

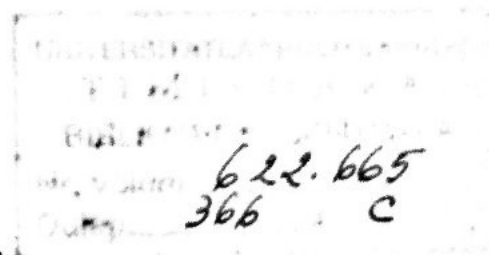
UNIVERSITATEA TEHNICĂ DIN TIMIȘOARA
Facultatea de Automatică și Calculatoare

Autor:
Ing. IOAN Z. MIHU

**SOLUȚII DE IMPLEMENTARE A
REȚELELOR NEURONALE PE
ARHITECTURI SISTOLICE**

TEZĂ DE DOCTORAT

BIBLIOTECA CENTRALĂ
UNIVERSITATEA "POLITEHNICA"
TIMIȘOARA



Conducător științific:
Prof. Dr. Ing. Nicolae Budișan

1998

Cuvânt înainte

Ființele vii evoluate își coordonează acțiunile punând în joc o cantitate considerabilă de celule nervoase, interconectate printr-un număr și mai considerabil de conexiuni sinaptice. Acest sistem puternic interconectat de celule nervoase, este sediul unei intense activități electro-chimice, și capacitatea sa de adaptare și acțiune constituie un subiect de cercetare abordat cu deosebit interes.

Așadar, principiul conexionismului se bazează pe faptul că activitatea colectivă a unor unități elementare simple, care se influențează reciproc, poate conduce la un comportament variat și complex.

Începând cu anii 1960, pe baza acestui principiu, au fost propuse mai multe modele care au cunoscut anumite succese în rezolvarea problemelor complexe și insuficient fundamentate teoretic pentru a putea fi soluționate prin metode tradiționale.

Aceste modele fac în general apel la două faze de calcul: o fază de adaptare sau de învățare care servește la stabilirea funcției pe care trebuie să o realizeze rețeaua neuronală pornind de la un set de exemple, apoi o fază de recunoaștere care exploatează funcția învățată pe un ansamblu de stimuli, în general mai vast decât setul de antrenament. Cele două faze pun în evidență caracteristicile de bază ale rețelelor conexioniste: învățare din exemple și capacitate de generalizare.

Pentru validarea acestor modele și pentru aplicarea lor în problematica reală, inginerul face apel la mijloacele de procesare a informației de care dispune. Prin urmare, simularea acestor structuri cu grad fin și înalt de paralelism pe calculatoare clasice, dotate cu unul sau mai multe procesoare, nu permite atingerea unui nivel de performanță suficient pentru abordarea aplicațiilor reale.

Numai problemele care reclamă rețele neuronale de mici dimensiuni pot fi tratate cu succes pe sistemele de calcul clasice. Realizarea unor circuite VLSI, cu grad înalt de paralelism intern, specializate pentru emularea rețelelor conexioniste, și care să poată funcționa autonom sau integrate într-un neurosistem (calculator gazdă dotat cu un *hardware* accelerator dedicat algoritmilor neuronali), reprezintă deci un obiectiv important pentru înțelegerea și exploatarea domeniului rețelelor neuronale artificiale.

Obiectivul acestei teze este de a propune un principiu de concepție a circuitelor VLSI digitale, capabile să integreze un anumit număr de funcții necesare emulării mai multor modele de rețele conexioniste. Ambele funcții (de învățare și de recunoaștere) vor fi integrate pentru a permite realizarea unor eventuale sisteme (VLSI-uri) autonome.

Primul capitol descrie principalele modele neuronale în termenii resurselor revendicate și a operațiilor efectuate pe aceste resurse. Această etapă ne-a permis evidențierea punctelor comune aferente acestor modele, în termeni structurali și funcționali (ai operațiilor aritmetico - logice pe care le implică).

Al doilea capitol realizează o sinteză a celor mai importante realizări în domeniul neuroprocesoarelor digitale și analogice și stabilește principiile de proiectare a neuroprocesoarelor sistolice.

În al treilea capitol se elaborează o arhitectură sistolică generală, dedicată rețelelor neuronale, din care vor deriva trei arhitecturi exploatabile.

În capitolul al patrulea se definește arhitectura sistolică optimă și se realizează proiectarea, testarea funcțională și evaluarea acestei arhitecturi.

Capitolul cinci este dedicat elaborării unei strategii de perfecționare a rețelei sistolice dedicate algoritmilor conexioniști, prin *pipeline*-izarea fluxurilor de date și de instrucțiuni în

interiorul rețelei sistolice, prin perfecționarea structurii celulei sistolice și prin implementarea unui mecanism performant (*pipeline*-izat) de interschimbare a matricilor sinaptice.

Capitolul șase face analiza unei propuneri alternative: Rețeaua Sistolice Liniară cu integrare de neuroni și nu de sinapse.

Capitolul șapte este dedicat elaborării unor versiuni alternative pentru algoritmi neuronali clasici, versiuni care vor permite un grad mai înalt de paralelism la nivelul implementării *hardware* pe arhitectura sistolică anterior propusă.

Capitolul opt prezintă concluziile finale, contribuțiile originale ale autorului și inventariază posibilele dezvoltări viitoare pornind de la studiul prezentat în această lucrare.

Am avut șansa ca pe parcursul activității mele profesionale și științifice și a pregătirii acestei teze să întâlnesc o serie de oameni deosebiți al căror sprijin mi-a fost de mare ajutor pe un drum nu totdeauna ușor.

În primul rând aduc un omagiu reputaților profesori Alexandru Rogoian și Vasile Pop, întemeietori ai școlii timișorene de calculatoare și dascăli care mi-au călăuzit primii pași în acest domeniu.

Domnului profesor Nicolae Budișan îi sunt profund recunoscător pentru curajul care mi l-a insuflat în momentele cele mai dificile, pentru răbdare și pentru încrederea cu care m-a onorat atunci când i-am devenit doctorand. Îndrumarea și ajutorul de care m-am bucurat în toți acești ani și-au pus amprenta în mod decisiv asupra tezei pe care o prezint astăzi.

Le sunt recunoscător bunilor mei prieteni prof. Ioan Mureșan, prof. Nicolae Robu, prof. Florin Breabăn, prof. Octavian Proștean și conf. Daniel - Ioan Curiac din cadrul școlii timișorene de automată pentru ajutorul și încurajările de care m-am bucurat în toți acești ani.

Nu pot să nu amintesc excelența colaborare, căldura și prietenia cu care am fost primit în câteva instituții universitare de prestigiu din Belgia, unde am desfășurat o bună parte din activitatea de cercetare prezentată în această lucrare. Le sunt profund recunoscător domnilor profesori Georges Hennuy și Jean - Benoît Cabo de la *ECAM Bruxelles*, Dragomir Mirojevic de la *Univerité Libre Bruxelles* și Michel Verleysen de la *Université Catholique de Louvain - la - Neuve* pentru șansa de a lucra într-un mediu de înalt profesionalism și de a obține rezultate care nu ar fi fost posibile fără această colaborare.

Mulțumesc de asemenea și colegilor mei de la Universitatea "Lucian Blaga" din Sibiu care, într-un fel sau altul, m-au ajutat în munca de elaborare a acestei lucrări; tuturor le sunt profund recunoscător.

Cuvintele mele sunt prea modeste pentru a putea exprima recunoștința ce o datorez familiei. Părinților mei le sunt profund recunoscător pentru grija, atenția și căldura cu care mi-au îndrumat pașii în viață. Lor le datorez și le dedic realizările mele de astăzi. Soției mele îi mulțumesc pentru fericirea deplină pe care a adus-o în viața mea. Iar în ceea ce privește pe copiii noștri, Diana și Lucian, totul devine foarte simplu: abia prin ei totul capătă un sens.

Sibiu, 25.06.1998

Autorul

CUPRINS

INTRODUCERE.....	1
Rețele neuronale în Automatică.....	1
A. Identificarea proceselor.....	2
B. Controlul proceselor.....	3
C. Rețele neuronale în Robotică.....	5
1. MODELE NEURONALE.....	7
1.1. Principiile conexiunismului.....	7
1.2. Scurt istoric.....	8
1.3. Modele conexiuniste.....	9
1.3.1. Modelul McCulloch - Pitts.....	9
1.3.2. Funcția neuronului.....	11
1.3.3. Regula de învățare (modelul lui Hebb).....	15
1.4. Modelul <i>Adaline</i>	16
1.4.1. Introducere.....	16
1.4.2. Metode utilizate.....	17
1.4.3. Notatii.....	18
1.4.4. Căutarea minimului prin metoda gradientului.....	20
1.4.5. Regula Delta.....	20
1.4.6. Metoda gradientului stohastic (regula Delta).....	21
1.4.7. Aplicație.....	21
1.4.8. Extensia modelului <i>Adaline</i>	23
1.4.9. Rețeaua <i>Madaline (Many Adalines)</i>	23
1.4.10. Concluzii.....	25
1.5. Modelul Perceptronului monostrat.....	25
1.5.1. Principiul algoritmului.....	26
1.5.2. Algoritmul de învățare.....	26
1.5.3. Exemplu de funcționare.....	27
1.5.4. Concluzii.....	29
1.6. Modelul retro - propagării gradientului.....	30
1.6.1. Introducere.....	30
1.6.2. Problema Sau - Exclusiv.....	30
1.6.3. Arhitectura RN multistrat.....	32
1.6.4. Algoritmul retropropagării erorii (<i>Error Back - Propagation</i>).....	33
1.6.5. Aplicarea algoritmului EBP.....	38
A. Termenul " <i>momentum</i> ".....	38
B. Utilizarea erorii medii sau instantanee.....	39
C. Inițializarea ponderilor.....	39
D. Panta funcției de activare.....	40
1.6.6. Aplicație demonstrativă.....	41
1.6.7. Concluzii.....	43
1.7. Modelul <i>Hopfield</i>	44
1.7.1. Introducere.....	44
1.7.2. Dinamica rețelelor <i>Hopfield</i>	45
1.7.3. Problema dinamicii inverse.....	47
1.7.4. Regula lui <i>Hebb</i>	47
1.7.4.1. Interpretare.....	48
1.7.4.2. Demonstrație.....	48
1.7.5. Proprietățile rețelei <i>Hopfield</i>	49

1.7.6. Regula proiecției	50
1.7.7. Algoritm de proiecție iterativ	51
1.7.8. Metoda gradientului (Proiecția Delta).....	51
1.7.9. Aplicațiile rețelelor recurente.....	53
1.7.10. Concluzii	55
1.8. Algoritm de auto - organizare al lui Kohonen	55
1.8.1. Observații fiziologice și funcționale	55
1.8.2. Principiul algoritmului	56
1.8.3. Funcția primului strat.....	56
1.8.4. Funcția celui de-al doilea strat	57
1.8.5. Algoritm SOFM simplificat	59
1.8.6. Exemple de utilizare a algoritmului SOFM	60
1.8.7. Concluzii	64
1.9. Modele avansate.....	64
2. IMPLEMENTĂRILE VLSI	67
2.1. Introducere	67
2.2. Soluții posibile pentru implementarea rețelelor neuronale	68
2.3. Structuri candidate la integrare	69
2.3.1. Integrarea neuronilor	70
2.3.2. Integrarea sinapselor	70
2.4. Implementările VLSI analogice	71
2.4.1. Fotoreceptorul	72
2.4.2. Diferențiatorul (celulele bipolare).....	72
2.4.3. Rețeaua rezistivă	73
2.4.4. Rețeaua completă	73
2.5. Implementările VLSI digitale	74
2.5.1. Sistemul CNAPS.....	77
A. Arhitectura CNAPS.....	78
B. Circuitul CNAPS - 1064.....	79
C. Circuitul de secvențiere	80
2.5.2. Sistemul multiprocesor MA 16.....	81
2.5.3. Circuitul NI 1000	83
2.5.4. Alte neuroprocesoare digitale	85
2.6. Principii de proiectare a neuroprocesoarelor sistolice	85
2.6.1. Arhitectura sistemului	86
2.6.2. Proiectarea elementului de procesare (EP)	87
3. ANALIZA IMPLEMENTĂRILOR SISTOLICE	91
3.1. Principiul procesării sistolice	91
3.1.1. Graful de dependență (DG).....	92
3.1.2. Graful de mișcare a datelor (SFG).....	94
3.1.3. Prima soluție sistolică	97
3.1.4. A doua soluție sistolică	98
3.1.5. A treia soluție sistolică	98
3.1.6. A patra soluție sistolică	99
3.1.7. Implantabilitatea rețelelor sistolice	100
3.2. Arhitecturi de RS destinate algoritmilor conexioniști	101
3.2.1. Arhitectura generală	102
3.2.2. Inel sistolic cu componentele vectorului fixate.....	103
3.2.3. Inel sistolic cu sume parțiale fixe.....	106
3.2.4. Arhitectura cu coeficienți sinapțici ficși	107

3.3. Concluzii	109
4. PROIECTAREA RS DEDICATE ALGORITMILOR CONEXIONIȘTI.....	111
4.1. Soluția cu valori de activare și sume parțiale mobile.....	111
4.2. O soluție care implică transmisii locale de date	112
4.2.1. Derularea algoritmului.....	113
A. Prima etapă	113
B. A doua etapă	114
C. A treia etapă.....	115
D. A patra etapă	115
E. A cincea etapă.....	116
F. A șasea etapă.....	117
4.2.2. Detecția convergenței	118
4.2.3. Definierea unei arhitecturi logice	119
4.2.4. Specificațiile celulei sistolice	120
4.2.5. Descrierea Sistemului de Ecuații Recurente Uniforme	123
4.3. Extensia arhitecturii pentru învățare.....	127
4.3.1. Resursele comune diferitelor modele	127
4.3.2. Regula lui Hebb.....	128
4.3.3. Funcționarea RS pentru faza de învățare Hebb	130
4.3.4. Regula Perceptronului	134
4.3.5. Regula lui Kohonen	137
A. Construcția matricii de conexiuni laterale.....	140
B. Normalizarea datelor	142
4.4. Evaluarea Rețelei Sistolice propuse	142
4.5. Arhitectura RS.....	143
4.5.1. Specificațiile celulei sistolice	143
4.5.2. Formatul operanzilor	144
4.5.3. Unitatea de calcul	147
4.5.4. Unitatea de comunicație	149
4.5.5. Adaptarea arhitecturii pentru învățare	150
4.6. Evaluarea performanțelor	150
4.7. Implementarea RS în structuri FPGA.....	152
4.7.1. Testarea funcțională a CS	153
4.7.2. Implementarea CS	161
4.8. Concluzii	161
5. CONSIDERAȚII PRIVIND PERFECTIONAREA RS DEDICATE ALGORITMILOR CONEXIONIȘTI.....	163
5.1. Introducere.....	163
5.2. Arhitectura RS.....	165
5.2.1. Fluxul <i>feed forward</i> de date.....	165
5.2.2. Fluxul de date impus de ajustarea ponderilor sinaptice.....	167
5.3. Fluxul sistolic al instrucțiunilor.....	168
5.4. Structura celulei sistolice.....	169
5.4.1. Dublarea registrului W	170
5.4.2. Precizia calculelor.....	172
5.4.3. Unitatea ALU.....	173
A. Notății și simboluri utilizate.....	173
B. Algoritmi fundamentali	174
C. Înmulțitorul serie - serie	176
5.5. Concluzii	180

6. O PROPUNERE ALTERNATIVĂ: REȚEAUA SISTOLICĂ LINIARĂ	181
6.1. Arhitectura RSL dedicată RN <i>feed forward</i> multistrat	183
6.1.1. Corespondența CS - neuron	183
6.1.2. Implementarea algoritmului EBP pe RSL.....	184
A. Faza <i>feed forward</i>	184
B. Învățarea.....	186
6.2. Analiza performanțelor	189
6.2.1. Faza <i>feed forward</i>	190
6.2.2. Învățarea.....	194
6.3. Considerații privind optimizarea implementării	197
6.3.1. Faza <i>feed forward</i>	197
6.3.2. Învățarea.....	198
6.4. Optimizarea performanțelor aplicațiilor	199
6.5. Concluzii	200
7. ALGORITMI NEURONALI PERFORMANȚI ADAPTAȚI IMPLEMENTĂRII PE REȚELE SISTOLICE	201
7.1. Algoritmul SOFM clasic și algoritmi SOFM modificați	201
7.1.1. Algoritmul SOFM clasic	202
7.1.2. Algoritmul SOFM <i>Batch</i>	202
7.1.3. Algoritmul <i>Batch - M</i>	203
7.1.4. Implementarea algoritmilor clasic, <i>batch</i> și <i>batch - M</i>	203
7.2. Simulări experimentale	206
7.2.1. Inițializare	206
A. Setul de antrenament	206
B. Criteriul de convergență	206
C. Parametrii de învățare.....	207
D. Limitarea ponderilor sinaptice	207
7.2.2. Rezultate obținute	208
A. Evoluția erorii.....	208
B. Viteza de convergență	208
C. Robustețea	211
D. Concluzii	212
8. CONCLUZII.....	213
8.1. Obiectivele lucrării.....	213
8.2. Contribuții	215
8.3. Dezvoltări viitoare	217
Prescurtări utilizate	219
Bibliografie	221

INTRODUCERE

Tehnicile concrete sau punctuale au fost întotdeauna inspirate de vise eterne. Zborul lui Icar, cucerirea spațiului cosmic și a planetelor, înțelegerea universului au constituit din totdeauna motoare puternice care au împins știința spre ungherele cele mai adânci ale cunoașterii.

Iată că, după aspirația ființei sale spre imensitatea externă, omul încearcă descoperirea abisului său interior dând științei un nou pretext: mașinile de procesare a informației nu sunt adaptabile. Ori, în sfera lumii cunoscute, omul pare a fi un rezultat al adaptării. El nu se adaptează numai din necesitate ci și prin voință. În mod esențial, el datorează această aptitudine creierului său și, într-un sens mai larg, sistemului său nervos. Ce poate fi deci mai natural decât a-l observa, a-l înțelege și a extrage principiile directoare pentru construcția calculatoarelor de mâine.

Conexionismul a cunoscut în ultimii ani un interes deosebit, în principal datorită tehnicilor de investigare dezvoltate cu ajutorul informaticii tradiționale. Dar, puterea de calcul a mașinilor tradiționale, care au permis de altfel revenirea la principiul conexiunii, a început să devină o frână în dezvoltarea acestuia.

Bazate pe arhitecturi cu grad înalt și fin de paralelism, modelele conexiunii sunt dezvoltate de regulă pe mașini secvențiale, vectoriale, cu grad grosier de paralelism. Ori, datorită potențialului pus în evidență prin simulare, rețelele conexiunii nu pot fi pe deplin exploatate decât pe mașini care reflectă concret construcția internă a acestora.

Trebuie deci construite mașini optime cu un grad înalt și fin de paralelism; acestea trebuie să fie optime în sensul eficienței și a vitezei de execuție. Scopul acestei lucrări este de a da un răspuns parțial la această problemă. Demersul nostru este canalizat pe 4 direcții:

- studiul modelelor conexiunii;
- elaborarea unei strategii de implementare a rețelelor neuronale pe arhitecturi sistolice;
- perfecționarea rețelelor sistolice dedicate algoritmilor neuronalilor;
- perfecționarea algoritmilor neuronalilor dedicați implementării rețelelor neuronale pe rețele sistolice.

Rețele Neuronale în Automatică

În câmpul cunoașterii suntem în prezent martorii declanșării unor energii creatoare nebănuite prin fenomenul de fuziune al nucleelor unor discipline științifice. Considerate inițial de sine stătătoare, pe o anumită treaptă a dezvoltării acestea își descoperă surprinzătoare rădăcini și fundamentări comune, ceea ce pune în evidență unitatea substanțială și structurală a lumii materiale.

Teoria sistemelor reprezintă un astfel de produs al interacțiunii interdisciplinare care, într-o viziune structuralistă integratoare adâncită până la nivelul cantitativ, reușește să impună perenitatea conceptului de sistem și a celor asociate la nivel de categorii fundamentale în analiza și sinteza fenomenelor lumii materiale, la diverse nivele de organizare a acestora.

Specialistul în automatică este un om cu o pregătire științifică generală, cu preocupări și experiență în multe ramuri bine delimitate ale științei. Pentru el sistemul poate îmbrăca cele mai variate forme fizice: electric, mecanic, biologic, economic, chimic, neuronal sau de orice altă natură. Prima sa sarcină este de a modela matematic sistemul și de a-i extrage caracteristicile fundamentale. Pentru el prezintă mare importanță tipul sistemului sau procesului studiat: liniar sau neliniar, discret sau continuu în raport cu timpul, concentrat sau distribuit, determinist sau stohastic.

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

cu stări continue sau discrete, pasiv sau activ. Automatistul va utiliza toate aceste instrumente în munca sa de concepție și proiectare a sistemelor automate destinate conducerii proceselor.

Sistemele neuronale artificiale, denumite succint rețele neuronale, sunt sisteme fizice celulare capabile să dobândească, să memoreze și să utilizeze cunoștințe. Cunoștințele, stocate sub forma unui ansamblu de stări stabile sau a unor funcții de transfer încapsulate în rețea, pot fi apelate prin aplicarea unor vectori pe intrarea rețelei. Răspunsul acesteia la stimulii aplicați pe intrare se obține pe baza cunoștințelor acumulate sau a comportamentului anterior învățat.

