorași cuvinte produce redundanța necesară algoritmilor de instruire care pot extrage trăsăturile dominante ale vocii fiecărui vorbitor. Răspunsurile dorite ale vectorilor de instruire au fost ca și la recunoașterea cuvintelor secvențe binare cu o componentă care are valoarea 1 și care arată clasa din care face parte vectorul, iar restul au valoarea 0. Fiind vorba de 25 de clase, stratul de ieșire al rețelelor neurale a fost redimensionat la această valoare. Stratul ascuns a crescut la 150 de noduri deoarece am constatat experimental că este soluția optimă.

În cazul instruirii SOM am optat pentru o rețea de 30x30 de noduri care a fost calibrată cu vectori din cele 25 de clase.

Sinteza rezultatelor obținute pentru acest set de experimente se găsește în tabelul 3.2.

Tabelul 3.2: Rezultatele experimentelor de recunoaștere a vorbitorilor.

|            | MLP   | TDNN   | SOM   |
|------------|-------|--------|-------|
| Speechdata | 71,2% | 57,2%  | 52,8% |
| Telephdata | 6,7%  | 32,33% | 44,7% |

## 3.4 Concluzii

Rezultatele experimentale prezentate în acest capitol creează o imagine a capacității unor modele clasice de rețele neurale de a rezolva probleme de recunoaștere a vorbirii.

Recunoașterea cuvintelor izolate din Speechdata a condus la rezultate satisfăcătoare. La recunoașterea cuvintelor din Telephdata, perceptronul multistrat a avut o comportare foarte slabă. Modul de construire a bazei de date a permis MLP obținerea unui rezultat apropiat de TDNN, chiar dacă structura rețelei neurale nu permite extragerea unor informații legate de secvența în timp a cadrelor ce alcătuiesc un cuvânt. Performanțele SOM arată potențialul crescut

al rețelelor neurale cu autoorganizare. Studiind ratele de recunoaștere, putem trage concluzia că modelul standard MLP nu este potrivit pentru recunoașterea eficientă a cuvintelor izolate.

Pe de altă parte, rezultatele obținute la testele de recunoașterea vorbitorilor arată că modelele clasice sunt insuficiente. Ratele de succes de până la 50% dovedesc faptul că ele nu sunt capabile să extragă caracteristicile vocale ale indivizilor într-o măsură suficient de bună și este nevoie de implementarea unor tehnici care să răspundă mai bine cerințelor. Pentru a îmbunătăți performanțele rețelelor neurale în recunoașterea vorbitorilor, în capitolul următor vom introduce rețelele neurale concurente, un nou model de recunoaștere care folosește ideea de modularitate și principiul *divide et impera*.

## Capitolul 4

# Rețele neurale concurente

*Propunem în acest capitol un nou model neural de recunoaștere numit Concurrent Neural Networks (CNN) care constă dintr-o colecție de rețele neurale, iar decizia la clasificare se ia prin regula câștigătorul ia tot. Fiecare rețea neurală componentă a sistemului este instruită separat pentru a răspunde corect intrărilor dintr-o singură clasă. Am aplicat acest model la recunoașterea vorbitorilor și a vocalelor folosind perceptroni multistrat, rețele neurale cu întârziere în timp și hărți cu autoorganizare ca module ale sistemului concurent.*

### 4.1 Introducere

Experimentele de recunoaștere a vorbitorilor descrise în capitolul anterior au demonstrat că modelele clasice pe care le-am testat nu oferă întotdeauna rezultate mulțumitoare. Perceptronul multistrat conduce uneori la rezultate mai slabe datorită faptului că arhitectura sa are un caracter general și nu exploatează caracteristicile semnalului vocal.

Performanțele acestor rețele scad foarte mult pentru recunoașterea vorbitorilor. În experimentele noastre, numărul claselor folosite la recunoașterea vorbitorilor este dublu față de cel folosit pentru recunoașterea cuvintelor. Reducând numărul de vorbitori, se constată, însă, că performanțele rămân inferioare celor stabilite la recunoașterea cuvintelor, ceea ce arată că problema este una mai complexă.

Pe de altă parte, este evident că reducerea numărului de vorbitori provoacă îmbunătățirea comportamentului rețelelor neurale prin creșterea ratei de succes pentru că problema pe care o implementează devine mai puțin dificilă. Problema globală poate fi redusă la o sumă de subprobleme prin separarea vorbitorilor în grupuri mai mici în interiorul cărora recunoașterea este mai ușor de realizat.

*Rețelele neurale modulare* [30] au la bază ideea că, pe lângă nivelele ierarhice de organizare a rețelelor neurale artificiale: sinapse, neuroni, straturi de neuroni și rețeaua însăși, se poate crea un nou nivel care combină mai multe rețele neurale. Modularitatea are o justificare neurobiologică, fiind un principiu important în arhitectura sistemelor nervoase ale vertebratelor. De altfel, sistemul vizual oferă un bun exemplu de calcul modular deoarece imaginile percepute sunt descompuse în imagini de dimensiuni reduse care sunt analizate separat, iar rezultatele individuale sunt apoi combinate.

Haykin [30] arată că o rețea neurală se numește modulară atunci când calculul realizat de rețea poate fi descompus în două sau mai multe module care operează pe seturi distincte de intrări, fără a comunica între ele. Ieșirile modulelor sunt mediate de o unitate de integrare care nu permite transmiterea informației înapoi către module. Unitatea de integrare trebuie să decidă cum se combină ieșirile modulelor pentru a forma ieșirea finală a sistemului și trebuie să decidă ce seturi de vectori trebuie învățate de fiecare modul. Modularitatea poate fi privită ca o formă de implementare a principiului *divide et impera* care permite rezolvarea problemelor complexe prin descompunerea acestora în subprobleme și combinarea soluțiilor individuale.

Rețelele neurale modulare au fost studiate în diverse contexte. Rueckl, Cave și Kosslyn pun în [74] problema identificării unor obiecte cunoscute și a găsirii locației acestora în spațiu. S-a constatat că timpul de instruire a fost mai scurt atunci când s-au folosit două rețele neurale în locul uneia singure. Tehnica numită *connectionist glue* propusă de Waibel, Sawai și Shikano [91] presupune instruirea individuală a rețelelor neurale pentru diferite probleme de recunoaștere a vorbirii urmată de regruparea lor. Jacobs, Jordan și Barto [34] instruiesc un sistem modular în care modulele concurează pentru a învăța diverși vectori particulari. Anand, Mehora, Mohan și Ranka [1] propun o metodă prin care problema clasificării vectorilor în  $K$  clase este transformată în  $K$  probleme de clasificare a vectorilor în două clase. Modulele decid dacă un vector aparține unei clase sau nu. Algoritmul backpropagation este modificat pentru a suplini diferența mare dintre exemplele pozitive și cele negative pe care o presupune această soluție.

Modelul propus de noi [11] numit *Concurrent Neural Networks (CNN)* introduce o nouă tehnică neurală de recunoaștere care se bazează pe ideea competiției între mai multe rețele neurale care lucrează în paralel, decizia finală luându-se pe baza regulii *câștigătorul ia tot*. Numărul de rețele folosit este egal cu numărul de clase în care sunt grupați vectorii, iar instruirea este supervizată. Fiecare rețea este proiectată să recunoască corect vectori dintr-o singură clasă, astfel că răspunsurile cele mai bune apar doar atunci când îi sunt prezentați vectori

din clasa cu care a fost instruită. Acest model este de fapt un cadru care oferă flexibilitate arhitecturii fiindcă modulele pot fi reprezentate de diverse tipuri de rețele neurale.

Pornind de la modelul CNN din [11], Neagoe [60] a dezvoltat și experimentat modelul *Concurrent Self-Organizing Maps (CSOM)* care se detașează ca o tehnică nouă cu performanțe excelente. A fost demonstrată importanța CSOM pentru cazul general al vectorilor cu repartiții normale multimodale și s-a arătat că poate fi folosit cu succes pentru clase largi de probleme de recunoaștere a formelor. CSOM a fost aplicat cu succes de Neagoe [60] pentru recunoașterea fețelor umane și de Neagoe și Ropot în clasificarea imaginilor satelitare [61], în aplicații de clasificare a formelor [62], pentru recunoașterea fețelor umane și la recunoașterea vorbitorilor [64].

Ideea competiției dintre hărțile cu autoorganizare a fost utilizată ulterior și de Lapidot, Guterman și Cohen în [51] pentru recunoașterea nesupervizată a vorbitorilor. Fiecare vorbitor este modelat printr-o hartă, în timp ce un algoritm iterativ permite ajustarea hărților.

Acest capitol prezintă noul model al rețelelor neurale concurente și analizează performanțele sale pentru recunoașterea vorbitorilor și a vocalelor.

## 4.2 Rețele neurale concurente

Schema generală folosită pentru instruirea rețelelor neurale concurente este cea din figura 4.1.

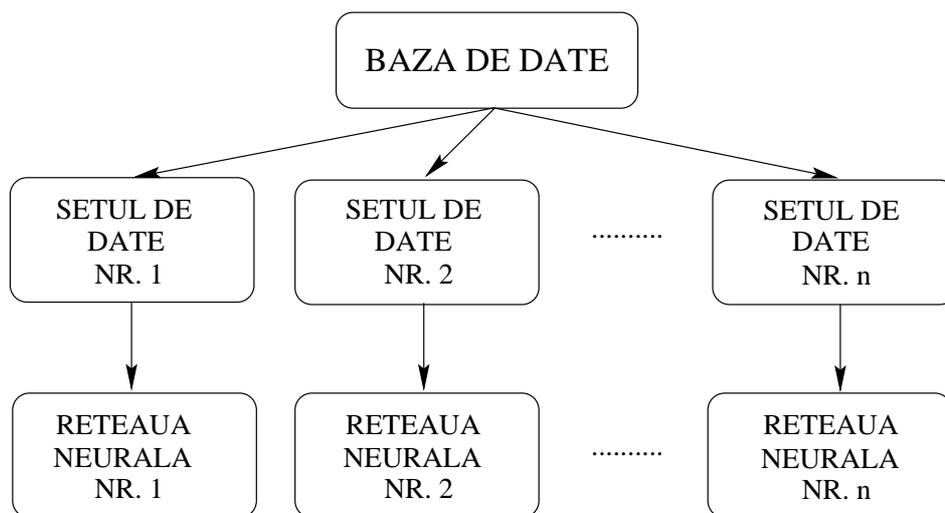


Figura 4.1: Modelul folosit la instruirea rețelelor neurale concurente.

În această schemă,  $n$  reprezintă numărul de rețele neurale care lucrează în paralel, dar este egal totodată și cu numărul de clase în care sunt grupați vectorii de instruire.

Considerăm mulțimea  $X \in R^p$  formată din  $M$  vectori  $p$ -dimensionali, adică:

$$X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M\}, \mathbf{x}_i \in R^p, 1 \leq i \leq M.$$

Presupunem că pentru fiecare vector din mulțimea  $X$  se cunoaște apartenența la una dintre cele  $n$  clase, unde:

$$X = X_1 \cup X_2 \cup \dots \cup X_n$$

și

$$X_1 \cap X_2 \cap \dots \cap X_n = \Phi,$$

iar  $X_1, X_2, \dots, X_n$  sunt submulțimi ale mulțimii  $X$  și conțin vectorii de intrare pentru fiecare dintre cele  $n$  rețele neurale.

Mulțimea  $X$  de vectori este obținută în urma preprocesării semnalelor vocale achiziționate în scopul instruirii rețelelor. Din această mulțime sunt extrase seturile de vectori  $X_j$ ,  $j = 1, 2 \dots n$  cu care vor fi instruite cele  $n$  rețele neurale. În urma aplicării procedurii de învățare, fiecare rețea neurală va trebui să răspundă pozitiv unei singure clase de vectori și să dea răspunsuri negative pentru toți ceilalți vectori. Algoritmul de instruire pentru rețeaua concurentă este următorul:

**Etapa 1.** *Se creează baza de date. Aceasta conține vectorii de instruire obținuți în urma preprocesării semnalului vocal.*

**Etapa 2.** *Se extrag din baza de date seturile de vectori specifice fiecărei rețele neurale în parte. Dacă este cazul, se stabilesc ieșirile dorite.*

**Etapa 3.** *Se aplică algoritmul de instruire fiecărei rețele neurale folosind seturile de vectori create la pasul 2.*

Etapa de recunoaștere se desfășoară în paralel, conform schemei din figura 4.2.

Cele  $n$  rețele neurale se presupune că au fost instruite prin algoritmul descris mai înainte. La aplicarea vectorului de test, rețelele generează câte un răspuns individual, iar selecția constă în alegerea rețelei care a generat răspunsul cel mai puternic. Rețeaua selectată prin regula *câștigătorul ia tot* este declarată câștigătoare. Indicele rețelei câștigătoare va fi indicele clasei în care este plasat vectorul de test. Această metodă de recunoaștere a formelor presupune, așadar,

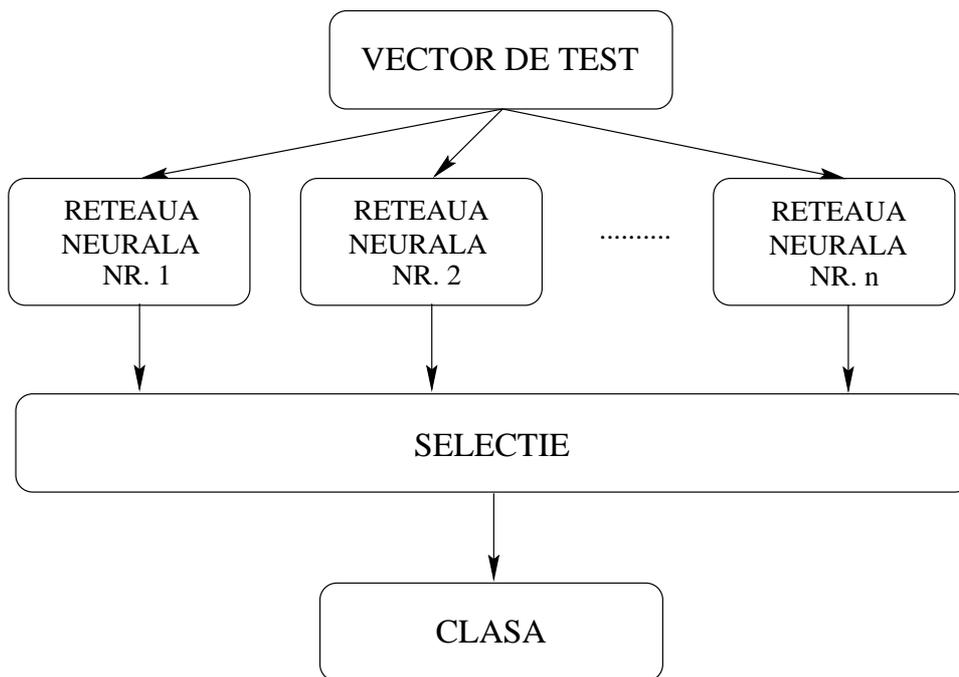


Figura 4.2: Modelul folosit la recunoaștere de rețelele neurale concurente.

că este *a priori* cunoscut numărul de clase cu care va lucra rețeaua concurentă și că pentru fiecare clasă există suficienți vectori de instruire. Algoritmul de recunoaștere este cel de mai jos:

- Etapa 1.** *Este creat vectorul de test prin preprocesarea semnalului vocal.*
- Etapa 2.** *Vectorul de test este transmis în paralel tuturor celor  $n$  rețele neurale instruite în prealabil.*
- Etapa 3.** *Blocul de selecție stabilește indicele rețelei cu răspunsul cel mai bun. Acesta va fi indicele clasei în care este încadrat vectorul.*

Modelele din figurile 4.1 și 4.2 pot fi particularizate plasând în locul celor  $n$  rețele neurale diferite arhitecturi. În experimentele prezentate în această lucrare am folosit perceptroni multistrat (MLP), rețele neurale cu întârziere în timp (TDNN) și hărți Kohonen (SOM), obținând astfel trei tipuri diferite de rețele neurale concurente.

Pentru cele care folosesc MLP și TDNN se utilizează algoritmul de instruire *backpropagation of error*. Atunci când implementăm rețele neurale concurente cu SOM, abordăm instruirea ușor diferit din cauză că instruirea modulelor

este nesupervizată. De aceea, diferă modul în care alcătuim seturile de vectori pentru instruire, dar și maniera în care se face selecția rețelei câștigătoare.

Pentru primele două tipuri de rețea concurentă, MLP-CNN și TDNN-CNN, am folosit seturi de vectori de instruire cu exemple pozitive și negative. Un *exemplu pozitiv* pentru o rețea componentă a CNN este un vector din clasa asociată ei. Un *exemplu negativ* este un vector care aparține altei clase. Pentru a ne asigura că rețelele vor fi capabile să facă diferența între exemplele pozitive și cele negative, am echilibrat numeric seturile de vectori: am inclus un număr egal de vectori pozitivi și negativi, iar exemplele negative au fost selectate din toate clasele în mod egal.

MLP și TDNN sunt rețele neurale care sunt instruite supervizat, motiv pentru care fiecărui exemplu pozitiv folosit la instruire i se asociază câte un răspuns dorit. Acesta este un vector  $n$ -dimensional,  $n$  fiind și numărul de clase. Fiecare componentă a sa va avea valoarea 1 sau 0. Există o singură componentă 1, plasată pe poziția care are ca indice numărul clasei din care face parte exemplul pozitiv. Pentru exemplele negative, răspunsul dorit este complementul binar al vectorului corespunzător exemplului pozitiv.

Decizia la recunoaștere se ia pe baza regulii *câștigătorul ia tot*. Fiecare dintre cele  $n$  rețele concurente produce în mod independent câte un răspuns la stimulul de test. Se selectează vectorul care conține exact o componentă apropiată de valoarea 1 și toate celelalte componente apropiate de valoarea 0, în limitele unei erori acceptabile. Pentru ca răspunsul să fie considerat valid, poziția componentei maxime în vectorul de ieșire trebuie să coincidă cu indicele rețelei care l-a generat.

Rețeaua concurentă care are la bază SOM este instruită cu seturi de vectori alcătuite doar din exemple pozitive. Ca și în celelalte cazuri, fiecare rețea componentă este instruită separat, cu setul de vectori special alcătuit pentru ea.

Clasificarea se face prin selecția hărții care generează o eroare de cuantizare minimă.

### 4.3 Recunoașterea vorbitorilor cu rețele neurale concurente. Experimente

Experimentele descrise în această secțiune [11], [13] compară capacitatea de recunoaștere a MLP, TDNN, SOM și a rețelelor neurale concurente care folosesc aceste tipuri de rețele neurale drept module de bază. Rezultatele prezentate aici au fost obținute în urma testelor de recunoaștere a vorbitorilor folosind mai întâi

baza de date Speechdata, iar apoi baza de date Telephdata care conțin cuvinte pronunțate de 25 de vorbitori pe care i-am organizat în 25 de clase.

Scopul testelor este acela de a recunoaște cei 25 de vorbitori pe baza diverselor pronunții ale aceluiași cuvinte, diferite de cele folosite la instruire. Am adoptat ca bază de comparație rezultatele obținute cu MLP, TDNN și SOM descrise în secțiunea 3.3.2 și care se găsesc și în tabelele 4.1, 4.2 și 4.3 pe coloanele marcate cu litera B (model neural de bază).

Modelul concurent ales de noi grupează 25 de rețele clasice, câte una pentru fiecare vorbitor. Metoda de instruire este descrisă în continuare, fiind supervizată pentru MLP și TDNN și nesupervizată pentru SOM. Etapa de recunoaștere impune alegerea celei mai performante rețele care este declarată câștigătoare. Setul de experimente descrise în acest subcapitol testează modelul rețelelor neurale concurente pentru recunoașterea vorbitorilor.

Tabelul 4.1: Rezultatele experimentelor de recunoaștere a vorbitorilor cu MLP și MLP-CNN.

| MLP        |        |            |        |
|------------|--------|------------|--------|
| Speechdata |        | Telephdata |        |
| B          | CNN    | B          | CNN    |
| 71,25%     | 86,22% | 6,75%      | 72,86% |

Tabelul 4.2: Rezultatele experimentelor de recunoaștere a vorbitorilor cu TDNN și TDNN-CNN.

| TDNN       |        |            |        |
|------------|--------|------------|--------|
| Speechdata |        | Telephdata |        |
| B          | CNN    | B          | CNN    |
| 57,28%     | 89,31% | 32,33%     | 72,72% |

Tabelul 4.3: Rezultatele experimentelor de recunoaștere a vorbitorilor cu SOM și SOM-CNN.

| SOM        |        |            |        |
|------------|--------|------------|--------|
| Speechdata |        | Telephdata |        |
| B          | CNN    | B          | CNN    |
| 52,86%     | 90,42% | 44,72%     | 75,25% |

Am folosit mai întâi MLP și TDNN drept componente de bază ale CNN.

Fiecărei clase i-am asociat câte una din cele 25 de rețele neurale compuse dintr-un strat de intrare cu 270 de neuroni, două straturi ascunse și un strat de ieșire alcătuit din 25 de neuroni. Am folosit algoritmul *backpropagation of error* pentru instruirea acestor rețele neurale feedforward și pentru fiecare dintre ele am creat seturi distincte de vectori de instruire. Astfel, vectorii din clasa 1, cele extrase din cuvintele pronunțate de vorbitorul 1, au fost exemple pozitive pentru rețeaua neurală cu indicele 1 și exemple negative pentru celelalte 24 de rețele neurale. Răspunsurile dorite asociate exemplilor pozitive au același format cu cel descris în secțiunea 3.3 unde am discutat despre instruirea MLP și TDNN ca rețele complexe care recunosc toți vorbitorii. Am asociat răspunsul dorit  $[1 \ 0 \ 0 \ \dots \ 0]^t$  exemplilor pozitive pentru rețeaua neurală 1 și răspunsul dorit  $[0 \ 1 \ 1 \ \dots \ 1]^t$  exemplilor negative pentru aceeași rețea. Pentru a echilibra numărul exemplilor pozitive și negative, am repetat de 24 de ori exemplele pozitive, făcând această operație pentru fiecare set de vectori. Urmărind această idee, am creat seturile de vectori pentru fiecare rețea neurală componentă a MLP-CNN și TDNN-CNN. În faza de recunoaștere, un vector de test trebuie să activeze doar o singură rețea neurală, cu alte cuvinte, acea rețea neurală activă trebuie să producă un răspuns apropiat de cel al unui exemplu pozitiv. Toate celelalte rețele neurale componente ale CNN trebuie să genereze un răspuns apropiat de răspunsul unui exemplu negativ. Luăm o decizie asupra recunoașterii numai atunci când exact o rețea neurală dintre cele 25 generează un răspuns pozitiv și toate celelalte generează răspunsuri negative. În caz contrar, nu se poate stabili apartenența vectorului de intrare la nici o clasă.

Figura 4.3 prezintă mediile ratelor de recunoaștere ale fiecărei rețele neurale componentă a CNN pentru 25 de vorbitori în două experimente care ilustrează folosirea TDNN și TDNN-CNN cu baza de date Telephdata. Este evidentă diferența dintre mediile ratelor de recunoaștere deși, așa cum se observă, este posibil ca TDNN să dea rezultate mai bune pentru anumiți vorbitori.

Instruirea SOM-CNN a fost mult mai simplă. Am creat 25 de seturi de vectori și am folosit algoritmul nesupervizat de instruire a rețelelor Kohonen pentru instruirea independentă a fiecăreia dintre cele 25 de rețele Kohonen. De această dată, am folosit doar exemple pozitive. La recunoaștere, vectorul de test este aplicat în paralel tuturor SOM componente ale SOM-CNN. Rețeaua în care apare cea mai mică eroare de cuantizare este declarată câștigătoare, iar indicele său este indicele clasei din care se presupune că face parte vectorul de test.

În tabelele 4.1, 4.2 și 4.3 sunt prezentate ratele de recunoaștere pentru MLP-CNN, TDNN-CNN și SOM-CNN, în coloanele marcate cu CNN.

Rezultatele obținute la recunoașterea vorbitorilor cu rețele neurale con-

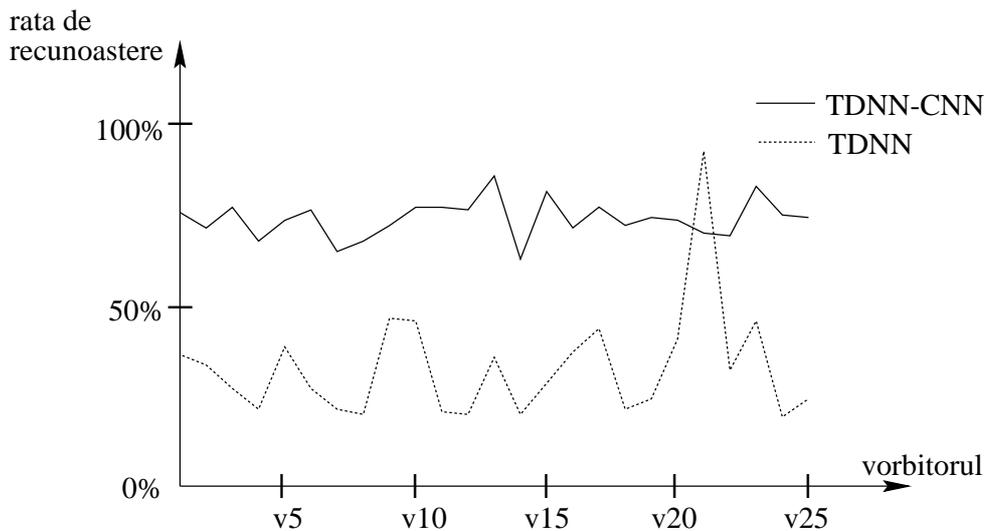


Figura 4.3: Mediile ratelor de recunoaștere pentru cei 25 de vorbitori în două experimente care ilustrează folosirea TDNN și TDNN-CNN cu Telephdata.

curente sunt mult îmbunătățite față de modelele clasice. MLP s-a dovedit a fi sensibil mai bun decât TDNN și SOM în testele cu modelul clasic folosind Speechdata, dar mult mai slab pentru Telephdata. Folosirea rețelelor concurente a îmbunătățit semnificativ ratele de recunoaștere. Modelul SOM-CNN s-a dovedit a fi cel mai performant, rezultatele fiind cu mai mult de 4% mai bune decât modelul MLP-CNN și cu aproximativ 1% mai bune decât modelul TDNN-CNN. În general, s-a păstrat o diferență de aproximativ 15% între testele cu Speechdata și cele cu Telephdata. Experimentele au arătat că folosind câte o rețea specializată pentru fiecare cuvânt se obțin rezultate mai bune la recunoaștere, chiar dacă faza de instruire durează în mod evident mai mult.

Este cunoscut faptul că algoritmul *backpropagation of error* are o convergență lentă, iar atunci când condiția de oprire înseamnă o eroare de clasificare mică, instruirea poate să fie foarte costisitoare din punct de vedere al timpului alocat. Când instruiem un set de rețele neurale concurente, timpul de instruire este proporțional cu numărul de module. Odată instruite aceste rețele neurale, timpul de calcul pentru stabilirea apartenenței unui vector de intrare la o clasă devine foarte scurt deoarece constă dintr-o singură propagare a informație dinspre intrări înspre ieșiri. În experimentele noastre, am alocat perioade până la 24 de ore pentru o instruire foarte fină a tuturor modulelor. Algoritmul de instruire a hărților cu autoorganizare este mai rapid deoarece condiția de terminare este dată de numărul de pași de execuție. Faza de recunoaștere este, însă, mai lentă decât pentru rețelele neurale feedforward. Pentru stabilirea erorii minime de cuantizare, aceasta presupune compararea fiecărui vector de

test cu toate prototipurile.

Pentru experimentele de recunoaștere a vorbitorilor cu MLP-CNN și TDNN-CNN am folosit module instruite prin programul Stuttgart Neural Network Simulator (SNNS). Rezultatele modulelor au fost combinate printr-o aplicație pe care am scris-o în limbajul de programare oferit de SNNS. Pentru testele de recunoaștere cu SOM-CNN am instruit hărțile cu ajutorul programelor din pachetul SOM\_PAK, iar pentru combinarea rezultatelor am realizat o aplicație în limbajul de programare C++.

## 4.4 Recunoașterea vocalelor cu rețele neurale concurente

După ce am văzut că modelul CNN aduce îmbunătățiri remarcabile ratelor de recunoaștere a vorbitorilor în comparație cu modelele clasice, în această secțiune ne vom concentra asupra recunoașterii unor vocale din limba română având la dispoziție diverse variante de pronunție a cifrelor de la 0 la 9. Înainte de a analiza metoda implementată și rezultatele acestor experimente, vom face o prezentare a principalelor noțiuni de fonetică și fonologie specifice limbii române [3].

### 4.4.1 Fonetică și fonologie în limba română

Vorbirea este un mijloc de transmitere orală a mesajelor și presupune existența unui sistem de sunete și a unui set de reguli de asociere a lor. Sunetele limbii alcătuiesc unitățile de bază ale acesteia, fiind acceptate în limitele unor variante și varietăți admisibile.

#### Elemente de fonetică

*Sunetele fonice articulate* sunt unitățile minimale purtătoare de informație ale unui limbaj vorbit, stând la baza comunicării dintre oameni. Ele constau în vibrații ale aerului, modulate de organele vorbirii. Sunt caracterizate prin tărie, înălțime, timbru, durată, trăsături perceptibile de urechea umană. *Fonetica* este disciplina care studiază producerea, transmiterea și receptarea sunetelor limbajului articulat, adică a sunetelor vorbite.