Din punctul de vedere al automatistului rețeaua neuronală este un sistem sau un subsistem destinat integrării într-un sistem. Rețeaua *feedforward* este un sistem capabil să învețe un anumit comportament, cu alte cuvinte, să-și însușească o anumită funcție de transfer. Rețelele neuronale recurente nu sunt altceva decât sisteme dinamice, continue sau discrete. Aplicarea unui vector pe intrarea rețelei înseamnă de fapt forțarea acesteia într-o anumită stare. Rețeaua va evolua apoi liber spre o stare de echilibru numită *atractor*. Fiecare atractor are propriul său set de condiții inițiale, care inițiază evoluții ale rețelei ce sfârșesc în atractorul respectiv. Acest set de condiții inițiale aferente unui atractor este denumit *bazin de atracție*. Dacă atractorul este reprezentat de un punct unic în spațiul stărilor, atunci el se numește atractor fix. Atractorul poate fi reprezentat și de o secvență periodică de stări, caz în care el va fi denumit atractor ciclic, și poate avea și o structură mai complicată. Drept consecință, toate metodele matematice de studiu aplicate în domeniul sistemelor dinamice pot fi aplicate și în domeniul rețelelor neuronale recurente. Criteriile de stabilitate aplicate în domeniul sistemelor dinamice continue sau discrete sunt valabile și în domeniul rețelelor neuronale recurente. Iată deci un prim grad de apropiere, materializat într-un amplu set de similitudini care se cantonează în sfera fundamentelor teoretice.

Cel de-al doilea set de similitudini se cantonează în sfera aplicațiilor. Rețelele neuronale sunt aplicate cu succes la modelarea, identificarea, controlul și optimizarea proceselor; ele pot fi de asemenea utilizate în construcția dispozitivelor de comandă aferente roboților industriali capabili să învețe din experiență.

A. Identificarea Proceselor

În general, prin identificarea unui proces trebuie să se obțină coeficienții ecuației diferențiale aferentă procesului respectiv sau coeficienții funcției sale de transfer. Identificarea proceselor este de foarte mare importanță în cazul sistemelor de control adaptive, deoarece parametrii procesului sunt variabili în timp și controlul adaptiv trebuie ajustat în funcție de variația parametrilor procesului.

Schema de principiu utilizată pentru identificarea funcției de transfer directe aferente unui proces este redată în figura I.1. a).

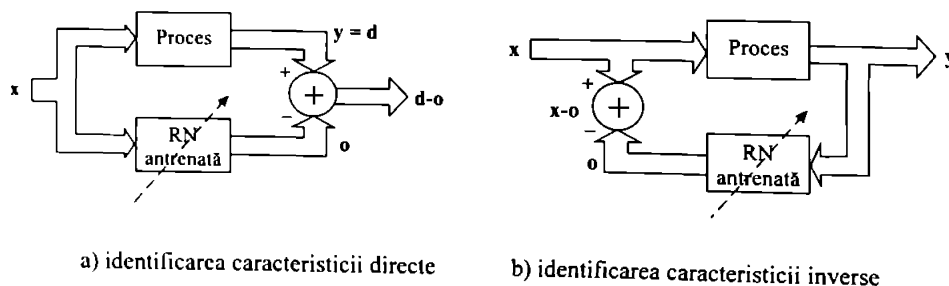


Fig. I.1. Configurarea rețelei neuronale pentru identificarea proceselor

Rețeaua neuronală neliniară recepționează aceeași intrare x ca și procesul identificat. În faza de antrenare a rețelei, procesul va furniza ieșirea dorită d . Scopul identificării constă în a găsi rețeaua neuronală care generează răspunsul o , identic cu răspunsul procesului y , pentru un set dat de vectori de intrare x . Pe durata fazei de identificare, norma vectorului de eroare $\|d - o\|$ va fi minimizată prin ajustarea ponderilor sinaptice (învățare). Să notăm că pot fi utilizate și scheme mai complexe pentru identificarea proceselor. Pentru a lua în considerare dinamica proceselor, pe intrarea rețelei neuronale se va aplica un set de vectori de intrare x întârziați cu diverse cuante de timp. Întârzierile se pot realiza prin inserarea unei cascade de elemente de întârziere pe intrarea procesului.

În multe situații, identificarea caracteristicii inverse a procesului, oferă o alternativă viabilă pentru proiectarea sistemului de conducere aferent procesului respectiv. Pentru identificarea caracteristicii inverse (fig. I.1. b)), ieșirea procesului y va fi utilizată drept intrare în rețeaua neuronală. Vectorul de eroare va fi $x - o$, unde x reprezintă intrarea procesului. Învățarea rețelei are drept scop minimizarea normei vectorului de eroare $\|x - o\|$. După antrenare rețeaua neuronală va găzdui caracteristica inversă aferentă procesului identificat.

B. Controlul Proceselor

În figura I.2. se prezintă un sistem de conducere *feedforward* bazat pe o rețea neuronală.

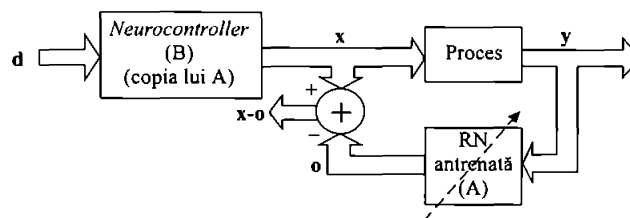


Fig. I.2. Arhitectură de conducere cu învățarea inversei caracteristicii procesului

Neurocontroller-ul B este o copie exactă a rețelei neuronale A, care suportă învățarea. Rețeaua A este astfel antrenată încât să realizeze caracteristica inversă a procesului condus. Vectorul de intrare în *neurocontroller* va fi chiar răspunsul dorit la ieșirea procesului (d). Răspunsul real al procesului va fi y și se datorează intrării x generate de către *neurocontroller*-ul B. Eroarea utilizată la antrenarea rețelei neuronale va fi diferența dintre semnalele de ieșire aferente rețelelor A și B ($x - o$). Deoarece rețeaua B urmărește (copiază) rețeaua A după fiecare pas de învățare, y trebuie să fie identic cu d pentru ca eroarea să fie nulă. Să notăm că această configurație de conducere devine inutilizabilă atunci când caracteristica inversă aferentă procesului nu este unic definită. Această arhitectură de conducere se mai numește și arhitectură cu învățare indirectă și are câteva avantaje evidente. Rețeaua neuronală A poate fi antrenată *on-line* în timp ce *neurocontroller*-ul realizează conducerea procesului; deci nu este necesară o fază de învățare separată. Suplimentar, vectorii de intrare în *neurocontroller* sunt chiar vectorii dorii la ieșirea procesului. Deci, antrenarea rețelei se va realiza chiar în regiunea de interes din domeniul vectorilor de ieșire. O astfel de învățare, realizată în regiunea de interes, se numește învățare specializată. Mai mult, rețeaua neuronală învață continuu și este prin urmare adaptivă.

O arhitectură derivată din cea prezentată în fig. I.2. și care realizează simultan conducerea procesului și învățarea specializată a domeniului de ieșire aferent unui proces static este prezentată în figura I.3. a). Această arhitectură posedă toate avantajele celei descrise anterior, dar utilizează o singură copie a rețelei neuronale. Obiectivul acestei arhitecturi de conducere cu învățare specializată este de a obține vectorul dorit d la ieșirea procesului, în condițiile în care intrarea x care

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

generează această ieșire este necunoscută. Pentru procese dinamice, pe intrarea rețelei antrenate trebuie aplicat și setul de ieșiri aferente procesului, întârziate cu diferite cuante de timp (fig. I.3. b)). Aceasta va permite reconstrucția stării curente aferente procesului, pe baza celei mai recente istorii derulate la ieșirea acestuia. În figura I.3. b) se redă arhitectura de conducere cu învățare specializată pentru un proces dinamic de ordinul k . În acest caz, pe intrarea *neurocontroller*-ului se vor aplica, pe lângă vectorul d , și ultimii k vectori obținuți la ieșirea procesului. Cuantele de timp Δ plasate în buclele de reacție, specifică timpul scurs între eșantioanele (prototipurile) aplicate la intrare în timpul învățării; ele specifică de asemenea că procesul este controlat în timp real.

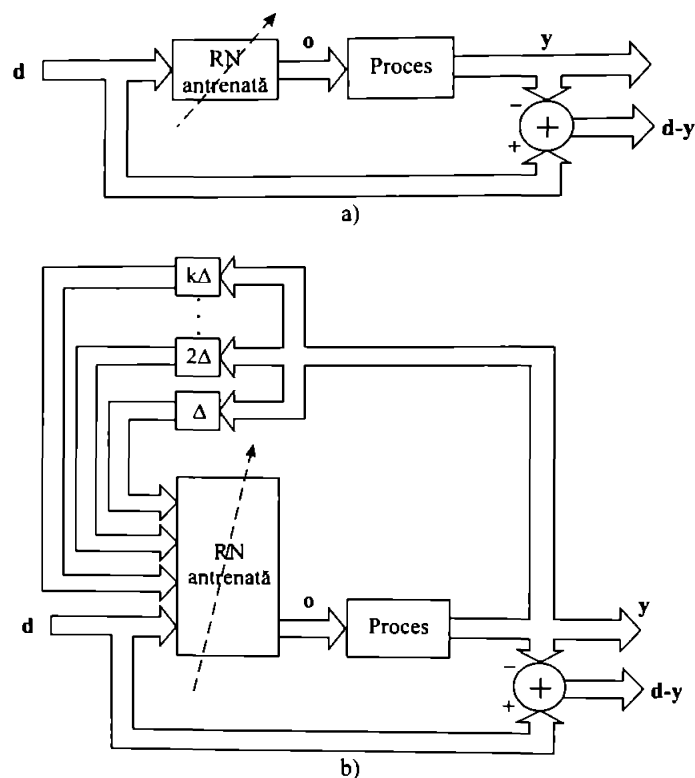


Fig. I.3. Arhitectura de conducere cu învățare *on-line* specializată
a) procese statice b) procese dinamice

Arhitectura de conducere din fig. I.3. are anumite inconveniente care pot conduce la o învățare lentă dacă procesul nu este cunoscut decât aproximativ. Algoritmii de învățare supervizată (ca de exemplu EBP - *Error Back Propagation*) nu pot fi utilizați direct deoarece numai eroarea dintre ieșirea reală y și cea dorită d este cunoscută utilizatorului; eroarea $d - y$ trebuie mai întâi propagată și transpusă la intrarea procesului (furnizată algoritmului EBP).

Să notăm că în literatura de specialitate sunt propuse și alte arhitecturi de conducere; unele dintre acestea combină sistemele de conducere clasice cu cele bazate pe *neurocontroller*-e în ideea obținerii unor performanțe superioare.

C. Rețele Neuronale în Robotică

Deși rețelele neuronale utilizate în robotică sunt, de fapt, *neurocontroller*-e, acestea sunt specializate în rezolvarea problemelor pe care le implică mișcarea brațului unui robot. În vederea obținerii unor performanțe deosebite, cinematica joacă un rol esențial în proiectarea și controlul unui robot. Controlul clasic al traiectoriei brațului robotic constă în succedarea unei secvențe programate de mișcări. Controlul unui robot revendică generarea unor semnale de comandă aplicate dispozitivelor de acționare a încheieturilor robotului, atunci când se specifică o traiectorie dorită sau o secvență de poziții dorită. Brațele robotului operează în plan sau spațiu. Pentru a mișca brațul, coordonatele finale (x , y , z) sunt aplicate *controller*-ului și acesta va genera unghiurile (θ_1 , θ_2 , θ_3) pentru comanda motoarelor care acționează elementele brațului. Pentru a controla mișcarea brațului spre o poziție finală, trebuie rezolvate două probleme:

1. problema cinematicii inverse: Fiind date coordonatele carteziene, specificate fie ca punct final al traiectoriei fie ca o succesiune de puncte de pe o traiectorie, trebuie generate unghiurile, respectiv succesiunea de unghiuri, aferente încheieturilor;
2. controlul poziției finale: Fiind dată poziția finală a brațului robotului, trebuie generată o secvență corespunzătoare de unghiuri necesară atingerii acestei poziții.

Problemele cinematicii directe și inverse pot fi soluționate eficient cu ajutorul rețelelor neuronale *feedforward* multistrat. Modelarea traiectoriilor dorite precum și transformările cinematice pot fi rezolvate cu un înalt grad de acuratețe cu ajutorul *neurocontroller*-elor. Valorile dorite la ieșire (d), în cadrul învățării supervizate, pot fi furnizate fie prin măsurători geometrice fie prin rezolvarea ecuațiilor cinematice directe sau inverse (pentru zona de lucru a robotului).

Rețelele neuronale pot fi utilizate în robotică și pentru învățarea unei secvențe de transformări. Procesul poate fi privit ca o învățare a unei secvențe de transformări în spațiu. Prin urmare, învățarea unei poziții finale poate fi privită ca o învățare a unui control dinamic. Brațele robotului pot fi tratate ca un proces dinamic și pot fi comandate să execute o secvență de pași, cu scopul de a atinge poziția finală dorită. Rețelele neuronale sunt capabile să memoreze diversele secvențe de semnale de comandă necesare atingerii unor poziții finale specificate.

Din cele prezentate până aici, se poate constata un grad înalt de întrepătrundere între domeniul automaticii și cel al rețelelor neuronale. Aceste întrepătrunderi se cantonează atât în sfera fundamentelor teoretice cât și în sfera aplicațiilor. Teoria sistemelor furnizează un set larg de concepte și de instrumente matematice necesare fundamentării teoretice a rețelelor neuronale. Pe de altă parte, rețelele neuronale rezolvă cu succes multe dintre obiectivele vizate în cadrul automaticii: modelarea și identificarea proceselor, conducerea și optimizarea proceselor, robotică etc. Iată de ce considerăm justificată și oportună încadrarea acestei lucrări în sfera mai largă a sistemelor automate.

1. MODELE NEURONALE

Din toate timpurile, creierul a constituit un subiect de reflecție pentru om. Încă din antichitate medicii acordau o atenție particulară acestui organ în care ei au descoperit sursa funcțiilor de comandă. Conexionismul reprezintă astăzi versiunea tehnologică a acestei cercetări. Dacă obiectivul este ambițios, dezvoltările actuale sunt încă sărace. Cunoașterea neuro-psihicului a cunoscut o dezvoltare considerabilă în ultimii 15 ani, dar modelele ingineresti datează de cel puțin 50 ani. Ecartul între observația clinică a creierului și modelele ingineresti este deci destul de mare.

Totuși, cu toată sărăcia rezultatelor obținute, impasul în care se găsește actualmente calculul formal nu ne oferă o altă alternativă. Trebuie căutat, în toate direcțiile posibile, un nou principiu de abordare capabil să asigure uneltelor de tratare a informației autonomia și adaptabilitatea pe care încă nu le au. Conexionismul este încă mai mult speranță decât tehnică.

1.1. Principiile conexionismului

Toată informatica contemporană este bazată pe principiul "mașinii *Von Neumann*". Acest concept cvasi-universal a condus la un model unic de programare a mașinilor:

- stabilirea unui model pentru fenomenul de simulat;
- transcrierea acestui model într-o secvență de instrucțiuni executate pe mașina sau mașinile suport.

Acest principiu operează excelent atunci când se pot extrage din realitatea observată modelele care satisfac exigențele exprimate. Dar modelele devin tot mai complexe și, după toate evidențele, noțiunea de modelare este limitată. Cu tot progresul tehnologic și arhitectural pe care-l înregistrează, calculatoarele "*Von Neumann*" nu fac față cerințelor care devin tot mai complexe.

O soluție alternativă o propune conexionismul care nu apelează la modele și care propune o soluție de "învățare din exemple". Sursa de inspirație a conexionismului este celula nervoasă. O rețea neuronală (RN) este formală dintr-un număr mare de elemente de procesare (EP) simple, numite neuroni, conectate între ele printr-o rețea complexă de conexiuni locale. În sistemele tradiționale (*Von Neumann*), existența unei unități de control unice și a unui tact global conduce spre un comportament individual. Resursele implicate în transferurile de date sunt puține și trebuie gestionate sever. Suplimentar, transferurile de date trebuie sincronizate iar resursele eliberate la sfârșitul fiecărui transfer. La antipod, sistemele conexioniste exploatează reprezentarea distribuită a informațiilor, asociată cu o evoluție complet asincronă care le îndepărtează de noțiunea de control global. Elementele EP sunt numeroase precum și legăturile (transferurile de informații) dintre acestea. Obținerea unei soluții nu va fi un proces previzibil și cuantificabil în manieră deterministă, ci va fi mai mult efectul de manifestare a unei tendințe pe care l-am putea considera ca o satisfacere aproximativă a soluției.

Calculul local simplifică complexitatea elementelor EP (neuronilor) dar crește densitatea căilor de comunicație; conexiunea își pierde statutul de operator ideal pe care îl avea în cadrul sistemelor tradiționale. Mii de neuroni trebuie conectați între ei iar legăturile vor ocupa în mod esențial spațiul disponibil. Reprezentând în sistemele clasice volume și costuri aproape nule, interconexiunile par a deveni "tendonul lui Achile" în cazul rețelelor conexioniste. Este de dorit deci localizarea calculelor și a schimburilor de date în scopul reducerii interconexiunilor la nivelul unei simple comunicații între procesoarele vecine.

Învățarea în cazul RN constă în ajustarea parametrilor rețelei. Regulile și algoritmiile utilizați în procesul de învățare (antrenare) înlocuiesc programarea tradițională din cazul calculatoarelor

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

convenționale. Procesarea în cadrul RN se situează undeva între inginerie și inteligență artificială. Tehnicile matematice utilizate sunt ingineresti dar aplicarea unui principiu experimental adhoc se dovedește adesea necesară.

RN au atras atenția oamenilor de știință dintr-un număr relativ mare de domenii și sunt pe cale să devină un instrument de lucru în inginerie. Inginerii sunt de asemenea interesați în construcția circuitelor electronice neuronale destinate mașinilor inteligente. RN sunt capabile să învețe din experimente, să memoreze, să generalizeze și să utilizeze cunoștințele dobândite.

1.2. Scurt istoric

Geneza ideii de conexiunism datează din anii 1940 și are la bază observațiile neurofiziologice ale lui *Warren McCulloch* [McC 43]. Au urmat apoi lucrările de definire a perceptronului elaborate de *Franck Rosenblatt* [ROS 58], care au pus bazele teoretice pentru dezvoltările actuale. În paralel, și într-un alt cadru de cercetare *Bernard Widrow* propune un model foarte asemănător, modelul *Adaline*. Perceptronul original constituia ieșirea unei rețele artificiale formată din 400 fotoreceptori și era capabil să recunoască motive simple. Această muncă, inițial artizanală, a fost profund reanalizată în lucrările teoretice ale lui *Block* [BLO 62]. Din nefericire entuziasmul lui *Rosenblatt* a fost un pic cam excesiv. El nu a sesizat posibilitățile limitate ale perceptronului; puterea acestuia, lui i se părea extraordinară. În același timp, el era de asemenea conștient de câteva dintre slăbiciunile inerente ale modelului său. El compara perceptronul cu modelul unui creier defect. Analogia era ambițioasă dar conținea problematica esențială a conexiunismului: se poate simula un comportament inteligent pornind de la un ansamblu de elemente simple?

Problema enunțată de *Rosenblatt* a constituit un subiect de controversă în perioada respectivă. Controversa a fost stopată brutal de lucrările lui *Minsky și Papert: "Perceptrons"* [MIN 88]. Ceea ce îi miră încă pe cercetătorii de astăzi, este brutalitatea cu care a fost abandonat perceptronul. Opera lui *Minsky și Papert* conține abordări matematice și vorbește de limitele de generalizare ale modelului. Ea arată explicit ceea ce perceptronul nu poate să facă, dar nu-l respinge.

Din nefericire, în lucrarea "*Perceptrons*", *Minsky și Papert* comit aceeași eroare ca și *Rosenblatt*, dar excelează de această dată prin pesimism. Ei explică faptul că perceptronul monostrat are limite și că aceste limite pot fi depășite dacă se construiesc perceptroni multistrat; cu condiția determinării experimentale (manuale) a parametrilor corecți deoarece nu este posibilă determinarea automată și deci găsirea unui algoritm de învățare pentru o astfel de structură. A trebuit să se aștepte anul 1986, cu lucrările lui *Rumelhart* [RUM 86] și *Le Chun* [LeC 85] pentru a găsi un răspuns la această problemă: da este posibilă găsirea unui algoritm și a unei reguli de învățare pentru rețeaua de perceptroni multistrat; este algoritmul retro-propagării gradientului sau erorii (EBP - *Error Back-Propagation*). Diversele variante ale acestui algoritm sunt încă și astăzi studiate.

Au fost propuse și alte arhitecturi de RN. De exemplu, la începutul anilor 1970, *Kohonen* preia ideea lui *Hebb* și a lui *Willshaw* și propune o rețea cu 2 straturi, dintre care unul are o structură de conexiuni laterale inhibitorii recurente. Este o arhitectură care se detașează de celelalte prin capacitatea ei de învățare nesupervizată (fără profesor).

J.J. Hopfield în anii 1980, relansează interesul pentru rețelele conexiuniste propunând modelul de rețea care îi poartă astăzi numele. Total rebucată, această structură recurentă are o dinamică proprie care se încearcă a fi exploatată prin forțarea stărilor stabile (atractorilor) pe anumite valori cărora li se pot atribui anumite semnificații.

Astăzi cercetătorii se întrec între ei încercând să regăsească filonul inspirației din neurofiziologie. *Herault și Jutten* au propus un algoritm de separare a surselor, direct inspirat de captorii proprio-receptori. *Grossberg* lucrează în strânsă colaborare cu cercetătorii psihologi la elaborarea

modelelor ART (*Adaptive Resonance Theory*). Mead copiază structura retinei pentru construcția circuitelor VLSI analogice.

Domeniul este deci în plină evoluție; el trebuie să-și regăsească sursele de inspirație originale, printr-o strânsă colaborare între biologi, ingineri și matematicieni și trebuie găsite rapid aplicații convingătoare pentru a-i asigura perenitatea și pentru a evita o nouă cădere în unul sau mai multe decenii de uitare. În capitolele următoare vom defini structurile de bază aferente modelelor conexioniste precum și diversele reguli de învățare exploatabile pe aceste structuri. Vom studia de asemenea și modelarea matematică a algoritmilor de învățare.

1.3. Modele conexioniste

1.3.1. Modelul McCulloch - Pitts

Modelul *McCulloch - Pitts* [McC 43] este cu certitudine cel mai vechi model și constituie elementul de bază al rețelelor conexioniste. Acesta descrie neuronul ca pe un element sumator care realizează o sumă ponderată a intrărilor; intrările traversează conexiunile sinaptice iar coeficienții de ponderare vor fi ponderile sinaptice (fig. 1.1).

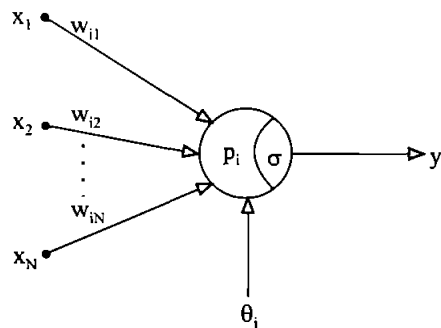


Fig. 1.1. Modelul general al neuronului

Dacă notăm:

x_j - valoarea stimulului de la intrare

w_{ij} - valoarea ponderii sinaptice dintre intrarea j și neuronul i

θ_i - o valoare numită valoare de prag a neuronului i

atunci, potențialul neuronului i va fi:

$$p_i = \sum_{j=1}^N w_{ij} \cdot x_j \quad (1.1)$$

Ieșirea efectivă a neuronului se calculează comparând potențialul p_i cu pragul θ_i . Dacă potențialul depășește valoarea de prag atunci ieșirea y_i va fi +1, dacă nu va fi 0.

$$y_i = 1 \text{ dacă } p_i > \theta_i$$

$$y_i = 0 \text{ dacă } p_i \leq \theta_i$$

Funcția de transfer (activare) σ este o funcție treaptă care generează valori binare la ieșirea neuronului. Aici funcția σ este definită ca în fig. 1.2. a). Această funcție reprezintă o aproximare a răspunsului intrare - ieșire aferent neuronului real.

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

În general, se consideră că funcția neuronală reprezintă o variație de frecvență, medie sau instantanee, în jurul valorii de răspuns aferentă neuronului real. Astfel, aceste variații vor fi pozitive în momentele de excitație și respectiv negative în momentele de inhibiție. Să notăm că funcțiile de transfer de tip neliniar, cum este cea prezentată în fig. 1.3, sunt destul de aproape de caracteristica intrare - ieșire a neuronului real.

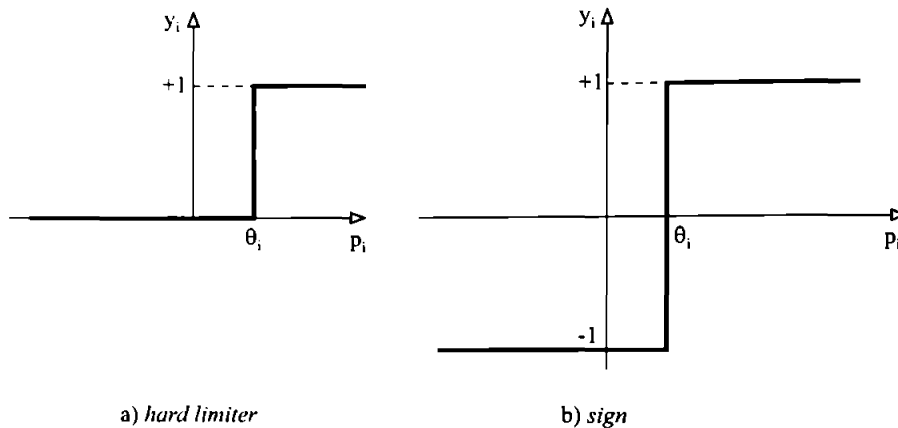


Fig. 1.2. Reprezentarea funcției de transfer (activare)

Funcția de transfer poate lua mai multe forme, care depind de rețea, de natura continuă sau discretă a valorilor potențialului, de aplicația vizată etc. O altă variantă simplă a funcției de activare este funcția *sign*, prezentată în fig. 1.2. b). Ecuația generală a acestei funcții este de forma:

$$y_i(p) = \text{sign}(p_i - \theta_i) \quad (1.2)$$

Accastă funcție generează valori binare pe ieșirea neuronului.

Adesea utilizată în modelele complexe, funcția sigmoidală este destul de aproape de caracteristica intrare - ieșire a neuronilor reali. Forma analitică a acestei funcții este:

a) funcția sigmoidală unipolară:

$$\sigma(p) = \frac{1}{1 + e^{-\frac{(p - \theta)}{T}}} \quad (1.3)$$

b) funcția sigmoidală bipolară:

$$\sigma(p) = \frac{2}{1 + e^{-\frac{(p - \theta)}{T}}} - 1 \quad (1.4)$$

unde:

θ reprezintă valoarea de prag
T controlează panta în origine

Să notăm că în anumite lucrări se notează $\lambda = \frac{1}{T}$.

În fig. 1.3. sunt reprezentate cele 2 variante ale funcției sigmoidale.

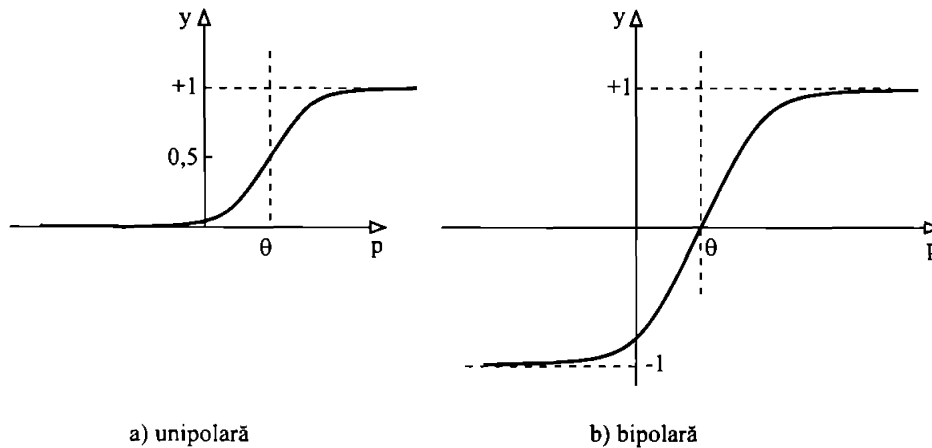


Fig. 1.3. Funcția de activare sigmoidală

1.3.2. Funcția neuronului

Funcția unui neuron *McCulloch - Pitts* poate fi interpretată ușor în termeni geometrici. Să considerăm spațiul R^N înzestrat cu produsul scalar și norma euclidiană. Vom considera valorile de intrare x_1, x_2, \dots, x_N ca și componente ale vectorului x și ansamblul de coeficienți sinaptici $w_{i1}, w_{i2}, \dots, w_{iN}$ ca și componente ale vectorului W_i . În acest caz potențialul neuronului i va fi:

$$p_i = \langle W_i, x \rangle = \|W_i\| \cdot \|x\| \cdot \cos(W_i, x) \quad (1.5)$$

Dacă $\|W_i\|=1$ și $\|x\|=1$ atunci:

$$p_i = \cos(W_i, x) \quad (1.6)$$

Funcția principală a neuronului *McCulloch - Pitts* este deci de a determina gradul de "similaritate de direcție" între 2 vectori. Așadar, neuronul care va genera valoarea cea mai mare pe ieșire (cel mai mare răspuns) va determina care dintre vectorii de ponderi sinaptice corespunde cel mai bine vectorului stimul aplicat pe intrare. Funcția de transfer va avea deci rolul de a selecta unitățile (neuronii) care răspund la un anumit criteriu, adică acelea al căror potențial depășește un anumit nivel (valoarea de prag). Procesul de selecție este descris prin figurile 1.4 până la 1.8; se consideră o RN în care toți neuronii au 2 intrări (spațiul R^2) și ai căror vectori de ponderi sinaptice sunt inițial orientați așa cum indică fig. 1.4.

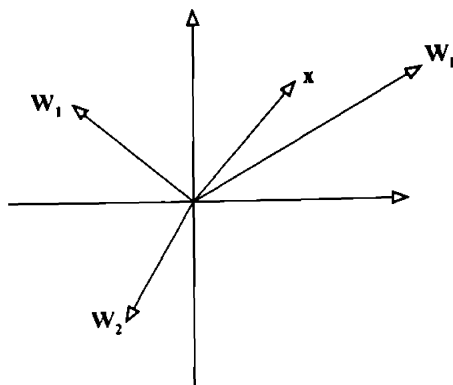


Fig. 1.4. Reprezentarea geometrică a vectorilor de intrare și de ponderi sinaptice

Prima etapă constă în normalizarea vectorilor, inclusiv a vectorului x aplicat la intrare. Pentru normalizare este suficient să se calculeze (fig. 1.5):

$$x' = \frac{x}{\|x\|} \quad \text{și} \quad w'_i = \frac{w_i}{\|w_i\|}$$

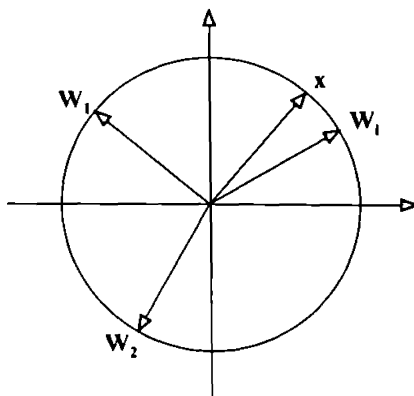


Fig. 1.5. Efectul normalizării vectorilor

Potrivit modelului *McCulloch - Pitts*, vom calcula cosinusul unghiurilor dintre vectorul aplicat la intrare și vectorii ponderilor sinaptice. Vom obține astfel un vector de potențiale p care va conține potențialele tuturor neuronilor. Se aplică apoi funcția de transfer care va avea ca efect selecția unui subsamblu de neuroni în funcție de valorile potențialelor acestora. Dacă se utilizează, de exemplu, funcția "hard limiter", (fig. 1.2.a)), și pragurile θ_i sunt nule, vom clasifica neuronii așa cum indică figura 1.6.

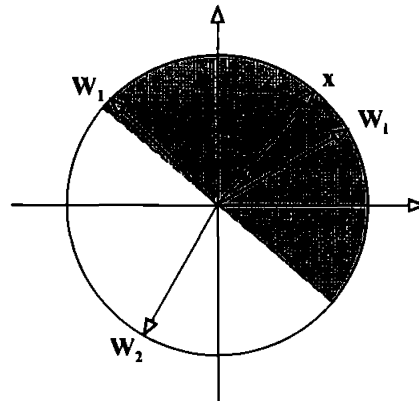


Fig. 1.6. Cele 2 zone din spațiul ponderilor sinaptice separate prin valoarea de prag

Toți neuronii ai căror vectori de ponderi sinaptice generează un cosinus pozitiv în raport cu vectorul de intrare x vor avea ieșirea egală cu 1, toți ceilalți vor avea ieșirea -1. În spațiul ponderilor sinaptice se realizează o selecție a unui subspațiu (format din 2 cadrane) centrat pe vectorul x . Cu $\theta = 0$ spațiul ponderilor sinaptice va fi împărțit în două părți egale. Dacă se fixează o valoare de prag $\theta \neq 0$ identică pentru toți neuronii, dimensiunea spațiului selectat se va reduce proporțional. De exemplu, dacă se utilizează aceeași funcție *hard limiter*, și $\theta = 1/2$ vor fi selectați toți neuronii pentru care $\cos(x, W_i) \leq 1/2$; toți neuronii ai căror vectori sinaptici W_i fac un unghi cu vectorul de intrare x mai mic de 60° vor avea ieșirea în 1, toți ceilalți vor avea ieșirea în 0 (fig. 1.7).

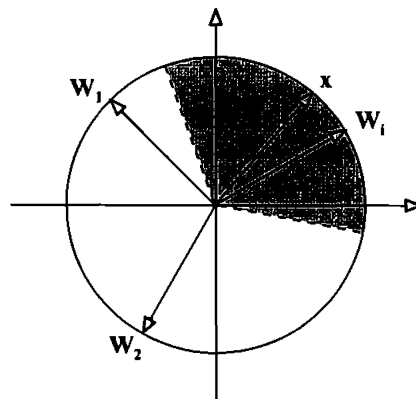


Fig. 1.7. Selecția unui subspațiu în funcție de valoarea de prag

Când datele nu pot fi normalizate, reprezentarea funcției neuronului *McCulloch - Pitts* în coordonate polare devine mai dificilă deoarece, pe lângă direcția vectorilor, va interveni și norma acestora. Se poate considera funcționarea neuronului în coordonate carteziane constatând astfel că el delimitează regiuni separate prin drepte. Să considerăm o RN elementară constituită dintr-un singur neuron cu 2 intrări (fig. 1.8).

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

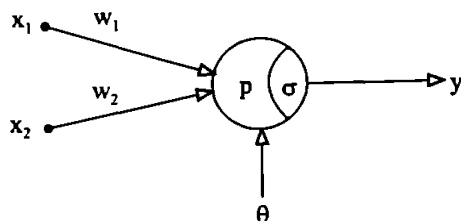


Fig. 1.8. RN elementară constituită dintr-un singur neuron

Dacă se aplică vectorul de intrare $x = (x_1, x_2)$, potențialul neuronului va fi:

$$p = w_1x_1 + w_2x_2$$

și va fi comparat cu intrarea de prag θ .

Dacă $w_1x_1 + w_2x_2 > \theta$ neuronul generează ieșirea $y = 1$, iar dacă $w_1x_1 + w_2x_2 \leq \theta$ neuronul generează ieșirea $y = 0$.

Geometric, aceasta înseamnă că s-a divizat planul (spațiul vectorilor stimuli) în 2 regiuni separate printr-o dreaptă de ecuație:

$$w_1x_1 + w_2x_2 - \theta = 0$$

$$x_2 = -\frac{w_1}{w_2} \cdot x_1 + \frac{\theta}{w_2} \quad (1.7)$$

Fig. 1.9 prezintă funcția rețelei prezentate:

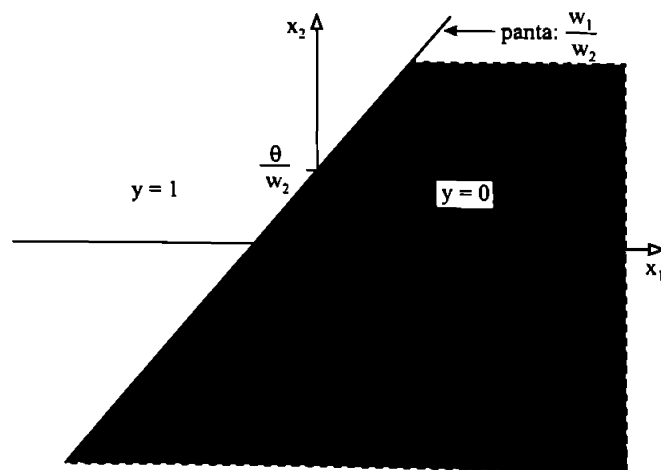


Fig. 1.9. Interpretarea geometrică a funcției neuronului din fig. 1.8.

Aceste interpretări geometrice sunt revelatoare în studiul rețelelor conexioniste.

1.3.3. Regula de învățare (modelul lui Hebb)

Donald Hebb este un neurofiziolog, autorul unei modelări al cărui succes se datorează în mod cert aparentei sale simplități. Studiul său [HEB 49] este axat pe ipoteza existenței unui ansamblu de celule, puternic interconectate, care se excită reciproc și care formează reprezentări ale informației în sistemul nervos. Este vorba de o reprezentare distribuită despre care *Hebb* afirmă că se regăsește atât la nivel funcțional (psihologic) cât și la nivel anatomic. Ideile sale au constituit o bogată sursă de inspirație pentru justificarea modelelor dezvoltate, printre alții, de către *Hopfield* [HOP 82], *Anderson* [AND 77] și *Grossberg* [GRO 76].

Din punct de vedere informatic, esența studiului lui *Hebb* se rezumă la un rezultat experimental calitativ, expus în cartea

"The first stage of perception: growth and assembly":

When an axon of cell A is near enough to excite a cell B and repeatedly or persistently takes part in firing it, some growth process or metabolic change takes place in one or both cells such A's efficiency, as one of the cells firing B, is increased.

Această afirmație a fost intens exploatată pentru construcția modelelor conexioniste, căci ea constituie un enunț cantitativ care pune în evidență plasticitatea sinaptică. Practic, are loc o intensificare sinaptică până la nivelul unei activități comune (unite) a neuronului pre-sinaptic și post-sinaptic. Altfel spus, cauza și efectul intervin pentru modelarea mediului. Se poate astfel elabora o tabelă de adevăr pentru regula lui *Hebb*:

Activitate pre - sinaptică	Activitate post - sinaptică	Regula lui <i>Hebb</i>
Activă	Activă	Intensificare +
Inactivă	Activă	Stare neutră
Activă	Inactivă	Stare neutră
Inactivă	Inactivă	Stare neutră

Să remarcăm că regula lui *Hebb* nu prevede decât o intensificare sinaptică astfel încât ponderile sinaptice pot să crească și să conducă la saturația rețelei. Un amendament [RUM 86] la regula lui *Hebb* a permis ulterior centrarea semnalelor:

Adjust the strength of the connection between cell A and B in proportion of the product of their simultaneous activity.

Această formulare este mai adecvată unui tratament informatic deoarece ea autorizează atât o variație pozitivă cât și una negativă a ponderilor sinaptice. Această regulă este rezumată în tabelul următor:

Activitate pre - sinaptică	Activitate post - sinaptică	Regula de activitate sinaptică
Activă	Activă	Intensificare
Inactivă	Activă	Reducere
Activă	Inactivă	Reducere
Inactivă	Inactivă	Intensificare

Trebuie să notăm că informația necesară calculului variației coeficientului sinaptic aferent unei anume sinapse este disponibilă local, în fiecare conexiune considerată. Variația coeficientului

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

sinaptic depinde doar de activitatea măsurată la cele 2 extremități ale sinapsei respective. În consecință, într-o rețea complexă, fiecare interconexiune își poate evalua propria variație și realizează acest lucru în paralel cu celelalte interconexiuni.

Această regulă de ajustare a ponderilor sinaptice în funcție de activitatea comună (unită) a neuronilor interconectați suferă totuși de anumite limitări, care au condus la dezvoltarea unor reguli mai elaborate (de exemplu regula *Delta*). Trebuie totuși să considerăm modelul lui *Hebb* ca stând la baza tuturor modelelor actuale, care nu fac decât o variație a structurii rețelei sau adaugă un "profesor" pentru procesul de învățare [McG 87]. Din punct de vedere formal, descrierea modificării ponderilor sinaptice în funcție de activitatea comună a neuronilor poate fi înțeleasă ca un calcul iterativ a unei măsuri de corelație între activitatea neuronului pre - sinaptic și potențialul neuronului post - sinaptic.

Reconsiderând reprezentarea neuronului din fig. 1.1, putem elabora prima regulă de ajustare a ponderilor sinaptice:

$$W_{ij}(t + 1) = W_{ij}(t) + \alpha \cdot x_j(t) \cdot y_i(t) \quad (1.8)$$

unde:

W - matricea ponderilor sinaptice

x_j - semnalul pre - sinaptic (activitatea neuronului j)

y_i - semnalul post - sinaptic (activitatea neuronului i)

t - timpul

α - constantă de învățare, mică, cuprinsă în general între 0 și 1.

Toate rețelele conexioniste actuale sunt construite pornind de la diversele variante ale neuronului elementar *McCulloch - Pitts* și de la diversele variante ale regulii lui *Hebb* de ajustare a ponderilor sinaptice.

1.4 Modelul Adaline

1.4.1. Introducere

Modelul *Adaline* (*ADaptive LINear Element*) reprezintă rezultatul cercetărilor lui *Bernard Widrow* și *Marcian Hoff* [WID 60] în domeniul filtrelor adaptive. Un *Adaline* este compus dintr-un singur neuron și dintr-un strat de interconexiuni care conține N ponderi sinaptice (N intrări și o ieșire) așa cum indică fig. 1.10.

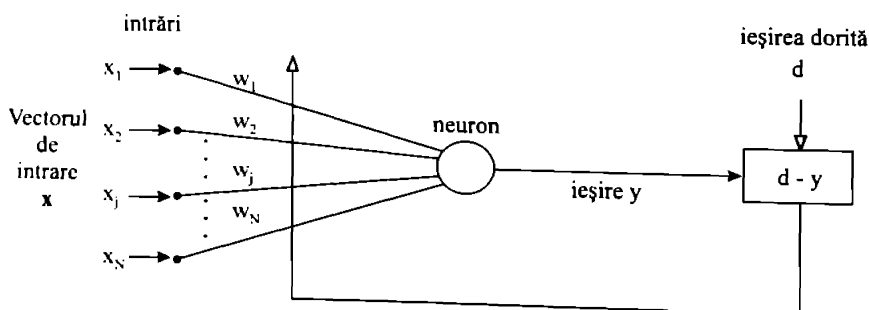


Fig. 1.10. Modelul *Adaline*

Funcția neuronului este identică cu cea a neuronului *McCulloch - Pitts*; potențialul se obține printr-o însumare a intrărilor ponderate cu ponderile sinaptice aferente conexiunilor respective. Într-o primă interpretare vom considera că funcția de transfer a neuronului este funcția identică (ieșirea neuronului este egală cu potențialul acestuia).