Sunetele fonice articulate sunt produse la intrarea în vibrație a aerului expirat. La acest proces participă plămânii, traheea, laringele, coardele vocale, glota, epiglota, faringele, maxilarele, gura, buzele, dinții, limba, fosele nazale. Ansamblul organelor care participă la producerea sunetelor se numește *aparat fonoarticulator*.

Articularea este produsă în timpul expirației, adică în momentul evacuării aerului din plămâni. În acest fel ia naștere fluxul de aer care produce unda fonică articulată. Această cantitate de aer trece printr-o serie de organe și cavități care, în principal, participă la funcția biologică a respirației și care formează *aparatură vorbirii*. În interiorul laringelui se află două lamele formate din țesuturi elastice și mușchi numite *coarde vocale* care delimitează în interiorul lor un spațiu numit *glota*. Coardele vocale realizează mișcări de închidere și deschidere a glotei. În funcție de acest grad de închidere, dar și de viteza cu care aerul este evacuat din plămâni, apar sunete mai înalte sau mai joase și de amplitudini diferite. Glota este un organ complex al vorbirii care produce prima diferențiere a sunetelor. Astfel, sunetele vor fi *vocoide* sau *contoide sonore* ori *surde*.

*Faringele, cavitatea bucală și cea nazală* sunt rezonatori în care sunetul este articulat, modulat și amplificat. *Uvula* dirijează fluxul fonic înspre cavitatea bucală sau nazală determinând *rezonanța bucală* sau *rezonanța nazală*. *Limba* este organul cu rolul cel mai important în articulare, iar *buzele* formează un ultim rezonator de îngustare sau ocluziune a fluxului fonic.

*Articularea* reprezintă totalitatea mișcărilor organelor vorbirii care determină forma și mișcările rezonatorilor, cu influențe asupra formării, modificării și amplificării fluxului fonic. Articularea are trei componente:

- *modul de articulare*, determinat de forma și mișcările rezonatorilor;
- *locul de articulare*, adică organul articulării;
- *prezența sau absența vibrațiilor* glotale.

Caracteristica acestor componente este că ele pot apărea doar în combinație, și nu izolate una de alta. Prin combinația lor se pot defini tipurile articulatorii care pot apărea în limba română vorbită.

*Vocalele (vocoidele)* sunt sunete produse de coloana de aer vibrat prin trecerea ei liberă prin traiectul vocal. Diversele vocale rezultă prin gradul de deschidere a rezonatorului bucal combinat cu mișcările verticale și orizontale ale limbii. Ca mod de articulare, vocalele pot fi *deschise* (*a*), *semî închise* (*e*, *ă*, *o*) și *închise* (*i*, *î*, *u*), în funcție de nivelul de apropiere a limbii de palatul moale și de gradul de deschidere a maxilarelor. După locul de articulare, ele pot fi *posterioare* (*u*, *o*), *centrale* (*ă*, *î*) sau *anterioare* (*e*, *i*), iar *a* este o vocală neutră. Pe lângă mușchiul lingual, un rol important îl joacă poziția buzelor, producând o vocală *rotunjită* (*o*, *u*) sau *nerotunjită* (*a*, *e*, *i*, *ă*, *î*). Conform acestor criterii, sistemul vocalic românesc se poate sintetiza ca în tabelul 4.4.

*Consoanele* sunt produse de fluxul de aer care întâlnește în calea sa diferite obstacole. Diferențele între ele se fac încă de la nivelul rezonatorului glotal,

Tabelul 4.4: Sistemul vocalic al limbii române.

| Poziția buzelor   |             | nerotunjite |          | rotunjite   |
|-------------------|-------------|-------------|----------|-------------|
| Poziția limbii    |             | anterioare  | centrale | posterioare |
| Deschiderea gurii | deschise    | a           |          |             |
|                   | semiînchise | e           | ă        | o           |
|                   | închise     | i           | î        | u           |

consoanele fiind sonore sau surde. Perechile de consoane clasificate după criteriul sonore – surde pot fi grupate în felul următor:  $b-p$ ,  $d-t$ ,  $v-f$ ,  $z-s$ ,  $ž-ts$ ,  $ğ-č$ ,  $g'-k'$ ,  $g-k$ . Toate aceste consoane sunt *nesonante*. În prezența vibrațiilor faringiene, apar consoanele *sonante*,  $l$ ,  $m$ ,  $n$ ,  $r$ . După locul de articulare, consoanele se clasifică astfel: *bilabiale*:  $b$ ,  $p$ ,  $m$  (buze), *labiodentale*:  $v$ ,  $f$  (buza inferioară și dinții superiori), *dentale*:  $d$ ,  $t$ ,  $z$ ,  $s$ ,  $ts$ ,  $n$ ,  $l$ ,  $r$  (dinții), *prepalatale*:  $ž$ ,  $š$ ,  $ğ$ ,  $č$  (palatul anterior), *palatale*:  $g'$ ,  $k'$  (palatul mediu), *velare*:  $g$ ,  $k$  (palatul moale) și *laringala*  $h$  (laringele). La producerea unor constante, palatul moale permite trecerea liberă a unei părți din fluxul de aer fonator prin fosele nazale, producându-se consoanele *nazale*:  $m$ ,  $n$ . Alteori, traiectul vocal este închis complet. Prin deschiderea lui bruscă se produc consoanele *occlusive (explozive)*:  $b$ ,  $p$ ,  $m$ ,  $d$ ,  $t$ ,  $n$ ,  $g'$ ,  $k'$ ,  $g$ ,  $k$ . Când traiectul vocal este doar îngustat și aerul îl parcurge tot timpul cât durează emisia, se produc consoanele *constrictive*. Acestea sunt *fricative*:  $v$ ,  $f$ ,  $z$ ,  $s$ ,  $ž$ ,  $š$ , când îngustarea este mai pronunțată, sau *lichide*:  $l$ ,  $r$ .

Consoana  $r$  este *vibrantă* pentru că vârful limbii vibrează, iar  $l$  este *laterală* deoarece aerul se scurge de-a lungul marginilor limbii. Când traiectul vocal este inițial închis complet și se deschide treptat, se produc consoanele *semioclusive (africate)*:  $ts$ ,  $č$ ,  $ğ$ . Sistemul consonantic al limbii române este alcătuit din 22 de consoane care se clasifică așa cum se vede în tabelul 4.5.

*Semivocalele* sunt articulate ca și vocalele, dar durata lor este mai mică și efortul respirator mai mare. Limba română are următoarele semivocale:  $e$ ,  $i$ ,  $o$ ,  $u$ .

## Elemente de fonologie

*Fonemele* sunt unități minimale ale sistemului limbii, caracterizate prin trăsături distinctive a căror modificare ar altera chiar sistemul limbii [33].

Tabelul 4.5: Sistemul consonantic al limbii române.

| Vibrațiile glotale  |              | sonante         |              |          | nesonante    |              |              |
|---------------------|--------------|-----------------|--------------|----------|--------------|--------------|--------------|
|                     |              |                 |              |          | surde/sonore | surde/sonore | surde/sonore |
| Modul de articulare |              | oclusive nazale | constrictive |          | explozive    | fricative    | africate     |
|                     |              |                 | vibrante     | laterale |              |              |              |
| Locul de articulare | bilabiale    | m               | -            | -        | b p          | --           | --           |
|                     | labiodentale | -               | -            | -        | --           | v f          | --           |
|                     | dentale      | n               | r            | l        | d t          | z s          | - ts         |
|                     | prepalatale  | -               | -            | -        | --           | ž š          | ğ ċ          |
|                     | palatale     | -               | -            | -        | g' k'        | --           | --           |
|                     | velare       | -               | -            | -        | g k          | --           | --           |
|                     | laringale    | -               | -            | -        | --           | - h          | --           |

Fonemul este o entitate care nu poate fi analizată în unități mai mici și are funcția de a diferenția între ele cuvinte diferite sau forme gramaticale ale aceluiași cuvânt. *Fonologia* studiază funcția fonemelor într-o limbă și stabilește sistemul de foneme al acelei limbi, realizând inventarierea lor.

Fonemul este o clasă de sunete asemănătoare, cu aceeași funcție. În limba vorbită, un fonem poate fi realizat prin mai multe *variante* care sunt sunete asemănătoare din punct de vedere al pronunției și al percepției, dar diferă datorită contextelor fonetice în care se găsesc. De asemenea, un fonem poate avea mai multe *varietăți* care sunt pronunții ale aceluiași sunet de către vorbitori diferiți [3]. Așadar, fonemul este un element abstract și cu caracter general, iar sunetul este o particularizare, o concretizare a fonemului.

Fonemul are două funcții: *distinctivă* și *contrastivă*. Prima funcție, cea distinctivă, realizează distingerea cuvintelor și a formelor lor (**dar**–**far**–**var**–**zar**). Există mai multe grupe de foneme care *comută*: *d, f, v, z; p, r, t; a, o, u; u, i; k, ȇ*. Fiecare fonem are propriile trăsături relevante legate de modul și locul de articulare, dar și de prezența sau absența vibrațiilor glotale. Aceste trăsături le deosebesc unele de altele și creează opozițiile între foneme. Cea de-a doua funcție, cea contrastivă permite stabilirea structurii fonematice a unui cuvânt, facilitând înțelegerea lui.

Cele șapte vocale din limba română pot constitui centrul de sonoritate al unei silabe sau pot alcătui singure o silabă, adică sunt *silabice*. Semivocalele se definesc drept vocale nesilabice și apar în structura *diftongilor* sau a *triftongilor*. Diftongii sunt formați dintr-o vocală și o semivocală, iar triftongii din două semivocale și o vocală. Diftongii și triftongii sunt segmente vocalice constituite din două, respectiv trei vocale din care doar una se realizează ca vocală. Consoanele din sistemul consonantic al limbii române sunt *nesilabice*, fiind dependente de o vocală.

*Silaba* este construcția superioară fonemului și este caracterizată printr-un singur accent. Ea este formată dintr-un singur fonem care trebuie să fie neapărat o vocală sau din mai multe foneme între care se găsește o vocală.

*Accentul* evidențiază o silabă printr-o pronunțare mai intensă și este în contrast cu silabele neaccentuate.

#### 4.4.2 Experimente de recunoaștere a vocalelor

Abordăm problema recunoașterii vocalelor prin studierea similarității dintre vectorii prototip și vectorii obținuți prin analiza cepstrală a regiunilor vocalizate din cuvintele țintă. Prezentăm în această secțiune un sistem alcătuit din rețele neurale concurente pentru recunoașterea vocalelor care apar în pronunțiile în

limba română ale cifrelor de la 0 la 9.

### **Analiza cepstrală a semnalelor vocale**

Cepstrul se obține printr-o transformare complexă care se poate aplica și asupra semnalului vocal. Această transformare produce separarea unor componente ale semnalului care apoi sunt combinate liniar. Pentru semnalul:

$$s(n), n = 0, 1, 2, \dots,$$

transformata sa va fi de forma:

$$c_s(n) = c_e(n) + c_{\Theta}(n), n = 0, 1, 2, \dots$$

În această relație,  $c_e(n)$  este partea din cepstru care reprezintă contribuția excitației periodice, iar  $c_{\Theta}(n)$  este partea corespunzătoare funcției de transfer a traiectului vocal. Ca urmare a acestei transformări, partea inferioară a cepstrului unui semnal vocal, pentru care  $n$  are valori mici, este asociată semnalului  $c_{\Theta}(n)$ , iar partea superioară, în care  $n$  este mare, corespunde semnalului  $c_e(n)$ .

### **Baza de date**

Am preprocesat prin analiza cepstrală bazată pe scala mel pronunțiile în limba română ale celor 10 cifre. Fiecare dintre cei 4 subiecți a pronunțat aceste cuvinte de câte 10 ori. De fiecare dată am extras vectorii de instruire a sistemului de recunoaștere din primele 5 repetiții, păstrând ultimele 5 repetiții pentru teste. Am avut astfel la dispoziție 200 de repetiții pentru antrenarea sistemului și alte 200 pentru verificarea sa. Pentru instruire, am segmentat manual fiecare repetiție și am selectat regiunile din semnal corespunzătoare vocalelor și semivocalelor. Extragerea vectorilor de caracteristici a fost realizată prin parcurgerea semnalelor selectate cu o fereastră de 256 de eșantioane și un pas de 50 de eșantioane. Din fiecare fereastră s-au calculat câte 12 coeficienți cepstrali care formează vectorul asociat ei. Numărul de ferestre și, implicit, numărul de vectori extrași din semnalul asociat unui fonem este influențat de lungimea acestuia în timp, adică de caracterul vocalic sau semivocalic și de debitul verbal, ca și de modul de pronunție specific fiecărui vorbitor. La nivelul întregii colecții de pronunții, numărul de vectori asociați unui fonem este direct proporțional cu frecvența de apariție a vocalelor în cele 10 cuvinte. Cuvintele au fost înregistrate în condiții normale, fără o amenajare specială a incintei, cu un microfon de uz general. Tabelul 4.6 conține o statistică a numărului de vectori pe care i-am obținut pentru fiecare fonem din colecția de semnale vocale descrisă mai sus. Figura 4.4 prezintă câte o secvență tipică de coeficienți

cepstrali calculați pentru fiecare dintre cele șase vocale care alcătuiesc baza de date pentru un vorbitor, iar în anexa A am listat valorile numerice pe baza cărora am trasat aceste grafice.

Se observă că în partea inferioară a cepstrului apar diferențe vizibile între secvențele asociate acestor foneme, confirmând astfel că traiectul vocal contribuie într-o măsură importantă la producerea fonemelor. În a doua jumătate a cepstrului, diferențele nu mai sunt atât de evidente, însă se remarcă un vârf datorat frecvenței fundamentale a semnalului excitator.

### **Energia semnalului vocal**

Analiza oricărui semnal se face într-un timp finit. Pentru aceasta trebuie să îi cunoaștem caracteristicile pe toată durata și proprietatea pe care o căutăm trebuie să rămână invariantă. Semnalul vocal se caracterizează prin faptul că nu este staționar pentru un timp infinit de lung. Proprietatea de staționaritate rămâne valabilă doar pentru câteva milisecunde. Acesta este motivul pentru care metodele de analiză "pe termen lung" a semnalelor sunt aplicate în contextul nostru pe o fereastră. Analiza pe fereastră presupune "decuparea" unei regiuni din semnal și prelucrarea sa în mod independent. Vorbirea este un proces dinamic, de aceea un singur cadru nu este întotdeauna suficient. Din necesitatea de a surprinde caracteristicile temporale ale semnalului vocal, obișnuim să folosim cadre extrase cu ajutorul unei ferestre care se deplasează de-a lungul acestuia.

Atunci când dorim să extragem informația pe care ne-o oferă cadrul care cuprinde porțiunea de semnal vocal din intervalul  $n = m - N + 1, \dots, m$  folosim o versiune adaptată a metodei care se folosește de regulă pentru semnale infinite în timp [16].

Să presupunem că dorim să știm dacă o secvență de semnal este vocalizată sau nevocalizată pe porțiunea care se termină la momentul  $n = m$ . Se cunoaște că semnalul vocalizat are, în general, o energie mai mare decât semnalul nevocalizat. Folosim această idee pentru a ne asista în luarea deciziei. Având două semnale  $s_1(n)$  vocalizat și  $s_2(n)$  nevocalizat de durată infinită, energia lor va fi în relația:

$$E_{s_1} > E_{s_2},$$

unde

$$E_s \triangleq \sum_{n=-\infty}^{\infty} s^2(n).$$

Pentru o fereastră de semnal care conține punctele din jurul lui  $m$  avem inega-

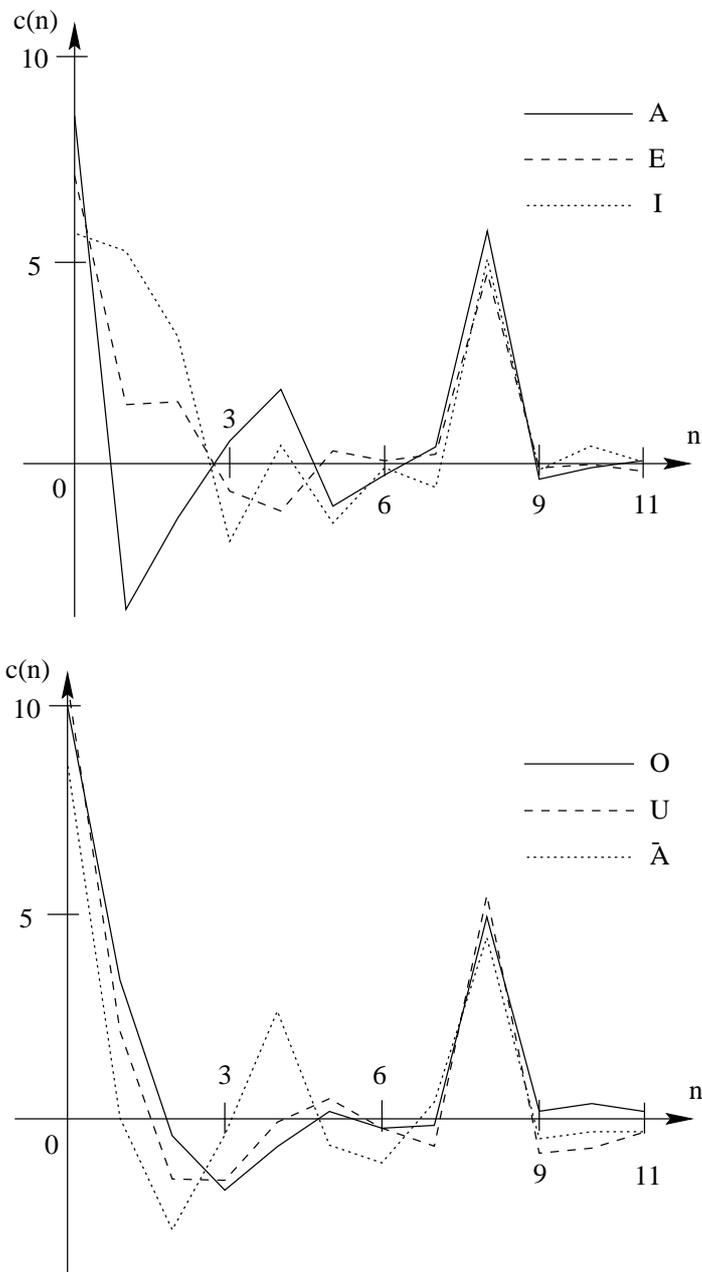


Figura 4.4: Secvențe tipice de coeficienți cepstrali asociate fonemelor A, E, I, O, U, ă.

litatea:

$$E_{s_1}(m) > E_{s_2}(m),$$

unde

$$E_s(m) = \sum_{n=m-N+1}^m s^2(n).$$

În această relație,  $N$  este numărul eșantioanelor din semnal care se regăsesc în interiorul ferestrei. Ca să putem automatiza luarea unei decizii privind caracterul vocalizat sau nevocalizat al semnalului, raportăm valoarea lui  $E_s(m)$  la un prag: dacă valoarea energiei depășește pragul semnalul este vocalizat, iar dacă valoarea sa este sub prag, semnalul este nevocalizat.

Un element important în extragerea unor informații pe baza energiei este fixarea numărului  $N$  de eșantioane cuprinse în fereastră. Atunci când fereastra are dimensiune mare, graficul energiei obținut prin baleierea semnalului este mai "neted". Aceasta înseamnă că în vecinătatea pragului decizia nu este influențată de "vârfuri" ale energiei. Pe de altă parte, o fereastră de dimensiuni mari poate provoca "netezirea" prea pronunțată a formei energiei. Drept consecință, fie zona vocalizată determinată prin calcul depășește dimensiunile reale, fie nu se mai pot detecta zonele nevocalizate de durată scurtă inserate între două zone vocalizate (fig. 4.5).

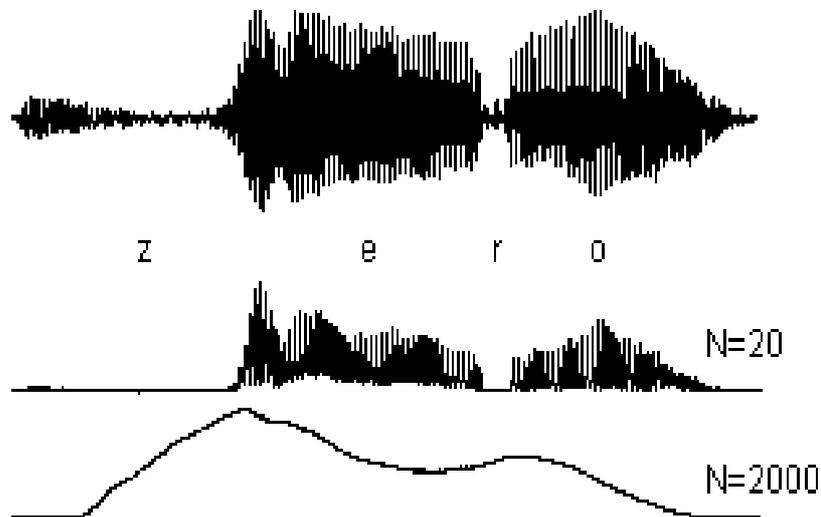


Figura 4.5: Funcția de energie a semnalului corespunzător unei pronunții a cifrei *zero* pentru  $N=20$  și  $N=2000$ .

Datele pe care le obținem prin studierea energiei au un caracter preponderent informativ și nu influențează în mod fundamental corectitudinea recunoașterii vocalelor din semnal. De aceea, în scopul simplificării calculului

folosim o fereastră dreptunghiulară a cărei dimensiune am stabilit-o experimental la  $N = 300$ . Deoarece am eşantionat semnalul vocal cu o frecvență de 8000 Hz, fereastra folosită la calculul energiei selectează o porțiune de 37,5 ms. Am ales o fereastră dreptunghiulară pentru că ne permite calculul recurent al energiei:

$$E_s(m) = E_s(m - 1) + s^2(m) - s^2(m - N).$$

Astfel, timpul de estimare a energiei semnalului este liniar, fiind proporțional doar cu dimensiunea semnalului vocal și nu depinde de dimensiunea ferestrei.

Pentru a lua o decizie în legătură cu caracterul vocalizat sau nevocalizat al semnalului într-o anumită regiune, am stabilit două praguri pentru energie: unul superior și unul inferior. Atunci când curba ascendentă a energiei depășește pragul superior, considerăm că începe o zonă vocalizată. Când pentru un segment vocalizat energia scade sub pragul inferior, considerăm că începe o zonă nevocalizată. Avantajele folosirii celor două praguri sunt, pe de o parte, o mai bună delimitare a începutului zonei vocalizate și, pe de altă parte, evitarea plasării în zona nevocalizată a vocalelor mai puțin accentuate de la sfârșitul cuvintelor.

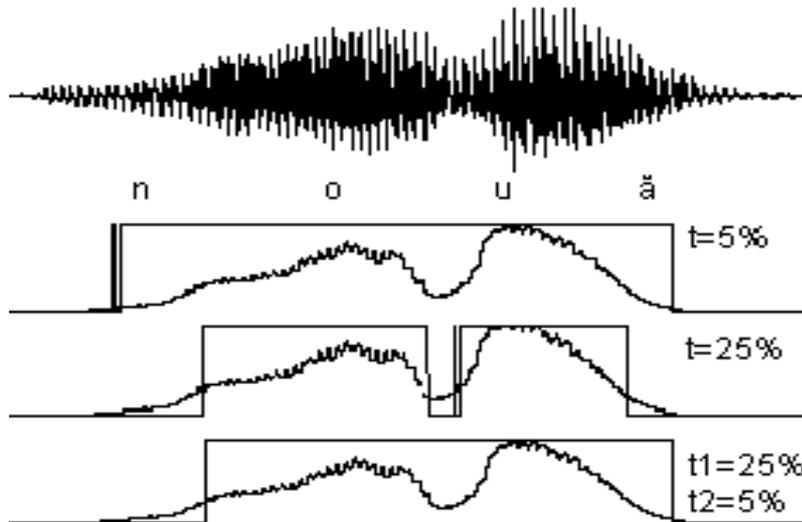


Figura 4.6: Stabilirea zonelor vocalizate și a celor nevocalizate pentru pronunția cifrei *nouă*. Valorile pragurilor sunt raportate la nivelul maxim atins de energie.

Figura 4.6 prezintă modul în care se decide dacă o zonă este vocalizată sau nu atunci când se folosește un prag și când se folosesc două praguri, pentru o pronunție a cifrei *nouă*. Atunci când pragul unic are valoare mică, regiunea vocalizată cuprinde și fonemul *ă*. Această zonă cuprinde în mod eronat și fonemul *n* care, deși este nevocalizat, are amplitudine mare în comparație cu cea a unei vocale. Un prag mare elimină fonemul *n*, dar ignoră și fonemul *ă*

care datorită poziției sale finale are o amplitudine mică. Pragul unic are în plus dezavantajul că zona vocalizată poate debuta și se poate încheia cu o oscilație, comportament eliminat de folosirea a două praguri.

Analiza semnalului pe baza energiei dă o primă informație legată de prezența sau absența zonelor vocalizate. Pe baza ei restrângem domeniul de căutare. A doua etapă constă în stabilirea identității fiecărei vocale în interiorul regiunilor vocalizate.

## Recunoașterea vocalelor

Ca și la recunoașterea vorbitorilor, am realizat un set de experimente folosind mai întâi o singură rețea neurală pentru recunoașterea tuturor vocalelor, iar apoi CNN. De data aceasta am testat doar rețeaua Kohonen, date fiind rezultatele pe care le-am obținut la recunoașterea vorbitorilor și i-am comparat comportamentul cu cel al unui SOM-CNN.

Pentru varianta care folosește un singur SOM am instruit rețeaua cu întreaga secvență de vectori obținuți după preprocesarea semnalelor vocale selectate. Am instruit în două etape o rețea Kohonen cu 10x15 noduri prin algoritmul SOFM descris în paragraful 2.9.2. Prima etapă, cea de organizare a clusterelor, s-a desfășurat de-a lungul a 1000 de pași și vecinătatea a scăzut treptat până la un singur neuron. În a doua etapă, instruirea s-a făcut în 10000 de pași, iar vecinătatea a rămas fixată la dimensiunea minimă. În urma instruirii și a calibrării rețelei neurale cu vectorii de instruire am obținut un set de prototipuri etichetate a cărui structură este cea din tabelul 4.6.

Tabelul 4.6: Numărul de vectori de instruire a SOM și SOM-CNN, numărul de prototipuri obținute în urma instruirii SOM și numărul de prototipuri obținute după instruirea SOM-CNN.

|           | A   | E   | I   | O   | U   | Ă  |
|-----------|-----|-----|-----|-----|-----|----|
| Instruire | 372 | 516 | 339 | 525 | 453 | 52 |
| SOM       | 42  | 37  | 16  | 21  | 21  | 2  |
| SOM-CNN   | 25  | 25  | 25  | 25  | 25  | 25 |

Se observă că fonemul *ă* este cel mai slab reprezentat din cauza volumului redus al datelor de instruire disponibile. Numărul total de prototipuri este mai mic decât cel al neuronilor din rețea pentru că nu toți neuronii au fost etichetați la calibrare. În figura 4.7 este prezentată proiecția U-matrix pentru această rețea Kohonen. *Unified distance matrix* (U-matrix) [47] este o metodă de vizualizare printr-o imagine care folosește nivele de gri a clusterelor generate

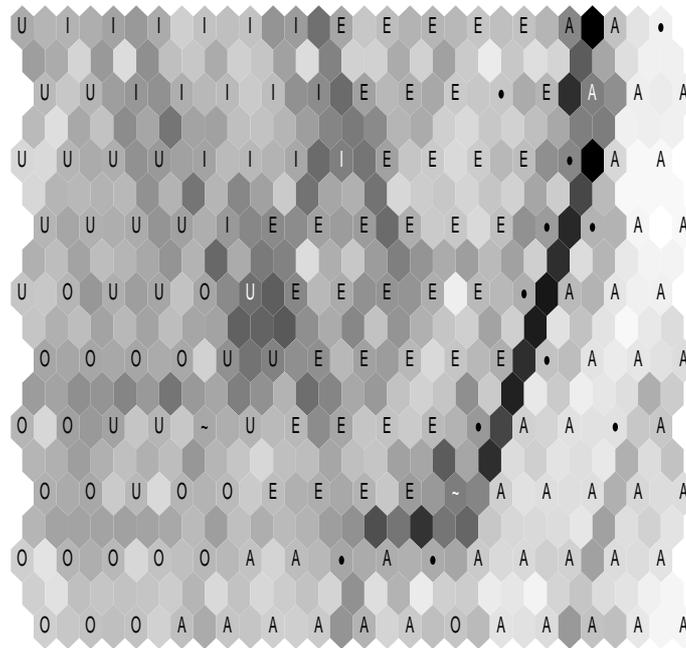


Figura 4.7: U-matrix pentru rețeaua Kohonen instruită să recunoască fonemele *a*, *e*, *i*, *o*, *u*, *ă*. Pentru fonemul *ă* am folosit simbolul  $\sim$ .

de o rețea Kohonen.

Și aici se observă slaba reprezentare a prototipurilor fonemului *ă*. Prototipurile celorlalte foneme formează grupuri compacte și poziția clusterelor reflectă calitățile acustico-fonetice ale fonemelor.