Noutatea, în cazul modelului *Adaline*, rezidă în noțiunea de învățare sau de adaptare în sens ingineresc. Vom face presupunerea că sistemul are posibilitatea de a cunoaște ieșirea ideală, d , asociată unui vector de intrare x . Devine posibilă definirea unui termen de eroare (ϵ) pornind de la diferența dintre ieșirea dorită (d) și ieșirea calculată (y). Cunoscând această eroare, scopul devine de a adapta ponderile sinaptice astfel încât eroarea să fie redusă la un minim. Mai teoretic, ideea constă în modelarea unei mărimi (funcții) printr-o combinație liniară de variabile și de a găsi automat cei mai buni coeficienți de ponderare pentru combinația liniară respectivă.

Această metodă de adaptare (învățare) este denumită învățare supervizată deoarece presupune existența unui "profesor" capabil să furnizeze ieșirea dorită (d) atunci când sistemul învață. Pentru a elabora un algoritm adaptiv avem la dispoziție 3 mărimi asupra cărora se va concentra întreaga analiză:

- x - vectorul de intrare (format din N componente reale)
- W - vectorul ponderilor sinaptice (N componente reale)
- d - ieșirea dorită furnizată de profesor (mărime scalară)

$$x = \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_N \end{pmatrix} \quad W = \begin{pmatrix} w_1 \\ w_2 \\ \dots \\ w_N \end{pmatrix}$$

Analiza care urmează este inspirată din lucrarea lui *B. Windrow* [WID 85] în care sunt descrise mai multe metode adaptive pentru procesarea de semnal.

1.4.2. Metode utilizate

Intuitiv, pentru ajustarea ponderilor sinaptice, trebuie căutată o măsură care să indice dacă ajustarea efectuată este justă sau nu. Ori, această măsură o avem deja la dispoziție; este termenul de eroare ($d-y$). Dacă, la o anumită ajustare a ponderilor, eroarea crește, atunci ajustarea nu este corectă. Dimpotrivă, dacă eroarea scade, putem considera că ajustarea ponderilor este corectă. Pentru automatizarea procesului de ajustare, trebuie determinată automat valoarea ce trebuie adăugată ponderilor pentru a produce scăderea erorii. Acest calcul poate fi realizat determinând gradientul erorii în raport cu ponderile sinaptice. Metoda care utilizează această procedură se numește metoda gradientului.

În cadrul metodei gradientului se determină:

- un termen de eroare pătratică;
- gradientul termenului de eroare în raport cu variația ponderilor sinaptice;
- regula de ajustare a ponderilor în funcție de gradientul calculat.

Dacă termenul de eroare este un termen de eroare medie, avem de-a face cu metoda gradientului. Dacă termenul de eroare este eroarea instantanee, avem de-a face cu metoda gradientului stohastic. În cele ce urmează vom dezvolta ambele metode.

Ajustarea ponderilor se va face după următoarea regulă generală:

$$\Delta W(t) = \alpha (-\text{grad}_w \xi(t)) \quad (1.9)$$

622.465/2000

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

unde:

$\Delta \mathbf{W}(t)$ reprezintă ajustarea vectorului \mathbf{W} calculată la momentul t .

α este constanta de învățare cu valoarea cuprinsă între 0 și 1.

$\text{grad}_w \xi(t)$ este gradientul erorii ξ la momentul t în raport cu ponderile sinaptice \mathbf{W} .

1.4.3. Notății

Considerăm rețeaua cu N intrări și o ieșire (fig. 1.10) care operează cu o funcție de activare liniară. Funcția modelului *Adaline* se rezumă la un simplu produs scalar între vectorul de intrare \mathbf{x} și vectorul ponderilor sinaptice \mathbf{W} care conectează intrările la unitatea de ieșire.

Metoda coborării pe direcția gradientului negativ se aplică unui ansamblu de vectori de intrare aplicați la diferite momente de timp. Ieșirea calculată / ieșirea dorită / vectorul ponderilor sinaptice / vectorul prezentat la intrare vor fi notate cu $y(t) / d(t) / \mathbf{W}(t) / \mathbf{x}(t)$. Ieșirea calculată va fi:

$$y(t) = \mathbf{W}^T(t) \cdot \mathbf{x}(t) \quad (1.10)$$

unde \mathbf{W}^T este vectorul \mathbf{W} transpus.

Pentru vectorul de intrare $\mathbf{x}(t)$ vom dispune de ieșirea dorită $d(t)$. Problema este de a ajusta ponderile \mathbf{W} astfel încât, în medie, diferența ($y - d$) să fie nulă. Vom defini un termen de eroare, numit termen de eroare pătratică, cu ajutorul căruia metoda va permite calculul valorilor de ajustare a ponderilor (ajustare adaptivă).

Fixăm:

$$\varepsilon(t) = (d(t) - \mathbf{W}^T(t) \cdot \mathbf{x}(t)) \quad (1.11)$$

de unde rezultă:

$$\varepsilon^2(t) = (d(t) - \mathbf{W}^T(t) \cdot \mathbf{x}(t))^2 = d^2(t) - 2d(t)(\mathbf{x}^T(t) \cdot \mathbf{W}(t)) + (\mathbf{W}^T(t) \cdot \mathbf{x}(t))(\mathbf{x}^T(t) \cdot \mathbf{W}(t)) \quad (1.12)$$

Relația (1.12) reprezintă eroarea pătratică instantanee asociată vectorului $\mathbf{x}(t)$. Dar, noi dorim să calculăm eroarea medie, adică eroarea după aplicarea întregului ansamblu de vectori $\mathbf{x}(t)$. Pentru aceasta vom face ipoteza că $\mathbf{x}(t)$, $d(t)$ și $\varepsilon(t)$ sunt staționare (variabile aleatoare independente și echidistribuite). Putem utiliza media empirică sau experimentală a acestor variabile pentru a calcula eroarea medie pentru întregul ansamblu de vectori de intrare. Aceasta prezintă avantajul că permite previziunea minimumului global precum și a traseului pe suprafața de eroare. Definim deci:

$$\xi = E[\varepsilon^2(t)] \quad (1.13)$$

unde ξ reprezintă eroarea medie pătratică.

În calculele ce urmează vom nota:

$$\text{a) } \mathbf{R} = E[\mathbf{x}(t) \cdot \mathbf{x}^T(t)] \quad (1.14)$$

Matricea \mathbf{R} , simetrică prin definiție, se numește matrice de covarianță a intrărilor (dacă acestea sunt centrate). Analiza detaliată a acesteia permite obținerea de informații importante în legătură cu orientarea suprafeței de eroare. Două elemente au o semnificație geometrică direct exploatabilă pentru a înțelege algoritmul: vectorii proprii ai matricii \mathbf{R} definesc axele principale ale suprafeței de

eroare și valorile proprii ale lui \mathbf{R} dau derivatele secunde ale suprafeței de eroare după aceste axe principale.

$$b) \quad \mathbf{q} = E[d(t) \cdot \mathbf{x}(t)] \quad (1.15)$$

Vectorul \mathbf{q} reprezintă ansamblul de croscorelații între răspunsurile dorite și vectorii de intrare. Utilizând aceste notații, eroarea medie pătratică se va putea scrie:

$$\xi = E[\varepsilon^2(t)] = E[d^2(t)] - 2\mathbf{q}^T \mathbf{W}(t) + \mathbf{W}^T(t) \cdot \mathbf{R} \cdot \mathbf{W}(t) \quad (1.16)$$

Fără a afecta generalitatea acestui calcul, vom utiliza ca măsură a erorii jumătatea mediei pătratice, pentru a obține o expresie simplă pentru gradient la finele calculelor.

$$\xi = \frac{1}{2} \cdot E[\varepsilon^2(t)] \quad (1.17)$$

Pornind de la această expresie a erorii medii pătratice, putem reprezenta forma de variație a erorii (suprafața de eroare) într-un spațiu cu $N + 1$ dimensiuni (primele N dimensiuni corespund celor N ponderi sinaptice w_i preluat ca variabile iar pe a $(N + 1)$ -a axă se reprezintă ξ). Un exemplu de suprafață de eroare în spațiul tridimensional este redat în fig. 1.11. Trebuie să notăm că, în absența neliniarităților, această suprafață va prezenta un singur minim global care va reprezenta soluția problemei noastre. Vom încerca să găsim o metodă iterativă care să conducă spre acest minim având ca scop final determinarea vectorului \mathbf{W}^* care reprezintă soluția dorită.

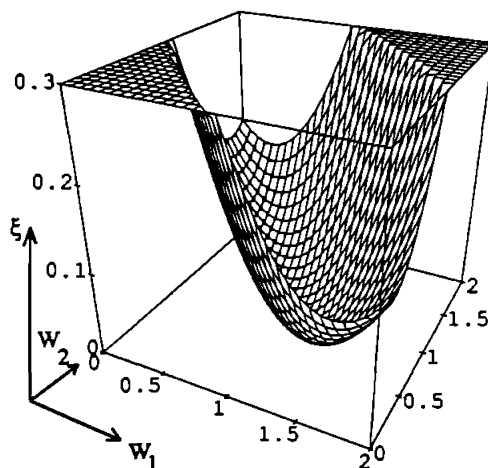


Fig. 1.11. Variația erorii pătratice (suprafața de eroare)

Această metodă a gradientului va permite ajustarea ponderilor sinaptice astfel încât să obținem o evoluție spre minimul global. Metoda constă în a evalua valoarea gradientului erorii într-un punct \mathbf{W} dat pe suprafața de eroare. Acest gradient ne indică clar direcția în care termenul de eroare crește odată cu variația ponderilor sinaptice. Cum noi dorim reducerea termenului de eroare, este suficient să evoluăm în direcție inversă.

1.4.4. Căutarea minimumului prin metoda gradientului

Dacă se impune:

$$\mathbf{W}(t+1) = \mathbf{W}(t) + \alpha \cdot (-\text{grad}_{\mathbf{w}} \xi) = \mathbf{W}(t) - \alpha \cdot (\mathbf{R} \cdot \mathbf{W}(t) - \mathbf{q}) \quad (1.18)$$

obținem un algoritm iterativ în care $\mathbf{W}(t)$ converge spre \mathbf{W}^* , dacă α este mic. Valorile proprii ale matricii $(\mathbf{I} - \alpha \cdot \mathbf{R})$ trebuie să fie toate cuprinse între -1 și +1.

Vom modifica așadar ponderile sinaptice cu o fracțiune din gradientul calculat în punctul respectiv. Parametrul α controlează pasul cu care se coboară spre punctul de minim al suprafeței de eroare.

Această metodă prezintă mai multe inconveniente. Pentru o singură ajustare a ponderilor, trebuie examinați toți vectorii de intrare în cadrul unui ciclu de test; apoi se va minimiza suma erorilor pătratice. Nu este deci o metodă locală din punct de vedere temporal și în afara unor cazuri particulare, va impune restricții greu de respectat în cazul aplicării în timp real. Mai mult, a fost arătat că pentru a obține o corecție care reduce eroarea [GUI 89] trebuie respectată condiția:

$$0 < \alpha < \alpha_{\text{opt}} \quad \text{unde} \quad \alpha_{\text{opt}} = \frac{\text{grad}_{\mathbf{w}} \xi}{\sum_i (\text{grad}_{\mathbf{w}}^T \xi \cdot \mathbf{x}(t))^2} \quad (1.19)$$

Corecția optimală revendică deci două faze succesive în cadrul fiecărui ciclu de baleiere a intrărilor:

1. acumularea vectorului gradient
2. calculul pătratelor produselor scalare între vectorul gradient și vectorii de intrare.

În scopul eliminării acestor inconveniente, *B. Windrow* [WID 60] a propus o metodă de ajustare a ponderilor care nu ține cont de media semnalelor ci de valoarea instantanee a acestora: metoda gradientului stohastic. După cum vom constata, această metodă conduce la același rezultat în ciuda unei convergențe mai haotice.

1.4.5. Regula Delta

Principalul inconvenient al metodei anterioare rezidă din non-localizarea temporală și spațială a calculelor. Pentru a tinde spre soluția optimală, trebuie cunoscută speranța matematică a semnalelor de intrare, ceea ce ne obligă să considerăm o estimare a mediei pornind de la o medie pe termen scurt a semnalelor de intrare.

Pentru a reduce cantitatea calculelor și pentru a propune un algoritm mult mai local, *B. Windrow* a dezvoltat o metodă de ajustare a ponderilor sinaptice în care traiectoria de coborâre spre punctul de minim se estimează la fiecare prezentare a unui vector de intrare. Convergența este mai complexă dar calculele nu au nevoie de media semnalelor de intrare.

Ideea revine la a considera că, pentru fiecare vector de intrare prezentat, vom genera o suprafață de eroare asemenea celei anterior prezentate. Pentru fiecare suprafață nouă, se estimează gradientul erorii în funcție de ponderi și se ajustează ponderile în sensul descreșterii erorii.

Vor exista atâtea suprafețe de eroare câți vectori se aplică la intrare și, la fiecare aplicare, se va pointa spre un punct de minim instantaneu. Convergența spre minimumul global teoretic este asigurată, dar traiectoria este mai haotică.

1.4.6. Metoda gradientului stohastic (regula Delta)

Ajustarea ponderilor, după fiecare prezentare a intrării $x(t)$ se va calcula astfel:

$$W(t+1) = W(t) - \alpha \cdot \text{grad}_w \xi(t) \tag{1.20}$$

Regăsim astfel expresia de ajustare a ponderilor sinaptice identică cu cea prezentată în introducere:

$$W(t+1) = W(t) + \alpha \cdot (d(t) - y(t)) \cdot x(t) \tag{1.21}$$

Și aici parametrul α reprezintă constanta de învățare al cărei rol este de a regla viteza și stabilitatea procesului de învățare. El poate fi menținut constant; în anumite versiuni ale algoritmului, valoarea sa variază în funcție de timp.

Utilizând regula *Delta*, procesul de ajustare a ponderilor spre minimul teoretic se face utilizând informații strict locale în timp. Derularea procesului de învățare va comporta mai mult zgomot pe parcursul primelor iterații dar acest zgomot va fi compensat prin numărul de eșantioane extrase pe durata învățării. Redăm un exemplu de convergență spre un minim global în figura 1.12. Punctul inițial este $[w_1, w_2] = [-5, 3]$ și punctul de minim care trebuie atins este $[w_1^*, w_2^*] = [1, 1]$.

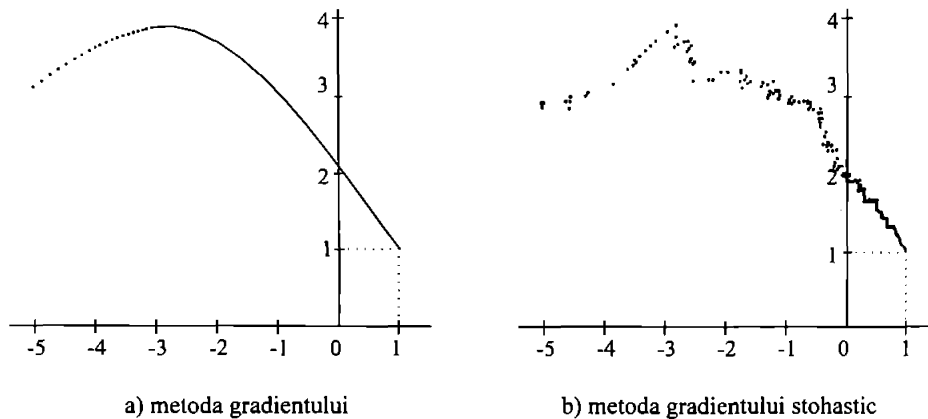


Fig. 1.12. Traectoria ponderilor sinaptice în timpul învățării

În fig. 1.12. a), traectoria ponderilor a fost obținută prin metoda gradientului. Această metodă revendică calculul mediei tuturor semnalelor de intrare. Observăm că traectoria obținută este regulată și atinge direct punctul de minim. În figura 1.12. b) se prezintă traectoria obținută prin metoda gradientului stohastic. Punctul de minim este același dar traectoria urmată este mai haotică. Ea corespunde de fapt diferitelor evaluări ale gradientului, asociate fiecărei curbe de eroare generată de fiecare vector aplicat la intrare.

1.4.7. Aplicație

Funcția de clasificare este una dintre aplicațiile directe ale modelului *Adaline*. Să presupunem că dorim împărțirea în 2 clase a unei mulțimi de puncte (vectori de intrare). Dacă cele 2 clase sunt liniar separabile, este posibilă utilizarea unei rețele *Adaline* pentru a construi automat clasificatorul.

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

În exemplul următor, care este un caz tipic, au fost generate aleator două mulțimi de puncte. Fiecare dintre acestea este centrată pe un punct de referință, celelalte puncte fiind generate cu o varianță comună și o medie distinctă.

Vom utiliza o rețea *Adaline* cu 3 intrări (coordonatele punctelor și valoarea de prag) și o ieșire. Ieșirea dorită este furnizată de un supervisor care indică, pentru fiecare vector de intrare careia dintre cele 2 clase aparține.

Printre cele 3 intrări avem una a cărei valoare este fixată permanent la 1, dar a cărei pondere sinaptică poate varia. Această intrare va juca rolul de prag pentru funcția de activare. În interpretarea geometrică, valoarea de prag va servi la ajustarea ordonatei în origine aferentă dreptei de separare. Rezultatele simulărilor sunt redată în figura 1.13, în 3 momente succesive: la începutul simulării, după 20 și respectiv după 100 vectori aplicați la intrare.

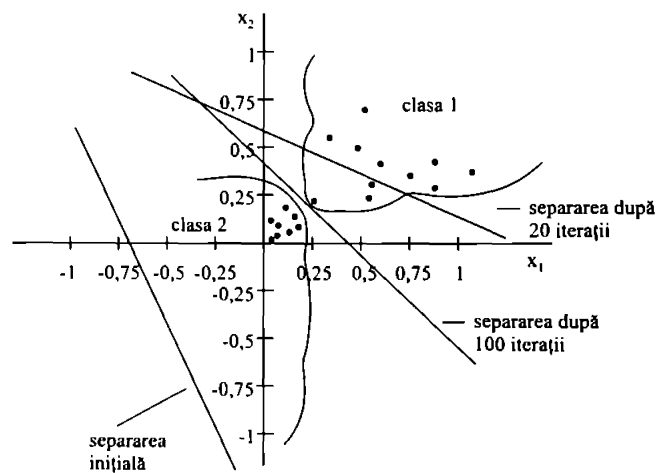


Fig. 1.13. Evoluția dreptei de separare pe parcursul algoritmului

Valoarea erorii evoluează pe durata învățării (fig. 1.14) și descrește asimptotic spre un punct de minim. Să notăm că pentru modelul *Adaline* eroarea este o valoare reală.

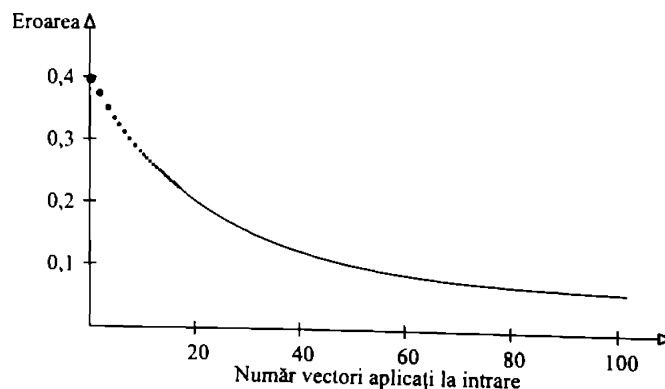


Fig. 1.14. Evoluția erorii pe durata procesului de învățare

În rezumat, metoda de ajustare a ponderilor sinaptice W decurge după următorul algoritm:

1. Inițializarea ponderilor W cu valori mici aleatoare
2. Aplicarea unui vector de intrare $x(t)$ și a unei ieșiri dorite $d(t)$
3. Calculul ieșirii reale:

$$y(t) = \sum_{i=1}^N W_i(t) \cdot x_i(t)$$
4. Ajustarea ponderilor sinaptice după următoarea regulă:

$$W(t+1) = W(t) + \alpha \cdot (d(t) - y(t)) \cdot x(t)$$

$W(t)$ este matricea sinaptică la momentul t .
 α este constanta de învățare, în general cuprinsă între 0 și 1.

1.4.8. Extensia modelului Adaline

Modelul pe care-l vom analiza poate fi utilizat pentru implementarea funcțiilor booleene. Intrările și ieșirea nu vor mai fi valori reale ci valori binare, ceea ce impune utilizarea unei funcții de transfer care transformă valorile continue în valori discrete. O primă soluție propusă constă în utilizarea unei funcții de transfer de tip "hard limiter", asociată cu o funcție de prag.

În articolul din 1960 [WID 60], *Windrow* utiliza un *Adaline* de acest tip pentru clasificarea modelelor; modelele exersate erau literele T, G și F. Valorile asociate la ieșire au fost: -60 pentru T, 0 pentru G și +60 pentru F. Să notăm că, în cazul *Adaline*, evaluarea erorii se face luând în considerare potențialul neuronului și nu ieșirea sa (ieșirea se obține prin aplicarea funcției de transfer asupra valorii potențialului).

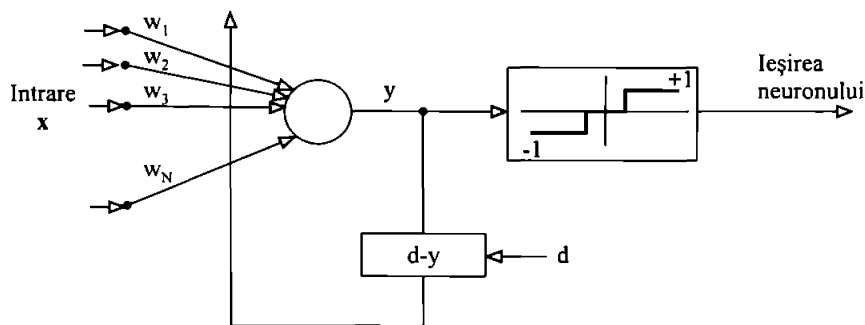


Fig. 1.15. Extensia modelului *Adaline*

1.4.9. Rețeaua Madaline (Many Adalines)

O rețea monostrat înzestrată cu o regulă de învățare de tip *Adaline* nu poate realiza decât o separare liniară. În problemele reale această funcție nu este suficientă. Este însă posibilă extinderea posibilităților rețelei prin adăugarea unui nou strat care va realiza o combinație de separatori liniari.

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

Astfel se vor putea genera operatori capabili să separe două regiuni dintre care cel puțin una este convexă sau închisă. Termenul convexă înseamnă că toate segmentele de dreaptă care unesc două puncte arbitrare situate pe marginea regiunii sunt integral conținute în regiunea respectivă. Această separație poate fi realizată cu o rețea *Madaline* a cărei structură este prezentată în fig. 1.16.

Primul strat de interconexiuni va separa spațiul de intrare prin hiperplane. Fiecare neuron din stratul intern va genera un hiperplan care separă 2 semispații; ieșirea lui va fi activă pentru punctele (vectorii de intrare) dintr-un semispațiu și inactivă pentru punctele din celălalt semispațiu.

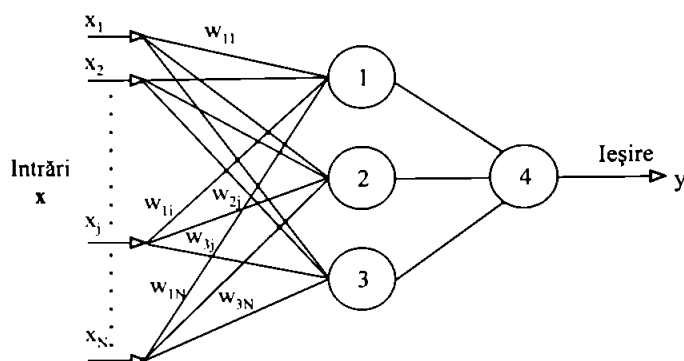


Fig. 1.16. Rețeaua *Madaline*

Al doilea strat de interconexiuni nu realizează decât o conexiune directă (neponderată). Putem considera că este un strat în care toate ponderile sinaptice sunt egale cu 1. Dacă se fixează pragul neuronului de ieșire la valoarea $(P - \epsilon)$, unde P reprezintă numărul de neuroni din stratul intern și ϵ eroarea, atunci funcția neuronului de ieșire va fi un ȘI logic între toate ieșirile stratului intern.

Regiunea de decizie a neuronului de ieșire va fi astfel o intersecție de semispații generate în primul strat. Se va forma deci o regiune convexă care va avea exact atâtea laturi câți neuroni există în stratul intern. În figura 1.17 se prezintă un exemplu de regiuni separabile și un exemplu de regiuni neseperabile printr-o RN de acest tip.

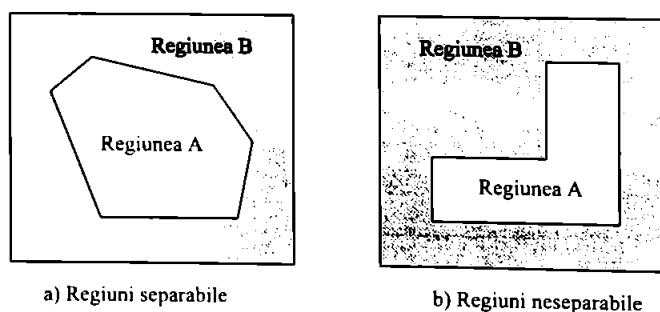


Fig. 1.17. Regiuni de decizie

RN cu 2 straturi de interconexiuni, cu toate posibilitățile lor interesante, sunt neutilizabile în majoritatea problemelor reale. Trebuie deci să se asigure rețelei un grad de libertate suplimentar pentru a putea separa și regiuni non-convexe.

1.4.10. Concluzii

Am studiat în acest paragraf fundamentele RN cu un singur strat de interconexiuni și am enunțat un algoritm de învățare care conduce la o rețea capabilă să facă doar separații liniare.

Descrierea intuitivă a funcționării rețelelor cu 2 straturi de interconexiuni a arătat că este posibilă combinarea separatorilor liniari în scopul construirii unor regiuni convexe. Totuși, a rămas deschisă problema unui algoritm care să permită, în cazul rețelelor multistrat, calcularea automată a ponderilor sinaptice (învățarea). Această problemă a rămas fără răspuns până în anul 1985, dată la care mai mulți autori [LeC 85], [RUM 86] au propus o soluție bazată pe generalizarea algoritmului *Windrow - Hoff*. Vom studia acest algoritm în secțiunea consacrată Perceptronilor multistrat.

1.5. Modelul Perceptronului monostrat

Modelul perceptronului se bazează pe neuronul *McCulloch - Pitts* și descrie un algoritm capabil să ajusteze ponderile sinaptice în funcție de valorile pe care dorim să le învețe. Termenul perceptron acoperă actualmente o clasă mult mai largă de modele decât perceptronul original introdus de *Rosenblatt* [ROS 58]. În figura 1.18 este rezumată structura perceptronului lui *Rosenblatt*.

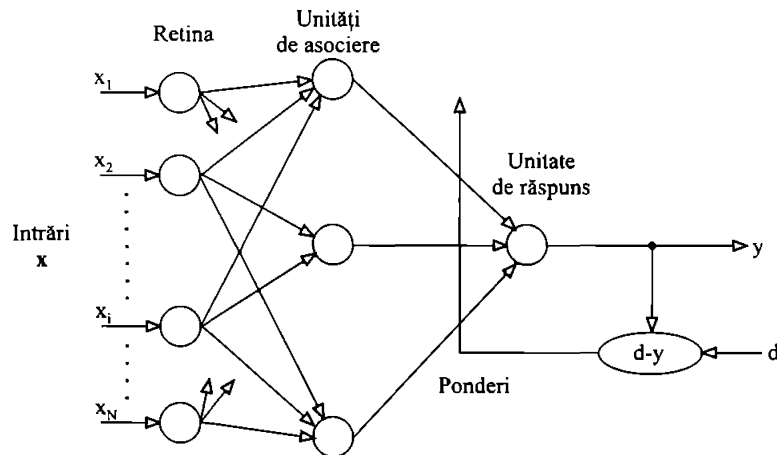


Fig. 1.18. Reprezentarea originală a perceptronului

Perceptronul este compus dintr-un prim strat de conexiuni (numite "*demons*") situat între unitățile de intrare (retina) și unitățile de asociere. Acest prim strat nu face decât o pre - procesare fixă.

Al doilea strat va fi constituit dintr-un singur neuron, conectat la unitățile de asociere prin ponderi sinaptice variabile. Formal, acest perceptron este format dintr-un singur strat de conexiuni plastice (ajustabile). După cum am concluzionat deja în cazul modelului *McCulloch - Pitts*, această structură nu poate funcționa decât ca separator liniar. Problema care se pune rezidă în elaborarea unui algoritm capabil să calculeze ponderile astfel încât să se realizeze o separare convenabilă a claselor.

Termenul perceptron denotă astăzi o structură mult mai simplă decât cea prezentată mai sus. El este mult mai apropiat de modelul *Adaline*, atât prin structură cât și prin regula de învățare.

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

Diferențele apar în 2 puncte:

1. măsura erorii
2. valorile de ieșire

Aceste diferențe au totuși o influență considerabilă asupra algoritmului de învățare. În continuare vom studia perceptronul elementar a cărui structură este prezentată în fig. 1.19. Constatăm că stratul "demons" a dispărut; n-a rămas decât stratul adaptiv asupra căruia se va aplica regula de învățare.

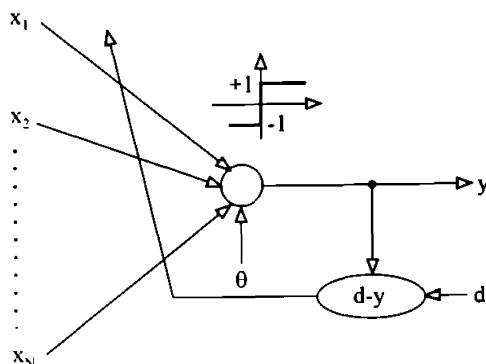


Fig. 1.19. Separator liniar în 2 clase

1.5.1. Principiul algoritmului

Algoritmul [ROS 58] trebuie să determine dreapta care va separa cele 2 clase. Construcția dreptei se va realiza minimizând un termen de eroare prin metoda gradientului stohastic. Pentru aceasta, vom calcula mai întâi potențialul neuronului pornind de la produsul scalar dintre vectorul de intrare x și vectorul de ponderi sinaptice W , vom însuma valoarea de prag și apoi vom aplica funcția de transfer neliniară.

$$p = W^T \cdot x \quad (1.22)$$

$$y = \sigma(p + \theta) \quad (1.23)$$

Dacă vectorii de ponderi sinaptice și respectiv intrare sunt normalizați, potențialul neuronului (p) va conține informații referitoare la gradul de similaritate dintre cei doi vectori (într-o interpretare geometrică unghiul dintre cei 2 vectori). În sfârșit, funcția de transfer (σ) va indica dacă acest potențial depășește sau nu valoarea de prag. În funcție de aceasta neuronul va răspunde cu +1 sau -1.

1.5.2. Algoritmul de învățare

Algoritmul de învățare vizează dotarea perceptronului cu capacitatea de clasificare. Această capacitate poate fi însușită prin compararea răspunsului dorit d (care va fi +1 sau -1), furnizat sistemului de către supervisor (profesor), și ieșirea reală y (care va fi tot +1 sau -1). Este motivul pentru care această tehnică de ajustare a ponderilor sinaptice se numește învățare supervizată.

Perceptronul trebuie să separe liniar 2 clase pe care le vom nota A și B. Pentru fiecare clasă vom asocia o ieșire dorită, $d = 1$ pentru clasa A și $d = -1$ pentru clasa B. Valoarea reală de la ieșirea y va fi tot binară (-1 sau $+1$). Pentru a ști dacă vectorului de intrare x i se atribuie clasa corectă, este suficient să comparăm ieșirea dorită d cu ieșirea reală y :

dacă $(d - y) = 0$,

vectorul x este bine clasat. În acest caz nu se întreprinde nimic.

dacă $(d - y) = 2$ sau $(d - y) = -2$,

vectorul x nu este clasat corect. În acest caz trebuie modificate ponderile sinaptice în scopul minimizării erorii. Soluția constă în a face o corecție după următorul algoritm [BLO 62]:

1. Se inițializează ponderile W cu valori aleatoare mici.
2. Se aplică o intrare $x(t)$ și ieșirea dorită asociată $d(t)$.
3. Se calculează ieșirea reală:

$$y(t) = \sigma(W(t) \cdot x(t))$$

unde σ este funcția "hard limiter"

4. Se ajustează ponderile sinaptice după următoarea regulă:

$$W(t+1) = W(t) + \alpha \cdot (d(t) - y(t)) \cdot x(t)$$

$W(t)$ este matricea (vectorul) de conexiuni la momentul t
 α este constanta de învățare cuprinsă între 0 și 1.

Principalul inconvenient al acestei reguli de învățare provine din ignorarea mărimii erorii atunci când se ajustează ponderile sinaptice (spre deosebire de cazul *Adaline*). Indiferent dacă eroarea este mică sau mare, ponderile sunt ajustate cu o valoare constantă. În general, procedura de ajustare de tip perceptron se bazează pe minimizarea unei distanțe în spațiul ponderilor și nu pe minimizarea unui criteriu de eroare pătratică (cum am constatat în cazul *Adaline* și în alte proceduri mai elaborate). În general, pentru RN multistrat (capabile să realizeze funcții complexe), se aplică criteriul minimizării erorii pătratice.

Problema care se pune acum este de a afla dacă această procedură de ajustare a ponderilor este capabilă să atingă o soluție; cu alte cuvinte, dacă algoritmul este convergent sau nu.

1.5.3. Exemplu de funcționare

Considerăm 2 mulțimi de puncte, reprezentate în planul cartezian în fig. 1.20.

Se va atribui un cod de clasă fiecăruia dintre puncte. Coordonatele punctelor reprezintă vectorii de intrare iar codul de clasă atribuit reprezintă ieșirea dorită.

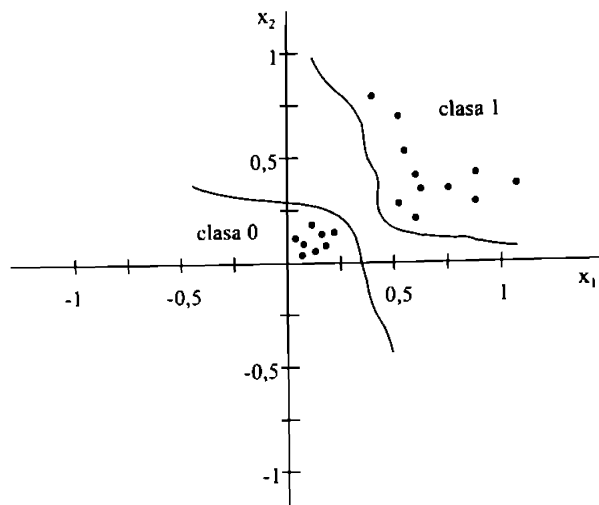


Fig. 1.20. Repartiția punctelor ce urmează a fi clasificate

Pentru a realiza clasificarea se poate utiliza un perceptron cu 2 intrări, o valoare de prag variabilă și o ieșire (fig. 1.21). Valoarea de prag poate fi asimilată cu o intrare de valoare fixă, +1, (polarizare) și cu o pondere sinaptică variabilă (θ). Prin acest artificiu, calculul valorii de prag revine la calculul unei ponderi suplimentare.

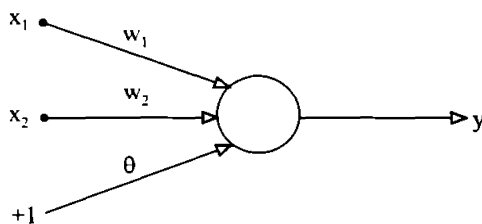


Fig. 1.21. Perceptronul utilizat pentru separarea celor 2 mulțimi de puncte

Se începe prin fixarea ponderilor sinaptice la valori mici aleatoare. Se aplică apoi perechile (vector intrare, ieșirea dorită). Se efectuează produsul scalar între vectorul de intrare aplicat ($x_1, x_2, 1$) și vectorul ponderilor sinaptice (w_1, w_2, θ). Se obține un rezultat care este transformat de funcția de activare a neuronului (în cazul perceptronului funcția *sign*). Rețeaua răspunde cu +1 sau -1.

Această ieșire se compară cu ieșirea dorită și se măsoară (calculează) diferența dintre ieșirea dorită și ieșirea reală. Dacă diferența este nulă, ponderile sinaptice nu vor fi ajustate. Dacă diferența este nenulă, ponderile se ajustează în funcție de intrare (\mathbf{x}), constanta de învățare (α) și semnul diferenței respective.

Această operație se repetă până când nu se va mai înregistra nici o diferență între ieșirea dorită și cea reală (pentru toate perechile de intrare). Trebuie să subliniem că, dacă separarea liniară a celor 2 clase este posibilă, teorema convergenței ne asigură că numărul de corecții aplicate ponderilor sinaptice este finit.

Dreapta care va separa cele 2 clase, în momentul atingerii soluției finale, va avea ecuația:

$$x_2 = -\frac{w_1}{w_2} \cdot x_1 - \frac{\theta}{w_2} \quad (1.24)$$

În momentul $t = 0$ dreapta de separare se va afla într-o poziție inițială aleatoare (fig. 1.22). Dreapta evoluează apoi, ca pantă și ca ordonată în origine, pentru a veni și a se plasa între cele 2 mulțimi de puncte. Din momentul plasării între cele 2 mulțimi de puncte termenul de eroare rămâne pe 0 și convergența este atinsă. Figura 1.22 indică evoluția dreptei separatoare pe parcursul aplicării perechilor de antrenament (x, d) precum și evoluția termenului de eroare în funcție de numărul de cicluri de antrenament efectuate (într-un ciclu de învățare se aplică succesiv toate perechile de antrenament).

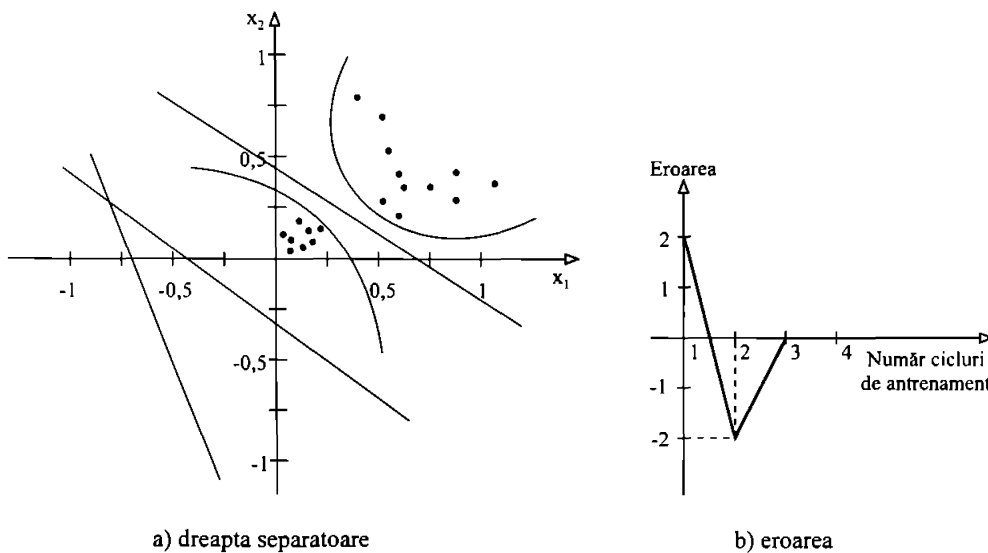


Fig. 1.22. Evoluția procesului de învățare

Cum se poate constata din acest exemplu, din momentul în care termenul de eroare se anulează, el rămâne definitiv nul și nu va mai avea loc nici o corecție. Deoarece convergența algoritmului a fost demonstrată atunci când există o soluție [MIN 88], aceasta este atinsă și ponderile sinaptice rămân într-o configurație stabilă.

1.5.4. Concluzii

Algoritmul perceptronului este un algoritm de învățare de bază care se limitează la probleme care admit soluții pe bază de clasificatori liniari. Aceasta reprezintă o limitare severă care nu permite decât tratarea cazurilor simple. Cu toate acestea, metoda a fost larg utilizată în aplicațiile de recunoaștere a formelor [DUD 73]. În majoritatea aplicațiilor trebuie mers dincolo de modelul perceptronului. Vom vedea în cele ce urmează că există algoritmi de ajustare a ponderilor (învățare) care permit separarea spațiului vectorilor de intrare în regiuni cu contururi oarecare și care depășesc deci limitele acestui model simplu.

1.6. Modelul retro-propagării gradientului

1.6.1. Introducere

Strategiile de învățare aplicate rețelelor monostrat adaptive oferă o funcționalitate redusă la nivelul unui simplu clasificator liniar. Capacitatea de procesare a informației este foarte limitată și nu permite decât rezolvarea problemelor foarte simple. Pentru depășirea acestor limite trebuie utilizate RN multistrat care vor fi apte să soluționeze probleme mult mai complexe. Din nefericire, algoritmul perceptronului nu permite ajustarea ponderilor unui strat care nu-și cunoaște obiectivul pe care trebuie să-l atingă și aceasta este problema care apare în momentul conectării în cascadă a 2 sau mai multe straturi.

Această problemă a limitat dezvoltarea rețelelor de perceptroni multistrat și o soluție a fost propusă doar relativ recent de *Yann Le Cun* [LeC 85], *Rumelhart, Hinton și Williams* [RUM 86] și *Paul Werbos* [WER 88], pentru a nu cita decât autorii cei mai cunoscuți. Vom analiza mai întâi carențele perceptronilor monostrat, pentru a arăta mai apoi care sunt posibilitățile oferite de perceptronii multistrat.

1.6.2. Problema SAU - Exclusiv

Pentru a sesiza limitele intrinseci ale modelului perceptron, vom studia următoarea problemă: implementarea funcției SAU - Exclusiv utilizând un singur perceptron. În figura 1.23 s-a reprezentat funcția SAU - Exclusiv în spațiul vectorilor de intrare.

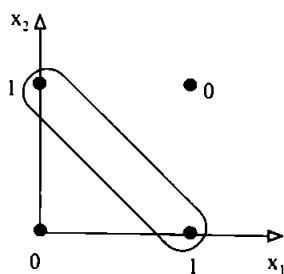


Fig. 1.23. Reprezentarea problemei SAU - Exclusiv

Fiecărui punct din fig. 1.23 îi este asociată ieșirea, care denotă clasa la care el trebuie să aparțină. După cum se observă, punctele din aceeași clasă (0 sau 1) nu pot fi separate cu o dreaptă. Această problemă relativ simplă nu poate fi rezolvată cu ajutorul unui singur perceptron. Dacă utilizăm o rețea bazată pe unități binare cu prag, cu o unitate ascunsă, de forma celei din fig. 1.24, problema poate fi rezolvată [McC 43].