Instruirea componentelor SOM-CNN s-a făcut după gruparea manuală a vectorilor de instruire în șase clase. Fiecare rețea neurală componentă s-a specializat în recunoașterea unei singure clase. Am plasat în același grup atât realizările ca vocale cât și ca semivocale ale fonemelor. Am folosit șase rețele Kohonen cu câte 5x5 neuroni instruite de asemenea în două etape. Numărul total de neuroni al SOM-CNN este egal cu numărul de neuroni al SOM folosit la experimentul cu o singură rețea neurală. Aceasta înseamnă că numărul prototipurilor ar trebui să fie egal pentru ambele tipuri de experimente: 150. Diferența este minoră și apare datorită faptului că în primul caz nu toți neuronii au fost etichetați la calibrare, așa cum am arătat în tabelul 4.6.

Forma semnalului corespunzător uneia dintre cele 20 de pronunții disponibile ale cifrei *nouă* este desenată în figura 4.8. Am delimitat regiunea din semnal corespunzătoare fonemului *ă* pe care am folosit-o la instruire și din care am extras 4 vectori. În vorbirea curentă pronunția cuvântului *nouă* este uneori distorsionată și acesta se transformă în *noo*. În baza noastră de date am avut un astfel de exemplu pe care l-am eliminat și din acest motiv am folosit doar



în acea zonă se găsește doar de fonemul  $o$ . În vorbirea curentă, trecerea de la  $u$  la  $ă$  în cuvântul *nouă* se face prin ușoara inserare între ele a semivocalei  $o$ , fapt sesizat atât de SOM cât și de SOM-CNN.

Tabelul 4.7: Rezultatele experimentelor de recunoaștere a vocalelor și semivocalelor cu SOM și SOM-CNN.

|      | SOM    | SOM-CNN |
|------|--------|---------|
| E4.1 | 73,68% | 84,21%  |
| E4.2 | 62,47% | 73,93%  |
| E4.3 | 50%    | 56,25%  |
| E4.4 | 5,88%  | 5,48%   |
| E4.5 | 13,07% | 8,61%   |

## Rezultate experimentale

Pentru a compara performanțele SOM și SOM-CNN, în tabelul 4.7 sunt descrise rezultatele următoarelor experimente [12]:

- E4.1: stabilirea procentului vocalelor și semivocalelor recunoscute de cele două modele;
- E4.2: proporția în care dimensiunea reală a zonei vocalei sau a semivocalei este sesizată de SOM și SOM-CNN;
- E4.3: numărul de succesiuni vocală-semivocală recunoscute corect;
- E4.4: procentul anunțurilor false în legătură cu existența unei vocale sau semivocale într-o regiune din semnal în care în realitate sunt consoane;
- E4.5: procentul confuziilor dintre două vocale.

Experimentul 4.1 ne ajută să ne formăm o imagine corectă asupra capacității de a recunoaște noi vectori de către SOM și SOM-CNN. Progresul obținut cu ajutorul modelului nostru este de peste 10%. Am considerat că o vocală este recunoscută corect dacă secvența de etichete este cea reală pe zona centrală a segmentului vocalic, așa cum se poate vedea și în figura 4.9.

Fenomenul coarticulației influențează recunoașterea în zonele de tranziție dintre două foneme. Prin experimentul 4.2 stabilim în ce măsură sistemele testate de noi determină dimensiunile corecte ale vocalelor și semivocalelor.

Atunci când apar în diftongi, vocalele sunt pronunțate cu mai multă forță și tind să domine semivocalele, așa cum se întâmplă și pentru cuvântul *doi*.

Experimentul 4.3 stabilește procentul de astfel de succesiuni pentru care sunt recunoscute atât vocala cât și semivocala.

Experimentele 4.4 și 4.5 determină capacitatea SOM și a SOM-CNN de a evita confuziile. Energia semnalului vocal dă o primă informație referitoare la caracterul vocalizat sau nevocalizat al unei porțiuni din semnal. Rețeaua neurală stabilește pentru segmentele vocalizate care este succesiunea fonemelor. Avem, astfel, două nivele de decizie a existenței vocalelor. În experimentul 4.4 am calculat rata confuziilor între vocale și consoane.

În experimentul 4.5 am stabilit rata confuziilor între vocale, determinând care este procentul în care o vocală este înlocuită de sistemul de recunoaștere cu altă vocală.

## 4.5 Concluzii

Rezultatele obținute în urma experimentelor de recunoaștere a vorbitorilor folosind CNN au fost de fiecare dată superioare celor obținute folosind modelele clasice de rețele neurale. Pe de altă parte, rezultatele SOM-CNN au fost mai bune cu aproximativ 4% față de cele ale MLP-CNN și cu 1% față de TDNN-CNN.

Am testat SOM și SOM-CNN pentru recunoașterea unor vocale din limba română și am constatat că procentul vocalelor și semivocalelor recunoscute de SOM-CNN este superior celui obținut cu SOM cu mai mult de 10%. Regiunile din semnal corespunzătoare vocalelor și semivocalelor sunt mai cuprinzătoare atunci când se folosește SOM-CNN decât atunci când se folosește SOM.

Am văzut în acest capitol că rețelele neurale concurente implementează conceptul de concurență la nivelul unei colecții de rețele neurale. Experimentele demonstrează că modelul rețelelor neurale concurente este superior în majoritatea cazurilor folosirii unei singure rețele neurale pentru rezolvarea aceleiași probleme.

În următoarele două capitole vom atinge a doua problemă importantă a tezei și vom introduce o serie de algoritmi prin care se poate determina importanța intrărilor pentru rețelele neurale care folosesc algoritmi de instruire tip LVQ.

## Capitolul 5

# Recunoaștere prin calculul relevanțelor intrărilor folosind operatori OWA

*Algoritmul Relevance Learning Vector Quantization (RLVQ) [10] este o variantă a algoritmului LVQ care permite determinarea printr-o metodă euristică a relevanțelor pentru componentele intrărilor. Metoda se bazează pe învățarea hebbiană și introduce factorii de ponderare a intrărilor, adaptați automat la problema de clasificare care trebuie implementată. Pe lângă stabilirea importanței intrărilor, această metodă conduce la îmbunătățirea performanțelor algoritmului LVQ dar, în același timp, poate fi folosită pentru reducerea numărului de intrări prin renunțarea la cele mai puțin influente.*

*În acest capitol propunem o nouă metodă de calcul al relevanțelor intrărilor în RLVQ. Relevanțele sunt calculate on-line ca ponderi ale unui operator Ordered Weighted Aggregation (OWA). Operatorii OWA reprezintă o familie de operatori de agregare și au fost definiți de Yager [93]. Principalul avantaj al acestui algoritm numit OWA-RLVQ este acela că face legătura dintre RLVQ și consistentul model matematic OWA.*

### 5.1 Introducere

Kohonen a introdus LVQ ca o metodă care definește o clasificare pe baza unui număr de vectori grupați într-un set de instruire. Așa cum am văzut deja, cuantizarea vectorială înseamnă înlocuirea vectorilor de instruire cu un număr redus de prototipuri care formează un *codebook*. Fiecare vector din *codebook* are atașată eticheta clasei din care face parte. Algoritmul LVQ adaptează iterativ

acești vectori prin optimizarea unui criteriu global bazat pe distanța euclidiană.

Algoritmul LVQ standard nu face distincția între componentele vectoriale mai influente și cele mai puțin influente, importanța lor fiind considerată egală. Spre deosebire de această abordare, algoritmul *Distinction Sensitive Learning Vector Quantizer (DSLQV)* introdus în [68] folosește o pondere modificabilă pentru fiecare componentă vectorială și folosește o funcție de distanță ponderată pentru clasificare. Se utilizează un algoritm euristic iterativ pentru adaptarea acestor ponderi la specificul problemei: influența componentelor care în mod frecvent generează o clasificare greșită este diminuată, în timp ce influența componentelor care contribuie la clasificarea corectă este amplificată. Procesul de ponderare poate fi privit ca o transformare prin scalare a spațiului vectorilor de instruire într-un spațiu bazat pe o distanță ponderată: crește distanța dintre valorile componentelor cele mai puternic discriminative și o micșorează pe cea dintre valorile componentelor care aduc mai puțină informație. Aceasta facilitează distincția dintre clase și face sistemul mai puțin sensibil la zgomot [69].

Algoritmul DSLQV folosește o funcție de distanță ponderată care se bazează pe distanța euclidiană ponderată. Se pot folosi, însă, și alte distanțe, după cum se precizează în [68]. Modificarea algoritmului LVQ este minimală, fiind înlocuită funcția distanței. Ponderile și vectorii codebook se modifică în paralel în timpul instruirii, actualizarea având loc on-line, după procesarea fiecărui vector de intrare. Ponderile rezultate pot fi interpretate ca ranguri ale componentelor vectoriale și folosite pentru reducerea dimensiunii spațiului de intrare.

RLVQ este o variantă modificată a LVQ, similară cu DSLQV, și introduce factorii de relevanță pentru pentru fiecare dimensiune a vectorilor de intrare. Metrica euclidiană este înlocuită cu una modificată prin asocierea unei ponderi fiecărei componente vectoriale. Adicional algoritmului standard, factorii de scalare sunt adaptați iterativ printr-o tehnică bazată pe învățarea hebbiană, oferind o imagine a influenței fiecărei componente vectoriale asupra procesului de clasificare.

Operatorii OWA reprezintă o clasă de operatori de agregare care realizează o agregare bazată pe reordonarea criteriilor care trebuie satisfăcute. Ponderile operatorului OWA sunt asociate pozițiilor acestor criterii obținute în urma reordonării, și nu unui criteriu anume.

Noul algoritm de calcul al relevanțelor pe care îl introducem în acest capitol calculează importanțele componentelor de intrare considerându-le ponderi ale unui operator OWA. Ponderile sunt adaptate on-line, în paralel cu instruirea prototipurilor. Metrica pe care o folosim este un caz particular al mediei generalizate de tip OWA [37]. Pe lângă recunoașterea eficientă a formelor, metoda

noastră face și o conexiune interesantă între două abordări diferite: RLVQ și OWA.

Experimentele prezentate în finalul capitolului demonstrează capacitatea tehnicii propuse de noi de a rezolva probleme dificile, iar rezultatele confirmă acest lucru. Am folosit trei baze de date care reprezintă standarde în clasificarea formelor. Cu bazele de date Iris și Ionosphere am testat performanțele generale în recunoaștere ale algoritmului nostru. Pentru testele de recunoaștere a vorbirii am utilizat baza de date Deterding care conține un set de vectori de referință în domeniu, obținuți din pronunții ale vocalelor în diferite contexte.

## 5.2 Algoritmul Relevance LVQ (RLVQ)

Uneori nu toate componentele vectorilor de intrare au aceeași importanță în decizia pe care o ia un sistem de clasificare. Unele se pot dovedi mai influente decât altele, astfel încât putem deveni interesați în a calcula ponderile lor.

RLVQ este un algoritm euristic, supervizat, care amplifică valoarea ponderilor acelor componente care contribuie cel mai mult la clasificarea corectă a vectorilor de instruire și slăbește ponderile componentelor care au o influență negativă. Clustering-ul este realizat de un set de prototipuri care sunt adaptate vectorilor de instruire printr-un algoritm LVQ standard. Figura 5.1 este o reprezentare schematică a modului în care lucrează acest algoritm.

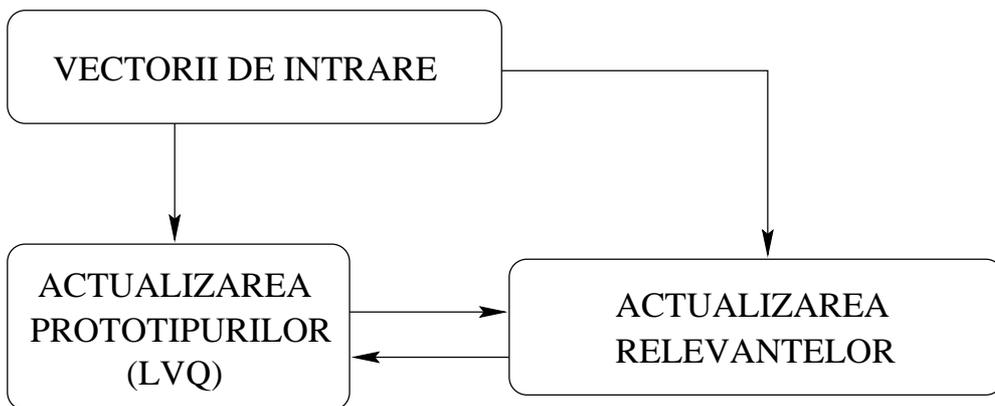


Figura 5.1: Reprezentarea schematică a modului în care se actualizează relevanțele și prototipurile prin algoritmul RLVQ.

Presupunem că  $\mathbf{x}_i = [x_{i1}, \dots, x_{in}]^t$ ,  $i = 1, \dots, N$  sunt vectorii de intrare. Un număr de  $P$  prototipuri sunt adaptate prin LVQ1 astfel încât să se obțină cea mai mică eroare de cuantizare posibilă pentru toți vectorii de instruire:

**Pasul 1.** Pentru o intrare  $\mathbf{x}_i$  se găsește cel mai apropiat prototip  $\mathbf{w}_j$ , câștigătorul, pentru care distanța  $\|\mathbf{x}_i - \mathbf{w}_j\|$  este minimă. Dacă  $\mathbf{x}_i$  și  $\mathbf{w}_j$  fac parte din aceeași clasă, atunci vectorul de intrare este clasificat corect.

**Pasul 2.** Se actualizează prototipul câștigător:

$$\mathbf{w}_j = \begin{cases} \mathbf{w}_j + \eta(\mathbf{x}_i - \mathbf{w}_j) & \text{dacă } \mathbf{x}_i \text{ este clasificat corect} \\ \mathbf{w}_j - \eta(\mathbf{x}_i - \mathbf{w}_j) & \text{în caz contrar,} \end{cases}$$

unde  $\eta > 0$  este rata de instruire.

Algoritmul RLVQ folosește o metrică modificată:

$$\|\mathbf{x}_i - \mathbf{w}_j\|_\lambda = \sqrt{\sum_{k=1}^n \lambda_k (x_{ik} - w_{jk})^2},$$

unde  $\lambda = [\lambda_1, \dots, \lambda_n]^t$  este vectorul relevanțelor, iar  $\sum_{k=1}^n \lambda_k = 1$ . Urmând o regulă similară algoritmului LVQ, factorii de ponderare sunt adaptați iterativ astfel [10]:

**Pasul 1.** Pentru fiecare  $k = 1, \dots, n$

$$\lambda_k = \begin{cases} \max\{\lambda_k - \alpha|x_{ik} - w_{jk}|, 0\} & \text{când } \mathbf{x}_i \text{ este clasificat corect} \\ \lambda_k + \alpha|x_{ik} - w_{jk}| & \text{în rest,} \end{cases}$$

unde  $\alpha > 0$  este un parametru de instruire.

**Pasul 2.** Ponderile  $\lambda_k$  sunt normalizate pentru fiecare  $k = 1, \dots, n$ .

Determinarea relevanțelor poate fi realizată după instruirea prototipurilor cu LVQ sau simultan cu aceasta, cea de-a doua variantă conducând la un algoritm on-line. Rezultatele experimentale din [10] demonstrează că rata de recunoaștere a RLVQ este mai bună în comparație cu algoritmul LVQ standard.

Această regulă de actualizare a relevanțelor folosește principiul învățării hebbiene, așa cum se demonstrează în [10].

Atunci când vectorul  $\mathbf{x}_i$  este clasificat corect, ponderile care corespund diferențelor mari între componentele lui  $\mathbf{x}_i$  și ale lui  $\mathbf{w}_j$  sunt diminuate cu o valoare mai mare. Ponderile corespunzătoare unor diferențe pe componente mai mici sunt diminuate mai puțin, exploatând astfel similaritatea componentelor vectoriale. Prin aplicarea normalizării, vor fi amplificate relevanțele componentelor care contribuie cel mai mult la clasificarea corectă.

Atunci când clasificarea este incorectă, relevanțele corespunzătoare diferențelor mici trebuie să piardă din importanță, crescând cu o valoare mică. Cele

corespunzătoare diferențelor mari sunt asociate componentelor cu o pondere mai mare în clasificare și sunt modificate tot în sens crescător, dar cu o valoare mai mare. Prin normalizare, acele relevanțe care contribuie la clasificarea incorectă sunt diminuate cel mai mult.

## 5.3 Operatorii OWA

În această secțiune vom descrie conceptul de operator OWA [93], proprietățile fundamentale ale acestuia [93] și vom relua o metodă propusă de Yager în [21] pentru instruirea operatorilor OWA pe baza unei colecții de eşantioane.

### 5.3.1 Problema generală a agregării

Presupunem că  $C_1, \dots, C_n$  sunt  $n$  criterii care definesc o problemă multicriterială. Notăm cu  $X$  domeniul în care iau valori argumentele celor  $n$  criterii și cu  $I$  intervalul  $[0, 1]$ . Problema constă în formularea unei funcții globale de decizie  $D$  în așa fel încât pentru orice alternativă  $x \in X$ , funcția  $D(x) \in I$  indică gradul în care  $x$  îndeplinește condițiile cerute în raport cu cele  $n$  criterii, atunci când  $C_i(x) \in I$ ,  $i = 1, \dots, n$  indică gradul în care  $x$  satisface criteriul  $C_i$ . Putem scrie, așadar:

$$D(x) = F(C_1(x), \dots, C_n(x)),$$

unde funcția  $F$  este un operator de agregare care îndeplinește condițiile de:

- a) Monotonie. Cu cât un criteriu individual este satisfăcut într-o măsură mai mare, cu atât valoarea funcției globale de decizie crește:

$$\text{dacă } C_i(x) \geq C_i(y) \forall i = 1, \dots, n \text{ și } x, y \in X \Rightarrow D(x) \geq D(y);$$

- b) Simetrie. Ordinea în care sunt tratate criteriile de către funcția globală de decizie nu este importantă, cu alte cuvinte un argument particular al funcției  $F$  poate apărea pe orice poziție în lista de argumente ale lui  $F$ .

Analiza unui operator de agregare presupune studiul relației care există între criteriile care descriu problema. O situație extremă este cea în care o alternativă  $x$  trebuie să satisfacă toate criteriile  $C_1, \dots, C_n$ . În acest caz este vorba despre un *anding* aplicat celor  $n$  valori. Un exemplu de astfel de operator este  $D(x) = \text{Min}(C_1(x), \dots, C_n(x))$ . Cealaltă extremă este situația în care trebuie satisfăcut cel puțin un criteriu, cerință care presupune un *oring* al valorilor, operatorul  $D(x) = \text{Max}(C_1(x), \dots, C_n(x))$  făcând parte din această categorie. Majoritatea cazurilor se găsesc între cele două extreme, adică între cerința de a îndeplini toate criteriile și cea de a îndeplini cel puțin unul dintre ele.

### 5.3.2 Conceptul de operator OWA

O categorie specială de operatori de agregare sunt operatorii OWA [93].

Un operator OWA de dimensiune  $n$  este o funcție

$$F : R^n \rightarrow R$$

$$F(a_1, \dots, a_n) = \sum_{i=1}^n w_i a_i^*,$$

unde  $\mathbf{W} = [w_1 \dots w_n]^t$  este un vector de ponderi asociat operatorului, cu  $w_i \in [0, 1]$ ,  $\sum_{i=1}^n w_i = 1$ , iar  $\mathbf{A}^* = [a_1^* \dots a_n^*]^t$  este o reordonare descrescătoare a argumentelor lui  $\mathbf{A} = [a_1 \dots a_n]^t$  în așa fel încât  $a_i^*$  este al  $i$ -lea cel mai mare dintre  $a_i$ ,  $a_1^* \geq \dots \geq a_n^*$ .

Elementul central în utilizarea unui operator OWA este pasul de reordonare a argumentelor în urma căruia un argument particular  $a_i$  nu mai este asociat unei ponderi particulare  $w_i$ , ci ponderea  $w_i$  este asociată valorii de pe poziția  $i$  obținută în urma reordonării.

Se demonstrează [93] că operatorii OWA satisfac proprietățile de comutativitate, monotonie și idempotență. O proprietate importantă a acestor operatori este că ei includ operatorii *Min*, *Max* și media aritmetică pentru anumite valori particulare ale matricii  $\mathbf{W}$ :

- a) Pentru  $\mathbf{W} = [0 \dots 0 \ 1]^t$ ,  $F(a_1, \dots, a_n) = \text{Min}(a_1, \dots, a_n)$ ;
- b) Pentru  $\mathbf{W} = [1 \ 0 \dots 0]^t$ ,  $F(a_1, \dots, a_n) = \text{Max}(a_1, \dots, a_n)$ ;
- c) Pentru  $\mathbf{W} = [1/n \dots 1/n]^t$ ,  $F(a_1, \dots, a_n) = \frac{1}{n} \sum_{i=1}^n a_i$ .

Operatorii OWA au, de asemenea, proprietatea de mărginire:

$$\text{Min}(a_1, \dots, a_n) \leq F(a_1, \dots, a_n) \leq \text{Max}(a_1, \dots, a_n).$$

Pornind de la această inegalitate, a fost introdusă în [93] o măsură care caracterizează tipul de agregare realizat de o anumită valoare a vectorului de ponderi  $\mathbf{W}$ . Această măsură se numește *orness* și se definește prin relația:

$$O(\mathbf{W}) = \frac{1}{n-1} \sum_{i=1}^n (n-i) w_i.$$

Domeniul său de valori este  $[0, 1]$  și se apropie de 1 cu atât mai mult cu cât operatorul OWA are caracter mai pronunțat al unui *or*, adică al funcției *Max*. Pentru setul de valori ale lui  $\mathbf{W}$  care corespund operatorilor *Min*, *Max* și media aritmetică, valorile lui  $O(\mathbf{W})$  sunt:

- a)  $O([0 \dots 0 1]^t) = 0$ ;
- b)  $O([1 0 \dots 0]^t) = 1$ ;
- c)  $O([1/n \dots 1/n]^t) = 0,5$ .

O altă măsură introdusă în [93] este *dispersia*:

$$D(\mathbf{W}) = \sum_{i=1}^n w_i \ln w_i$$

care a fost propusă pentru calculul cantității de informație înglobată de ponderile  $\mathbf{W}$  ale unui operator OWA.

### 5.3.3 Instruirea operatorilor OWA

Au fost propuse diverse metode de alegere a ponderilor operatorilor OWA. O'Hagan [21] a introdus o tehnică bazată pe un *orness* prestabilit și o maximizare a dispersiei. Torra [85], [86] a utilizat un operator de agregare particular numit Weighted OWA și i-a determinat parametrii printr-o procedură care generează câte o ieșire ideală pentru fiecare vector de instruire. Karayiannis [37] a folosit două familii de ponderi OWA care cuprind, prima, un set de ponderi egale completat cu un număr de ponderi cu valoarea 0, iar a doua, ponderi care descresc liniar complete, de asemenea, cu un număr de valori de 0. Beliakov [6] a aproximat operatorii OWA rezolvând o problemă de tip Least Squares with Equality and Inequality (LSEI).

Filev și Yager [21] au calculat ponderile operatorilor OWA prin metoda gradientului descrescător folosind o valoare țintă pentru valoarea agregată. Prezentăm pe scurt această metodă.

Algoritmul presupune existența a  $m$  eșantioane alcătuite din valori ale celor  $n$  criterii,  $[a_{i1} \dots a_{in}]^t$ ,  $i = 1, \dots, m$  și a unei valori  $d_i$  asociate fiecărui eșantion, numită *valoare agregată*. Metoda prin care se fixează valoarea agregată nu are relevanță pentru acest algoritm. De exemplu, ea poate fi apreciată de un expert care, studiind scorul acordat fiecărui criteriu, stabilește o valoare globală.

Scopul algoritmului este de a obține ponderile unui operator OWA care modelează procesul de agregare ce a generat setul de date de instruire.

Presupunem că  $[a_{i1}^* \dots a_{in}^*]^t$  reprezintă o reordonare descrescătoare a criteriilor  $[a_{i1} \dots a_{in}]^t$ . Atunci valoarea agregată va fi:

$$d_i = a_{i1}^* w_1 + \dots + a_{in}^* w_n, \quad \forall i = 1, \dots, m.$$

Rescriem această relație în sensul unei erori pătratice medii pe care vom căuta să o minimizăm:

$$e_i = \frac{1}{2} (a_{i1}^* w_1 + \dots + a_{in}^* w_n - d_i)^2,$$

cu  $w_j \in [0, 1]$ ,  $j = 1, \dots, n$  și  $\sum_{j=1}^n w_j = 1$ . Pentru a elimina aceste două constrângeri care definesc un operator OWA, presupunem că ponderile sunt definite astfel:

$$w_j = \frac{e^{\lambda_j}}{\sum_{k=1}^n e^{\lambda_k}}, \quad j = 1, \dots, n.$$

Astfel, problema de minimizare cu constrângeri enunțată mai sus se rescrie în felul următor:

$$e_i = \frac{1}{2} \left( a_{i1}^* \frac{e^{\lambda_1}}{\sum_{k=1}^n e^{\lambda_k}} + \dots + a_{in}^* \frac{e^{\lambda_n}}{\sum_{k=1}^n e^{\lambda_k}} - d_i \right)^2. \quad (5.1)$$

Eroarea pătratică medie descrisă de relația (5.1) trebuie minimizată în funcție de parametrii  $\lambda_1, \dots, \lambda_n$ .

Pentru rezolvarea acestei probleme de optimizare se folosește metoda gradientului descrescător și se obține regula de actualizare a parametrilor  $\lambda_j$ ,  $j = 1, \dots, n$ :

$$\lambda_j^{(t+1)} = \lambda_j^{(t)} - \eta \frac{\partial e_i}{\partial \lambda_j},$$

unde  $t$  este pasul de instruire și  $\eta$  este rata de instruire. Notând cu  $\hat{d}_i$  valoarea estimată a lui  $d_i$ :

$$\hat{d}_i = a_{i1}^* \frac{e^{\lambda_1}}{\sum_{k=1}^n e^{\lambda_k}} + \dots + a_{in}^* \frac{e^{\lambda_n}}{\sum_{k=1}^n e^{\lambda_k}}$$

rezultă:

$$\frac{\partial e_i}{\partial \lambda_j} = w_j (a_{ij} - \hat{d}_i) (\hat{d}_i - d_i), \quad i = 1, \dots, n$$

sau

$$\lambda_j^{(t+1)} = \lambda_j^{(t)} - \eta w_j (a_{ij} - \hat{d}_i) (\hat{d}_i - d_i)$$

cu

$$w_j = \frac{e^{\lambda_j^{(t)}}}{\sum_{k=1}^n e^{\lambda_k^{(t)}}}, \quad j = 1, \dots, n.$$

Parametrii  $\lambda_j$ ,  $j = 1, \dots, n$  sunt actualizați prin propagarea erorii  $\hat{d}_i - d_i$  dintre valoare agregată  $\hat{d}_i$  estimată de ponderile operatorului OWA obținute până la momentul  $t$  și valoarea agregată reală  $d_i$ .

Această metodă calculează ponderile unui operator OWA având la dispoziție un set de eşantioane care constau din valori ale celor  $n$  criterii și câte o valoare agregată pentru fiecare eşantion.

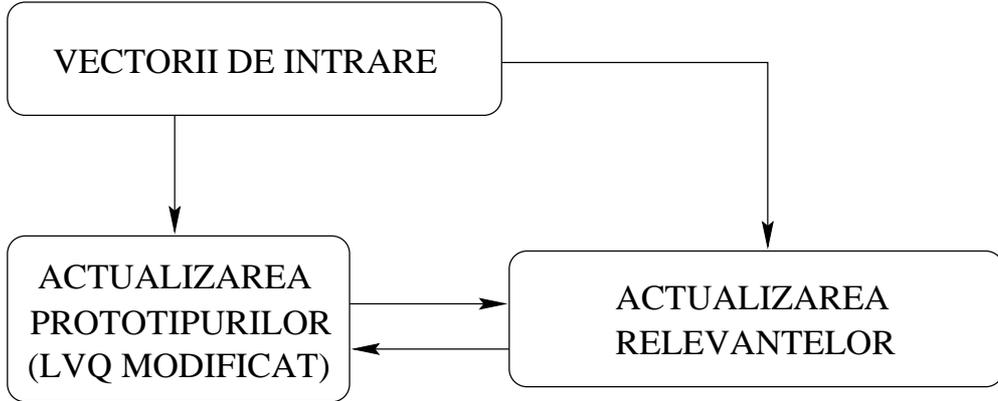


Figura 5.2: Reprezentarea schematică a modului în care se actualizează relevanțele și prototipurile prin algoritmul OWA-RLVQ. Algoritmul LVQ este modificat prin înlocuirea distanței euclidiene cu distanța  $D_{ij}^*$ , iar relevanțele sunt calculate ca ponderi ale unui operator OWA.

## 5.4 Algoritmul OWA - Relevance LVQ (OWA-RLVQ)

În această secțiune introducem un nou algoritm [14] care calculează ponderile unui operator OWA ca relevanțe ale componentelor vectorilor de instruire pentru probleme de clustering supervizat, așa cum se poate vedea în figura 5.2.

Fie  $\mathbf{x}_i \in \mathbf{R}^n$  un vector de instruire,  $i = 1, \dots, N$ ,  $N$  fiind numărul de vectori de instruire. Considerăm  $\mathbf{w}_k$ ,  $k = 1, \dots, P$  prototipurile care pot fi calculate printr-un algoritm LVQ. Pentru a calcula ponderile OWA, vom considera ca operator de agregare următoarea distanță modificată:

$$D_{ij}^* = \sqrt{\sum_{k=1}^n \lambda_k (|x_{ik} - w_{jk}|^*)^2},$$

unde  $\sum_{k=1}^n \lambda_k = 1$  și  $|x_{ik} - w_{jk}|^*$  este cea de-a  $k$ -a cea mai mare distanță dintre componentele corespondente ale vectorului de instruire și vectorului prototip.