Acest comportament se explică simplu, constatând că adăugarea unei unități ascunse permite rescrierea ecuației funcției SAU - Exclusiv adăugând un termen redundant:

$$y = x_1 \oplus x_2 = (x_1 + x_2) \cdot \overline{(x_1 \cdot x_2)}$$

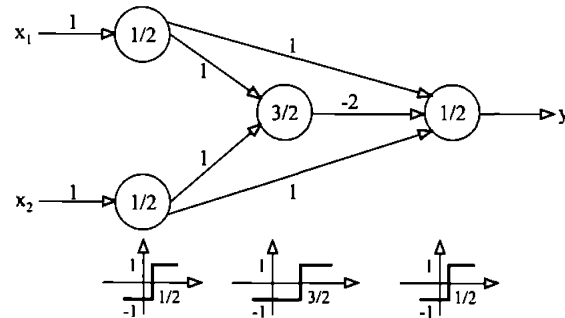


Fig. 1.24. Rețea de unități binare pentru rezolvarea problemei SAU - Exclusiv

Reprezentarea geometrică se va face acum utilizând o dimensiune suplimentară; pe axa z se reprezintă termenul $x_1 \cdot x_2$ (fig. 1.25).

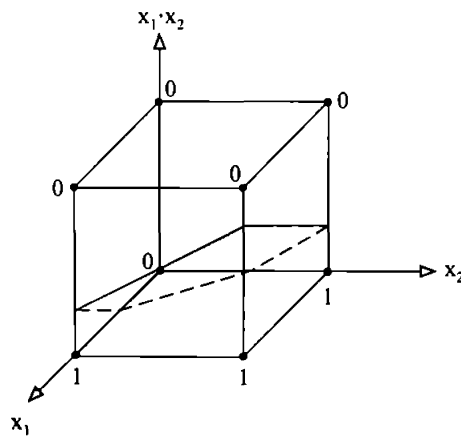


Fig. 1.25. Planul de separare a celor 2 clase

Am găsit astfel o soluție pentru problema SAU - Exclusiv, bazată pe o rețea de celule liniare cu valori de prag asimetrice și care generează un plan de separare capabil să claseze convenabil vectorii de intrare. Totuși, ceea ce am făcut manual relevă o anumită complexitate care dă naștere la o întrebare: cum vom alege numărul de unități ascunse (pentru o problemă dată) și cum își va ajusta rețeaua ponderile și valorile de prag pentru a rezolva problema dată?

Pentru RN multistrat, *Y. Le Chun* [LeC 85] a propus minimizarea unui termen de eroare pătratică printr-o metodă similară celei propuse de *Windrow* la modelul *Adaline*. Metoda presupune bineînțeles existența unui supervisor capabil să genereze ieșirea dorită pentru fiecare vector aplicat la intrare. Inovația soluției lui *Le Chun* rezidă în posibilitatea de a comunica recursiv neuronilor din straturile ascunse termenul de eroare care le revine (contribuția fiecăruia la eroarea globală de la ieșire). Cunoscându-și starea precum și obiectivul care trebuie atins (minimizarea propriului termen de eroare), neuronii din straturile ascunse își vor putea ajusta ponderile sinaptice în acord cu acest obiectiv.

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

Totuși, metoda propusă de *Le Chun* nu oferă răspuns la întrebarea: care este arhitectura optimă a RN chemată să rezolve o anumită problemă? Metoda permite doar calculul ponderilor pentru o rețea cu straturi ascunse (neuroni care nu aparțin stratului de ieșire).

Singurele rezultate exploatabile, referitoare la optimizarea arhitecturii RN, se referă la capacitatea de clasificare în funcție de numărul de straturi. Lucrări relativ recente au arătat că o rețea cu 2 straturi permite realizarea oricărei funcții booleene în raport cu variabilele de pe intrare [HAM 87] și că o rețea cu 3 straturi este suficientă pentru rezolvarea tuturor problemelor de clasificare [LIP 87]. Noi presupunem că există de asemenea și un număr optim de neuroni pe fiecare strat dar acesta nu putem încă să-l determinăm cu precizie. După cum vom vedea, acest parametru are o importanță crucială asupra vitezei de convergență a algoritmului de învățare și chiar, pur și simplu, asupra posibilității de a găsi o soluție la o problemă.

În general, o RN care conține în structura sa neuroni ascunși va putea genera o suprafață de decizie complexă și, prin urmare, va putea rezolva probleme de clasificare complexe. Suplimentar, utilizând o funcție de transfer neliniară la nivelul neuronilor, capacitatea RN crește mult putând atinge o capacitate vecină metodelor de interpolare neliniară.

1.6.3. Arhitectura RN multistrat

Într-o rețea multistrat, toate celulele stratului de intrare sunt urmate de o succesiune de straturi, referite cu termenul straturi ascunse, pentru ca în final să se ajungă la stratul de ieșire. Când rețeaua este simplu etajată, fiecare neuron ce aparține unui strat nu va primi pe intrare decât ieșirile neuronilor din stratul anterior. În sens larg, într-o RN multistrat, nimic nu ne împiedică să luăm în considerare și interacțiuni între straturi îndepărtate. În general se face uz doar de interconexiuni locale între straturi adiacente așa cum arată figura 1.26.

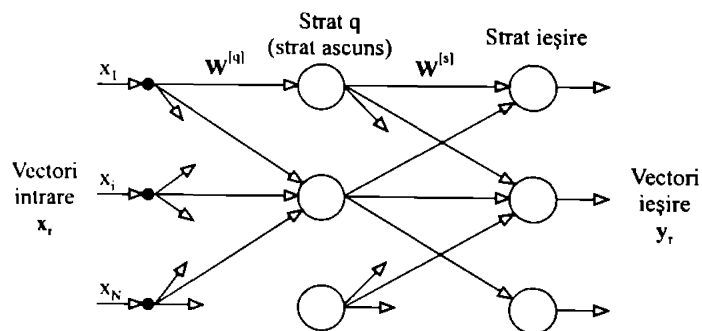


Fig. 1.26. Arhitectura generală a rețelei multistrat

Prezența stratului intermediar (straturilor intermediare) aduce un plus de bogăție parametrică care trebuie să conducă la creșterea capabilităților rețelei. Trebuie să menționăm faptul că neuronii interni nu au nici o conexiune predefinită; ei vor contribui doar la obținerea rezultatelor dorite la ieșire.

O altă diferență, esențială, se referă la funcțiile de transfer aferente neuronilor. Vor fi funcții neliniare, numite și "funcții de compresie" pentru că ele comprimă potențialul neuronului în intervalul $[-1, +1]$. Funcțiile de transfer utilizate trebuie să fie continue și derivabile (vom vedea de ce). Prin urmare, neuronii emit valori continue și nu discrete pe ieșire.

Pentru un neuron i din stratul $[q]$, considerând că $x_j^{[q]}$ sunt intrările în acest strat, potențialul $p_i^{[q]}$ va fi:

$$p_i^{[q]} = \sum_j W_{ij}^{[q]} \cdot x_j^{[q]} \quad (1.25)$$

iar ieșirea $y_i^{[q]}$ se va obține pornind de la acest potențial:

$$y_i^{[q]} = \sigma(p_i^{[q]}) \quad (1.26)$$

Funcția de transfer σ trebuie să fie continuă și derivabilă; se utilizează de regulă funcția sigmoidală unipolară sau bipolară (relațiile (1.3) și (1.4)) sau alte funcții continue și derivabile cum ar fi:

$$\sigma(p) = \text{th}(p) = \frac{e^p - e^{-p}}{e^p + e^{-p}} \quad (1.27)$$

Arhitectura pe care o descriem comportă mai mulți parametri precum ar fi: forma funcției de transfer, numărul de straturi, numărul de neuroni per strat. Nu există încă o metodă analitică pentru determinarea acestor parametri structurali. Este nevoie de experiență în domeniu pentru a ajunge la cel mai bun model pentru o anumite aplicație.

1.6.4. Algoritmul retropropagării erorii (Error Back - Propagation)

Rețeaua multistrat ("feedforward") are o structură formată din straturi succesive, traversate de semnalele de intrare care se propagă până la ieșire. Dacă dorim să transpunem algoritmul de învățare dezvoltat pentru rețelele monostrat, va trebui să definim un termen de eroare. Învățarea (ajustarea ponderilor sinaptice) va urmări minimizarea acestui termen de eroare. Termenul de eroare poate fi definit dacă există un set de vectori de intrare pentru care se cunosc valorile dorite la ieșirea RN. Învățarea va fi deci supervizată și se bazează pe un set de perechi de învățare sau de antrenament (x_r, d_r) . Termenul de eroare va fi definit ca în fig. 1.27.

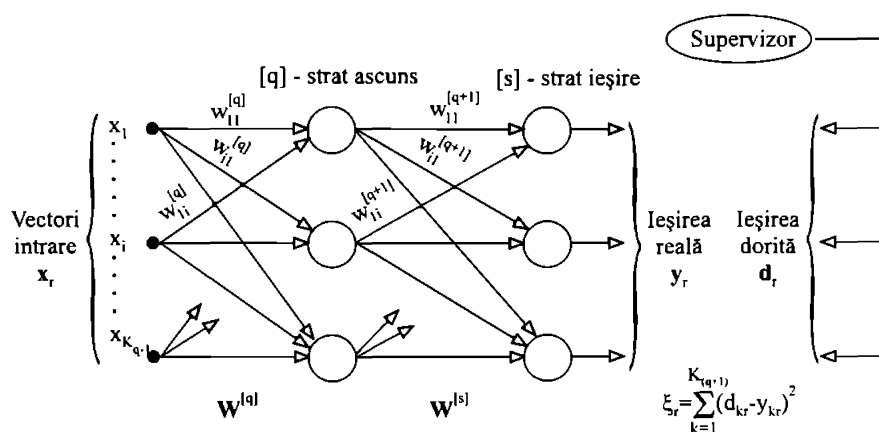


Fig. 1.27. Arhitectura rețelei feedforward și definirea termenului de eroare

Se utilizează strategia de minimizare a erorii ξ_r bazată pe gradientul stohastic. Dar, așa cum se observă din fig. 1.27, numai stratul de interconexiuni asociat stratului de ieșire cunoaște direct

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

valoarea de minimizat. Straturile interne, asociate neuronilor ascunși, nu dispun de un termen de eroare pentru operația de minimizare. Rezultă că, pentru straturile ascunse, nu se vor putea calcula direct valorile de ajustare a ponderilor sinaptice.

Principiul algoritmului *Error Back - Propagation* (EBP) constă în ideea de a impune o funcție obiectiv neuronilor din straturile ascunse. Pe baza acestei funcții obiectiv se vor putea calcula și valorile de ajustare a ponderilor sinaptice din straturile ascunse. Funcția obiectiv trebuie să rezulte din eroarea ξ_r , măsurată la ieșirea rețelei. Un neuron situat într-un strat intern nu poate să vadă eroarea calculată la ieșire decât traversând toate ponderile sinaptice care-l leagă de stratul de ieșire. Așadar, algoritmul ar putea modifica ponderile aferente conexiunilor care intră în neuronul respectiv cu scopul de a reduce participarea acestuia la eroarea calculată la ieșire. Aceasta înseamnă că eroarea calculată la ieșirea rețelei este propagată în sens invers cu aceleași ponderi sinaptice pe conexiuni (identice cu cele utilizate în faza *feedforward* în care se generează ieșirea y_r). În acest mod, pe ieșirea neuronului din stratul ascuns, ajunge un termen de eroare care reprezintă contribuția sa la eroarea globală calculată la ieșirea rețelei; pe baza acestui termen de eroare vor putea fi ajustate ponderile sinaptice aferente conexiunilor de intrare în neuronul respectiv. În consecință:

- ieșirile neuronilor din faza *feedforward* devin puncte de însumare a erorilor retro-propagate;
- punctele de însumare din faza *feedforward* (intrările) devin ieșiri pentru erori;
- funcțiile de transfer ale neuronilor devin propriile lor derivate.

Procedura de ajustare a ponderilor sinaptice prevăzută în cadrul algoritmului EBP va fi:

Date inițiale: Setul de antrenament format din R perechi

$(\mathbf{x}_r, \mathbf{d}_r)$; $r = 1, 2, \dots, R$

\mathbf{x}_r - vectorul aplicat la intrare

\mathbf{d}_r - vectorul dorit la ieșire (aferent vectorului de intrare \mathbf{x}_r)

1. Se inițializează toate ponderile sinaptice (matricile $\mathbf{W}^{(q)}$, unde $q = 1, 2, \dots, s$) cu valori mici aleatoare.
2. Se aplică o pereche de antrenament $(\mathbf{x}_r, \mathbf{d}_r)$;
 $r = 1, 2, \dots, R$

NOTA: Pentru o convergență cât mai bună, perechile de antrenament $(\mathbf{x}_r, \mathbf{d}_r)$ se selectează în ordine aleatoare din setul de antrenament (minimizarea erorii se face pe baza gradientului stohastic!)

3. Se calculează ieșirea reală prin propagarea directă (*feedforward*) a vectorului de intrare prin toate straturile RN:

$$y_{jr}^{(q)} = \sigma\left(\sum_{i=1}^{k_q} w_{ji}^{(q)} \cdot y_{ir}^{(q-1)}\right) = \sigma\left(\sum_{i=1}^{k_q} w_{ji}^{(q)} \cdot x_{ir}^{(q)}\right) \quad (1.28)$$

unde σ este funcția de transfer și $q = 1, 2, \dots, s$.
Pentru stratul de intrare $q = 1$, $y_{ir}^{(q-1)} = x_{ir}$.

4. Se calculează eroarea la ieșire:

$$\xi_r = \frac{1}{2} \sum_{k=1}^{k_s} (d_{kr} - y_{kr}^{(s)})^2 \quad (1.29)$$

unde: \mathbf{d}_r - ieșirea dorită asociată vectorului de intrare \mathbf{x}_r
 $y_r^{[s]}$ - ieșirea reală din RN obținută în urma prezentării vectorului de intrare \mathbf{x}_r .
 ξ_r - eroarea la ieșire calculată pentru cuplul $(\mathbf{x}_r, \mathbf{d}_r)$

5. Retro - propagarea semnalului de eroare (δ_{ir}) dinspre stratul de ieșire spre stratul de intrare.

a) pentru neuronii din stratul de ieșire:

$$\delta_{ir}^{[s]} = \sigma'(p_{ir}^{[s]}) \cdot (d_{ir} - y_{ir}^{[s]}) \quad (1.30)$$

b) pentru neuronii din straturile ascunse:

$$\delta_{ir}^{[q]} = \sum_{k=1}^{K_q} \delta_{kr}^{[q+1]} \cdot w_{ki}^{[q+1]} \cdot \sigma'(p_{ir}^{[q]}) \quad (1.31)$$

unde $\sigma'(p)$ reprezintă derivata funcției de activare.

6. Ajustarea ponderilor sinaptice după regula:

$$\Delta w_{ij}^{[q]} = \alpha \cdot \delta_{ir}^{[q]} \cdot x_{jr}^{[q]} \quad (1.32)$$

unde α este constanta de învățare cuprinsă între 0 și 1.

7. Se revine la pasul 2 atâta timp cât mai există perechi $(\mathbf{x}_r, \mathbf{d}_r)$ de aplicat.

8. Un ciclu de învățare se încheie atunci când sunt epuizate (aplicate) toate perechile de antrenament. Cu ponderile sinaptice $\mathbf{W}^{[q]}$ obținute în ciclul de învățare curent se va starta un nou ciclu. Învățarea se încheie când se constată convergența algoritmului (eroarea ξ_r atinge un minim și rămâne staționară).

Nu vom prezenta în această lucrare demonstrația algoritmului EBP. Această demonstrație poate fi găsită în [LeC 85], [RUM 86], [WER 88] precum și în alte lucrări care au preluat această problematică.

Schema bloc a algoritmului EBP este prezentată în figura 1.28.

Pe schema bloc din figura 1.28 se disting 2 faze:

1. faza directă (*feedforward*) în care sunt implicate blocurile (1) și (2) și în care RN, pe baza vectorului aplicat la intrare (\mathbf{x}_r), generează ieșirea reală y_r . Blocurile (1) realizează produsul între o matrice și un vector:

$$\mathbf{p}^{[q]} = \mathbf{W}^{[q]} \cdot \mathbf{y}^{[q-1]} \quad (1.34)$$

unde: $\mathbf{p}^{[q]}$ - vectorul potențialelor neuronilor din stratul [q]
 $\mathbf{W}^{[q]}$ - matricea ponderilor sinaptice aferentă neuronilor din stratul [q]
 $\mathbf{y}^{[q-1]}$ - vectorul de ieșire din stratul anterior [q-1].

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

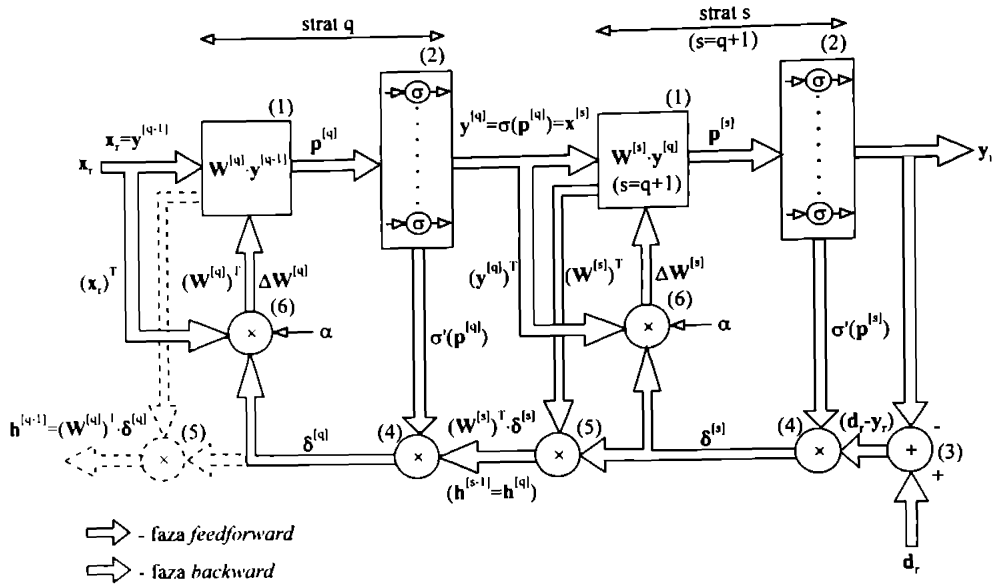


Fig. 1.28. Schema bloc a algoritmului EBP

$$p^{[q]} = \begin{bmatrix} p_1 \\ p_2 \\ \dots \\ p_{K[q]} \end{bmatrix}^{[q]}; \quad W^{[q]} = \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1K[q-1]} \\ w_{21} & w_{22} & \dots & w_{2K[q-1]} \\ \dots & \dots & \dots & \dots \\ w_{K[q]1} & w_{K[q]2} & \dots & w_{K[q]K[q-1]} \end{bmatrix}^{[q]}; \quad y^{[q-1]} = \begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_{K[q-1]} \end{bmatrix}^{[q-1]}$$

- w_{ij} - reprezintă ponderea sinaptică care leagă ieșirea neuronului j din stratul $[q-1]$ cu intrarea neuronului i din stratul $[q]$
- p_i - reprezintă potențialul neuronului i din stratul $[q]$
- y_j - reprezintă ieșirea neuronului j din stratul $[q-1]$
- $K[q]$ - numărul de neuroni din stratul $[q]$
- $K[q-1]$ - numărul de neuroni din stratul $[q-1]$

Linia i a matricii $W^{[q]}$ conține ponderile conexiunilor de intrare în neuronul i din stratul $[q]$. Coloana j a matricii $W^{[q]}$ conține ponderile conexiunilor de ieșire aferente neuronului j din stratul $[q-1]$.

Blocurile (2) implementează funcția de transfer σ :

$$y^{[q]} = \sigma(p^{[q]}) \tag{1.35}$$

Să notăm că în cadrul aplicațiilor faza 1 este unica fază parcursă. În procesul de învățare faza 1 este urmată de faza 2 (ajustarea ponderilor).

2. faza de propagare a erorii înapoi (backward) în care sunt implicate blocurile (3), (4), (5) și (6). Blocul (3) calculează diferența dintre vectorul dorit la ieșire \mathbf{d}_r și vectorul de ieșire real \mathbf{y}_r . Blocurile (5) realizează produsul dintre o matrice și un vector:

$$\mathbf{h}^{[q]} = (\mathbf{W}^{[q+1]})^T \cdot \delta^{[q+1]} \quad (1.36)$$

unde:

$(\mathbf{W}^{[q+1]})^T$ - transpusa matricii $\mathbf{W}^{[q+1]}$
 $\delta^{[q+1]}$ - vectorul de eroare asociat stratului $[q+1]$

$$\delta^{[q+1]} = \begin{bmatrix} \delta_1 \\ \delta_2 \\ \dots \\ \delta_{K[q+1]} \end{bmatrix}^{[q+1]}; \quad (\mathbf{W}^{[q+1]})^T = \begin{bmatrix} w_{11} & w_{21} & \dots & w_{K[q+1]1} \\ \dots & \dots & \dots & \dots \\ w_{1K[q]} & w_{2K[q]} & \dots & w_{K[q+1]K[q]} \end{bmatrix}^{[q+1]}; \quad \mathbf{h}^{[q]} = \begin{bmatrix} h_1 \\ h_2 \\ \dots \\ h_{K[q]} \end{bmatrix}^{[q]}$$

Să notăm că, în faza *feedforward*, în calcule se utilizează matricea sinaptică \mathbf{W} , iar în faza *backward* (retro - propagarea erorii), se utilizează transpusa acesteia \mathbf{W}^T .

Blocurile (4) calculează vectorii de eroare $\delta^{[q]}$. Vectorul de eroare asociat stratului de ieșire ($\delta^{[s]}$) se obține:

$$\delta^{[s]} = (\mathbf{d}_r - \mathbf{y}_r) \otimes \sigma'(\mathbf{p}^{[s]}) \quad (1.37)$$

unde:

\mathbf{d}_r - vectorul dorit la ieșire (pentru intrarea \mathbf{x}_r)
 \mathbf{y}_r - vectorul real obținut la ieșire (în urma aplicării intrării \mathbf{x}_r)
 $\sigma'(\mathbf{p}^{[s]})$ - vectorul derivatelor funcției de transfer în raport cu potențialele neuronilor din stratul de ieșire.

$$\sigma'(\mathbf{p}^{[s]}) = \begin{bmatrix} \sigma'(p_1) \\ \sigma'(p_2) \\ \dots \\ \sigma'(p_{K[s]}) \end{bmatrix}^{[s]} \quad \text{unde: } \sigma'(p_i^{[s]}) = \frac{\partial y_i^{[s]}}{\partial p_i^{[s]}}; i = 1, 2, \dots, K[s] \quad \mathbf{d}_r = \begin{bmatrix} d_1 \\ d_2 \\ \dots \\ d_{K[s]} \end{bmatrix}^{[s]}; \quad \mathbf{y}_r = \begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_{K[s]} \end{bmatrix}^{[s]}$$

$K[s]$ - numărul de neuroni din stratul de ieșire $[s]$

În continuare, vectorii de eroare se obțin prin retro - propagarea vectorilor prin blocurile (5) și (4). Fiecare strat are asociat o pereche de blocuri (5) - (4). Vectorii de eroare se obțin printr-un pseudoprodus a 2 vectori.

$$\delta^{[q]} = \mathbf{h}^{[q]} \otimes \sigma'(\mathbf{p}^{[q]}) \quad (1.38)$$

unde $\sigma'(\mathbf{p}^{[q]})$ reprezintă vectorul derivatelor funcției de transfer (derivate parțiale) în raport cu potențialele neuronilor din stratul $[q]$.

$$\sigma'(\mathbf{p}^{[q]}) = \begin{pmatrix} \sigma'(p_1) \\ \sigma'(p_2) \\ \vdots \\ \sigma'(p_{K[q]}) \end{pmatrix}^{[q]} \quad \text{unde : } \sigma'(p_i^{[q]}) = \frac{\partial y_i^{[q]}}{\partial p_i^{[q]}}; i = 1, 2, \dots, K[q]$$

Fiecare componentă a vectorului de eroare $\delta^{[q]}$ (și $\delta^{[s]}$) se obține din produsul a 2 componente:

$$\delta_i^{[q]} = h_i^{[q]} \cdot \sigma'(p_i^{[q]}); \quad i = 1, 2, \dots, K[q] \quad (\text{indicele neuronului})$$

Am notat această operație cu \otimes și am denumit-o pseudoprodus deoarece nu este vorba de produsul (scalar) a 2 vectori (în sens matematic); este o operație de înmulțire a 2 vectori în sensul în care aceasta este implementată în limbajele de programare.

Să notăm că operația cea mai complexă și mai costisitoare din punct de vedere al timpului de execuție, atât pentru faza *feedforward* cât și pentru faza *backward*, este înmulțirea unei matrici cu un vector. Astfel de operații se implementează foarte eficient pe rețele sistolice (RS); acesta este de fapt obiectivul acestei lucrări.

1.6.5. Aplicarea algoritmului EBP

Aplicarea algoritmului EBP în aplicațiile practice reale reclamă stabilirea unui context de lucru. Acest context introduce parametrii suplimentari care trebuie stăpâniți și ajustați (pe parcursul învățării) astfel încât să se obțină o funcționare corectă și rapid convergentă a algoritmului. Redăm în continuare câțiva dintre cei mai importanți parametrii de context, rezultați din experimentarea algoritmului pe aplicații practice concrete.

A. Termenul "momentum"

Versiunea teoretică a algoritmului EBP reclamă în mod normal o variație infinezimală a ponderilor pe parcursul învățării. Această variație este controlată de parametrul α numit constantă (parametru) de învățare. Cu cât α este mai mare cu atât variația ponderilor va fi mai mare. În acest caz, procesul de învățare este rapid dar poate conduce la un comportament oscilant. O variație bruscă a ponderilor poate avea o acțiune semnificativă asupra termenului de eroare. Corecția, care ar trebui să fie semnificativă, va avea tendința să genereze o nouă variație bruscă. Va apărea în consecință un fenomen de oscilație care va împiedica stabilizarea procesului de învățare. Este posibilă limitarea acestui comportament oscilant (fără a micșora α) prin introducerea în corecția din pasul curent a unui procent din corecția pasului anterior. Formula de ajustare a ponderilor sinaptice devine:

$$w_{ij}^{[q]}(t+1) = w_{ij}^{[q]}(t) + \alpha \cdot \delta_i^{[q]} \cdot x_j^{[q]}(t) + \beta \cdot (\Delta w_{ij}^{[q]}(t)) \quad (1.39)$$

Valoarea parametrului β trebuie determinată experimental. Se ia în general între 0,6 și 0,9. Acest nou termen impune memorarea corecției precedente (pentru fiecare pondere sinaptică) ceea ce înseamnă consum suplimentar de resurse de memorie.

B. Utilizarea erorii medii sau instantanee

Ca și în cazul *Adaline*, versiunea exactă a algoritmului EBP revendică cunoașterea mediei semnalelor de intrare și a mediei semnalelor dorite la ieșire. Este totuși posibilă învățarea pornind de la valori instantanee, dacă se menține o constantă de învățare (α) mică. Observăm că apare un antagonism între viteza algoritmului de învățare și procedura de învățare instantanee (care este mult mai ușor de implementat în practică).

Se poate totuși implementa un algoritm care realizează o formă de mediere pe termen scurt, acumulând variațiile gradientului pentru mai mulți vectori de intrare aplicați succesiv (fără a ajusta ponderile), apoi făcând o ajustare a ponderilor care ține cont de toate variațiile de gradient cumulate. Este vorba de metoda ajustării cumulative spre deosebire de metoda ajustării incrementale (care realizează ajustarea ponderilor după fiecare prototip aplicat la intrare). Ambele tipuri de ajustări pot conduce la rezultate bune.

Trebuie subliniat însă că învățarea decurge bine în condiții aleatoare. Pare, deci, mai convenabilă ajustarea incrementală (pentru că este mai ușor de implementat), dar cu preluarea prototipurilor din setul de antrenament și aplicarea acestora pe intrarea RN într-o secvență aleatoare. Acest caracter aleator al secvenței de prototipuri introduce acel "zgomot" necesar rezolvării problemei de mediere.

C. Inițializarea ponderilor

Ponderile nu trebuie inițializate cu valori egale. În acest caz, pentru că termenul de eroare este perceput de către neuronii unui strat după ce acesta a traversat ponderile sinaptice, stratul intern care precede stratul de ieșire va percepe erori identice pentru toți neuronii săi (fig. 1.29).

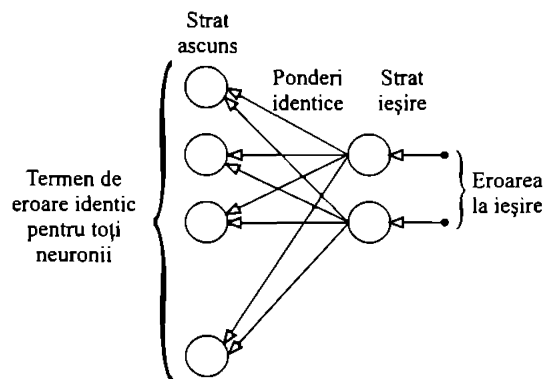


Fig. 1.29. Problema ponderilor inițiale egale

Ponderile stratului ascuns din fig. 1.29 nu vor varia independent creând un fel de stare stabilă parazită care împiedică un proces de învățare corect și care reduce artificial numărul de straturi din componența RN. Soluția la această problemă rezidă în inițializarea ponderilor sinaptice cu valori aleatoare mici.

Faza de inițializare a ponderilor este importantă și din alt punct de vedere; vom ilustra acest lucru printr-un exemplu inspirat din [McC 86].

Considerăm o RN cu 2 straturi, în care fiecare strat conține un singur neuron (fig. 1.30).

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

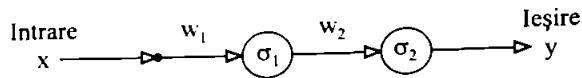


Fig. 1.30. Exemplu relativ la inițializarea ponderilor sinaptice

Cele 2 funcții de transfer σ_1 și σ_2 sunt identice (funcția sigmoidală unipolară exprimată prin relația (1.3)). Valoarea de prag a ambilor neuroni este 0. Dorim ca această RN să învețe funcția identică pe două perechi de antrenament:

x	d
0	0
+1	+1

Dacă trasăm suprafața de eroare $(y - d)^2$ în funcție de ponderile w_1 și w_2 obținem suprafața din fig. 1.31.

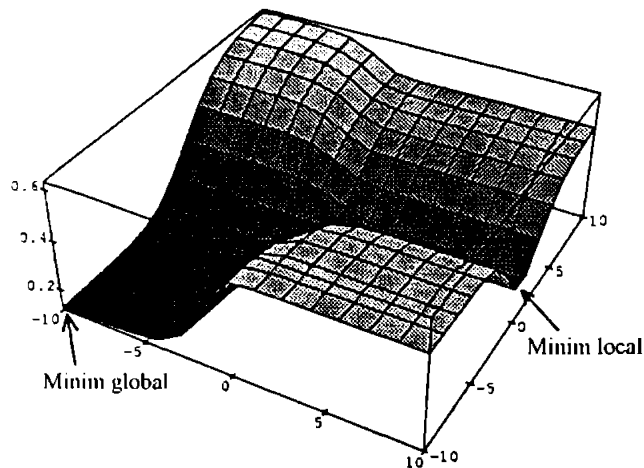


Fig. 1.31. Suprafață de eroare cu 2 puncte de minim

Suprafața de eroare din fig. 1.31 indică faptul că sistemul admite două puncte de minim: un minim global cu o eroare minimă de aproximativ 0,025 și un minim local cu o eroare asociată de 0,45. Atingerea unuia sau a altuia dintre aceste 2 minime depinde esențial de inițializarea ponderilor sinaptice. În acest caz simplu, vizualizarea suprafeței de eroare a fost posibilă, dar nu va fi posibilă și în cazul general. Așadar, chiar și inițializarea cu valori aleatoare poate conduce RN (în timpul învățării) spre minime locale. Pentru a evita această capcană învățarea trebuie restartată cu alte ponderi inițiale (aleatoare).

D. Panta funcției de activare

Funcțiile de activare (transfer) sigmoidale (relațiile (1.3) și (1.4) și figura 1.3) se caracterizează prin factorul de pantă $\lambda = 1/T$. Derivata funcției de activare atinge valoarea maximă ($\lambda/2$) în punctul $p = 0$. Algoritmul EBP ajustează ponderile sinaptice proporțional cu factorul de scară $\sigma'(p)$. Drept consecință, ponderile aferente neuronilor nedeciși vor fi puternic ajustate în timp

ce ponderile aferente neuronilor supraexcitați sau suprainhibați vor fi ajustate mult mai puțin în procesul de învățare. Deoarece semnalele locale de eroare $\delta_w^{(q)}$ se obțin utilizând $\sigma'(p)$ ca factor multiplicator, componentele de eroare transmise înapoi spre straturile anterioare, vor fi mari numai pentru neuronii aflați în zona de indecizie. Pentru o anumită constantă de învățare α fixată, ajustarea ponderilor va fi proporțională cu factorul de pantă $\lambda = 1/T$. Deci, utilizând funcții de transfer cu pantă mare vom obține rezultate similare cu cele obținute în cazul utilizării constantelor de învățare (α) mari. Sugestia noastră este de a menține panta funcției de activare la valoarea standard $\lambda = 1/T = 1$ și de a controla viteza de convergență a procesului de învățare numai din constanta α .

1.6.6. Aplicație demonstrativă

Vom considera 2 mulțimi de puncte în planul cartezian (fig. 1.32), care definesc două regiuni dintre care una este convexă.

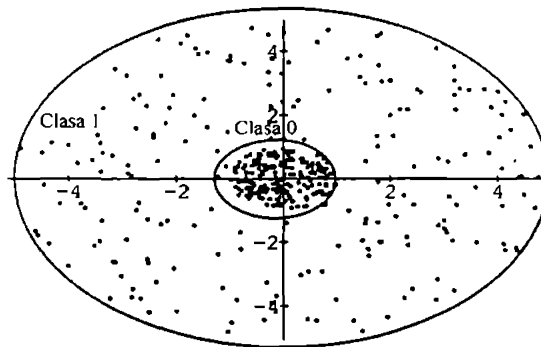


Fig. 1.32. Repartiția punctelor în planul cartezian

Dacă dorim să clasificăm aceste puncte în două clase (asociate celor 2 regiuni), în funcție de coordonatele carteziene ale acestor puncte, putem utiliza un perceptron cu 2 straturi (fig. 1.33). Vectorii de intrare x vor conține coordonatele carteziene aferente punctelor de clasificat iar ieșirea y va specifica clasa căreia îi aparține vectorul de intrare x . Aplicând algoritmul EBP, cu următorii parametri:

- $\alpha = 0,7$
- fără termen *momentum*
- ajustare incrementală a ponderilor (învățare instantanee)
- valoarea de prag a fiecărui neuron este realizată printr-o intrare fixată la +1 și un coeficient sinaptic variabil
- variabile reale
- obiectivul vizat este clasificarea punctelor în cele 2 clase cărora li s-au asociat valorile $y = +1$ și respectiv $y = 0$
- funcția de transfer utilizată:

$$\sigma(x) = \text{th}(x)$$

$$\sigma'(x) = 1 - \text{th}^2(x)$$

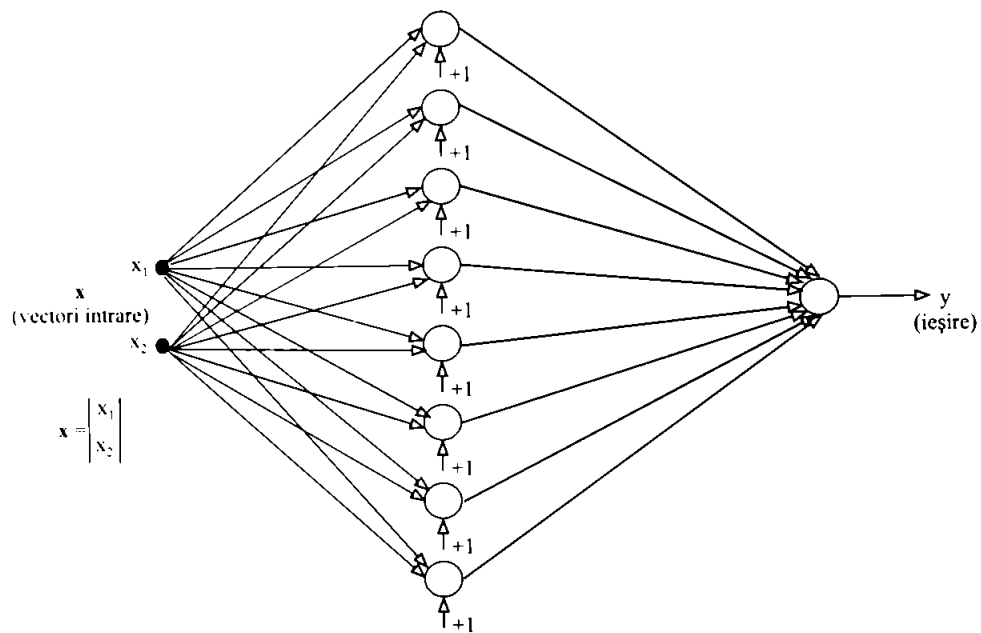
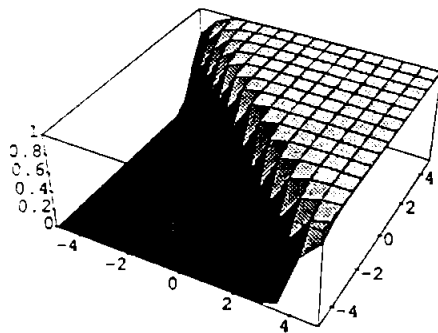
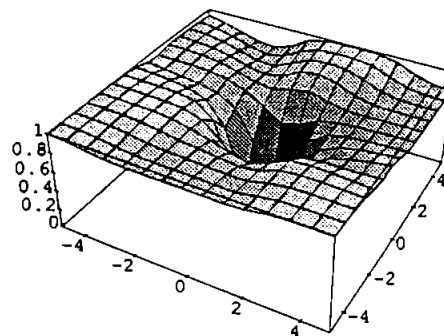


Fig. 1.33. Structura rețelei

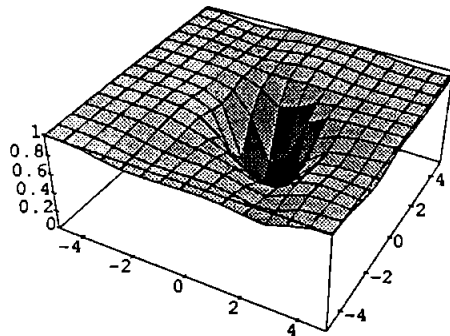
S-a obținut separarea în clase reprezentată în fig. 1.34. Setul de antrenament (prototipuri) a fost format din 500 perechi (x, d) . În figura 1.34 se redau suprafețele de separare $(y = f(x))$ obținute. După cum se poate constata rezultatul nu este perfect dar constituie un separator convenabil. Trebuie să notăm că numărul de iterații poate fi mărit și că parametrii utilizați n-au fost aleși în scopul obținerii celei mai bune soluții posibile (parametrii optimali nu sunt cunoscuți).



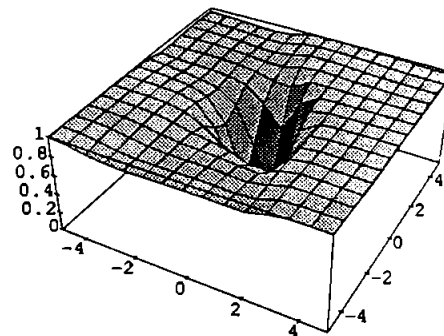
a) suprafața de separare inițială



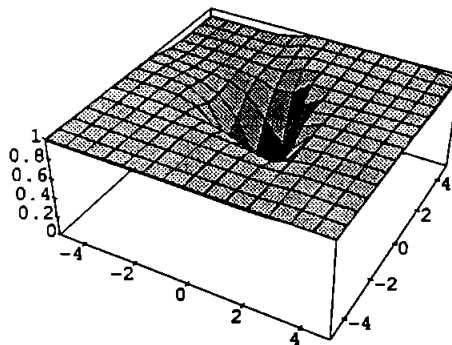
b) după 200 prezentări



c) după 400 prezentări



d) după 1000 prezentări



e) după 2000 prezentări

Fig. 1.34. Evoluția suprafeței de separare în timpul învățării

1.6.7. Concluzii

După cum am putut constata, aplicarea algoritmului EBP este relativ complicată și calculele efectuate sunt destul de departe de ceea ce se întâmplă într-un sistem nervos. În plus, algoritmul suferă de anumite limitări:

- incapacitatea de a determina numărul optim de neuroni din fiecare strat în funcție de aplicație;
- viteza de convergență lentă care depinde de mărimea setului de antrenament și de parametrii rețelei;
- riscul deraierii spre minime locale.

Trebuie să-i recunoaștem totuși o certă aplicabilitate industrială, care l-a transformat într-un motor al cercetării în domeniul conexiunilor. Dar RN nu trebuie limitate la arhitectura *feedforward*; cercetările au evoluat și trebuie să evolueze și mai departe pe arhitecturi mai complexe (rețele recurente, rețele probabiliste etc.).

1.7. Modelul Hopfield

1.7.1. Introducere

Modelul *Hopfield* este un model recurent sau rebusat. Într-o rețea rebusată fiecare neuron este în relație directă cu toți ceilalți. Nu se disting deci unități (neuroni) de intrare și unități de ieșire. Fiecare neuron este în același timp și intrare și ieșire, și mărimea care-l caracterizează este starea sa (fig. 1.35).

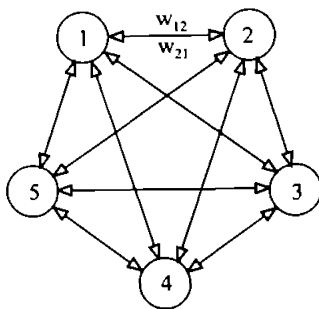


Fig. 1.35. Structura unei rețele *Hopfield*

O rețea recurentă se caracterizează prin dinamica sa, ceea ce înseamnă evoluția stărilor unităților sale în timp. Această evoluție nu este controlată și nivelul de activitate al neuronilor va evolua evident datorită interacțiunilor acestora cu ceilalți neuroni din componența rețelei. În practică, pentru a crea sisteme exploatabile, se caută ca aceste rețele să aibă o dinamică care să conducă spre un punct de stabilitate, numit atractor. Această dinamică va interesa atunci când aplicația revendică verdicte (răspunsuri) unice. Nu trebuie exclusă utilitatea unor atractori circulari. În acest caz starea de echilibru a rețelei nu este punctuală ci este ciclică (un ciclu limitat de stări); aceste cicluri limitate pot fi utilizate pentru memorarea unor secvențe.

Într-o rețea *Hopfield* există o dinamică naturală. Obiectivul demersului nostru este de a exploata această dinamică naturală pentru a determina rețeaua să funcționeze ca o memorie asociativă. Problema trebuie abordată sub următoarele 4 aspecte:

1. În anumite condiții se poate asocia o energie globală întregului sistem dinamic (rețelei). Reducerea acestei energii reprezintă o condiție de stabilitate pentru sistemul dinamic. Cum va putea acționa o variație locală (la nivelul unui neuron) pentru reducerea energiei globale?
2. În general, care sunt condițiile care conduc la un sistem a cărui energie scade? Răspunsul la această întrebare ne va permite să controlăm rețeaua.
3. Cum se poate construi un sistem ale cărui stări stabile reprezintă cu precizie informația pe care dorim s-o memorăm? Dacă vom putea găsi o regulă care să răspundă la această întrebare, atunci vom fi capabili să construim memorii asociative a căror dinamică naturală de sistem va conduce spre informațiile pe care le-am stocat.
4. Există reguli care răspund la întrebarea precedentă: regula lui *Hebb*, regula proiecției. Vom face analiza acestora.