Această distanță este un caz special al mediei generalizate de tip OWA [37]:

$$M = \left( \sum_{k=1}^n \lambda_k a_k^{*p} \right)^{\frac{1}{p}}, \quad (5.2)$$

unde  $p \in \mathbf{R}^*$ . Pentru  $p = 2$ , (5.2) se reduce la distanța modificată  $D_{ij}^*$  propusă de noi dacă facem următoarele înlocuiri:

$$a_k^* = |x_{ik} - w_{jk}|^*$$

și

$$|x_{i1} - w_{j1}|^* \geq \dots \geq |x_{in} - w_{jn}|^*.$$

Algoritmul LVQ1 trebuie reformulat pentru a minimiza o funcție obiectiv bazată pe această distanță modificată. Calculăm:

$$\Delta \mathbf{w}_j^* = \eta \lambda \mathbf{I}(\mathbf{x}_i - \mathbf{w}_j)^*$$

dacă  $\mathbf{x}_i$  este clasificat corect conform distanței  $D_{ij}^*$  și:

$$\Delta \mathbf{w}_j^* = -\eta \lambda \mathbf{I}(\mathbf{x}_i - \mathbf{w}_j)^*$$

dacă  $\mathbf{x}_i$  nu este clasificat corect. În relațiile de mai sus am notat cu  $(\mathbf{x}_i - \mathbf{w}_j)^*$  vectorul diferență ale cărui componente sunt reordonate, iar relevanțele se vor aplica după reordonare. Din această cauză, setul de prototipuri obținut în urma acestei proceduri nu va fi același cu cel obținut în urma algoritmului LVQ standard pentru că reflectă două abordări diferite ale instruirii.

Un vector prototip  $\mathbf{w}_j$  va eticheta cu numele clasei sale toți vectorii de intrare care îndeplinesc următoarea inegalitate:

$$|\mathbf{x}_i - \mathbf{w}_j|^* < |\mathbf{x}_i - \mathbf{w}_l|^*, \forall j \neq l.$$

Fie

$$\lambda = [\lambda_1 \dots \lambda_n]^t$$

$$\mathbf{x}_i = [x_{i1} \dots x_{in}]^t$$

$$\mathbf{w}_j = [w_{j1} \dots w_{jn}]^t$$

și

$$\mathbf{d} = [d_1 \dots d_n]^t,$$

unde  $d_k = x_{ik} - w_{jk}$ ,  $k = 1, \dots, n$ . Dacă  $\mathbf{x}_i$  este clasificat corect conform distanței  $D_{ij}^*$ , o valoare mică a lui  $|d_k|^*$  ar trebui să inducă o valoare mare pentru  $\Delta \lambda_k$ , iar o valoare mare a lui  $|d_k|^*$  ar trebui să conducă la o valoare mică pentru  $\Delta \lambda_k$ . Astfel, pentru o clasificare corectă conform criteriului formulat anterior,  $|d_k|^* \leq |d_{k'}|^*$  implică  $\Delta \lambda_k \geq \Delta \lambda_{k'}$ , adică  $-|d_k|^* \geq -|d_{k'}|^*$  înseamnă  $\Delta \lambda_k \geq \Delta \lambda_{k'}$  pentru  $k, k' = 1, \dots, n$ . Putem, astfel, să luăm:

$$\Delta \lambda_k = -\alpha |d_k|^*$$

sau:

$$\Delta \lambda_k = -\alpha |x_{ik} - w_{jk}|^*,$$

unde  $\alpha > 0$  este un parametru de instruire. Atunci când clasificarea conform aceluiași criteriu este incorectă,  $|d_k|^* \leq |d_{k'}|^*$  implică  $\Delta \lambda_k \leq \Delta \lambda_{k'}$ , pentru  $k, k' = 1, \dots, n$ . Putem lua:

$$\Delta \lambda_k = \alpha |d_k|^*$$

sau:

$$\Delta\lambda_k = \alpha|x_{ik} - w_{jk}|^*.$$

Folosim transformarea

$$\lambda_k = \frac{e^{\lambda_k}}{\sum_{i=1}^n e^{\lambda_i}}$$

pentru  $i = 1, \dots, n$  pentru a ne asigura că  $\sum_{k=1}^n \lambda_k = 1$  și  $\lambda_k \in [0, 1]$  care sunt constrângerile impuse ponderilor operatorilor OWA.

Următoarea procedură actualizează on-line atât relevanțele cât și rangul fiecărei componente a vectorilor de intrare.

**Pasul 1.** Inițializează  $\eta$  și  $\alpha$ . Inițializează  $\lambda_k = \frac{1}{n}$ ,  $k = 1, \dots, n$ , unde  $n$  este numărul de componente ale vectorului de intrare.

**Pasul 2.** Inițializează prototipurile.

**Pasul 3.** Actualizează prototipurile conform regulii LVQ modificate:

$$\mathbf{w}_j^* = \begin{cases} \mathbf{w}_j^* + \eta\lambda\mathbf{I}(\mathbf{x}_i - \mathbf{w}_j)^* & \text{dacă } \mathbf{x}_i \text{ este clasificat corect} \\ & \text{conform distanței } D_{ij}^* \\ \mathbf{w}_j^* - \eta\lambda\mathbf{I}(\mathbf{x}_i - \mathbf{w}_j)^* & \text{în rest.} \end{cases}$$

**Pasul 4.** Actualizează relevanțele:

$$\lambda_k = \begin{cases} \lambda_k - \alpha|x_{ik} - w_{jk}|^* & \text{dacă } \mathbf{x}_i \text{ este clasificat corect} \\ & \text{conform distanței } D_{ij}^* \\ \lambda_k + \alpha|x_{ik} - w_{jk}|^* & \text{în rest} \end{cases}$$

pentru  $k = 1, \dots, n$ .

**Pasul 5.** Transformă relevanțele:

$$\lambda_k = \frac{e^{\lambda_k}}{\sum_{i=1}^n e^{\lambda_i}}$$

pentru  $k = 1, \dots, n$ .

**Pasul 6.** Calculează ponderea fiecărei componente a vectorului de intrare ca medie a indicilor de poziție calculați înaintea pasului de ordonare, pentru toți pașii parcurși de la începutul instruirii.

**Pasul 7.** Repetă pașii 3–6 pentru fiecare vector de instruire.

Acest algoritm calculează ponderile unui operator OWA ca relevanțe prin minimizarea distanței modificate  $D_{ij}^*$  care este și valoarea agregată. Trebuie să facem distincția între relevanțe și rangul unei componente. Rangul este atașat unei componente. Relevanța este atașată unei poziții specifice din vectorul ordonat al distanței pe componente dintre un vector de intrare și un prototip. Aceasta explică semnificația pasului 6.

## 5.5 Experimente

În studiul nostru am testat algoritmul OWA-RLVQ pe baze de date standard pentru a vedea care este structura tipică a vectorilor de relevanțe și care sunt performanțele în recunoaștere. Programul care implementează acest algoritm a fost scris în limbajul C++. Instruirea este on-line, fiecare vector fiind prezentat o singură dată, astfel încât timpul de execuție este practic neglijabil.

Baza de date Iris [9] conține 3 clase alcătuite din câte 50 de vectori. Două dintre ele nu sunt liniar separabile. Problema este detectarea claselor pe baza celor 4 componente vectoriale. Am instruit OWA-RLVQ cu 6 prototipuri și am obținut o rată medie de recunoaștere de 96,6%. Un vector tipic al relevanțelor este  $[0,15 \ 0,21 \ 0,23 \ 0,38]^t$ . Am luat  $\eta=0,3$  și  $\alpha=0,2$ . Experimentele au arătat că a patra componentă vectorială are influența cea mai mare, iar prima componentă vectorială este cea mai puțin semnificativă, aceste rezultate fiind confirmate și de cele din [10]. Rangurile obținute pentru fiecare componentă după instruirea cu RLVQ și OWA-RLVQ folosind același set de prototipuri inițiale sunt prezentate în tabelul 5.1.

Tabelul 5.1: Rangul componentelor vectoriale pentru baza de date Iris.

| Rangul                           | 1    | 2    | 3    | 4    |
|----------------------------------|------|------|------|------|
| Componentele RLVQ                | 4    | 2    | 3    | 1    |
| Componentele OWA-RLVQ            | 4    | 3    | 2    | 1    |
| Ponderile componentelor OWA-RLVQ | 1,86 | 1,44 | 1,24 | 1,17 |

Baza de date Deterding (Vowel Recognition) [9] conține vectori extrași din pronunțiile vocalelor în 11 contexte. Fiecare pronunție este repetată de 6 ori de cei 15 vorbitori. Problema cere să se folosească pronunțiile primilor 8 vorbitori pentru instruire și pronunțiile ultimilor 7 vorbitori pentru recunoaștere. Un vector tipic al relevanțelor obținut prin algoritmul OWA-RLVQ este  $[0,032 \ 0,039 \ 0,043 \ 0,044 \ 0,077 \ 0,136 \ 0,137 \ 0,15 \ 0,163 \ 0,173]^t$ . Rata de recunoaștere pe care am obținut-o cu acest algoritm a fost de 46,75%. Rata medie de recunoaștere pentru LVQ1 la acest experiment a fost de 44,8%, iar pentru RLVQ a fost de 46,32%. Problema este cunoscută ca fiind una dificilă, de aceea am folosit 59 de prototipuri. În experimentele OWA-RLVQ am luat  $\eta=1,7$  și  $\alpha=1,9$ . Rangurile obținute pentru această bază de date sunt prezentate în tabelul 5.2. Componenta 2 este considerată cea mai importantă de RLVQ și pe poziția a doua de OWA-RLVQ, după cum se poate remarca studiind aceste rezultate.

Componenta 10 este cea mai puțin importantă pentru ambele modele.

Tabelul 5.2: Rangul componentelor vectoriale pentru baza de date Deterding.

|                                     |       |       |      |      |      |
|-------------------------------------|-------|-------|------|------|------|
| Rangul                              | 1     | 2     | 3    | 4    | 5    |
| Componentele RLVQ                   | 2     | 5     | 1    | 9    | 6    |
| Componentele OWA-RLVQ               | 8     | 2     | 4    | 5    | 6    |
| Ponderile componentelor<br>OWA-RLVQ | 5,209 | 5,209 | 5,20 | 5,17 | 5,16 |
| Rangul                              | 6     | 7     | 8    | 9    | 10   |
| Componentele RLVQ                   | 3     | 4     | 8    | 7    | 10   |
| Componentele OWA-RLVQ               | 9     | 3     | 7    | 1    | 10   |
| Ponderile componentelor<br>OWA-RLVQ | 5,15  | 5,15  | 5,14 | 5,14 | 5,13 |

Baza de date Ionosphere [9] constă din 351 de vectori de date cu câte 34 de caracteristici. Vectorilor le sunt asociate etichetele "bad" și "good", fiind vorba de o clasificare binară. Problema constă în folosirea primilor 200 de vectori, echilibrați numeric între exemple pozitive și negative, pentru instruire și ultimii 151 pentru test. Am folosit 8 prototipuri, iar valorile parametrilor pentru instruirea OWA-RLVQ au fost  $\eta=3,3$  și  $\alpha=3,5$ . Rata medie de recunoaștere a fost de 93,37%, în timp ce pentru LVQ1 am obținut 90,06%, iar pentru RLVQ am obținut 92,71%. Ratele uzuale de recunoaștere corectă pentru această bază de date încep, așa cum se arată în [9], de la 90%. În acest context, rezultatele obținute de noi sunt foarte bune. Rangurile celor mai importante 5 componente vectoriale obținute cu RLVQ și OWA-RLVQ sunt prezentate în tabelul 5.3.

Tabelul 5.3: Rangul componentelor vectoriale pentru baza de date Ionosphere. Sunt reprezentate cele mai importante 5 componente.

|                                     |       |       |       |       |       |
|-------------------------------------|-------|-------|-------|-------|-------|
| Rangul                              | 1     | 2     | 3     | 4     | 5     |
| Componentele RLVQ                   | 20    | 28    | 26    | 12    | 6     |
| Componentele OWA-RLVQ               | 14    | 12    | 1     | 3     | 28    |
| Ponderile componentelor<br>OWA-RLVQ | 19,20 | 19,16 | 19,10 | 19,07 | 19,06 |

Tabelul 5.4 compară ratele de recunoaștere obținute cu LVQ1, RLVQ și OWA-RLVQ folosind de fiecare dată același set inițial de prototipuri. În general am obținut rezultate mai bune cu algoritmul OWA-RLVQ față de cele obținute

cu RLVQ. Ambele metode sunt superioare algoritmului LVQ1, fapt care ilustrează utilitatea relevanțelor în clasificare. Trebuie accentuată observația că valorile relevanțelor obținute cu RLVQ și OWA-RLVQ nu se suprapun în mod obligatoriu pentru că reflectă două abordări diferite. În cazul metodei propuse de noi, relevanțele se referă la distanțele reordonate, în timp ce algoritmul RLVQ generează un set de valori asociate în mod precis anumitor componente ale vectorilor de intrare.

Tabelul 5.4: Comparație între ratele de recunoaștere obținute cu LVQ1, RLVQ și OWA-RLVQ.

| Baza de date | LVQ1   | RLVQ   | OWA-RLVQ |
|--------------|--------|--------|----------|
| Iris         | 91,33% | 95,33% | 96,60%   |
| Deterding    | 44,80% | 46,32% | 46,75%   |
| Ionosphere   | 90,06% | 92,71% | 93,37%   |

S-a constatat că, uneori, componentele care au valori neschimbate în tot setul de vectori de instruire pot conduce la falsa accentuare a relevanței lor, această problemă fiind cunoscută și pentru RLVQ. Din acest motiv, se poate recurge la o preprocesare în scopul eliminării lor.

## 5.6 Concluzii

În acest capitol am introdus o nouă metodă de calcul al relevanțelor dimensiunilor vectorilor de intrare ca ponderi OWA. Algoritmul propus calculează on-line relevanțele, oferind posibilitatea unei permanente adaptări la datele de intrare. Rezultatele experimentale au arătat că algoritmul propus de noi oferă o bună rată de recunoaștere. În același timp, factorii de relevanță pot fi utilizați pentru stabilirea on-line a rangului fiecărei dimensiuni a vectorilor de intrare.

Algoritmul OWA-RLVQ introdus de noi este diferit de algoritmul LVQ bazat pe operatorii OWA introdus de Karayiannis [37] deoarece noi modificăm metrica LVQ pentru ca să urmărească strategia RLVQ.

Este interesant de observat că abordarea noastră este o nouă metodă de calcul al ponderilor OWA. În comparație cu metoda descrisă în [21] și [94] care folosește tehnica gradientului descrescător, noi realizăm acest lucru în mod dinamic.

În capitolul următor vom dezvolta două noi metode care calculează factorii de relevanță prin maximizarea unor mărimi bazate pe informația mutuală și pe energia informațională.

## Capitolul 6

# Recunoaștere prin calculul relevanțelor intrărilor folosind noțiuni de teoria informației

*Selecția componentelor vectorilor de intrare într-un sistem de clasificare este o etapă de preprocesare absolut necesară, în special atunci când avem de a face cu cantități foarte mari de date. Introducem în acest capitol algoritmul Energy Relevance LVQ (ERLVQ) bazat pe energia informațională Onicescu [65] ca o metodă incrementală de instruire pentru clasificare supervizată și calcul al relevanțelor. Algoritmul Mutual Information Relevance LVQ (MIRLVQ) propune o tehnică alternativă de dezvoltare a acestei metode pornind de la informația mutuală.*

### 6.1 Introducere

Informația mutuală este o mărime care a fost utilizată pentru a implementa diverse metode de selecție a intrărilor [5], [49]. În acest caz, informația mutuală evaluează ”conținutul informațional” al fiecărei componente în mod individual, raportat la clasa de ieșire. Metoda de selecție a componentelor constă în căutarea unui subset de componente relevante din setul initial de componente disponibile. Acest subset se alege în așa fel încât să maximizeze informația mutuală. O importanță însemnată în acest algoritm o are estimarea informației mutuale. Estimarea informației mutuale este o problemă dificilă datorită funcțiilor de densitate de probabilitate condiționată și a complexității

compuționale ridicate. Mulți algoritmi de selecție a componentelor vectoriale care se bazează pe informația mutuală folosesc histograma ca estimator al densității de probabilitate [5], [49]. În spații cu multe dimensiuni, însă, histograma nu este practică și nici foarte precisă [32].

Informația mutuală bazată pe entropia Renyi pătratică [70] a fost utilizată recent într-o metodă eficientă de selecție a caracteristicilor [32]. Acest algoritm numit *toward-optimal feature selection methodology - mutual information (OFS-MI)* constă din două criterii bazate pe informația mutuală și un algoritm de căutare. Primul criteriu numit *feature relevancy criterion (FRC)* are scopul de a selecta componentele relevante, iar al doilea, *feature similarity criterion (FSC)*, este folosit pentru a reduce redundanța între componentele selectate. Considerăm mulțimea  $F$  de componente candidate și  $S$  componentele selectate. Pentru a determina cea mai semnificativă componentă din  $F$ , toate componentele candidate sunt ordonate folosind relația

$$FRC(f) = I(f, S; C),$$

unde  $C$  este mulțimea claselor,  $f$  este una dintre componentele din  $F$ , iar  $I$  este informația mutuală. Adăugarea unei componente de intrare cu un aport scăzut în creșterea ratei de clasificare conduce la scăderea  $FRC$ . Al doilea criteriu pentru care

$$FSC(f; S) = \operatorname{argmax} \left( \frac{I(f; f_i)}{H(f_i)} \right)$$

măsoară similaritatea dintre componenta  $f$  și submulțimea  $S$ , unde  $f \notin S$  iar  $H$  este entropia. Atunci când  $FSC(f; S)$  depășește un prag, se consideră că  $f$  este similar cu  $S$ .  $FSC$  este estimat în acest caz folosind informația mutuală pătratică [70] dintre două variabile aleatoare continue.

Torkkola [80], [81], [82], [83], [84] maximizează informația mutuală pătratică obținând o metodă de instruire a unui sistem pentru extragerea caracteristicilor. Un astfel de sistem transformă vectorii de intrare în vectori cu un număr mai mic de dimensiuni prin recalcularea componentelor, fără a afecta performanțele clasificatorului. Diferența dintre un sistem de selecție a caracteristicilor și unul de extragere (transformare, construcție, generare) a caracteristicilor este aceea că primul folosește mai departe un subset al componentelor vectoriale originale, în timp ce al doilea construiește vectori cu un număr redus de dimensiuni prin recalcularea acestora, fără a păstra în mod necesar valorile inițiale. Metoda lui Torkkola adaptează informația mutuală pătratică din [70] pentru variabile aleatoare discrete și o maximizează printr-o metodă de gradient ascendent pentru a obține o transformare convenabilă.

Algoritmul *Mutual Information Relevance LVQ (MIRLVQ)* propus de noi în acest capitol maximizează informația mutuală pătratică pentru a obține relevanțele care conduc la îmbunătățirea performanțelor unui clasificator bazat pe un algoritmul LVQ. Algoritmul *Energy Relevance LVQ (ERLVQ)* este similar lui MIRLVQ, însă maximizează o mărime bazată pe energia informațională. Experimental vom arăta că acești algoritmi conduc la rate de recunoaștere mai bune decât algoritmul OWA-RLVQ introdus de noi în capitolul 5.

## 6.2 Entropia Shannon

În această secțiune vom prezenta câteva noțiuni fundamentale de teoria informației [54] care sunt folosite în dezvoltarea algoritmilor MIRLVQ și ERLVQ.

Presupunem că  $X$  este o variabilă aleatoare, iar  $x$  o realizare a sa care poate lua valori din mulțimea

$$\{x_1, \dots, x_n\}.$$

Fiecare valoare poate să apară cu probabilitatea

$$\{p_1, \dots, p_n\}$$

unde

$$P(x = x_i) = p_i, \quad p_i \geq 0, \quad i = 1, \dots, n$$

și

$$\sum_{x_i \in \{x_1, \dots, x_n\}} P(x = x_i) = 1.$$

Pentru a simplifica notația, îl vom înlocui pe  $P(x = x_i)$  cu  $P(x_i)$  sau cu  $P(x)$ .

Dacă  $X$  și  $Y$  sunt două variabile aleatoare, realizările fiind perechi ordonate  $(x, y)$  de valori cu  $x \in \{x_1, \dots, x_n\}$  și  $y \in \{y_1, \dots, y_m\}$ , atunci  $P(x = x_i, y = y_j)$  este probabilitatea de a se realiza evenimentul pentru care  $x$  și  $y$  au valorile particulare  $x_i$  și  $y_j$ ,  $i \in \{1, \dots, n\}$  și  $j \in \{1, \dots, m\}$ . Pentru simplificarea notației, scriem  $P(x, y)$  în loc de  $P(x = x_i, y = y_j)$ . Este cunoscută relația:

$$P(x = x_i) = \sum_{y \in \{y_1, \dots, y_m\}} P(x = x_i, y)$$

sau, mai simplu:

$$P(x) = \sum_{y \in \{y_1, \dots, y_m\}} P(x, y). \quad (6.1)$$

Probabilitatea ca  $x$  să fie egal cu  $x_i$  atunci când  $y$  este egal cu  $y_j$  se numește probabilitatea condiționată și se notează cu

$$P(x = x_i | y = y_j)$$

unde

$$x_i \in \{x_1, \dots, x_n\}$$

și

$$y_j \in \{y_1, \dots, y_m\}.$$

Probabilitatea condiționată respectă următoarea egalitate:

$$P(x = x_i | y = y_j) = \frac{P(x = x_i, y = y_j)}{P(y = y_j)}, \quad \text{unde } P(y = y_j) \neq 0. \quad (6.2)$$

Atunci când  $P(y = y_j) = 0$ , probabilitatea condiționată nu este definită. Folosind scrierea simplificată, relația (6.2) este echivalentă cu:

$$P(x|y) = \frac{P(x, y)}{P(y)}. \quad (6.3)$$

Spunem că variabilele aleatoare  $X$  și  $Y$  sunt independente dacă și numai dacă:

$$P(x, y) = P(x)P(y).$$

Relațiile de mai sus pot fi rescrise în diverse forme și obținem formule utile în aplicațiile care folosesc probabilități. Pornind de la expresia (6.3), deducem următoarele egalități:

$$P(x, y) = P(x|y)P(y) = P(y|x)P(x),$$

din care formulăm teorema lui Bayes:

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)},$$

iar din (6.1) și (6.3) obținem:

$$P(x) = \sum_{y \in \{y_1, \dots, y_m\}} P(x, y) = \sum_{y \in \{y_1, \dots, y_m\}} P(x|y)P(y).$$

Shannon [78] a definit entropia ca o măsură a cantității de informație produsă de o sursă de informație reprezentată printr-un proces markovian.

Presupunem că avem o mulțime de evenimente posibile cu probabilitățile  $p_1, \dots, p_n$ . Aceste probabilități sunt singurele elemente pe care le cunoaștem. Măsura  $H(p_1, \dots, p_n)$  arată care este incertitudinea în legătură cu ieșirile generate de sursa de informație și îndeplinește următoarele cerințe [78]:

- a)  $H$  este continuă în  $p_i$ ;
- b) Dacă toate valorile lui  $p_i$ ,  $i = 1, \dots, n$  sunt egale,  $p_i = \frac{1}{n}$ , atunci incertitudinea este maximă;

- c) Dacă "alegerea" unui eveniment de către sursa de informație poate fi împărțit în două "alegeri" succesive, atunci valoarea originală a lui  $H$  este suma ponderată a valorilor individuale ale lui  $H$ .

Shannon a demonstrat că singura funcție  $H$  care satisface cerințele de mai sus este :

$$H(X) = -k \sum_{x \in \{x_1, \dots, x_n\}} P(x) \log P(x)$$

sau

$$H(X) = k \sum_{x \in \{x_1, \dots, x_n\}} P(x) \log \frac{1}{P(x)},$$

cu posibilitatea de a alege  $k = 1$  și cu convenția că pentru  $P(x) = 0$  avem:

$$0 \times \log \frac{1}{0} = 0$$

deoarece

$$\lim_{\theta \rightarrow 0^+} \theta \log \frac{1}{\theta} = 0.$$

Funcția  $H$  se numește entropie.

Conținutul informațional al realizării  $x$  a variabilei aleatoare  $X$  se definește prin:

$$h(x) = -\log_2 P(x).$$

Atât conținutul informațional cât și entropia sunt măsurate în biți.

Atunci când avem două variabile aleatoare  $X$  și  $Y$ , entropia este:

$$H(X, Y) = \sum_{(x, y) \in \{x_1, \dots, x_n\} \times \{y_1, \dots, y_m\}} -P(x, y) \log P(x, y).$$

Entropia este aditivă pentru variabile aleatoare independente:

$$H(X, Y) = H(X) + H(Y) \Leftrightarrow P(x, y) = P(x)P(y).$$

Dacă variabilele aleatoare sunt continue, atunci nu mai putem discuta de probabilitatea  $P$  a unei anumite realizări și folosim funcția de densitate de probabilitate pe care o notăm în continuare cu  $p$ .

### 6.3 Estimarea densităților de probabilitate cu ferestre Parzen

În această secțiune vom aborda problema modelării unei funcții de densitate de probabilitate [8], [17] pe baza unui număr de eșantioane  $\mathbf{y}_i$ ,  $i = 1, \dots, N$  care sunt realizări ale unei variabile aleatoare  $Y$ . Vom vedea cum se poate folosi acest

rezultat pentru estimarea entropiei Renyi pătratică ca prim pas în estimarea din eşantioane a informaţiei mutuale dintre două variabile aleatoare.

Presupunem că fiecare eşantion  $\mathbf{y}_i$  aparţine uneia dintre clasele  $c_1, \dots, c_P$  care sunt realizări ale variabilei aleatoare discrete  $C$ . Spunem că vectorii  $\mathbf{y}_i$  sunt etichetaţi cu numele clasei din care fac parte.

Există două importante categorii de metode pentru estimarea densităţii de probabilitate: *metodele parametrice* şi *metodele neparametrice*.

Metodele parametrice pornesc de la presupunerea că funcţia densităţii de probabilitate are o formă cunoscută. Se foloseşte un set de parametri care sunt apoi optimizaţi prin adaptarea modelului la setul de date disponibile. Dezavantajul acestei metode este că forma particulară aleasă pentru funcţia densităţii de probabilitate poate să nu modeleze foarte fidel adevărata reprezentare.

Metodele neparametrice de estimare nu fac vreo presupunere asupra modelului densităţii de probabilitate şi permit ca forma acesteia să fie determinată exclusiv pe baza setului de date de instruire.

Ne vom concentra asupra metodelor neparametrice deoarece dorim să construim un model de recunoaştere care se bazează exclusiv pe fluxul de date, fără a face nici o presupunere asupra proprietăţilor acestor date. Acest gen de aplicaţii sunt caracteristice sistemelor care prelucrează cantităţi mari de date în timp real, fiind capabile totodată să îşi adapteze dinamic parametrii la proprietăţile în schimbare ale intrărilor. Este vorba aşadar, de sisteme care sunt instruite on-line, pentru aplicaţii de data mining unde datele sosesc secvenţial, nefiind stocate din cauza volumului foarte mare şi a costurilor pe care le implică aceasta. Este important de spus că algoritmi on-line oferă posibilitatea de a se renunţa la datele care au fost folosite, deoarece nu se mai revine asupra lor.

### Estimarea densităţii de probabilitate prin eşantioane

Reluăm în această secţiune o metodă de estimare descrisă în [8] şi [17].

Fie un vector oarecare  $\mathbf{y}$  care face parte din setul de date de instruire şi care contribuie la definirea densităţii de probabilitate  $p(\mathbf{y})$ . Probabilitatea ca vectorul  $\mathbf{y}$  să se găsească în subspaţiul  $S$  se defineşte prin:

$$P = \int_S p(\mathbf{y}') d\mathbf{y}'.$$

Dacă presupunem că  $p(\mathbf{y})$  este continuă şi nu are variaţii foarte mari pe domeniul  $S$ , putem aproxima această probabilitate prin:

$$P = \int_S p(\mathbf{y}') d\mathbf{y}' \simeq p(\mathbf{y})V, \quad (6.4)$$

unde  $V$  este volumul lui  $S$  şi  $\mathbf{y}$  este un punct din  $S$ .

Pe de altă parte, considerăm că avem  $N$  eşantioane extrase independent din mulțimea caracterizată prin  $p(\mathbf{y})$  și atunci probabilitatea ca un număr  $K$  dintre ele să se găsească în subspațiul  $S$  va fi dată de legea binomială:

$$P(K) = C_N^K P^K (1 - P)^{N-K}.$$

Această egalitate se poate justifica astfel: Probabilitatea fiecăruia dintre cele  $K$  eşantioane de a se găsi în regiunea  $S$  este  $P$ , iar probabilitatea fiecăruia dintre celelalte  $N - K$  eşantioane de a se găsi în afara acestei regiuni este  $1 - P$ . Deoarece cu cele  $N$  eşantioane se pot realiza  $C_N^K$  grupuri de câte  $K$ , obținem relația de mai sus.

Cunoscând probabilitatea  $P$  ca un număr  $K$  de eşantioane din totalul celor  $N$  să se găsească în domeniul  $S$ , în mod evident valoarea medie a lui  $K$  este:

$$\epsilon[K] = NP,$$

relație ca se poate scrie și sub forma:

$$\epsilon\left[\frac{K}{N}\right] = P.$$

Dispersia în jurul valorii medii a lui  $\frac{K}{N}$  este:

$$\epsilon\left[\left(\frac{K}{N} - P\right)^2\right] = \frac{P(1 - P)}{N}$$

și

$$\epsilon\left[\left(\frac{K}{N} - P\right)^2\right] \rightarrow 0 \text{ când } N \rightarrow \infty.$$

Aceasta înseamnă că, pentru un număr suficient de mare de eşantioane, valoarea  $\frac{K}{N}$  este o estimare corectă a probabilității  $P$ :

$$P \simeq \frac{K}{N}. \quad (6.5)$$

Din (6.4) și (6.5) rezultă:

$$p(\mathbf{y})V \simeq \frac{K}{N}$$

sau

$$p(\mathbf{y}) \simeq \frac{K}{NV}. \quad (6.6)$$

Pentru a obține această estimare a densității de probabilitate am făcut două presupuneri referitoare la domeniul  $S$ . Am scris relația (6.4) presupunând că  $S$  este relativ mic astfel încât să putem considera  $p(\mathbf{y})$  constantă în interiorul domeniului. Am considerat că  $S$  este relativ mare astfel încât să cuprindă un

număr mare de eşantioane pentru a obține relația (6.5). În realitate ne vom găsi între cele două cazuri extreme, de aceea întotdeauna trebuie făcută o alegere care definește gradul de precizie al estimării și, pentru un anumit set de date, există în mod clar o valoare optimă a dimensiunii  $S$  astfel încât  $p(\mathbf{y})$  să fie estimat cât mai bine posibil.

Atunci când se aplică relația (6.6), există două importante opțiuni. Pe de o parte, putem alege o valoare fixă pentru  $K$ , determinând volumul  $V$  corespunzător pe baza datelor disponibile (*K-nearest neighborhood methods*). Alternativ, putem fixa volumul  $V$ , determinând valoarea lui  $K$  din setul de date. Această alternativă face parte din clasa tehnicilor de estimare a densității de probabilitate folosind nuclee (*kernel-based methods*) și o vom detalia mai departe.

În condițiile unui număr infinit de puncte, această procedură de estimare poate deveni exactă pentru că volumul  $S$  poate fi restrâns către zero. În felul acesta, ne asigurăm că relația (6.4) devine din ce în ce mai precisă în paralel cu creșterea preciziei relației (6.5) prin mărirea numărului de eşantioane din domeniul  $S$ . Duda, Hart și Stork [17] au arătat că ambele abordări descrise mai sus converg către valoarea reală a lui  $p(\mathbf{y})$  în condițiile în care  $N \rightarrow \infty$ , iar  $V$  scade odată cu  $N$  și  $K$  crește odată cu  $N$  într-un mod convenabil.

## Metoda ferestrelor Parzen

Prezentăm metoda de aproximare cu ferestre Parzen a lui  $p(\mathbf{y})$  folosind o tehnică bazată pe nuclee [8], [17].

Să presupunem că regiunea  $S$  are forma unui hipercub cu muchiile de lungime  $h$ , centrat în punctul  $\mathbf{y}$ . Volumul său va fi:

$$V = h^d. \quad (6.7)$$

Putem găsi o expresie pentru  $K$  folosind o funcție nucleu  $W(\mathbf{u})$  numită și fereastră Parzen care se definește astfel:

$$W(\mathbf{u}) = \begin{cases} 1, & \text{când } |u_j| < \frac{1}{2}, j = 1, \dots, d \\ 0, & \text{în rest} \end{cases}$$

astfel încât  $W(\mathbf{u})$  corespunde hipercubului unitate centrat în origine, iar  $\mathbf{u} = [u_1, \dots, u_d]^t$ . În felul acesta, pentru un punct  $\mathbf{y}_n$  putem scrie:

$$W\left(\frac{\mathbf{y} - \mathbf{y}_n}{h}\right) = \begin{cases} 1, & \text{dacă } \mathbf{y}_n \text{ se găsește în interiorul hipercubului} \\ & \text{de muchie } d, \text{ centrat în } \mathbf{y} \\ 0, & \text{în rest.} \end{cases}$$

Folosind această expresie, numărul total de puncte care se găsesc în interiorul hipercubului va fi:

$$K = \sum_{n=1}^N W\left(\frac{\mathbf{y} - \mathbf{y}_n}{h}\right). \quad (6.8)$$

Înlocuind (6.7) și (6.8) în (6.6) obținem:

$$\widetilde{p}(\mathbf{y}) = \frac{1}{N} \sum_{n=1}^N \frac{1}{h^d} W\left(\frac{\mathbf{y} - \mathbf{y}_n}{h}\right),$$

unde  $\widetilde{p}(\mathbf{y})$  este estimatul densității de probabilitate.

Deoarece  $W(\mathbf{u})$  este o funcție discontinuă,  $\widetilde{p}(\mathbf{y})$  este și ea o funcție discontinuă. Se poate evita acest lucru alegând diverse alte forme pentru nucleu.

O alegere foarte utilizată este nucleul normal pentru care:

$$\widetilde{p}(\mathbf{y}) = \frac{1}{N} \sum_{n=1}^N \frac{1}{(2\pi\sigma^2)^{\frac{d}{2}}} \cdot e^{-\frac{\|\mathbf{y} - \mathbf{y}_n\|^2}{2\sigma^2}}. \quad (6.9)$$

Se folosește în mod frecvent notația:

$$G(\mathbf{y}, \sigma^2 \mathbf{I}) = \frac{1}{(2\pi\sigma^2)^{\frac{d}{2}}} \cdot e^{-\frac{\mathbf{y}^T \mathbf{y}}{2\sigma^2}}, \quad (6.10)$$

unde  $G(\mathbf{y}, \sigma^2 \mathbf{I})$  este nucleul gaussian (normal)  $d$ -dimensional. Mai general, nucleul gaussian într-un spațiu  $d$ -dimensional este definit astfel:

$$G(\mathbf{y}, \Sigma) = \frac{1}{(2\pi)^{\frac{d}{2}} \|\Sigma\|^{\frac{1}{2}}} e^{(-\frac{1}{2} \mathbf{y}^T \Sigma^{-1} \mathbf{y})}.$$

Forma (6.10) a nucleului gaussian este una simplificată [8], considerându-se că vectorii de intrare au componente independente și că pentru fiecare componentă au aceeași dispersie. Gaussienele vor fi simetrice, centrate în câte un eșantion  $\mathbf{y}$ . Astfel, matricea de covarianță este  $\sigma^2 \mathbf{I}$ , unde  $\mathbf{I}$  este matricea unitate. Alegerea nucleului gaussian este motivată de faptul că în dezvoltările care urmează acesta trebuie să fie derivabil pe întregul domeniu.

Estimatorul Parzen al densităților de probabilitate înseamnă plasarea câte unui nucleu în locul fiecărui eșantion, iar densitatea este estimată ca medie a valorilor calculate cu aceste nuclee:

$$\widetilde{p}(\mathbf{y}) = \frac{1}{N} \sum_{n=1}^N G(\mathbf{y} - \mathbf{y}_n, \sigma^2 \mathbf{I}). \quad (6.11)$$

Se demonstrează că nucleul gaussian are următoarea proprietate:

$$\int_{-\infty}^{+\infty} G(\mathbf{y} - \mathbf{y}_i, \sigma^2 \mathbf{I}) G(\mathbf{y} - \mathbf{y}_j, \sigma^2 \mathbf{I}) d\mathbf{y} = G(\mathbf{y}_i - \mathbf{y}_j, 2\sigma^2 \mathbf{I}). \quad (6.12)$$

Această relație reprezintă convoluția a două gaussiene care generează tot o gaussiană. Remarcabil este faptul că rezultatul se exprimă ca o interacțiune dintre două eșantioane, iar covarianța este egală cu suma covarianțelor individuale ale celor două gaussiene. Demonstrația acestei egalități este dată în anexa B.

În general, pentru ca  $\widetilde{p}(\mathbf{y}) \geq 0$  și  $\int_S \widetilde{p}(\mathbf{y}) d\mathbf{y} = 1$ , nucleul trebuie să îndeplinească următoarele două condiții:

$$W(\mathbf{u}) \geq 0$$

$$\int_S W(\mathbf{u}) d\mathbf{u} = 1.$$

În relația (6.9),  $\sigma$  joacă rolul unui parametru care determină calitatea estimării. Atunci când  $\sigma$  este mare, netezirea lui  $\widetilde{p}(\mathbf{y})$  este foarte accentuată, pierzându-se eventualele maxime locale. La extrema cealaltă, când  $\sigma \rightarrow 0$  nucleul se apropie de funcția delta, graficul lui  $\widetilde{p}(\mathbf{y})$  conține variații foarte dese. Din acest motiv, trebuie găsit un compromis între cele două variante, iar valoarea lui  $\sigma$  poate fi stabilită experimental.

## Entropia Renyi

Estimarea densității de probabilitate cu ferestre Parzen oferă o metodă de estimare a entropiei cu avantaje certe din punct de vedere computațional atunci când se folosește împreună cu entropia Renyi [84]. Această măsură alternativă a entropiei este, de fapt, o familie de măsuri care include entropia Shannon ca un caz particular [26], [79]. Pentru o variabilă aleatoare discretă  $C$ , entropia Renyi se scrie:

$$H_{R_\alpha}(C) = \frac{1}{1-\alpha} \log \sum_c p^\alpha(c), \quad (6.13)$$

iar pentru o variabilă aleatoare continuă  $Y$  entropia Renyi se definește prin:

$$H_{R_\alpha}(Y) = \frac{1}{1-\alpha} \log \int_{\mathbf{y}} p^\alpha(\mathbf{y}) d\mathbf{y}, \quad (6.14)$$

unde  $\alpha > 0$  și  $\alpha \neq 1$ .

Este cunoscută relația care există între entropia Shannon și entropia Renyi:

$$\lim_{\alpha \rightarrow 1} H_{R_\alpha} = H.$$

Această egalitate se poate verifica imediat după cum urmează. Pentru cazul discret avem:

$$\sum_c p(c) = 1$$

și atunci putem scrie:

$$\lim_{\alpha \rightarrow 1} H_{R_\alpha}(C) = \lim_{\alpha \rightarrow 1} \frac{\log \sum_c p^\alpha(c)}{1 - \alpha}.$$

Deoarece avem cazul de nedeterminare  $\frac{0}{0}$ , aplicăm regula lui l'Hospital și obținem:

$$\begin{aligned} \lim_{\alpha \rightarrow 1} H_{R_\alpha}(C) &= \lim_{\alpha \rightarrow 1} \left( -\frac{1}{\sum_c p^\alpha(c)} \sum_c p^\alpha(c) \log p(c) \right) = \\ &= -\sum_c p(c) \log p(c), \end{aligned} \quad (6.15)$$

adică relația dintre entropia Renyi și entropia Shannon pentru cazul discret. Atunci când entropia Renyi se definește pentru o variabilă aleatoare continuă, demonstrația este similară. Avem:

$$\int_{\mathbf{y}} p(\mathbf{y}) d\mathbf{y} = 1$$

și atunci obținem:

$$\lim_{\alpha \rightarrow 1} H_{R_\alpha}(Y) = \lim_{\alpha \rightarrow 1} \frac{\log \int_{\mathbf{y}} p^\alpha(\mathbf{y}) d\mathbf{y}}{1 - \alpha}.$$

Aplicăm regula lui l'Hospital pentru cazul de nedeterminare  $\frac{0}{0}$  și obținem relația dintre entropia Renyi și entropia Shannon în cazul continuu:

$$\lim_{\alpha \rightarrow 1} H_{R_\alpha}(Y) = -\int_{\mathbf{y}} p(\mathbf{y}) \log p(\mathbf{y}).$$

Pentru  $\alpha = 2$  avem entropia Renyi pătratică. În cazul discret, relația (6.13) devine:

$$H_{R_2}(C) = -\log \sum_c p^2(c),$$

iar pentru cazul continuu relația (6.14) se scrie:

$$H_{R_2}(Y) = -\log \int_{\mathbf{y}} p^2(\mathbf{y}) d\mathbf{y}. \quad (6.16)$$

Pentru dezvoltarea algoritmilor MIRLVQ și ERLVQ vom fi interesați doar de cazul continuu. Din relația (6.11), proprietatea (6.12) și expresia (6.16) obținem:

$$\begin{aligned} H_{R_2}(Y) &= -\log \frac{1}{N^2} \int_{\mathbf{y}} \left( \sum_{k=1}^N \sum_{j=1}^N G(\mathbf{y} - \mathbf{y}_k, \sigma^2 \mathbf{I}) G(\mathbf{y} - \mathbf{y}_j, \sigma^2 \mathbf{I}) \right) d\mathbf{y} = \\ &= -\log \frac{1}{N^2} \sum_{k=1}^N \sum_{j=1}^N G(\mathbf{y}_k - \mathbf{y}_j, 2\sigma^2 \mathbf{I}) \end{aligned} \quad (6.17)$$

după ce integrala a comutat cu cele două sume.

Am arătat astfel că estimarea densității de probabilitate cu ferestre Parzen este utilă pentru estimarea entropiei Renyi pătratice ca sumă de interacțiuni locale între toate perechile de eşantioane. În secțiunea următoare vom folosi aceste rezultate pentru estimarea informației mutuale dintre două variabile aleatoare.

## 6.4 Informația mutuală pătratică

Informația mutuală este o măsură a dependenței dintre două variabile aleatoare și se definește prin următoarea diferență de entropii Shannon:

$$I(X, Y) = H(Y) - H(Y|X),$$

unde  $X$  și  $Y$  sunt două variabile aleatoare.

În [53] se fac o serie de considerații referitoare la informația mutuală și la funcția de corelație care măsoară, de asemenea, dependența dintre două variabile aleatoare. Reluăm în paragraful următor, pe scurt, aceste observații.

Dacă cele două variabile aleatoare sunt independente, atunci informația mutuală dintre ele este nulă. Când cele două variabile aleatoare sunt dependente, de exemplu una se poate calcula în funcție de cealaltă, informația mutuală dintre ele este mare. Informația mutuală arată care este gradul de predictibilitate a unei variabile aleatoare atunci când cealaltă este cunoscută. Interpretările de mai sus ale informației mutuale sunt proprii și corelației. Cu toate acestea, informația mutuală măsoară dependența generală dintre două variabile aleatoare, în timp ce corelația măsoară dependența liniară. O altă diferență majoră dintre cele două măsuri este aceea că informația mutuală poate fi aplicată atât secvențelor simbolice cât și secvențelor numerice, iar corelația poate fi aplicată doar secvențelor numerice. Problema pe care o discutăm în acest capitol asimilează colecția de clase de forme unei secvențe simbolice, de aceea informația mutuală este o alegere naturală.

Deoarece entropia Shannon nu poate fi estimată din datele de instruire cu un efort computațional rezonabil, folosim o soluție alternativă de estimare a informației mutuale cu ajutorul divergenței Kulback-Leibler (KL) care se mai numește și entropie relativă. Divergența KL dintre două funcții de densitate de probabilitate  $f(\mathbf{x})$  și  $g(\mathbf{x})$  se definește prin [17]:

$$K(f, g) = \int f(\mathbf{x}) \log \frac{f(\mathbf{x})}{g(\mathbf{x})} d\mathbf{x} \quad (6.18)$$

care folosește implicit entropia Shannon și satisface inegalitatea lui Gibbs [54]:

$$K(f, g) \geq 0,$$

egalitatea având loc pentru:

$$f(\mathbf{x}) = g(\mathbf{x}), \forall \mathbf{x} \in X.$$

Informația mutuală dintre două variabile aleatoare  $X$  și  $Y$  se estimează prin divergența KL astfel:

$$\begin{aligned} I_S(X, Y) &= K(f_{XY}(\mathbf{x}, \mathbf{y}), f_X(\mathbf{x})f_Y(\mathbf{y})) = \\ &= \int \int f_{XY}(\mathbf{x}, \mathbf{y}) \log \frac{f_{XY}(\mathbf{x}, \mathbf{y})}{f_X(\mathbf{x})f_Y(\mathbf{y})} d\mathbf{x}d\mathbf{y}. \end{aligned}$$

Dorim să estimăm informația mutuală dintre două variabile aleatoare într-o manieră similară celei în care am estimat entropia pătratică Renyi. Relația de mai sus nu poate fi, însă, integrată ușor cu metoda de estimare a densităților de probabilitate cu ferestre Parzen deoarece nu putem aplica proprietatea (6.12). Pentru aceasta ar trebui ca toți termenii să aibă gradul doi.

Din acest motiv, în [70] se propune o nouă măsură a distanței dintre două funcții de densitate de probabilitate care conține termeni de gradul doi și care păstrează proprietăți ale divergenței KL. Această măsură pornește de la formula estimării entropiei Renyi și este derivată euristic menținând termenii de gradul doi, similar entropiei Renyi pătratice.

Pe baza distanței euclidiene se poate scrie următoarea inegalitate:

$$(\mathbf{x} - \mathbf{y})^T(\mathbf{x} - \mathbf{y}) \geq 0 \Leftrightarrow \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 - 2\mathbf{x}^T\mathbf{y} \geq 0$$

și atunci rezultă următoarea estimare a divergenței dintre  $f(\mathbf{x})$  și  $g(\mathbf{x})$ :

$$I_{ED}(f, g) = \int f^2(\mathbf{x})d\mathbf{x} + \int g^2(\mathbf{x})d\mathbf{x} - 2 \int f(\mathbf{x})g(\mathbf{x})d\mathbf{x}.$$

Principe, Xu și Fisher [70] au arătat că această nouă măsură are aceleași minime și maxime ca și divergența KL, de aceea o poate înlocui pe cea din (6.18) în probleme care folosesc optimizarea informației mutuale.

Atunci informația mutuală exprimată pe baza acestei noi măsuri este divergența dintre  $p(\mathbf{x}, \mathbf{y})$  și  $p(\mathbf{x})p(\mathbf{y})$ :

$$\begin{aligned} I_{ED}(X, Y) &= \int_{\mathbf{x}} \int_{\mathbf{y}} p^2(\mathbf{x}, \mathbf{y})d\mathbf{x}d\mathbf{y} + \int_{\mathbf{x}} \int_{\mathbf{y}} p^2(\mathbf{x})p^2(\mathbf{y})d\mathbf{x}d\mathbf{y} - \\ &- 2 \int_{\mathbf{x}} \int_{\mathbf{y}} p(\mathbf{x}, \mathbf{y})p(\mathbf{x})p(\mathbf{y})d\mathbf{x}d\mathbf{y}. \end{aligned} \quad (6.19)$$

În mod evident,  $I_{ED}(X, Y) \geq 0$  și egalitatea are loc dacă și numai dacă  $X$  și  $Y$  sunt statistic independente, cu alte cuvinte aceste două proprietăți ale divergenței KL se păstrează și pentru  $I_{ED}(X, Y)$ .

Torkkola introduce în [84] aceeași formulă de estimare a informației mutuale folosind, însă, o argumentație diferită pe care o prezentăm în continuare.

Așa cum am văzut, informația mutuală poate fi privită ca o divergență KL dintre  $p(c, \mathbf{y})$  și produsul  $p(c)p(\mathbf{y})$ . Există o serie de măsuri alternative ale divergenței care pot fi utilizate pentru scopul propus în această secțiune. Una dintre acestea este divergența propusă de Kapur [84] pentru două distribuții discrete  $P$  și  $Q$ :

$$D(P, Q) = \frac{1}{\alpha(\alpha - 1)} \sum_{i=1}^n [p_i^\alpha - \alpha p_i q_i^{\alpha-1} + (\alpha - 1) q_i^\alpha], \quad \alpha \neq 0, \alpha \neq 1.$$

Pentru  $\alpha = 2$ , ignorând factorul  $\frac{1}{2}$  și extinzând măsura pentru densități continue, se obține:

$$D(f, g) = \int_{\mathbf{x}} (f(\mathbf{x}) - g(\mathbf{x}))^2 d\mathbf{x}. \quad (6.20)$$

Această măsură este întotdeauna pozitivă și ia valoarea 0 atunci când  $f(\mathbf{x}) = g(\mathbf{x})$ , acestea fiind și proprietățile divergenței KL. Kapur nu dă, însă, o justificare formală a condițiilor în care relația de mai sus este adevărată [84].

O altă justificare a folosirii măsurii pătratice este dată tot de Torkkola în [84], bazându-se următoarea divergență:

$$V(f, g) = \int_{\mathbf{x}} |f(\mathbf{x}) - g(\mathbf{x})| d\mathbf{x}.$$

Același autor arată că inegalitatea lui Pinsker dă o limită inferioară pentru  $K(f, g)$ :

$$K(f, g) \geq \frac{1}{2} V(f, g)^2. \quad (6.21)$$

Pe de altă parte, din cauză că  $f(\mathbf{x})$  și  $g(\mathbf{x})$  sunt densități de probabilitate, iau valori în intervalul  $[0, 1]$  și:

$$|f(\mathbf{x}) - g(\mathbf{x})| \geq (f(\mathbf{x}) - g(\mathbf{x}))^2,$$

de unde rezultă că:

$$V(f, g) \geq D(f, g). \quad (6.22)$$

Din inegalitățile (6.21) și (6.22) rezultă că a îl maximiza pe  $D(f, g)$  este echivalent cu a maximiza limita inferioară a lui  $K(f, g)$ .

Am văzut că pentru două variabile aleatoare continue  $Y_1$  și  $Y_2$  informația mutuală se exprimă ca divergență între  $p(\mathbf{y}_1, \mathbf{y}_2)$  și produsul  $p(\mathbf{y}_1)p(\mathbf{y}_2)$ . Introducând aceste valori în (6.20) rezultă:

$$I_T(Y_1, Y_2) = \int \int (p(\mathbf{y}_1, \mathbf{y}_2) - p(\mathbf{y}_1)p(\mathbf{y}_2))^2 d\mathbf{y}_1 d\mathbf{y}_2$$

care este echivalentă cu  $I_{ED}$ .

Torkkola [83] a arătat că această relație dezvoltată pentru două variabile aleatoare continue poate fi adaptată pentru a obține o estimare a informației mutuale dintre o variabilă aleatoare discretă  $C$  și o variabilă aleatoare continuă  $Y$ , caz în care (6.19) devine:

$$I_{ED}(C, Y) = \sum_c \int_{\mathbf{y}} p^2(c, \mathbf{y}) d\mathbf{y} + \sum_c \int_{\mathbf{y}} p^2(c) p^2(\mathbf{y}) d\mathbf{y} - 2 \sum_c \int_{\mathbf{y}} p(c, \mathbf{y}) p(c) p(\mathbf{y}) d\mathbf{y}. \quad (6.23)$$

Notăm cu  $M_p$  numărul de eșantioane din clasa  $p$ , cu  $\mathbf{y}_k$  eșantionul numărul  $k$  din întreg setul de eșantioane disponibile fără a ne interesa clasa careia îi aparține și cu  $\mathbf{y}_{pj}$  același eșantion arătând că el aparține clasei  $p$  și că are indicele  $j$  în clasă. Folosim următoarele egalități cunoscute:

$$p(c_p, \mathbf{y}) = P(c_p) p(\mathbf{y} | c_p)$$

și

$$p(\mathbf{y}) = \sum_{p=1}^M p(c_p, \mathbf{y}),$$

unde  $M$  este numărul de clase.

Rescriem (6.23) folosind aproximările densităților de probabilitate cu ferestre Parzen și rezultă:

$$\begin{aligned} I_{ED}(C, Y) &= \frac{1}{N^2} \sum_{p=1}^M \sum_{k=1}^{M_p} \sum_{l=1}^{M_p} G(\mathbf{y}_{pk} - \mathbf{y}_{pl}, 2\sigma^2 \mathbf{I}) + \\ &+ \frac{1}{N^2} \left( \sum_{p=1}^M \left( \frac{M_p}{N} \right)^2 \right) \sum_{k=1}^N \sum_{l=1}^N G(\mathbf{y}_k - \mathbf{y}_l, 2\sigma^2 \mathbf{I}) - \\ &- 2 \frac{1}{N^2} \sum_{p=1}^M \frac{M_p}{N} \sum_{j=1}^{M_p} \sum_{k=1}^N G(\mathbf{y}_{pj} - \mathbf{y}_k, 2\sigma^2 \mathbf{I}), \end{aligned} \quad (6.24)$$

unde  $N$  reprezintă numărul total de eșantioane.

## 6.5 Algoritmul Mutual Information Relevance LVQ (MIRLVQ)

În această secțiune introducem o nouă metodă de instruire a importanței componentelor vectorilor de intrare, obținută prin maximizarea informației mutuale dintre etichetele claselor și o transformare a vectorilor de intrare. Așa

cum am arătat în secțiunea precedentă, informația mutuală pătratică este o sumă de termeni de gradul doi, fiind dezvoltată pornind de la similaritatea cu entropia pătratică Renyi. Informația mutuală pătratică estimată cu ferestre Parzen oferă avantajul că se calculează ca sumă de gaussiene, fiecare dintre ele putând fi evaluată folosind perechi de eşantioane din setul de instruire.

Considerăm că avem o variabilă aleatoare  $X$ , iar  $\mathbf{x}_i$ ,  $i = 1, \dots, N$  este un set de realizări ale acestei variabile aleatoare. Realizările vor fi chiar vectorii de intrare. Considerăm o variabilă aleatoare discretă  $C$  și  $c_j$ ,  $j = 1, \dots, M$  sunt realizările sale care reprezintă etichetele claselor din care fac parte vectorii  $\mathbf{x}_i$ . Pentru ca instruirea să fie eficientă, trebuie ca fiecărei clase să i se asocieze mai mulți vectori de intrare.

### Algoritmul Weighted LVQ

Am arătat în capitolul 5 că actualizarea factorilor de relevanță prin algoritmi DSLVQ și RLVQ este motivată euristic. Din acest motiv, am introdus algoritmul OWA-RLVQ care calculează relevanțele ca ponderi ale unui operator OWA prin minimizarea distanței modificate

$$D_{ij}^* = \sqrt{\sum_{k=1}^n \lambda_k (|x_{ik} - w_{jk}|^*)^2},$$

unde notațiile sunt cele din secțiunea 5.4. Algoritmul MIRLVQ va folosi metrica ponderată

$$D_{ij} = \|\mathbf{x}_i - \mathbf{w}_j\|_\lambda = \sqrt{\sum_{k=1}^n \lambda_k (x_{ik} - w_{jk})^2}$$

care a fost introdusă în algoritmul RLVQ.

Includem aceste variații ale lui LVQ într-o clasă generală de algoritmi pe care o numim *Weighted LVQ*. Un algoritm *weighted LVQ* are următoarele caracteristici:

- a) Folosește o distanță ponderată;
- b) Relevanțele și prototipurile sunt actualizate simultan în timpul fazei de instruire. Actualizarea poate fi realizată incremental după prelucrarea fiecărui vector de instruire;
- c) Relevanțele care se obțin după încheierea algoritmului de instruire pot fi privite ca ranguri ale componentelor vectorilor de intrare și se pot utiliza pentru reducerea spațiului intrărilor (*feature selection*).

Un algoritm *weighted LVQ* constă din următorii pași:

**Pasul 1.** Găsește cel mai apropiat prototip  $\mathbf{w}_j$ , câștigătorul, pentru care valoarea lui  $\|\mathbf{x}_i - \mathbf{w}_j\|_\lambda$  este minimă pentru întreaga colecția de vectori prototip.

**Pasul 2.** Actualizează prototipul  $\mathbf{w}_j$  câștigător după un algoritm LVQ.

**Pasul 3.** Actualizează vectorul relevanțelor  $\lambda$ .

**Pasul 4.** Actualizează rangul fiecărei componente a vectorilor de intrare ca o medie a pașilor anteriori.

Ne propunem să transpunem algoritmi de tip *weighted LVQ* într-un cadru teoretic diferit. Introducem algoritmul MIRLVQ ca un algoritm din această clasă bazat pe maximizarea informației mutuale. O problemă majoră pe care o ridică această abordare este marea complexitate computațională pe care o implică estimarea informației mutuale. Pentru a evita acest dezavantaj, vom folosi metoda de estimare a informației mutuale cu ferestre Parzen.

### Adaptarea factorilor de relevanță

Notăm

$$\mathbf{y}_i = \lambda \mathbf{I}(\mathbf{x}_i - \mathbf{w}_j), \quad (6.25)$$

unde  $\mathbf{x}_i$  este vectorul de intrare  $i$ ,  $\mathbf{w}_j$  este prototipul  $j$  obținut cu un algoritm din categoria *weighted LVQ*,  $\lambda$  este vectorul relevanțelor și  $\mathbf{y}_i$  este vectorul de ieșire. Aceasta este transformarea pe care dorim să o implementăm,  $\lambda$  fiind un vector global care caracterizează toate perechile intrare-prototip.

Considerăm informația mutuală

$$I(c_j; \mathbf{y}_i) = I(c_j; \mathbf{x}_i, \mathbf{w}_j, \lambda),$$

unde  $\mathbf{y}_i$  sunt realizări ale variabilei aleatoare continue  $Y$ . În relația anterioară, am folosit două notații care au aceeași semnificație cu  $I(C, Y)$ . Această egalitate înglobează ipoteza formulată în relația (6.25), conform căreia ieșirea sistemului nostru se calculează folosind un vector de intrare  $\mathbf{x}_i$ , un prototip  $\mathbf{w}_j$  asociat lui calculat printr-un algoritm tip *weighted LVQ* și vectorul relevanțelor  $\lambda$  care este comun tuturor vectorilor din setul de instruire. Atunci când vectorul de intrare  $\mathbf{x}_i$  este asemănător prototipului  $\mathbf{w}_j$ , componentele vectorului  $\lambda$  trebuie întărite. Putem interpreta această similaritate prin faptul că informația mutuală dintre clasa  $c_j$  și valoarea agregată  $\mathbf{y}_i$  este mare. Pornind de la această observație, putem calcula relevanțele prin maximizarea informației mutuale și

folosim tehnica gradientului crescător. La fiecare pas de instruire  $t + 1$  vom avea:

$$\begin{aligned}\lambda^{(t+1)} &= \lambda^{(t)} + \alpha \frac{\partial I}{\partial \lambda} \Leftrightarrow \\ \lambda^{(t+1)} &= \lambda^{(t)} + \alpha \sum_{i=1}^n \frac{\partial I}{\partial \mathbf{y}_{ij}} \mathbf{I} \frac{\partial \mathbf{y}_{ij}}{\partial \lambda}.\end{aligned}\quad (6.26)$$

Vectorii  $\mathbf{x}_i$  și  $\mathbf{w}_j$  nu sunt dependenți de  $\lambda$ , de aceea:

$$\frac{\partial \mathbf{y}_i}{\partial \lambda} = \mathbf{x}_i - \mathbf{w}_j$$

și (6.26) devine:

$$\lambda^{(t+1)} = \lambda^{(t)} + \alpha \sum_{i=1}^n \frac{\partial I}{\partial \mathbf{y}_i} \mathbf{I} (\mathbf{x}_i - \mathbf{w}_j).\quad (6.27)$$

Trebuie să estimăm acum

$$\frac{\partial I}{\partial \mathbf{y}_i}$$

și pentru aceasta folosim măsura informației mutuale dintre o variabilă aleatoare continuă și una discretă dată de relația (6.23):

$$\begin{aligned}I(C, Y) &= \sum_{j=1}^M \int_{\mathbf{y}} p^2(c_j, \mathbf{y}) d\mathbf{y} + \sum_{j=1}^M \int_{\mathbf{y}} p^2(c_j) p^2(\mathbf{y}) d\mathbf{y} - \\ &\quad - 2 \sum_{j=1}^M \int_{\mathbf{y}} p(c_j, \mathbf{y}) p(c_j) p(\mathbf{y}) d\mathbf{y}.\end{aligned}$$

Utilizăm relația probabilităților condiționate:

$$p(\mathbf{y}, c_j) = P(c_j) p(\mathbf{y}|c_j),$$

unde  $P(c_j)$  este probabilitatea *a priori* a clasei  $c_j$ ,  $p(\mathbf{y}|c_j)$  este densitatea de probabilitate condiționată a vectorului  $\mathbf{y}$  care provine dintr-un vector de intrare ce aparține clasei  $c_j$ , iar

$$p(\mathbf{y}) = \sum_{j=1}^M p(c_j, \mathbf{y})$$

este densitatea de probabilitate a eșantionului  $\mathbf{y}$ . Analizând formula anterioară, trebuie să facem observația că un  $\mathbf{y}$  particular poate fi generat de diverși vectori de intrare  $\mathbf{x}_i$  care pot aparține oricărei clase. Probabilitatea condiționată arată clasa din care a provenit vectorul de intrare care a generat respectiva ieșire, fiind elementul care stabilește legătura dintre clasă și vectorul de ieșire.

Calculul gradientului așa cum este descris în formulele anterioare poate fi înlocuit cu un algoritm mult mai simplu, de gradient stohastic. Se folosesc două eșantioane aleatoare  $\mathbf{x}_1$  și  $\mathbf{x}_2$  reprezentative pentru întreaga bază de date [83].

Prin estimarea informației mutuale cu ferestre Parzen, din (6.24) se pot face ajustări ale lui  $\lambda$  în direcția gradientului pozitiv.

Vom trata separat următoarele două cazuri: cel în care eşantioanele sunt din aceeași clasă și cel în care eşantioanele sunt din clase diferite.

- a) În prima situație, putem considera că pentru clasa reprezentată de cele două eşantioane avem probabilitatea *a priori*  $P(c_1) = 1$ , iar pentru celelalte avem  $P(c_k) = 0$ ,  $k \neq 1$ . În relația (6.24) vom avea  $N = 2$ ,  $M = 2$ ,  $M_1 = 2$ ,  $M_2 = 0$  și vom scrie:

$$\begin{aligned}
I(C, Y) &= \frac{1}{2^2} \sum_{p=1}^2 \sum_{k=1}^{M_p} \sum_{l=1}^{M_p} G(\mathbf{y}_{pk} - \mathbf{y}_{pl}, 2\sigma^2 \mathbf{I}) + \\
&+ \frac{1}{2^2} \left( \sum_{p=1}^2 \left( \frac{M_p^2}{2} \right) \right) \sum_{k=1}^2 \sum_{l=1}^2 G(\mathbf{y}_k - \mathbf{y}_l, 2\sigma^2 \mathbf{I}) - \\
&- 2 \cdot \frac{1}{2^2} \cdot \sum_{p=1}^2 \frac{M_p}{2} \sum_{j=1}^{M_p} \sum_{k=1}^2 G(\mathbf{y}_{pj} - \mathbf{y}_k, 2\sigma^2 \mathbf{I}) = \\
&= \frac{1}{4} [G(\mathbf{y}_{11} - \mathbf{y}_{11}, 2\sigma^2 \mathbf{I}) + G(\mathbf{y}_{11} - \mathbf{y}_{12}, 2\sigma^2 \mathbf{I}) + \\
&\quad + G(\mathbf{y}_{12} - \mathbf{y}_{11}, 2\sigma^2 \mathbf{I}) + G(\mathbf{y}_{12} - \mathbf{y}_{12}, 2\sigma^2 \mathbf{I})] + \\
&+ \frac{1}{4} \cdot 1 \cdot [G(\mathbf{y}_1 - \mathbf{y}_1, 2\sigma^2 \mathbf{I}) + G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2 \mathbf{I}) + \\
&\quad + G(\mathbf{y}_2 - \mathbf{y}_1, 2\sigma^2 \mathbf{I}) + G(\mathbf{y}_2 - \mathbf{y}_2, 2\sigma^2 \mathbf{I})] - \\
&- 2 \cdot \frac{1}{4} \cdot [G(\mathbf{y}_{11} - \mathbf{y}_1, 2\sigma^2 \mathbf{I}) + G(\mathbf{y}_{11} - \mathbf{y}_2, 2\sigma^2 \mathbf{I}) + \\
&\quad + G(\mathbf{y}_{12} - \mathbf{y}_1, 2\sigma^2 \mathbf{I}) + G(\mathbf{y}_{12} - \mathbf{y}_2, 2\sigma^2 \mathbf{I})].
\end{aligned}$$

Această relație se poate simplifica ținând cont de observația că

$$\mathbf{y}_{11} = \mathbf{y}_1$$

și

$$\mathbf{y}_{12} = \mathbf{y}_2,$$

justificate prin faptul că toate eşantioanele fac parte din aceeași clasă, cea cu indicele 1. În aceste condiții, informația mutuală aproximată cu ferestre Parzen conform relației (6.24) devine:

$$\begin{aligned}
I(C, Y) &= \frac{1}{2} [G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2 \mathbf{I}) + G(0, 2\sigma^2 \mathbf{I})] + \\
&+ \frac{1}{2} [G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2 \mathbf{I}) + G(0, 2\sigma^2 \mathbf{I})] - \\
&- 2 \cdot \frac{1}{2} \cdot [G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2 \mathbf{I}) + G(0, 2\sigma^2 \mathbf{I})],
\end{aligned}$$

adică

$$I(C, Y) = 0.$$

Așadar, atunci când cele două eșantioane aparțin aceleiași clase, informația mutuală calculată mai sus este nulă. Introducând acest rezultat în (6.27) rezultă că:

$$\lambda^{(t+1)} = \lambda^{(t)},$$

însemnând că asupra relevanțelor nu apare nici o modificare.

- b) A doua situație este cea în care cele două eșantioane sunt din clase diferite și atunci probabilitățile *a priori* ale claselor din care acestea fac parte sunt  $P(c_1) = \frac{1}{2}$  și  $P(c_2) = \frac{1}{2}$ , iar  $P(c_k) = 0$  pentru  $k \notin \{1, 2\}$ . În relația (6.24) vom avea  $N = 2$ ,  $M = 2$ ,  $M_1 = 1$ ,  $M_2 = 1$  și vom scrie:

$$\begin{aligned} I(C, Y) &= \frac{1}{2^2} \sum_{p=1}^2 \sum_{k=1}^{M_p} \sum_{l=1}^{M_p} G(\mathbf{y}_{pk} - \mathbf{y}_{pl}, 2\sigma^2 \mathbf{I}) + \\ &+ \frac{1}{2^2} \left( \sum_{p=1}^2 \left( \frac{M_p^2}{2} \right) \right) \sum_{k=1}^2 \sum_{l=1}^2 G(\mathbf{y}_k - \mathbf{y}_l, 2\sigma^2 \mathbf{I}) - \\ &- 2 \cdot \frac{1}{2^2} \cdot \sum_{p=1}^2 \frac{M_p}{2} \sum_{j=1}^{M_p} \sum_{k=1}^2 G(\mathbf{y}_{pj} - \mathbf{y}_k, 2\sigma^2 \mathbf{I}) = \\ &= \frac{1}{4} [G(\mathbf{y}_{11} - \mathbf{y}_{11}, 2\sigma^2 \mathbf{I}) + G(\mathbf{y}_{21} - \mathbf{y}_{21}, 2\sigma^2 \mathbf{I})] + \\ &+ \frac{1}{4} \cdot \frac{1}{4} \cdot 2 \cdot [G(\mathbf{y}_1 - \mathbf{y}_1, 2\sigma^2 \mathbf{I}) + G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2 \mathbf{I}) + \\ &+ G(\mathbf{y}_2 - \mathbf{y}_1, 2\sigma^2 \mathbf{I}) + G(\mathbf{y}_2 - \mathbf{y}_2, 2\sigma^2 \mathbf{I})] - \\ &- 2 \cdot \frac{1}{4} \cdot \frac{1}{2} [G(\mathbf{y}_{11} - \mathbf{y}_1, 2\sigma^2 \mathbf{I}) + G(\mathbf{y}_{11} - \mathbf{y}_2, 2\sigma^2 \mathbf{I}) + \\ &+ G(\mathbf{y}_{21} - \mathbf{y}_1, 2\sigma^2 \mathbf{I}) + G(\mathbf{y}_{21} - \mathbf{y}_2, 2\sigma^2 \mathbf{I})]. \end{aligned} \quad (6.28)$$

Remarcăm că  $\mathbf{y}_{11} = \mathbf{y}_{21} = \mathbf{y}_1$ , și de această dată relația (6.28) se va scrie:

$$\begin{aligned} I(C, Y) &= \frac{1}{2} G(0, 2\sigma^2 \mathbf{I}) + \frac{1}{4} [G(0, 2\sigma^2 \mathbf{I}) + G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2 \mathbf{I})] \\ &- 2 \frac{1}{4} [G(0, 2\sigma^2 \mathbf{I}) + G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2 \mathbf{I})] = \\ &= \frac{1}{4} [G(0, 2\sigma^2 \mathbf{I}) - G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2 \mathbf{I})]. \end{aligned} \quad (6.29)$$

Având acest rezultat, putem stabili acum care este relația de actualizare a factorilor de relevanță prin introducerea lui (6.29) în (6.27). Trebuie să calculăm mai întâi:

$$\frac{\partial I(C, Y)}{\partial \mathbf{y}_1} = \frac{\partial}{\partial \mathbf{y}_1} \left( \frac{1}{4} (G(0, 2\sigma^2 \mathbf{I}) - G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2 \mathbf{I})) \right) =$$

$$\begin{aligned}
&= -\frac{1}{4} \frac{\partial G(0, 2\sigma^2 \mathbf{I})}{\partial \mathbf{y}_1} - \frac{1}{4} \frac{\partial G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2 \mathbf{I})}{\partial \mathbf{y}_1} = \\
&= -\frac{1}{8\sigma^2} \frac{\partial G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2 \mathbf{I})}{\partial \mathbf{y}_1} \tag{6.30}
\end{aligned}$$

și

$$\begin{aligned}
\frac{\partial I(C, Y)}{\partial \mathbf{y}_2} &= \frac{\partial}{\partial \mathbf{y}_2} \left( \frac{1}{4} (G(0, 2\sigma^2 \mathbf{I}) - G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2 \mathbf{I})) \right) = \\
&= -\frac{1}{4} \frac{\partial G(0, 2\sigma^2 \mathbf{I})}{\partial \mathbf{y}_2} - \frac{1}{4} \frac{\partial G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2 \mathbf{I})}{\partial \mathbf{y}_2} = \\
&= \frac{1}{8\sigma^2} \frac{\partial G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2 \mathbf{I})}{\partial \mathbf{y}_1}, \tag{6.31}
\end{aligned}$$

unde am folosit faptul că

$$G(0, 2\sigma^2 \mathbf{I}) = \frac{1}{(2\pi)^{\frac{d}{2}} |\sigma|^{\frac{1}{2}}}$$

este o constantă și derivata ei este nulă. Din (6.27), (6.30) și (6.31) rezultă că actualizarea factorilor de relevanță se face astfel:

$$\begin{aligned}
\lambda^{(t+1)} &= \lambda^{(t)} - \frac{\alpha}{8\sigma^2} G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2 \mathbf{I}) (\mathbf{y}_2 - \mathbf{y}_1) \mathbf{I} [(\mathbf{x}_1 - \mathbf{w}_{j(1)}) - \\
&\quad - (\mathbf{x}_2 - \mathbf{w}_{j(2)})], \tag{6.32}
\end{aligned}$$

unde prin  $\mathbf{w}_{j(1)}$  și  $\mathbf{w}_{j(2)}$  am notat prototipurile câștigătoare pentru vectorii  $\mathbf{x}_1$  și  $\mathbf{x}_2$ .

Am stabilit această regulă de modificare iterativă a relevanțelor prin maximizarea informației mutuale  $I(C, Y)$ . Vom dezvolta în secțiunile 6.6 și 6.7 o metodă similară de calcul al factorilor de relevanță prin maximizarea unei mărimi bazate pe energia informațională a variabilelor aleatoare. În final, vom detalia algoritmul care calculează în paralel atât prototipurile  $\mathbf{w}_j$  cât și relevanțele  $\lambda$ .

## 6.6 Energia informațională

Prezentăm în această secțiune definițiile și cele mai importante proprietăți ale energiei informaționale care este o mărime introdusă de O. Onicescu în [65].

Pentru variabila aleatoare continuă  $X$ , energia informațională se definește astfel:

$$E(X) = \int_{-\infty}^{+\infty} p^2(\mathbf{x}) d\mathbf{x},$$

unde  $p(\mathbf{x})$  este funcția de densitate de probabilitate a variabilei aleatoare.

În cazul în care avem o variabilă aleatoare discretă  $C$ , energia informațională se definește prin:

$$E(C) = \sum_{i=1}^n p^2(c_i),$$

unde distribuția aleatoare este definită prin

$$p(c_k) \geq 0$$

și

$$\sum_{k=1}^n p(c_k) = 1.$$

Data fiind această definiție, următorul rezultat este adevărat [71]:

**Propoziție.** Pentru orice repartiție  $p(c_k) \geq 0$ ,  $\sum_{k=1}^n p(c_k) = 1$ , între energia informațională și entropia Renyi are loc relația

$$H_{R_2}(C) = -\log E(C)$$

sau

$$E(C) = 2^{-H_{R_2}(C)}.$$

În consecință, entropia pătratică Renyi și energia informațională se determină reciproc și în mod unic.

Vom discuta proprietățile energiei informaționale discrete, iar cazul continuu este similar.

Considerăm o altă variabilă aleatoare discretă  $T$  pentru care avem probabilitățile  $p(t_k) \geq 0$ ,  $k = 1, \dots, m$  corespunzătoare celor  $m$  realizări posibile și  $\sum_{k=1}^m p(t_k) = 1$ . Atunci:

$$p(c_i, t_k) = p(t_k)p(c_i|t_k) = p(c_i)p(t_k|c_i).$$

Variabilele aleatoare  $C$  și  $T$  sunt independente dacă și numai dacă  $p(c_i, t_k) = p(c_i)p(t_k)$ .

Energia informațională medie a variabilei  $C$  condiționată de  $T$  se definește astfel:

$$E(C|T) = \sum_{k=1}^m p(t_k)E(C|t_k),$$

unde

$$E(C|t_k) = \sum_{i=1}^n p^2(c_i|t_k).$$

Cu aceste relații, energia informațională condiționată se mai poate scrie:

$$\begin{aligned} E(C|T) &= \sum_{i=1}^n \sum_{k=1}^m p(t_k) p^2(c_i|t_k) = \\ &= \sum_{i=1}^n \sum_{k=1}^m p(c_i, t_k) p(c_i|t_k). \end{aligned}$$

În [2] se arată că există două strategii care se pot adopta atunci când se studiază relațiile dintre două sisteme care interacționează: se poate măsura interdependența acestora ca atribut mutual sau se poate cuantifica măsura în care un sistem depinde de celălalt.

Prima strategie a fost abordată de S. Watanabe care a folosit entropia Shannon pentru a dezvolta o măsură de interdependență sau de coeziune [26].

A doua strategie poate fi adoptată prin introducerea unei măsuri de dependență pe baza energiei informaționale [2]. Această măsură definește o dependență unilaterală între două variabile aleatoare:

$$o(C, T) = E(C|T) - E(C)$$

și are următoarele proprietăți:

- a)  $o$  nu este simetrică în raport cu argumentele sale;
- b)  $o(C, T) \geq 0$ , egalitatea îndeplinindu-se dacă și numai dacă  $C$  și  $T$  sunt independente;
- c)  $o(C, T) \leq 1 - E(C)$ , iar egalitatea are loc dacă și numai dacă  $C$  este complet dependentă de  $T$ .

Atunci când avem două variabile aleatoare  $C$  și  $T$  astfel încât  $C$  este complet dependentă de  $T$  și  $T$  este complet dependentă de  $C$ , se spune că  $C$  și  $T$  sunt mutual complet dependente. Dacă  $C$  și  $T$  sunt două variabile aleatoare astfel încât  $o(C, T) = o(T, C)$ , se spune că  $C$  și  $T$  sunt egal unilateral dependente. Valoarea  $o(C, T)$  corespunde cantității de informație despre  $C$  pe care o înglobează  $T$ , putând fi privită ca un indicator al dependenței unilaterale ce caracterizează variabila aleatoare  $C$  în raport cu variabila aleatoare  $T$ . Se poate constata, așadar, că  $o(C, T)$  și  $I(C, T)$  măsoară practic același fenomen [2].

## 6.7 Algoritm Energy Relevance LVQ (ERLVQ)

Am arătat în secțiunea anterioară că există o similaritate între informația mutuală  $I(Y, X) = H(Y) - H(Y|X)$  și valoarea  $o(Y, X) = E(Y|X) - E(Y)$ ,

unde  $E(Y)$  este energia informațională a variabilei aleatoare  $Y$ . Inițial,  $o(Y, X)$  a fost definită pentru două variabile aleatoare discrete. Ea poate fi însă extinsă pentru variabile aleatoare continue sau variabile aleatoare continue împreună cu variabile aleatoare discrete.

Maximizăm  $o(Y, C)$  pentru a obține valorile relevanțelor:

$$\begin{aligned}\lambda^{(t+1)} &= \lambda^{(t)} + \alpha \frac{\partial o(Y, C)}{\partial \lambda} \Leftrightarrow \\ \Leftrightarrow \lambda^{(t+1)} &= \lambda^{(t)} + \alpha \sum_{i=1}^N \frac{\partial o(Y, C)}{\partial \mathbf{y}_i} \mathbf{I}(\mathbf{x}_i - \mathbf{w}_j).\end{aligned}\quad (6.33)$$

Folosim următoarele notații:  $N$  - numărul total de eșantioane,  $M$  - numărul de clase,  $M_p$  - numărul de eșantioane din clasa  $c_p$ .

Conform definiției, putem scrie:

$$o(Y, C) = E(Y|C) - E(Y),$$

unde

$$\begin{aligned}E(Y|C) &= \int_{\mathbf{y}} \sum_{p=1}^M p(c_p) p^2(\mathbf{y}|c_p) d\mathbf{y} = \\ &= \sum_{p=1}^M p(c_p) \int_{\mathbf{y}} p^2(\mathbf{y}|c_p) d\mathbf{y} = \\ &= \sum_{p=1}^M \frac{1}{p(c_p)} \int_{\mathbf{y}} p^2(\mathbf{y}, c_p) d\mathbf{y}\end{aligned}$$

și

$$E(Y) = \int_{\mathbf{y}} p^2(\mathbf{y}) d\mathbf{y}.$$

Obținem

$$o(Y|C) = \sum_{p=1}^M \frac{1}{p(c_p)} \int_{\mathbf{y}} p^2(\mathbf{y}, c_p) d\mathbf{y} - \int_{\mathbf{y}} p^2(\mathbf{y}) d\mathbf{y}. \quad (6.34)$$

Calculul acestei expresii necesită un efort substanțial din punct de vedere computațional, de aceea putem recurge la aproximarea cu ferestre Parzen a densităților pătratice de probabilitate care apar sub integrale. Folosim expresia gaussienei multidimensionale dată de relația (6.10), iar densitatea de probabilitate se poate aproxima prin (6.11).

Folosind aproximările cu ferestre Parzen avem:

$$\int_{\mathbf{y}} p^2(\mathbf{y}, c_p) d\mathbf{y} = \frac{1}{N^2} \sum_{k=1}^{M_p} \sum_{l=1}^{M_p} G(\mathbf{y}_{pk} - \mathbf{y}_{pl}, 2\sigma^2 \mathbf{I})$$

și

$$\int_{\mathbf{y}} p^2(\mathbf{y}) d\mathbf{y} = \frac{1}{N^2} \sum_{k=1}^N \sum_{l=1}^N G(\mathbf{y}_k - \mathbf{y}_l, 2\sigma^2 \mathbf{I}).$$

Cu aceste relații, ecuația (6.34) se scrie:

$$\begin{aligned} o(Y, C) &= \frac{1}{N^2} \left( \sum_{p=1}^M \frac{1}{M_p} \right) \sum_{k=1}^{M_p} \sum_{l=1}^{M_p} G(\mathbf{y}_{pk} - \mathbf{y}_{pl}, 2\sigma^2 \mathbf{I}) - \\ &\quad - \frac{1}{N^2} \sum_{k=1}^N \sum_{l=1}^N G(\mathbf{y}_k - \mathbf{y}_l, 2\sigma^2 \mathbf{I}), \end{aligned}$$

adică

$$\begin{aligned} o(Y, C) &= \frac{1}{N} \left( \sum_{p=1}^M \frac{1}{M_p} \right) \sum_{k=1}^{M_p} \sum_{l=1}^{M_p} G(\mathbf{y}_{pk} - \mathbf{y}_{pl}, 2\sigma^2 \mathbf{I}) - \\ &\quad - \frac{1}{N^2} \sum_{k=1}^N \sum_{l=1}^N G(\mathbf{y}_k - \mathbf{y}_l, 2\sigma^2 \mathbf{I}). \end{aligned} \quad (6.35)$$

Pentru a dezvolta formula de actualizare a relevanțelor, folosim ideea din [83] conform căreia putem considera două eșantioane  $\mathbf{y}_1$  și  $\mathbf{y}_2$  consecutive ca reprezentante ale claselor. Avem următoarele două cazuri:

- a) Când cele două eșantioane aparțin aceleiași clase,  $N = 2$ ,  $M = 2$ ,  $M_1 = 2$  și  $M_2 = 0$ , deci primul termen al ecuației (6.35) nu se poate calcula datorită nedeterminării introduse de valoarea lui  $M_2$ . Acest caz este, așadar, ignorat.
- b) Atunci când cele două eșantioane sunt din clase diferite,  $N = 2$ ,  $M = 2$ ,  $M_1 = 1$  și  $M_2 = 1$ . Ecuația (6.35) se poate scrie astfel:

$$\begin{aligned} o(Y, C) &= \frac{1}{2} \left( \sum_{p=1}^2 \frac{1}{M_p} \right) \sum_{k=1}^{M_p} \sum_{l=1}^{M_p} G(\mathbf{y}_{pk} - \mathbf{y}_{pl}, 2\sigma^2 \mathbf{I}) - \\ &\quad - \frac{1}{4} \sum_{k=1}^2 \sum_{l=1}^2 G(\mathbf{y}_k - \mathbf{y}_l, 2\sigma^2 \mathbf{I}) = \\ &= \frac{1}{2} [G(\mathbf{y}_{11} - \mathbf{y}_{11}, 2\sigma^2 \mathbf{I}) + G(\mathbf{y}_{21} - \mathbf{y}_{21}, 2\sigma^2 \mathbf{I})] - \\ &\quad - \frac{1}{4} [G(\mathbf{y}_1 - \mathbf{y}_1, 2\sigma^2 \mathbf{I}) + G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2 \mathbf{I}) + \\ &\quad + G(\mathbf{y}_2 - \mathbf{y}_1, 2\sigma^2 \mathbf{I}) + G(\mathbf{y}_2 - \mathbf{y}_2, 2\sigma^2 \mathbf{I})] = \\ &= G(0, 2\sigma^2 \mathbf{I}) - \frac{1}{2} [G(0, 2\sigma^2 \mathbf{I}) + G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2 \mathbf{I})]. \end{aligned}$$

În final relația devine:

$$o(Y, C) = G(0, 2\sigma^2\mathbf{I}) - \frac{1}{2}G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2\mathbf{I}). \quad (6.36)$$

Conform egalității (6.34), va trebui să calculăm derivatele parțiale ale lui  $o(Y, C)$  în funcție de  $\mathbf{y}_1$  și apoi de  $\mathbf{y}_2$ . Pentru aceasta, folosim relația următoare:

$$\begin{aligned} \frac{\partial}{\partial \mathbf{y}_i} G(\mathbf{y}_i - \mathbf{y}_j, 2\sigma^2\mathbf{I}) &= G(\mathbf{y}_i - \mathbf{y}_j, 2\sigma^2\mathbf{I}) \frac{\mathbf{y}_j - \mathbf{y}_i}{2\sigma^2} = \\ &= -\frac{\partial}{\partial \mathbf{y}_j} G(\mathbf{y}_i - \mathbf{y}_j, 2\sigma^2\mathbf{I}). \end{aligned}$$

Atunci:

$$\begin{aligned} \frac{\partial o(Y, C)}{\partial \mathbf{y}_1} &= \frac{\partial}{\partial \mathbf{y}_1} \left( G(0, 2\sigma^2\mathbf{I}) - \frac{1}{2}G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2\mathbf{I}) \right) = \\ &= \frac{\partial G(0, 2\sigma^2\mathbf{I})}{\partial \mathbf{y}_1} - \frac{1}{2} \frac{\partial G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2\mathbf{I})}{\partial \mathbf{y}_1}. \end{aligned}$$

Deoarece

$$G(0, 2\sigma^2\mathbf{I}) = \frac{1}{(2\pi)^{\frac{d}{2}} |\sigma|^{\frac{1}{2}}}$$

este o constantă, avem:

$$\begin{aligned} \frac{\partial o(Y, C)}{\partial \mathbf{y}_1} &= -\frac{1}{2} \frac{\partial G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2\mathbf{I})}{\partial \mathbf{y}_1} = \\ &= -\frac{1}{2} G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2\mathbf{I}) \cdot \frac{\mathbf{y}_2 - \mathbf{y}_1}{2\sigma^2}. \end{aligned}$$

Similar obținem:

$$\begin{aligned} \frac{\partial o(Y, C)}{\partial \mathbf{y}_2} &= \frac{\partial}{\partial \mathbf{y}_2} \left( G(0, 2\sigma^2\mathbf{I}) - \frac{1}{2}G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2\mathbf{I}) \right) = \\ &= -\frac{1}{2} \frac{\partial G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2\mathbf{I})}{\partial \mathbf{y}_2} = \\ &= \frac{1}{2} G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2\mathbf{I}) \cdot \frac{\mathbf{y}_2 - \mathbf{y}_1}{2\sigma^2}. \end{aligned}$$

Înlocuind aceste rezultate în (6.33) rezultă:

$$\begin{aligned} \lambda^{(t+1)} &= \lambda^{(t)} + \alpha \left( -\frac{1}{2} G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2\mathbf{I}) \cdot \frac{\mathbf{y}_2 - \mathbf{y}_1}{2\sigma^2} \mathbf{I}(\mathbf{x}_1 - \mathbf{w}_{j(1)}) + \right. \\ &\quad \left. + \frac{1}{2} G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2\mathbf{I}) \cdot \frac{\mathbf{y}_2 - \mathbf{y}_1}{2\sigma^2} \mathbf{I}(\mathbf{x}_2 - \mathbf{w}_{j(2)}) \right) \Rightarrow \\ \Rightarrow \lambda^{(t+1)} &= \lambda^{(t)} - \alpha \left( \frac{1}{2} G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2\mathbf{I}) \cdot \frac{\mathbf{y}_2 - \mathbf{y}_1}{2\sigma^2} \mathbf{I}(\mathbf{x}_1 - \mathbf{w}_{j(1)}) - \right. \\ &\quad \left. - \frac{1}{2} G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2\mathbf{I}) \cdot \frac{\mathbf{y}_2 - \mathbf{y}_1}{2\sigma^2} \mathbf{I}(\mathbf{x}_2 - \mathbf{w}_{j(2)}) \right), \end{aligned}$$

sau

$$\lambda^{(t+1)} = \lambda^{(t)} - \alpha \frac{1}{4\sigma^2} G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2 \mathbf{I}) \cdot (\mathbf{y}_2 - \mathbf{y}_1) \mathbf{I} (\mathbf{x}_1 - \mathbf{w}_{j(1)} - \mathbf{x}_2 + \mathbf{w}_{j(2)}), \quad (6.37)$$

unde  $\mathbf{w}_{j(1)}$  și  $\mathbf{w}_{j(2)}$  sunt prototipurile câștigătoare pentru vectorii de intrare  $\mathbf{x}_1$  și  $\mathbf{x}_2$ .

Următoarea procedură actualizează simultan prototipurile și factorii de relevanță:

**Pasul 1.** Inițializează  $\eta$  și  $\alpha$ . Inițializează vectorul relevanțelor:

$$\lambda_k = \frac{1}{n}, \quad k = 1, \dots, n.$$

**Pasul 2.** Inițializează prototipurile.

**Pasul 3.** Actualizează vectorii prototip conform regulii LVQ modificate folosind distanța  $D_{ij}$ :

$$\mathbf{w}_j = \begin{cases} \mathbf{w}_j + \eta \lambda \mathbf{I} (\mathbf{x}_i - \mathbf{w}_j), & \text{dacă } \mathbf{x}_i \text{ a fost clasificat corect} \\ \mathbf{w}_j - \eta \lambda \mathbf{I} (\mathbf{x}_i - \mathbf{w}_j), & \text{în caz contrar} \end{cases}$$

**Pasul 4.** Actualizează relevanțele:

$$\lambda^{(t+1)} = \lambda^{(t)} - \alpha \frac{1}{4\sigma^2} G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2 \mathbf{I}) \cdot (\mathbf{y}_2 - \mathbf{y}_1) \mathbf{I} \cdot (\mathbf{x}_1 - \mathbf{w}_{j(1)} - \mathbf{x}_2 + \mathbf{w}_{j(2)}).$$

**Pasul 5.** Transformă relevanțele:

$$\lambda_k = \frac{e^{\lambda_k}}{\sum_{i=1}^n e^{\lambda_i}}$$

pentru  $k = 1, \dots, n$ .

**Pasul 6.** Calculează rangul fiecărei componente a vectorului de intrare.

**Pasul 7.** Repetă pașii 3-6 pentru fiecare vector de intrare.

Remarcăm asemănarea algoritmilor MIRLVQ și ERLVQ care, cu toate că sunt dezvoltați diferit, conduc la rezultate similare. Diferența este minoră și apare la pasul 4, formula (6.37) putând fi înlocuită cu (6.32). Rezultatele experimentale sunt identice deoarece constanta care face diferență între cele două variante introduce doar o scalare a ratei de instruire  $\alpha$ .

## 6.8 Experimente

Am testat algoritmi MIRLVQ și ERLVQ în aceleași condiții ca și OWA-RLVQ pentru a putea face o comparație obiectivă a rezultatelor și a performanțelor lor. Deoarece MIRLVQ și ERLVQ sunt similari, rezultatele înregistrate în tabelele 6.1, 6.2, 6.3 și 6.4 sunt cele pentru ERLVQ. Ele sunt identice cu cele obținute cu MIRLVQ, cu observația că în primul caz valoarea constantei  $\alpha$  trebuie înjumătățită. Importanța componentelor poate fi diferită de cea obținută pentru OWA-RLVQ din următorul motiv: OWA-RLVQ este focalizat pe o separare cât mai bună a centrelor clusterelor. ERLVQ încearcă, prin definiție, o distincție cât mai bună la limita dintre clustere deoarece folosește întotdeauna câte două eșantioane consecutive din clase diferite. Așadar, OWA-RLVQ găsește acele componente care sunt importante pentru accentuarea centrelor clusterelor, iar ERLVQ pe acelea care fac o discriminare cât mai bună între două clase adiacente. Aplicația C++ care implementează algoritmul ERLVQ este asemănătoare cu cea utilizată pentru testele din capitolul 5.

Am folosit mai întâi baza de date Iris [9] pentru care am instruit 6 prototipuri, identic cu testele pentru OWA-RLVQ, și am obținut o rată de recunoaștere de 97,33%, iar valorile parametrilor au fost  $\eta = 0,15$  și  $\alpha = 20$ . Pentru algoritmul MIRLVQ se folosește  $\alpha = 40$ . Rangurile componentelor vectoriale pe care le-am obținut pentru această bază de date cu algoritmi RLVQ, OWA-RLVQ și ERLVQ sunt în tabelul 6.1. Relevanțele componentelor vectoriale pentru ERLVQ sunt  $[0,714 \ 0,498 \ 0,475 \ 0,122]^t$ . Rezultatele arată că atunci când se utilizează baza de date Iris, ultima componentă vectorială este cea mai importantă dacă ne propunem să construim clase ale căror centroizi sunt foarte accentuați, în timp ce dacă dorim să facem o discriminare cât mai bună între cele trei clase, prima componentă vectorială este cea mai semnificativă. Această concluzie este asemănătoare celor prezentate în [15] unde se arată că atunci când se folosesc doar primele două componente vectoriale se obține o proiecție în planul bidimensional a claselor "1" și "2" cu centroizii relativ apropiați, dar cu clasele delimitate clar de o dreaptă, în timp ce atunci când se folosesc ultimele două componente, centroizii celor două clase proiectate în planul bidimensional sunt foarte depărtați, dar clasele se întrepătrund ușor, astfel că limita nu este atât de clară.

Pentru baza de date Deterding (Vowel Recognition) am păstrat același număr de 59 de prototipuri ca și în testele anterioare. De această dată, am obținut o rată de recunoaștere de 47,18% față de 46,75% pentru OWA-RLVQ și 46,32% pentru RLVQ. Parametrii de instruire au fost  $\eta = 0,35$ ,  $\alpha = 200$ . Pentru algoritmul MIRLVQ se folosește  $\alpha = 400$ . În tabelul 6.2 avem rangul

Tabelul 6.1: Rangul componentelor vectoriale pentru baza de date Iris.

| Rangul                | 1 | 2 | 3 | 4 |
|-----------------------|---|---|---|---|
| Componentele RLVQ     | 4 | 2 | 3 | 1 |
| Componentele OWA-RLVQ | 4 | 3 | 2 | 1 |
| Componentele ERLVQ    | 1 | 2 | 3 | 4 |

Tabelul 6.2: Rangul componentelor vectoriale pentru baza de date Deterding.

| Rangul                | 1 | 2 | 3 | 4  | 5  |
|-----------------------|---|---|---|----|----|
| Componentele RLVQ     | 2 | 5 | 1 | 9  | 6  |
| Componentele OWA-RLVQ | 8 | 2 | 4 | 5  | 6  |
| Componentele ERLVQ    | 2 | 1 | 3 | 4  | 6  |
| Rangul                | 6 | 7 | 8 | 9  | 10 |
| Componentele RLVQ     | 3 | 4 | 8 | 7  | 10 |
| Componentele OWA-RLVQ | 9 | 3 | 7 | 1  | 10 |
| Componentele ERLVQ    | 8 | 9 | 5 | 10 | 7  |

componentelor vectoriale pentru experimentele cu această bază de date și componenta 2 este cea mai importantă atât pentru RLVQ cât și pentru ERLVQ, în timp OWA-RLVQ a plasat-o pe locul 2. Componenta 7 are importanța cea mai mică pentru algoritmul ERLVQ, iar RLVQ și OWA-RLVQ o plasează pe penultimul, respectiv antepenultimul loc. Vectorul relevanțelor pe care l-am obținut la testul prezentat aici este  $[0,354 \ 0,439 \ 0,349 \ 0,310 \ 0,271 \ 0,299 \ 0,253 \ 0,287 \ 0,284 \ 0,268]^t$ .

În al treilea set de experimente am folosit baza de date Ionosphere [9] care este formată din vectori cu un număr mare de componente. Pentru a reprezenta cele două clase, am folosit 8 prototipuri așa cum am făcut și în experimentele cu OWA-RLVQ. Rata de recunoaștere pe care am obținut-o a fost de 94,03%, spre deosebire de 93,37% cu OWA-RLVQ și de 92,71% cu RLVQ. Valorile constantelor de instruire au fost  $\eta = 0,45$ ,  $\alpha = 0,9 \cdot 10^{13}$ , iar pentru algoritmul MIRLVQ se folosește  $\alpha = 1,8 \cdot 10^{13}$ . Componenta vectorială cea mai importantă stabilită de acest algoritm este cea cu numărul 8. Componenta numărul 12 este plasată între primele 5 atât de ERLVQ cât și de RLVQ și de OWA-RLVQ, după cum se poate constata din tabelul 6.3.

Comparația ratelor de recunoaștere pentru algoritmi LVQ1, RLVQ, OWA-RLVQ și ERLVQ din tabelul 6.4 arată că ultimul dintre ei oferă cele mai bune

Tabelul 6.3: Rangul componentelor vectoriale pentru baza de date Ionosphere. Sunt reprezentate cele mai importante 5 componente.

| Rangul                | 1  | 2  | 3  | 4  | 5  |
|-----------------------|----|----|----|----|----|
| Componentele RLVQ     | 20 | 28 | 26 | 12 | 6  |
| Componentele OWA-RLVQ | 14 | 12 | 1  | 3  | 28 |
| Componentele ERLVQ    | 8  | 24 | 16 | 12 | 6  |

rezultate pentru toate cele trei baze de date testate. Pentru baza de date Deterding care reprezintă o cunoscută problemă de recunoașterea vorbirii, creșterea de performanță obținută cu ERLVQ este de aproape 1%. Această metodă oferă rezultate satisfăcătoare și pentru probleme generale de recunoaștere a formelor cum sunt Iris și Ionosphere pentru care performanțele rămân superioare.

Tabelul 6.4: Comparație între ratele de recunoaștere obținute cu LVQ1, RLVQ, OWA-RLVQ și ERLVQ.

| Baza de date | LVQ1   | RLVQ   | OWA-RLVQ | ERLVQ  |
|--------------|--------|--------|----------|--------|
| Iris         | 91,33% | 95,33% | 96,66%   | 97,33% |
| Deterding    | 44,80% | 46,32% | 46,75%   | 47,18% |
| Ionosphere   | 90,06% | 92,71% | 93,37%   | 94,03% |

## 6.9 Concluzii

Am introdus în acest capitol o nouă metodă de calcul al relevanțelor pentru componentele vectorilor de intrare într-un sistem de recunoaștere a formelor. Cei doi algoritmi, MIRLVQ și ERLVQ, sunt similari dar primul este dezvoltat prin maximizarea entropiei pătratică Renyi, iar al doilea maximizează o măsură bazată pe energia informațională a lui Onicescu. MIRLVQ și ERLVQ reprezintă o tehnică alternativă algoritmului OWA-RLVQ pentru calculul relevanțelor.

Algoritmii propuși de noi oferă o altă perspectivă pentru recunoașterea formelor, calculând relevanțele prin folosirea ca suport matematic al puter-nicului aparat reprezentat de teoria informației. Metoda noastră stabilește on-line valorile relevanțelor prin continua adaptare la fluxul de date de instruire. Relevanțele pot fi folosite la ordonarea în mod dinamic a componentelor vectorilor de intrare în funcție de importanța lor, dar și la reducerea dimensiunii vectorilor prin renunțarea la componentele mai puțin importante.

# Capitolul 7

## Concluzii

### 7.1 Introducere

În acest capitol facem o sinteză a rezultatelor teoretice și experimentale obținute în lucrare. Secțiunea 7.2 prezintă contribuțiile pe care teza le aduce la dezvoltarea domeniului recunoașterii formelor folosind rețele neurale, iar secțiunea 7.3 definește posibilele direcții viitoare de dezvoltare a rezultatelor obținute aici.

### 7.2 Contribuții

Lucrarea abordează următoarele două probleme:

- Cum se poate transfera conceptul de concurență la nivelul unei colecții de rețele neurale pornind de la ideea de modularitate?
- Care este impactul contribuției individuale a componentelor intrărilor asupra performanțelor unor sisteme de recunoaștere a formelor bazate pe rețele neurale?

Modelul rețelelor neurale concurente (CNN) din capitolul 4 este o soluție consistentă pentru prima problemă, iar algoritmi OWA-RLVQ, MIRLVQ și ERLVQ din capitolele 5 și 6 reprezintă răspunsul nostru la a doua întrebare.

Contribuțiile pe care le aduce această lucrare se regăsesc în capitolele 3, 4, 5 și 6.

În capitolul 2 am descris modelele neurale fundamentale și principalele tehnici de instruire. Din punct de vedere istoric, neuronul McCulloch-Pitts este primul model neural care a fost descris în literatură și pe care l-am prezentat în secțiunea 2.2. Acest neuron cu o structură foarte simplă are un mare potențial de calcul, permițând implementarea funcțiilor logice cu complexitate

oricât de mare. Perceptronul lui Rosenblatt prezentat în secțiunea 2.3 realizează un pas important înainte pentru că are capacitatea de a învăța să separe vectori liniar separabili. Dacă perceptronii sunt organizați pe straturi, se poate realiza o structură complexă numită perceptron multistrat. Algoritmul de instruire *backpropagation of error* detaliat de noi în paragraful 2.7 este unul dintre cei mai cunoscuți și permite perceptronului multistrat să învețe supervizat asocieri între domeniul intrărilor și cel al ieșirilor chiar și pentru clase neseperabile liniar. Rețelele neurale cu întârziere în timp descrise în secțiunea 2.8 sunt o extensie a perceptronului multistrat și pot fi folosite pentru forme temporale. Aceste rețele neurale au fost introduse în special pentru aplicații de recunoaștere a vorbirii care, prin construcția lor, înglobează factorul timp. Rețelele neurale cu autoorganizare pe care le-am analizat în secțiunea 2.9 au capacitatea de a extrage și de a înmagazina trăsăturile esențiale ale formelor de instruire și, pe baza acestora, de a realiza clasificarea. Învățarea competitivă este suportul unei clase de algoritmi care promovează selecția celor mai performanți neuroni ai unei rețele pentru care se aplică instruirea. Hărțile cu autoorganizare fac nesupervizat o proiecție bidimensională a vectorilor de instruire. Algoritmii LVQ construiesc supervizat o mulțime de vectori care caracterizează întreg setul de instruire, fiind etichetați cu numele clasei pe care o reprezintă.

### **Recunoașterea vorbirii cu modele clasice de rețele neurale**

Este dificil să comparăm performanțele mai multor tipuri de rețele neurale atâta timp cât testele sunt realizate pe date diferite. Din acest motiv, nu putem trage o concluzie relevantă prin simpla alăturare a unor rezultate disponibile în diverse lucrări. Am preferat să facem un studiu comparativ al calității recunoașterii cu perceptroni multistrat, rețele neurale cu întârziere în timp și cu hărți cu autoorganizare a unui set de cuvinte izolate, pe de o parte, și a vorbitorilor dintr-un grup de subiecți, pe de altă parte. Capitolul 3 detaliază metoda de lucru și rezultatele pe care le-am obținut. Am folosit bazele de date Speechdata și Telephdata care conțin pronunții în limba engleză ale tuturor cifrelor, prima dată înregistrate într-un studio special amenajat și a doua oară înregistrate prin telefon. Datorită limitărilor pe care le produce transmisia pe liniile telefonice, a doua bază de date ridică dificultăți suplimentare în recunoaștere. Metoda de preprocesare a semnalelor vocale descrisă în secțiunea 3.2 este binecunoscută. Aceasta constă din calculul coeficienților cepstrali pe ferestre de semnal care sunt suprapuse pentru a surprinde evoluția temporală a semnalelor vocale. Experimental am stabilit că aceste modele neurale nu ne conduc la rezultate satisfăcătoare chiar dacă pentru testele de recunoaștere a

cuvintelor izolate avem rate de succes bune. Am putea argumenta că problemele de recunoaștere a vorbitorilor din Speechdata și Telephdata pe care ne-am propus să le rezolvăm sunt dificile însă, cu toate acestea, rezultatele pe care le-am grupat în tabelul 3.2 din paragraful 3.3 nu sunt suficiente.

## Rețele neurale concurente

Este evident că pe măsură ce o rețea neurală are de rezolvat o problemă din ce în ce mai puțin dificilă ea conduce la performanțe mai bune. Pentru a reduce complexitatea recunoașterii printr-o singură rețea neurală a întregii colecții de cuvinte izolate sau a tuturor vorbitorilor, am ales soluția unei rețele neurale modulare formată din rețele neurale specializate pe subprobleme ale problemei inițiale. În capitolul 4 am introdus noul model de clasificare al Rețelelor Neurale Concurente (*Concurrent Neural Networks - CNN*) ca o colecție de rețele neurale de dimensiuni reduse care lucrează în paralel, iar clasificarea se face conform regulii *câștigătorul ia tot*.

Instruirea rețelelor neurale concurente se face conform schemei din figura 4.1 care presupune că fiecare modul este instruit cu propriul set de date. Sistemul este alcătuit din rețele neurale cu diverse arhitecturi. Noi am folosit module de tip perceptron multistrat, de tip rețea neurală cu întârziere în timp și de tip hartă cu autoorganizare, însă se pot folosi și alte variante. Schema de recunoaștere din figura 4.2 constă dintr-o colecție de module instruite pe câte o subproblemă și un modul care selectează cel mai bun răspuns. Algoritmii de instruire și de recunoaștere din secțiunea 4.2 implementează aceste două tehnici care sunt particularizate în paragraful 4.3 pentru perceptroni multistrat, rețele neurale cu întârziere în timp și hărți cu autoorganizare. MLP-CNN și TDNN-CNN folosesc module instruite supervizat și seturile de vectori de instruire conțin atât exemple pozitive cât și negative. Spre deosebire de acestea, SOM-CNN constă din module care sunt instruite printr-un algoritm nesupervizat și datele constau doar din exemple pozitive. Secțiunea 4.3 oferă detaliile legate de testele de recunoaștere a vorbitorilor din bazele de date Speechdata și Telephdata cu rețele neurale concurente. În comparație cu rezultatele din capitolul 3, de data aceasta progresele sunt substanțiale iar performanțele sunt, în medie, cu 30%-35% mai bune.

Remarcăm comportamentul foarte bun al SOM-CNN pe care îl folosim mai departe în recunoașterea unor vocale din limba română. În acest scop, am realizat o bază de date care constă din secvențe vocalizate extrase din pronunții în limba română ale celor zece cifre. Am preprocesat aceste semnale vocale prin aceeași metodă folosită pentru Speechdata și Telephdata pe care am descri-

o în paragraful 3.2. Energia semnalului calculată pe ferestre oferă o primă informație legată de prezența sau absența unei vocale și ne permite să selectăm doar aceste zone. Identitatea vocalelor este stabilită apoi de o rețea neurală concurentă SOM-CNN.

Am realizat cinci tipuri de experimente diferite pentru a ne forma o imagine corectă asupra modelului propus de noi în acest capitol și pentru o analiză detaliată a sa. Primul experiment a fost cel de recunoaștere a celor șase vocale și semivocale cu care am instruit sistemul. Rezultatele experimentale au arătat în mod clar performanțele foarte bune ale SOM-CNN pentru care am obținut un progres de peste 10% față de SOM, până la 84,21%. Celelalte experimente au analizat elemente foarte importante care contribuie la calitatea recunoașterii: proporția în care dimensiunile reale ale zonelor vocalizate și nevocalizate sunt sesizate de SOM și SOM-CNN, numărul de succesiuni vocală-semivocală recunoscute corect, proporția în care sistemul sesizează prezența unei consoane în locul unei vocale sau a unei semivocale și procentul în care sunt confundate vocalele. De fiecare dată când am folosit rețelele neurale concurente rezultatele au fost superioare. Putem considera că modelul propus de noi în capitolul 4 este o bună alternativă pentru sistemele clasice de recunoaștere a formelor.

### **Recunoaștere prin calculul relevanțelor intrărilor folosind operatori OWA**

În capitolele 5 și 6 analizăm problema impactului contribuției individuale a componentelor intrărilor asupra performanțelor unor sisteme de recunoaștere a formelor care folosesc rețele neurale.

Importanțele componentelor vectoriale ale intrărilor într-un sistem de recunoaștere a formelor, pe care le-am numit relevanțe sau factori de relevanță, sunt calculate în capitolul 5 ca ponderi ale unor operatori speciali de agregare numiți operatori OWA. Am prezentat această clasă de operatori în secțiunea 5.3, am subliniat proprietățile lor matematice și am prezentat o tehnică de instruire supervizată a lor. Metoda introdusă de noi prin algoritmul OWA-RLVQ folosește relevanțele pentru a defini un nou algoritm LVQ în care prototipurile sunt modificate după o regulă ce înlocuiește distanța euclidiană cu una care înglobează ponderile OWA. În paralel cu actualizarea prototipurilor sunt calculați și factorii de relevanță după o regulă descrisă în detaliu în paragraful 5.4. Această regulă se bazează pe ideea că un vector clasificat corect induce o mărire a relevanțelor, în timp ce unul clasificat incorect provoacă diminuarea acestor factori.

Am experimentat algoritmul OWA-RLVQ pe baze de date standard, dispo-

nibile în rețeaua Internet, pentru recunoașterea formelor, după cum se poate vedea în paragraful 5.5. Am obținut rezultate superioare atât față de algoritmul LVQ standard cât și față de algoritmul RLVQ care calculează, de asemenea, relevanțele, însă după o metodă motivată euristic.

Algoritmul OWA-RLVQ este important deoarece face legătura între algoritmul RLVQ și consistentul model matematic al operatorilor OWA. Metoda noastră are rezultate bune atât pentru probleme de recunoaștere vocală cât și pentru probleme generale de recunoaștere a formelor. Relevanțele pot fi folosite și pentru stabilirea dinamică a rangului fiecărei dimensiuni a vectorilor de intrare.

### **Recunoaștere prin calculul relevanțelor intrărilor folosind noțiuni de teoria informației**

În capitolul 6 introducem o nouă metodă de recunoaștere a formelor care folosește factori de relevanță calculați prin maximizarea fie a informației mutuale dintre intrări și clase, fie a unei mărimi bazate pe energia informației. Algoritmul ERLVQ, care se dovedește a fi similar cu MIRLVQ, este o alternativă la OWA-RLVQ deoarece calculează relevanțele care conduc la o bună delimitare a claselor, în timp ce al doilea algoritm calculează relevanțele care reorientează prototipurile claselor către centroizii acestora.

Informația mutuală se calculează ca o diferență a două entropii Shannon. Este cunoscută, însă, complexitatea de calcul pe care o implică estimarea entropiei Shannon prin eșantioane. Pentru a evita acest dezavantaj, am folosit o metodă recentă de aproximare a informației mutuale prin estimarea densităților de probabilitate cu ferestre Parzen. Această metodă pe care am tratat-o în detaliu în secțiunea 6.3 folosește nuclee gaussiene multidimensionale. Una dintre proprietățile lor remarcabile este că integrala produsului a două nuclee gaussiene este tot o gaussiană. În anexa B am demonstrat această proprietate. Estimarea densității de probabilitate cu ferestre Parzen se folosește pentru estimarea entropiei Renyi care include și entropia Shannon ca un caz particular. Printr-o similaritate cu estimarea entropiei Renyi se poate estima informația mutuală sub forma informației mutuale pătratice, după cum arătăm în paragraful 6.4.

Odată stabilit acest cadru formal, am dezvoltat algoritmul MIRLVQ printr-o metodă de gradient crescător pentru maximizarea informației mutuale dintre variabila aleatoare continuă a intrărilor și variabila aleatoare discretă a claselor. Pentru aceasta, mai întâi am introdus în secțiunea 6.5 algoritmul general *Weighted LVQ* care sintetizează etapele parcurse de un algoritm ce folosește relevanțele pentru clasificarea prin LVQ. OWA-RLVQ este o particularizare a

acestui algoritm, iar MIRLVQ se bazează tot pe acest șablon.

Elementul cheie al lui MIRLVQ este faptul că el actualizează relevanțele prin compararea apartenenței la aceeași clasă a doi vectori de instruire consecutivi. Vectorii nu mai sunt apoi folosiți, lucru care îi conferă acestui algoritm atributul de on-line. El poate fi folosit pentru fluxuri de date de instruire și relevanțele se adaptează continuu, reflectând în mod dinamic schimbările care se petrec asupra intrărilor.

Algoritmul ERLVQ introdus în paragraful 6.7 este similar lui MIRLVQ, deși este dezvoltat prin maximizarea unei mărimi bazate pe energia informațională. Definiția energiei informaționale ne permite să folosim și de această dată estimarea densităților de probabilitate cu ferestre Parzen și obținem o formulă de actualizare a relevanțelor pe baza unor vectori consecutivi din setul de instruire. Regulile de modificare iterativă a relevanțelor pe care le-am dedus prin cele două metode au fost integrate în algoritmul weighted LVQ rezultând procedura din secțiunea 6.7.

Am experimentat algoritmul ERLVQ în condiții similare cu OWA-RLVQ pentru a putea face o comparație corectă a rezultatelor. Am constatat că abordarea prin prisma teoriei informației conduce la rate de recunoaștere mai bune, inclusiv pentru baza de date Deterding care este un cunoscut benchmark în recunoașterea vorbirii.

Cele două modele de recunoaștere a formelor sintetizate prin algoritmi OWA-RLVQ și ERLVQ reflectă două abordări diferite ale aceleiași probleme. Factorii de relevanță calculați fie prin prisma teoriei operatorilor OWA, fie prin cea a teoriei informației conduc la rezultate care îmbunătățesc ratele de recunoaștere în comparație cu algoritmi puternici care implementează aceeași gamă de aplicații. În plus, relevanțele sunt utile la ordonarea dinamică a componentelor vectorilor de intrare în funcție de importanța lor (*feature ranking*) sau la eliminarea componentelor puțin importante (*feature selection*).

Pentru orice tehnică de recunoaștere sau clasificare automată a formelor se poate face presupunerea că rezultatele obținute ar putea fi dependente de datele de instruire și de test. Acest lucru este adevărat în majoritatea cazurilor, motiv pentru care nu putem susține superioritatea absolută a unei metode. Putem, însă, să folosim date cât mai diverse pentru a testa propriile tehnici sau ne putem verifica metodele în diferite contexte și atunci spectrul cazurilor pentru care abordarea noastră este superioară altora se lărgeste. Din această cauză am testat modelul CNN atât pentru recunoașterea cuvintelor izolate cât și pentru recunoașterea vorbitorilor și pentru recunoașterea vocalelor. Dacă pentru primele două categorii de teste am folosit cuvinte în limba engleză, la

recunoașterea vocalelor am utilizat cuvinte pronunțate în limba română. Atunci când am testat algoritmi OWA-RLVQ și ERLVQ, pe lângă pronunții vocale am folosit și eșantioane care, în mod uzual, sunt folosite pentru sisteme generale de recunoaștere a formelor.

### 7.3 Direcții viitoare de cercetare

În lucrare au fost discutate diverse variante de abordare a problemelor care au fost enunțate încă o dată la începutul acestui capitol. Există, desigur, și alte idei care pot fi utile pentru extinderea acestor cercetări și doresc să menționez câteva dintre ele.

- Așa cum au fost introduse aici, rețelele neurale concurente sunt colecții de rețele neurale care sunt asociate câte unei singure clase de vectori. Este util de studiat cum poate fi modificat acest model în așa fel încât să se optimizeze numărul de rețele neurale componente prin unificarea unor module fără a afecta performanțele generale. De această dată, vom avea module asociate mai multor clase. Această abordare *bottom-up* poate fi înlocuită cu una *top-down* prin care pornim cu o rețea neurală globală pe care o înlocuim cu mai multe rețele neurale concurente. Această variantă a CNN poate conduce la o modelare a creierului pe baza ipotezei că el constă din zone distincte care răspund la stimuli bine determinați;
- Pornind de la rezultatele cercetărilor din neurobiologie care au stabilit existența unei memorii de scurtă durată și a unei memorii de lungă durată, se poate extinde modelul CNN prin transferarea unor vectori dintr-o clasă - memoria de scurtă durată - în altă clasă - memoria de lungă durată - atunci când se constată că ei generează erori de cuantizare minime;
- O continuare firească a rezultatelor materializate prin algoritmi OWA-RLVQ, MIRLVQ și ERLVQ este folosirea relevanțelor pentru eliminarea dinamică a componentelor de intrare ne semnificative, cu utilitate deosebită în *data-mining*.

Majoritatea modelelor neurale abordează recunoașterea formelor ca o problemă globală și fără a face distincția între aportul individual al intrărilor. Această teză demonstrează că tehnicile propuse de noi contribuie la sporirea acurateții recunoașterii și creează suportul pentru cercetări viitoare legate de extinderea conceptului de concurență la nivelul sistemelor neurale și de reducere a dimensiunii intrărilor prin eliminarea componentelor puțin importante.



# Anexa A

## Valori reprezentative ale coeficienților cepstrali pentru foneme din limba română

Tabelele B.1 și B.2 conțin valorile numerice pe baza cărora a fost realizată figura 4.4. Aceste date sunt valori reprezentative ale coeficienților cepstrali calculate pentru câte o pronunție a fiecărei vocale. Pronunțiile sunt ale primului dintre cei 4 vorbitori pe care i-am folosit la crearea bazei de date descrisă în secțiunea 4.4.2.

Tabelul B.1: Valori reprezentative ale coeficienților cepstrali calculate pentru câte o pronunție a vocalelor A, E, I.

| Coeficientul | A     | E     | I     |
|--------------|-------|-------|-------|
| 0            | 8,54  | 7,07  | 5,66  |
| 1            | -3,58 | 1,44  | 5,22  |
| 2            | -1,33 | 1,51  | 3,11  |
| 3            | 0,55  | -0,67 | -1,91 |
| 4            | 1,82  | -1,16 | 0,44  |
| 5            | -1,04 | 0,31  | -1,46 |
| 6            | -0,27 | 0,06  | -0,10 |
| 7            | 0,41  | 0,23  | -0,58 |
| 8            | 5,71  | 4,68  | 5,01  |
| 9            | -0,38 | -0,10 | -0,15 |
| 10           | -0,10 | -0,02 | 0,42  |
| 11           | 0,06  | -0,18 | 0,04  |

Tabelul B.2: Valori reprezentative ale coeficienților cepstrali calculate pentru câte o pronunție a vocalelor O, U, Ă.

| Coeficientul | O     | U     | Ă     |
|--------------|-------|-------|-------|
| 0            | 9,98  | 10,54 | 8,59  |
| 1            | 3,35  | 2,11  | 0,01  |
| 2            | -0,40 | -1,45 | -2,68 |
| 3            | -1,72 | -1,48 | -0,38 |
| 4            | -0,67 | -0,08 | 2,60  |
| 5            | 0,17  | 0,48  | -0,63 |
| 6            | -0,22 | -0,22 | -1,06 |
| 7            | -0,15 | -0,66 | -0,40 |
| 8            | 4,88  | 5,41  | 4,38  |
| 9            | 0,18  | -0,83 | -0,48 |
| 10           | 0,36  | -0,70 | -0,31 |
| 11           | 0,17  | -0,32 | -0,31 |

# Anexa B

## Estimarea densităților de probabilitate cu ferestre Parzen. Demonstrația proprietății (6.12)

Vom demonstra relația (6.12) conform căreia:

$$\int_{-\infty}^{+\infty} G(\mathbf{y} - \mathbf{y}_1, \sigma^2 \mathbf{I}) G(\mathbf{y} - \mathbf{y}_2, \sigma^2 \mathbf{I}) d\mathbf{y} = G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2 \mathbf{I}) \quad (\text{B.1})$$

unde

$$G(\mathbf{y}, \sigma^2 \mathbf{I}) = \frac{1}{(2\pi\sigma^2)^{\frac{d}{2}}} e^{-\frac{\|\mathbf{y}\|^2}{2\sigma^2}}, \quad (\text{B.2})$$

iar  $d$  este dimensiunea vectorilor  $\mathbf{y}$ ,  $\mathbf{y}_1$  și  $\mathbf{y}_2$ .

Facem următoarele notații:

$$\mathbf{y} = [y_1, \dots, y_d]^t$$

$$\mathbf{y}_1 = [y_{11}, \dots, y_{1d}]^t$$

$$\mathbf{y}_2 = [y_{21}, \dots, y_{2d}]^t.$$

Atunci, conform expresiei (B.2) obținem:

$$G(\mathbf{y} - \mathbf{y}_1, \sigma^2 \mathbf{I}) = \frac{1}{(2\pi\sigma^2)^{\frac{d}{2}}} e^{-\frac{\|\mathbf{y} - \mathbf{y}_1\|^2}{2\sigma^2}} \quad (\text{B.3})$$

și

$$G(\mathbf{y} - \mathbf{y}_2, \sigma^2 \mathbf{I}) = \frac{1}{(2\pi\sigma^2)^{\frac{d}{2}}} e^{-\frac{\|\mathbf{y} - \mathbf{y}_2\|^2}{2\sigma^2}}. \quad (\text{B.4})$$

Conform definiției normei scriem:

$$\|\mathbf{y} - \mathbf{y}_1\|^2 = (y_1 - y_{11})^2 + \dots + (y_d - y_{1d})^2 \quad (\text{B.5})$$

și

$$\|\mathbf{y} - \mathbf{y}_2\|^2 = (y_1 - y_{21})^2 + \dots + (y_d - y_{2d})^2. \quad (\text{B.6})$$

Introducem (B.5) și (B.6) în (B.3) și (B.4) și obținem:

$$\begin{aligned} G(\mathbf{y} - \mathbf{y}_1, \sigma^2 \mathbf{I}) &= \frac{1}{(2\pi\sigma^2)^{\frac{d}{2}}} e^{-\frac{(y_1 - y_{11})^2 + \dots + (y_d - y_{1d})^2}{2\sigma^2}} \\ &= \frac{1}{(2\pi\sigma^2)^{\frac{d}{2}}} e^{-\frac{(y_1 - y_{11})^2}{2\sigma^2}} \cdot \dots \cdot e^{-\frac{(y_d - y_{1d})^2}{2\sigma^2}} \end{aligned}$$

și

$$\begin{aligned} G(\mathbf{y} - \mathbf{y}_2, \sigma^2 \mathbf{I}) &= \frac{1}{(2\pi\sigma^2)^{\frac{d}{2}}} e^{-\frac{(y_2 - y_{21})^2 + \dots + (y_d - y_{2d})^2}{2\sigma^2}} \\ &= \frac{1}{(2\pi\sigma^2)^{\frac{d}{2}}} e^{-\frac{(y_1 - y_{21})^2}{2\sigma^2}} \cdot \dots \cdot e^{-\frac{(y_d - y_{2d})^2}{2\sigma^2}}. \end{aligned}$$

Deci:

$$\begin{aligned} G(\mathbf{y} - \mathbf{y}_1, \sigma^2 \mathbf{I}) G(\mathbf{y} - \mathbf{y}_2, \sigma^2 \mathbf{I}) &= \\ &= \frac{1}{(2\pi\sigma^2)^{\frac{d}{2}}} \cdot \frac{1}{(2\pi\sigma^2)^{\frac{d}{2}}} e^{-\frac{(y_1 - y_{11})^2 + (y_1 - y_{21})^2}{2\sigma^2}} \cdot \dots \cdot e^{-\frac{(y_d - y_{1d})^2 + (y_d - y_{2d})^2}{2\sigma^2}} \end{aligned}$$

de unde rezultă:

$$\begin{aligned} &\int_{-\infty}^{+\infty} G(\mathbf{y} - \mathbf{y}_1, \sigma^2 \mathbf{I}) G(\mathbf{y} - \mathbf{y}_2, \sigma^2 \mathbf{I}) = \\ &= \int_{-\infty}^{+\infty} \frac{1}{(2\pi\sigma^2)^{\frac{d}{2}}} \cdot \frac{1}{(2\pi\sigma^2)^{\frac{d}{2}}} e^{-\frac{(y_1 - y_{11})^2 + (y_1 - y_{21})^2}{2\sigma^2}} \cdot \dots \\ &\dots \cdot e^{-\frac{(y_d - y_{1d})^2 + (y_d - y_{2d})^2}{2\sigma^2}} dx_1 \dots dx_d = \\ &= \frac{1}{(2\pi\sigma^2)^d} \left[ \prod_{i=1}^d \int_{-\infty}^{+\infty} e^{-\frac{(y_i - y_{1i})^2 + (y_i - y_{2i})^2}{2\sigma^2}} dy_i \right] = \\ &= \frac{1}{(2\pi\sigma^2)^{\frac{d}{2}}} \left[ \prod_{i=1}^d \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{+\infty} e^{-\frac{(y_i - y_{1i})^2 + (y_i - y_{2i})^2}{2\sigma^2}} dy_i \right]. \quad (\text{B.7}) \end{aligned}$$

Calculăm:

$$\begin{aligned} &\frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{+\infty} e^{-\frac{(y_i - y_{1i})^2 + (y_i - y_{2i})^2}{2\sigma^2}} dy_i = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{+\infty} e^{-\frac{(x - y_{1i})^2 + (x - y_{2i})^2}{2\sigma^2}} dx = \\ &= \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{+\infty} e^{-\frac{1}{\sigma^2} \left[ x^2 - (y_{1i} + y_{2i})x + \left(\frac{y_{1i} + y_{2i}}{2}\right)^2 - \left(\frac{y_{1i} - y_{2i}}{2}\right)^2 + \frac{y_{1i}^2 + y_{2i}^2}{2} \right]} dx = \\ &= \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{+\infty} e^{-\frac{1}{\sigma^2} \left[ x - \frac{y_{1i} + y_{2i}}{2} \right]^2} \cdot e^{-\frac{1}{\sigma^2} \left[ \frac{y_{1i}^2 + y_{2i}^2}{2} - \frac{y_{1i}^2 + y_{2i}^2 + 2y_{1i}y_{2i}}{4} \right]} dx = \end{aligned}$$

$$\begin{aligned}
&= e^{-\frac{1}{4\sigma^2}[y_{1j}^2+y_{2j}^2-2y_{1j}y_{2j}]} \cdot \frac{1}{\sqrt{2}} \cdot \frac{1}{\frac{\sigma}{\sqrt{2}}\sqrt{2\pi}} \int_{-\infty}^{+\infty} e^{-\frac{\left(x-\frac{y_{1j}+y_{2j}}{2}\right)^2}{2\left(\frac{\sigma}{\sqrt{2}}\right)^2}} dx = \\
&= \frac{1}{\sqrt{2}} e^{-\frac{(y_{1j}+y_{2j})^2}{4\sigma^2}},
\end{aligned}$$

unde am trecut de la variabila  $y_i$  la  $x$  pentru o scriere mai clară. În final, am folosit un rezultat care este cunoscut:

$$\frac{1}{\frac{\sigma}{\sqrt{2}}\sqrt{2\pi}} \int_{-\infty}^{+\infty} e^{-\frac{\left(x-\frac{y_{1j}+y_{2j}}{2}\right)^2}{2\left(\frac{\sigma}{\sqrt{2}}\right)^2}} dx = 1.$$

În concluzie, (B.7) devine:

$$\begin{aligned}
&\int_{-\infty}^{+\infty} G(\mathbf{y} - \mathbf{y}_1, \sigma^2 \mathbf{I}) G(\mathbf{y} - \mathbf{y}_2, \sigma^2 \mathbf{I}) d\mathbf{y} = \\
&= \frac{1}{(2\pi\sigma^2)^{\frac{d}{2}}} \left( \prod_{j=1}^d \frac{1}{\sqrt{2}} e^{-\frac{(y_{1j}-y_{2j})^2}{4\sigma^2}} \right) = \\
&= \frac{1}{(2\pi\sigma^2)^{\frac{d}{2}}} \frac{1}{2^{\frac{d}{2}}} \left[ e^{-\frac{1}{4\sigma^2} \sum_{j=1}^d (y_{1j}-y_{2j})^2} \right] = \\
&= \frac{1}{(2\pi\sigma^2)^{\frac{d}{2}}} e^{-\frac{1}{4\sigma^2} \|\mathbf{y}_1 - \mathbf{y}_2\|^2} = \\
&= \frac{1}{[2\pi(2\sigma^2)]^{\frac{d}{2}}} e^{-\frac{\|\mathbf{y}_1 - \mathbf{y}_2\|^2}{2 \cdot (2\sigma^2)}} = \\
&= G(\mathbf{y}_1 - \mathbf{y}_2, 2\sigma^2 \mathbf{I}).
\end{aligned}$$

Am scris ultima egalitate ținând cont de definiția (B.2). Relația (B.1) este astfel demonstrată.



# Bibliografie

- [1] R. Anand, K. Mehrotra, C.K. Mohan, S. Ranka "Efficient classification for multiclass problems using modular neural networks". IEEE Transactions on Neural Networks 6, 117-124, 1995.
- [2] R. Andonie, F. Petrescu "Interacting systems and informational energy". Foundations of control and engineering 11, no. 2, 53-59, 1986.
- [3] M. Andrei, I. Ghiță "Limba română". Editura Didactică și Pedagogică, București, 1983.
- [4] J. Basak, R.K. De, S.K. Pal "Unsupervised feature selection using a neuro-fuzzy approach". Pattern Recognition Letters 19, 997-1006, 1998.
- [5] R. Battiti "Using mutual information for selecting features in supervised neural net learning". IEEE Transactions on Neural Networks 5, 537-550, 1994.
- [6] G. Beliakov "How to build aggregation operators from data". International Journal of Intelligent Systems 18, 903-923, 2003.
- [7] Y. Bengio, R. Cardin, R. de Mori, E. Merlo "Programable execution of multilayered networks for automatic speech recognition". Communications of the ACM 32, 195-199, 1989.
- [8] C.M. Bishop "Neural networks for pattern recognition". Oxford University Press, New York, 1995.
- [9] C.L. Blake, C.J. Merz "UCI repository of machine learning databases". University of California, Irvine, Department of Information and Computer Sciences, <http://www.ics.uci.edu/~mllearn/MLRepository.html>, 1998.
- [10] T. Bojer, B. Hammer, D. Schunk, K.T. von Toschanowitz "Relevance determination in learning vector quantization". In M. Verleysen (Ed.), Proceedings of the European Symposium on Artificial Neural Networks (ESANN 2001), D-Side publications, 271-276, 2001.

- [11] A. Cațaron, V.-E. Neagoe "Concurrent neural networks for speaker recognition". Proceedings of the International Conference on Telecommunications (ICT 2001), București, 2001.
- [12] A. Cațaron "Vowel recognition using concurrent neural networks". Proceedings of the 8th International Conference on Optimization of Electrical and Electronic Equipments (OPTIM 2002), Brașov, 779-782, 2002.
- [13] A. Cațaron, S. Mătase "Speaker recognition using multi-layer perceptron, time delay neural network, self-organizing map and concurrent neural networks". Proceedings of the 8th International Conference on Optimization of Electrical and Electronic Equipments (OPTIM 2002), Brașov, 775-778, 2002.
- [14] A. Cațaron, R. Andonie "RLVQ determination using OWA operators". In M.H. Hamza (Ed.), Proceedings of the 3rd IASTED International Conference on Artificial Intelligence and Applications (AIA 2003), September 8-10, Benalmadena, Spain, ACTA Press, 434-438, 2003.
- [15] R.K. De, N.R. Pal, S.K. Pal "Feature analysis: neural network and fuzzy set theoretic approaches". *Patterns Recognition* 10(30), 1579-1590, 1997.
- [16] J.R. Deller, J.G. Proakis, J.H.L. Hansen "Discrete-time processing of speech signals". Prentice Hall, Upper Saddle River, New Jersey, 1987.
- [17] R.O. Duda, P.E. Hart, D.G. Stork "Pattern classification". John Wiley & Sons Inc., New York, 2001.
- [18] J.L. Elman "Finding structure in time". *Cognitive Science* 14, 179-211, 1990.
- [19] D. Erdogmus, J.C. Principe "Generalized information potential criterion for adaptive system training". *IEEE Transactions on Neural Networks* 13, 1035-1044, 2002.
- [20] D. Filev, R.R. Yager "Learning OWA operator weights from data". Proceedings of the Third IEEE International Conference on Fuzzy Systems, Orlando, IEEE Press, 468-473, 1994.
- [21] D. Filev, R.R. Yager "On the issue of obtaining OWA operator weights". *Fuzzy Sets and Systems* 94, 157-169, 1998.
- [22] J.W. Fisher III, J.C. Principe "A methodology for information theoretic feature extraction". Proceedings of the IEEE World Congress on Computational Intelligence, Anchorage, Alaska, 1712-1716, 1998.

- [23] B. Freisleben, C.A. Bohn "Speaker-independent word recognition with backpropagation networks". In R.F. Albrecht, C.R. Reeves, N.C. Steele (Eds.), Artificial neural nets and genetic algorithms, Proceedings of the International Conference in Innsbruck, Springer, Wien, 243-248, 1993.
- [24] L.M. Fu "Neural networks in computer intelligence". McGraw-Hill, New York, 1994.
- [25] R.M. Golden "Mathematical methods for neural network analysis and design". The MIT Press, Cambridge, Massachusetts, 1996.
- [26] S. Guiasu "Information theory with applications". McGraw-Hill, New York, 1977.
- [27] B. Hammer, T. Villmann "Input pruning for neural gas architectures". In M. Verleysen (Ed.), Proceedings of the European Symposium on Artificial Neural Networks (ESANN 2001), D-Facto publications, 283-288, 2001.
- [28] B. Hammer, T. Villmann "Batch-RLVQ". In M. Verleysen (Ed.), Proceedings of the European Symposium on Artificial Neural Networks (ESANN 2002), D-Side publications, 295-300, 2002.
- [29] B. Hammer, T. Villmann "Generalized relevance learning vector quantization". Neural Networks 15, 1059-1068, 2002.
- [30] S. Haykin "Neural networks - a comprehensive foundation". Macmillan College Publishing Company, New York, 1994.
- [31] D.O. Hebb "The organization of behaviour: a neuropsychological theory". Wiley, New York, 1949.
- [32] D. Huang, T. Chow "Searching optimal feature subset using mutual information". In M. Verleysen (Ed.), Proceedings of the European Symposium on Artificial Neural Networks (ESANN 2003), D-Side publications, 161-166, 2003.
- [33] I. Iordan, V. Robu "Limba română contemporană". Editura Didactică și Pedagogică, București, 1978.
- [34] R.A. Jacobs, M.I. Jordan, A.G. Barto "Task decomposition through competition in a modular connectionist architecture: the what and where vision tasks". Cognitive Science 15, 219-250, 1991.
- [35] G. Jumarie "Relative information". Springer-Verlag, Berlin, 1990.

- [36] N.B. Karayiannis "An axiomatic approach to soft learning vector quantization and clustering". IEEE Transactions on Neural Networks 10, 1153-1165, 1999.
- [37] N.B. Karayiannis "Soft learning vector quantization and clustering algorithms based on ordered weighted aggregation operators". IEEE Transactions on Neural Networks 11, 1093-1105, 2000.
- [38] N.B. Karayiannis, M.M. Randolph-Gips "Soft learning vector quantization and clustering algorithms based on non-euclidian norms: multinorm algorithms". IEEE Transactions on Neural Networks 14, 89-102, 2003.
- [39] S.V. Kartalopoulos "Understanding neural networks and fuzzy logic". IEEE Press, Piscataway, New Jersey, 1996.
- [40] T. Kohonen "The 'neural' phonetic typewriter". IEEE Computer Magazine, March, 11-22, 1988.
- [41] T. Kohonen "Self-organization and associative memory". Series in Information Sciences 61, Springer-Verlag, Berlin, 1989.
- [42] T. Kohonen "The self-organizing map". Proceedings of the IEEE 78, 1464-1480, 1990.
- [43] T. Kohonen, J. Hynninen, J. Kangas, J. Laaksonen "SOM\_PAK: The self-organizing map program package". Laboratory of Computer and Information Science, Helsinki University of Technology, Finland, [http://www.cis.hut.fi/research/som\\_research/nncr\\_programs.html](http://www.cis.hut.fi/research/som_research/nncr_programs.html), 1995.
- [44] T. Kohonen, J. Hynninen, J. Kangas, J. Laaksonen, K. Torkkola "LVQ\_PAK: The learning vector quantization program package". Laboratory of Computer and Information Science, Helsinki University of Technology, Finland, [http://www.cis.hut.fi/research/som\\_research/nncr\\_programs.html](http://www.cis.hut.fi/research/som_research/nncr_programs.html), 1995.
- [45] T. Kohonen "Self-organizing maps". Springer-Verlag, 1997.
- [46] T. Kohonen, S. Kaski, K. Lagus, J. Salojarvi, J. Honkela, V. Paatero, A. Saarela "Self organization of a massive document collection". IEEE Transactions on Neural Networks 11, 574-585, 2000.
- [47] M.A. Kraaijveld, J. Mao, A.K. Jain "A nonlinear projection method based on Kohonen's topology preserving maps". IEEE Transactions on Neural Networks 6, 548-559, 1995.

- [48] N. Kwak, C.-H. Choi "Improved mutual information feature selector for neural networks in supervised learning". Proceedings of the International Joint Conference on Neural Networks (IJCNN 1999), vol. 2, 1313-1381, 1999.
- [49] N. Kwak, C.-H. Choi "Input feature selection for classification problems". IEEE Transactions on Neural Networks 13, 143-159, 2002.
- [50] N. Kwak, C.-H. Choi "Input feature selection by mutual information based on Parzen windows". IEEE Transactions on Pattern Analysis and Machine Intelligence 24, 1667-1671, 2002.
- [51] I. Lapidot (Voitovetsky), H. Guterman, A. Cohen "Unsupervised speaker recognition based on competition between self-organizing maps". IEEE Transactions on Neural Networks 13, 877-887, 2002.
- [52] M. Last, A. Kandel, O. Maimom "Informatic-theoretic algorithm for feature selection". Pattern Recognition Letters 22, 799-811, 2001.
- [53] W. Li "Mutual information functions versus correlation functions". Journal of Statistical Physics 60, 823-837, 1990.
- [54] D.J.C. MacKay "Information theory, inference, and learning algorithms". Cambridge University Press, 2003.
- [55] W.S. McCulloch, W Pitts "A logical calculus of the ideas immanent in nervous activity". Bulletin of Mathematical Biophysics 5, 115-133, 1943.
- [56] W.S. McCulloch, W. Pitts "How we know universals: the perception of auditory and visual forms". Bulletin of Mathematical Biophysics 9, 127-147, 1947.
- [57] P. Mitra, C.A. Murthy, S.K. Pal "Unsupervised feature selection using feature similarity". IEEE Transactions on Pattern Analysis and Machine Intelligence 24, 301-312, 2002.
- [58] V.-E. Neagoe, O. Cula "A fuzzy connectionist approach to vowel recognition". In B. Reusch and D. Dascalu (Eds.), Real World Applications of Intelligent Technologies (part II), Institutul Național de Cercetare și Dezvoltare în Microtehnologii, București, 1998.
- [59] V.-E Neagoe, O. Stănășilă "Recunoașterea formelor și rețele neurale - algoritmi fundamentali". Editura Matrix Rom, București, 1999.

- [60] V. Neagoe "Concurrent self-organizing maps for automatic face recognition". Proceedings of the 29th International Conference of the Romanian Technical Military Academy, Bucharest, Romania, 35-40, 2001.
- [61] V. Neagoe, A. Ropot "Concurrent self-organizing maps for classification of multispectral satellite imagery". Proceedings of the 6th World Multiconference on Systemics, Cybernetics, and Informatics (SCI 2002), Orlando, Florida, 126-131, 2002.
- [62] V. Neagoe, A. Ropot "Concurrent self-organizing maps for pattern classification". Proceedings of the First IEEE International Conference on Cognitive Informatics (ICCI 2002), Calgary, Alberta, Canada, 304-312, 2002.
- [63] V. Neagoe, I. Iatan "Face recognition using a fuzzy-gaussian neural network". Proceedings of the First IEEE International Conference on Cognitive Informatics (ICCI 2002), Calgary, Alberta, Canada, 361-368, 2002.
- [64] V. Neagoe, A. Ropot "Concurrent self-organizing maps - a powerful tool for face and speaker recognition". Proceedings of the International Conference COMMUNICATIONS 2002, Edited by the Technical Military Academy, Bucharest, Romania, 276-281, 2002.
- [65] O. Onicescu "Energia informațională". Studii și Cercetări Matematice 18, 1419-1420, 1966.
- [66] F.J. Owens, R. Andonie, G.H. Zheng, A. Cataron, S. Manciulea "A comparative study of the multy-layer perceptron, the multi-output layer perceptron, the time-delay neural network and the Kohonen self-organizing map in an automatic speech recognition task". Proceedings of the International ICSC Symposium on Engineering of Intelligent Systems (EIS 1998), ICSC Academic Press, 624-629, Tenerife, Spain, February 11-13, 1998.
- [67] S.K. Pal, R.K. De, J. Basak "Unsupervised feature evaluation: a neuro-fuzzy approach". IEEE Transactions on Neural Networks 11, 366-376, 2000.
- [68] M. Pregenzer, D. Flotzinger, G. Pfurtscheller "Distinction sensitive learning vector quantization - a noise-insensitive classification method". Proceedings of the IEEE International Joint Conference on Neural Networks (IJCNN 1994), Orlando, Florida, 2890-2894, 1994.
- [69] M. Pregenzer, G. Pfurtscheller, D. Flotzinger "Automated feature selection with a distinction sensitive learning vector quantizer". Neurocomputing 11, 19-29, 1996.

- [70] J.C. Principe, D. Xu, J.W. Fisher III "Information-theoretic learning". In Simon Haykin (Ed.), *Unsupervised adaptive filtering*, Wiley, New York, 2000.
- [71] I. Purcaru "Informație și corelație". Editura Științifică și enciclopedică, București, 1988.
- [72] L. Rabiner, B.-H. Juang "Fundamentals of speech recognition". Prentice Hall, Englewood Cliffs, New Jersey, 1993
- [73] A. Rosetti, A. Lăzăroiu "Introducere în fonetică". Editura Științifică și Enciclopedică, București, 1982.
- [74] J.G. Rueckl, K.R. Cave, S.M. Kosslyn "Why are 'what' and 'where' processed by separate cortical visual systems? A computational investigation". *Journal of Cognitive Neuroscience* 1, 171-186, 1989.
- [75] D.E. Rumelhart, J.L. McClelland (Eds.) "Parallel distributed processing: explorations in the microstructure of cognition". MIT Press, Cambridge, MA, 1986.
- [76] A. Sato, K. Yamada "Generalized learning vector quantization". In "Advances in neural information processing systems", G. Tesauro, D. Touretzky, T. Leen (Eds.), vol. 7, MIT Press, 423-429, 1995.
- [77] R. Setiono, H. Liu "Neural-network feature selector". *IEEE Transactions on Neural Networks* 8, 654-662, 1997.
- [78] C. Shannon "A mathematical theory of communication". *The Bell System Technical Journal* 27, 379-423, 623-656, 1948.
- [79] V. Ștefănescu "Aplicații ale energiei și corelației informaționale". Editura Academiei, București, 1979.
- [80] K. Torkkola, W.M. Campbell "Mutual information in learning feature transformations". *Proceedings of the International Conference of Machine Learning*, Stanford, CA, USA, 1015-1022, 2000.
- [81] K. Torkkola "Nonlinear feature transforms using maximum mutual information". *Proceedings of the IEEE International Joint Conference on Neural Networks (IJCNN 2001)*, Washington DC, USA, 2756-2761, 2001.
- [82] K. Torkkola "Learning discriminative feature transforms to low dimensions in low dimensions". In "Advances in neural information processing systems" 14 (NIPS 2001), Vancouver, BC, Canada, MIT Press, 2001.

- [83] K. Torkkola "On feature extraction by mutual information maximization". Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP 2002), Orlando, Florida, 2002.
- [84] K. Torkkola "Feature extraction by non-parametric mutual information maximization". Journal of Machine Learning Research 3, 1415-1438, 2003.
- [85] V. Torra "Learning weights for weighted OWA operators". Proceedings of the IEEE International Conference on Industrial Electronics, Control, and Instrumentation (IECON 2000), Nagoya, Japan, 2000.
- [86] V. Torra "Weighted OWA operators for synthesis of information". Proceedings of the 5th IEEE International Conference on Fuzzy Systems, New Orleans, USA, 996-971, 2000.
- [87] L.P.J. Veelenturf "Analysis and applications of artificial neural networks". Prentice Hall, London, 1995.
- [88] J. Vesanto, E. Alhoniemi "Clustering of the self-organizing map". IEEE Transactions on Neural Networks 11, 586-600, 2000.
- [89] X.L. Xie, G. Beni "A validity measure for fuzzy clustering". IEEE Transactions on Pattern Analysis and Machine Intelligence 13, 841-847, 1991.
- [90] A. Waibel, T. Hanazawa, G. Hinton, K. Shikano, K.J. Lang "Phoneme recognition using time-delay neural networks". IEEE Transactions on Acoustics, Speech, and Signal Processing 37, 328-339, 1989.
- [91] A. Waibel, H. Sawai, K. Shikano "Modularity and scaling in large phoneme neural networks". IEEE Transactions on Acoustics, Speech, and Signal Processing 37, 1989.
- [92] H. Watanabe "A new view of the formal entropy as a measure of interdependence and its application to pattern recognition". Proceedings of the IEEE International Conference on System, Man, and Cybernetics, Nashville, USA, 2827-2832, 2000.
- [93] R.R. Yager "On ordered weighted averaging aggregation operators in multicriteria decisionmaking". IEEE Transactions on System, Man, and Cybernetics 18, 188-190, 1988.
- [94] R.R. Yager, D. Filev "Induced ordered weighted averaging operators". IEEE Transactions on System, Man, and Cybernetics 29, 141-150, 1999.

- [95] R.R. Yager "Induced aggregation operators". Fuzzy Sets and Systems 137, 59-69, 2003.
- [96] D.S. Yeung, X.Z. Wang "Improving performance of similarity-based clustering by feature weight learning". IEEE Transactions on Pattern Analysis and Machine Intelligence 24, 556-561, 2002.
- [97] J.M. Zurada "Introduction to artificial neural systems". West Publishing Company, St. Paul, Minnessotta, 1992.
- [98] \*\*\* "Stuttgart Neural Network Simulator". SNNS Group, Institute for Parallel and Distributed High-Performance Systems (IPVR), University of Stuttgart, Germany, <http://www-ra.informatik.uni-tuebingen.de/SNNS/>, 1995.

#### **Apariții în anul 2004**

- [99] R. Andonie, A. Cațaron "An informational energy LVQ approach for feature ranking". Proceedings of the European Symposium on Artificial Neural Networks (ESANN 2004), Bruges, Belgium, 2004.
- [100] A. Cațaron, R. Andonie "Computing OWA weights as relevance factors". Proceedings of the 9th International Conference on Optimization of Electrical and Electronic Equipments (OPTIM 2004), Brașov, 2004.
- [101] A. Cațaron, R. Andonie "Energy generalized LVQ with relevance factors". Proceedings of the International Joint Conference on Neural Networks (IJCNN 2004), Budapest, Hungary, 2004.