1.7.2. Dinamica rețelelor Hopfield

Dinamica rețelelor *Hopfield* compuse din N neuroni este descrisă de ecuația clasică a neuronului *McCulloch-Pitts*. În formă discretă ecuația se scrie:

$$p_i(t) = \sum_{j=1}^N w_{ij} \cdot x_j(t-1) \quad (1.40)$$

$$\text{și } x_i(t) = \text{sign}(p_i(t) - 2\theta_i) \quad (1.41)$$

unde:

- p_i - potențialul neuronului i
- w_{ij} - ponderea conexiunii dinspre neuronul j spre neuronul i
- x_i - ieșirea (starea) neuronului i
- θ_i - valoarea de prag asociată neuronului i
- sign - funcția sign reprezentată în fig. 1.2. b).

În cazul rețelelor recurente se disting 2 metode de actualizare a stării neuronilor: secvențială aleatoare și respectiv paralelă. În cazul metodei secvențiale, la momentul t , se alege aleator un neuron i pe baza unei reguli de distribuție uniformă. Starea neuronului selectat se actualizează utilizând regula *McCulloch-Pitts* și se continuă selecția aleatoare. În cazul actualizării paralele, la momentul t , se calculează potențialele tuturor neuronilor și starea acestora se actualizează simultan.

Hopfield a arătat [HOP 82] că dinamica secvențială poate conduce la o stare stabilă în anumite condiții:

- matricea ponderilor sinaptice \mathbf{W} să fie simetrică ($w_{ij} = w_{ji}$);
- elementele diagonalei matricii \mathbf{W} să fie nenegative.

Demonstrația se bazează pe faptul că putem asocia o funcție de energie unui sistem dinamic astfel încât toate schimbările de stare aferente elementelor acestuia să nu poată conduce decât la scăderea funcției de energie. Suplimentar, această funcție de energie admite un minim absolut și descreșterea energiei se face în timp finit.

Prin definiție, funcția de energie se exprimă la momentul t prin relația:

$$E(t) = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N w_{ij} \cdot x_i(t) \cdot x_j(t) \quad (1.42)$$

unde N reprezintă numărul de componente din vectorul \mathbf{x} , cu alte cuvinte numărul de neuroni din rețea. Vom pune în evidență efectul pe care-l va avea asupra energiei sistemului, schimbarea stării unui singur neuron. Considerând timpul discret, vom avea:

$$E(t+1) - E(t) = -\frac{1}{2} \left[\sum_{i=1}^N \sum_{j=1}^N (w_{ij} \cdot x_j(t+1) - 2\theta_i) \cdot x_i(t+1) - \sum_{i=1}^N \sum_{j=1}^N (w_{ij} \cdot x_j(t) - 2\theta_i) \cdot x_i(t) \right] \quad (1.43)$$

Vom considera că actualizarea stării sistemului este de tip secvențială aleatoare. La fiecare moment t , se selectează aleator un neuron și i se calculează noua stare. Această metodă de selecție este denumită metoda *Monte Carlo*.

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

Pot să apară 2 situații:

- neuronul selectat nu-și schimbă starea și, în acest caz, dacă i este indicele neuronului extras, vom avea $x_i(t+1) = x_i(t)$;
- neuronul selectat își schimbă starea: $x_i(t+1) = -x_i(t)$.

Care va fi variația de energie dacă un singur neuron k își schimbă starea?

Trebuie dezvoltată expresia ΔE pentru a face să apară explicit sumele pentru $i \neq k, j \neq k, i = k, j = k$; apoi expresia trebuie prelucrată pentru a o face dependentă numai de variabila t . După câteva calcule simple [ZUR 92] se va obține:

$$\Delta E = \sum_{j=1}^N (w_{kj} \cdot x_j(t) - 2\theta_k) x_k(t) + \sum_{i=1}^N w_{ik} \cdot x_k(t) \cdot x_i(t) - 2w_{kk} \cdot x_k^2(t) \quad (1.44)$$

Prin ipoteză, matricea \mathbf{W} este simetrică, ceea ce conduce la:

$$\Delta E = 2x_k(t) \cdot \left(\sum_{j=1}^N w_{kj} \cdot x_j(t) - \theta_k \right) - 2w_{kk} \cdot x_k^2(t) \quad (1.45)$$

Dacă alegem valoarea de prag $\theta_i = 0$, și substituim valorile lui x_i , se obține:

$$\Delta E = -2 - 2w_{kk} \quad (1.46)$$

Dacă se respectă condiția $w_{kk} \geq -1$, variația de energie provocată de schimbarea stării unui singur neuron este întotdeauna negativă. Funcția de energie va fi monotonă descrescătoare și valoarea ei evoluează până la atingerea unui minim (cel puțin un minim local) după care va rămâne constantă.

Rămâne de arătat că funcția de energie are o valoare minimă finită. Din moment ce numărul de neuroni din rețea este finit și valorile de pe ieșirile acestora nu pot fi decât $+1$ sau -1 , este evident că minimul funcției de energie este finit. Funcția de energie este o sumă dublă în care intră stările tuturor neuronilor și, în consecință, va fi finită. Rezultă că, după un număr de iterații, va fi atins un punct de minim energetic care reprezintă un atractor al sistemului. În figura 1.36 se prezintă imaginea procesului de descreștere a funcției de energie în spațiul stărilor.

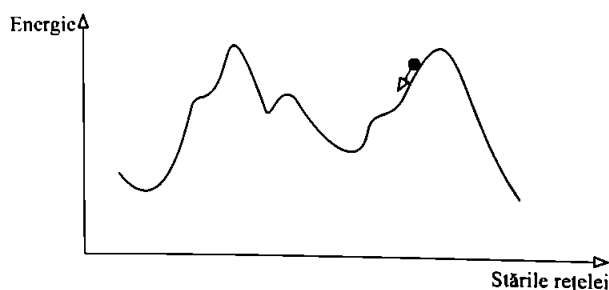


Fig. 1.36. Reprezentarea funcției de energie

Este important de subliniat că aceeași rețea (sistem) admite mai multe minime care pot avea energii asociate diferite și care sunt accesibile pornind de la configurații inițiale diferite. De asemenea, minimul absolut poate fi atins plecând din mai multe configurații inițiale.

Va trebui analizat cum dinamica unei rețele recurente va permite acesteia, în anumite condiții, să evolueze spre acești atractori. Dacă dorim să exploatăm această proprietate, pentru a construi o memorie asociativă de exemplu, trebuie să fim capabili să controlăm poziția acestor atractori. Această problemă se numește "problema dinamicii inverse" și se formulează astfel: "Este posibilă construcția unei rețele recurente care să admită ca atractori un ansamblu de stări alese arbitrar?"

1.7.3. Problema dinamicii inverse

Să considerăm că dispunem de un ansamblu de configurații (*patterns*) care dorim să constituie baza de referință (setul de stări stabile) pentru rețeaua recurentă. Această operație de forțaj este denumită fază de învățare.

Atunci când aplicăm o configurație (stare) deteriorată vom dori ca rețeaua să răspundă prin una dintre configurațiile (stările) de referință. Acest principiu de restaurare a unei configurații de referință pornind de la o configurație deteriorată se numește memorie asociativă. Principiul este descris în figura 1.37.

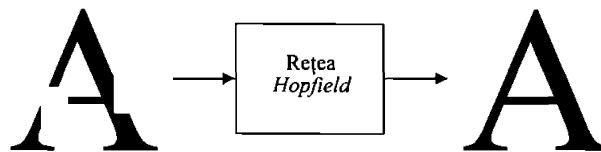


Fig. 1.37. Recunoașterea caracterelor cu ajutorul rețelei Hopfield

În acest exemplu se presupune că imaginea caracterului A este un prototip. La intrarea rețelei se aplică un A deteriorat. După câteva iterații, rețeaua trebuie să răspundă cu prototipul A.

Prima rezolvare (parțială) la care putem recurge este regula lui Hebb, care a fost descrisă anterior. Matricea sinaptică W controlează poziția atractorilor. Calculând convenabil valorile coeficienților sinaptici w_{ij} , atractorii vor deveni un codaj de prototipuri și rețeaua, prin dinamica ei, va juca rolul de decodor.

1.7.4. Regula lui Hebb

Regula lui Hebb pune la dispoziție o metodă foarte simplă pentru calculul ponderilor sinaptice. Formularea acesteia se rezumă la:

$$w_{ij} = \frac{1}{N} \left(\sum_{s=1}^P x_i(s) \cdot x_j(s) \right) \quad (1.47)$$

w_{ij} - coeficientul sinaptic ij din matricea W .

$x_i(s)$ - componenta i a stării prototip s

P - numărul de stări prototip

N - numărul de neuroni din rețea

Pentru fiecare prezentare a unui prototip $x(s)$, se ajustează matricea sinaptică cu valori rezultate din activitatea conjugată a componentelor x_i, x_j .

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

1.7.4.1. Interpretare

Din punct de vedere geometric, dacă stările prototip sunt ortogonale, matricea W poate fi interpretată ca matrice de proiecție ortogonală în subspațiul generat de vectorii prototip.

1.7.4.2. Demonstrație

Presupunem că stările prototip sunt reciproc ortogonale. Vom forma matricea A de dimensiune $N \times P$ și ale cărei coloane vor conține cele P stări prototip.

$$A = (\mathbf{x}(1), \mathbf{x}(2), \dots, \mathbf{x}(P))$$

Regula lui *Hebb* se va scrie:

$$W = \frac{1}{N} \cdot A \cdot A^T \quad (1.48)$$

Prin definiție, proiecția ortogonală a tuturor vectorilor \mathbf{x} pe subspațiul definit de coloanele matricii A va fi:

$$\text{Proj}_1 \mathbf{x} = A \cdot (A^T \cdot A)^{-1} \cdot A^T \cdot \mathbf{x} \quad (1.49)$$

unde Proj_1 reprezintă operatorul de proiecție ortogonală.
Matricea $A^T \cdot A$ va avea forma:

$$A^T \cdot A = \begin{vmatrix} \mathbf{x}(1)^2 & \mathbf{x}(2) \cdot \mathbf{x}(1) & \dots & \dots \\ \mathbf{x}(1) \cdot \mathbf{x}(2) & \mathbf{x}(2)^2 & \dots & \dots \\ \dots & \dots & \mathbf{x}(3)^2 & \dots \\ \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \mathbf{x}(P)^2 \end{vmatrix}$$

unde $\mathbf{x}(i)$ reprezintă starea prototip numărul i .

Am presupus prin ipoteză că vectorii stărilor prototip sunt ortogonali. Rezultă că produsele scalare $\mathbf{x}(i) \cdot \mathbf{x}(j)$ vor fi nule pentru $i \neq j$. Matricea $A^T \cdot A$ este deci o matrice diagonală, având pe diagonală elementele $\mathbf{x}(i)^2$. Dar componentele oricărui vector sunt egale cu $+1$ sau -1 ; rezultă că toate elementele de pe diagonală matricii $A^T \cdot A$ vor avea valoarea N . Deci:

$$A^T \cdot A = N \cdot \mathbf{I} \quad \text{iar} \quad (A^T \cdot A)^{-1} = \frac{1}{N} \cdot \mathbf{I} \quad (1.50)$$

unde \mathbf{I} reprezintă matricea unitate de dimensiune P .
Dacă prototipurile sunt ortogonale 2 câte 2, relația (1.49) devine:

$$\text{Proj}_1 \mathbf{x} = \frac{1}{N} \cdot A \cdot A^T \cdot \mathbf{x} \quad (1.51)$$

Am regăsit astfel formula lui *Hebb*. Din moment ce această operație este o proiecție ortogonală, va fi adevărată și ecuația:

$$\mathbf{W} \cdot \mathbf{x}(s) = \mathbf{x}(s) \text{ pentru orice prototip } \mathbf{x}(s); s = 1, 2, \dots, P \quad (1.52)$$

Dacă prototipurile sunt neortogonale relația (1.52) nu se verifică.

Rămân de valorificat, în final, și proprietățile funcției de transfer (neliniare) aplicată de către fiecare neuron propriului potențial. Reamintim că toate componentele vectorilor prototip au valorile ± 1 . Înseamnă că este suficient să fixăm o valoare de prag între 0 și +1, și stările prototip vor fi automat și stări stabile ale rețelei.

Operația de proiecție efectuată prin produsul dintre un vector de intrare (oarecare) și matricea ponderilor sinaptice nu conservă norma. Componentele vectorilor de intrare pot fi oarecare și pot să nu corespundă vectorilor prototip. Dar fiecare componentă care depășește valoarea de prag va fi restabilită la +1, iar cele care nu depășesc pragul vor fi restabilite la -1. Astfel, configurațiile prototip pot fi restaurate dacă se pornește de la o configurație vecină.

1.7.5. Proprietățile rețelei Hopfield

Rețeaua *Hopfield* este capabilă să realizeze funcția de memorie asociativă. Dacă inițializăm rețeaua într-o stare oarecare și această stare nu este prea îndepărtată de un prototip codificat (distanță *Hamming*), rețeaua va evolua spre prototip. În caz contrar, rețeaua va evolua spre o stare nedefinită numită stare parazită.

Atractivitatea se definește astfel: o configurație \mathbf{x} este un k -atractor ($1 \leq k \leq N$) dacă, plecând dintr-o stare alterată care prezintă k erori în raport cu configurația prototip \mathbf{x} , rețeaua va evolua într-un mod iterativ sincron și determinist stabilizându-se în starea \mathbf{x} .

Utilizând regula lui *Hebb*, dacă toate prototipurile sunt reciproc ortogonale, fiecare prototip va reprezenta un atractor cu un bazin de atracție până la distanța *Hamming*:

$$k = \left\lfloor \frac{N}{2P} \right\rfloor \quad (1.53)$$

Deci, se observă că, cu cât stocăm mai multe prototipuri cu atât atracția fiecărui prototip va scădea. Pentru a implementa o funcție de memorie asociativă performantă, va fi necesar un număr mare de neuroni corelat cu stocarea a cât mai puține prototipuri.

Această atractivitate reprezintă un fel de măsură globală a performanțelor rețelei dar principala constrângere provine din necesitatea de a avea vectori prototip ortogonali. Este evident posibilă aplicarea unor proceduri de ortogonalizare de tip *Gramm - Schmidt* (o preprocesare a vectorilor prototip pentru ortogonalizare), dar este de asemenea posibilă utilizarea altor metode de rezolvare mai puțin restrictive.

O altă limitare severă, în cazul utilizării regulii lui *Hebb*, se referă la capacitatea de stocare a configurațiilor prototip. Această capacitate este de ordinul 14% din numărul de neuroni dacă nu acceptăm o eroare mai mare de 3%, estimată prin distanța *Hamming* dintre vectorul prototip și vectorul aplicat. Pe măsură ce se depășește această capacitate, performanțele scad dramatic. Rețeaua începe să uite chiar și ce a învățat.

În cele ce urmează, vom prezenta câteva metode algebrice descrise în [COT 88] care permit creșterea capacității de stocare și, în același timp, eliminarea unor constrângeri foarte limitative. Primul principiu este cel enunțat în [KOH 86] (memoria optimală), regula reluată mai târziu sub denumirea de regula proiecției.

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

1.7.6. Regula proiecției

Vom studia posibilitatea construirii unei matrici sinaptice care garantează stabilitatea configurațiilor prototip și vom examina capacitatea de asociere aferentă rețelei obținute. Considerăm de această dată că vectorii prototip nu mai sunt ortogonali între ei, dar sunt încă liniar independenți. Formal, sistemul de rezolvat se reduce la:

$$\text{sign}(\mathbf{W} \cdot \mathbf{A}) = \text{sign}(\mathbf{A})$$

Într-o primă abordare, vom considera neuronii fără intrare de prag și fără iterații repetate. La nivel fundamental problema rămâne identică [COT 88] și poate fi rezolvată cu ajutorul sistemului de ecuații:

$$\mathbf{A} = \mathbf{W} \cdot \mathbf{A} \quad (1.54)$$

unde $\mathbf{A} = (\mathbf{x}(1), \mathbf{x}(2), \dots, \mathbf{x}(P))$. \mathbf{A} este o matrice ($N \times P$), unde N reprezintă numărul de neuroni și P numărul de prototipuri. Este cazul tipic de memorie asociativă.

Matricea \mathbf{A} nu este pătrată și nu admite inversă. Totuși, există o soluție non - unică, dacă rangul matricii \mathbf{A} (numărul de vectori prototip) este inferior lui N (numărul de componente dintr-un vector prototip = numărul de neuroni). Este vorba de generalizarea inversei pentru matrici dreptunghiulare.

$$\mathbf{W} = \mathbf{A} \cdot \mathbf{A}^+ \quad (1.55)$$

unde:

\mathbf{A}^+ este o matrice ($P \times N$) numită pseudo - inversa matricii \mathbf{A} .

\mathbf{A}^+ va fi:

$$\mathbf{A}^+ = (\mathbf{A}^T \cdot \mathbf{A})^{-1} \cdot \mathbf{A}^T \quad (1.56)$$

Condiția de liniar independență impusă coloanelor matricii \mathbf{A} , ne asigură că $\mathbf{A}^T \cdot \mathbf{A}$ are rangul P și este deci inversabilă, ceea ce asigură existența matricii \mathbf{A}^+ . Această metodă a fost tratată exhaustiv de [KOH 86], [PER 86], [PER 87], [GUY 88]. În cazul prototipurilor liniar independente, pentru problema auto - asocierii soluția va fi:

$$\mathbf{W} = \mathbf{A} \cdot (\mathbf{A}^T \cdot \mathbf{A})^{-1} \cdot \mathbf{A}^T \quad (1.57)$$

Și în acest caz, \mathbf{W} va fi o matrice de proiecție pe spațiul vectorial generat de vectorii prototip. Este evident că, în cazul în care configurațiile prototip sunt ortogonale, matricea pseudo - inversă devine:

$$\mathbf{A}^+ = \frac{1}{N} \cdot (\mathbf{A}^T) \quad (1.58)$$

În acest caz particular, matricea sinaptică devine:

$$\mathbf{W} = \frac{1}{N} \cdot (\mathbf{A} \cdot \mathbf{A}^T) \quad (1.59)$$

Am regăsit astfel regula lui *Hebb*.

Există mai multe metode de calcul a pseudo - inversei unei matrici. Prima se referă evident la calculul algebric cunoscut, dar acest principiu este relativ greoi deoarece el obligă la recalcularea completă a pseudo - inversei după aplicarea fiecărui prototip. Suplimentar, calculul pseudo - inversei necesită aplicarea unor algoritmi cum ar fi cel al lui *Choleski*. O altă metodă de calcul a pseudo - inversei constă în utilizarea algoritmului lui *Greville* [KOH 86] care simplifică mult procesul. Redăm acest algoritm în paragraful următor:

1.7.7. Algoritm de proiecție iterativ

1. Inițializare:

$$\begin{aligned} \mathbf{x}(0) &= 0 \\ \mathbf{W}(0) &= 0 \end{aligned}$$

2. Aplicarea unui prototip $\mathbf{x}(t)$ și calculul lui q :

$$q = \frac{1}{\|(\mathbf{x}(t) - \mathbf{W}(t-1)) \cdot \mathbf{x}(t)\|^2} \quad (1.60)$$

Dacă q este foarte mare se rejectează prototipul $\mathbf{x}(t)$ și se revine la pasul 2.
Dacă nu, se trece la pasul 3

3. Actualizarea matricii \mathbf{W} .

$$\mathbf{W}(t+1) = \mathbf{W}(t) + q \cdot [\mathbf{I} - \mathbf{W}(t)] \cdot \mathbf{x}(t) \cdot [\mathbf{I} - \mathbf{W}(t)] \cdot \mathbf{x}(t)^T \quad (1.61)$$

4. Revenire la pasul 2 atâta timp cât mai există prototipuri de învățat.

Este interesant de observat că, dacă q este foarte mare (tinde la infinit), avem de-a face cu un prototip $\mathbf{x}(t)$ linear dependent în raport cu coloanele (prototipurile) deja introduse. Acest nou prototip nu aduce informații suplimentare și din acest motiv va fi rejectat.

Pentru matricea de proiecție obținută, fiecare prototip stocat va fi un atractor până la o distanță *Hamming* [COT 88]:

$$k = \frac{1}{2 \cdot \max_i (w_{ii})}; \quad 1 \leq i \leq N \quad (1.62)$$

Inconvenientul acestei metode rezidă din faptul că, pentru fiecare pondere sinaptică w_{ij} , calculul face apel la un factor global reprezentat de termenul q . Există totuși o metodă locală [DIE 87] care, pentru un neuron dat, nu face apel decât la potențialul său și al neuronului la care acesta este conectat. Este metoda gradientului stohastic care se apropie de cea propusă de [WID 60] pentru modelul *Adaline*.

1.7.8. Metoda gradientului (Proiecția Delta)

Este posibilă utilizarea unei metode de gradient stohastic care minimizează o funcție de eroare (vezi modelul *Adaline*). În acest caz învățarea se va realiza printr-o prezentare repetată a tuturor vectorilor prototip. Această metodă este de asemenea costisitoare (în ceea ce privește

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

cantitatea de calcul) dar limitează aria calculului doar la interacțiuni strict locale și se înscrie ca o candidată preferabilă pentru implementările VLSI.

Metoda gradientului vizează minimizarea unui termen de eroare. Pentru o rețea recurentă, cum este cazul rețelei *Hopfield*, se propune utilizarea unui termen de eroare de forma [DIE 87]:

$$\varepsilon(t) = \sum_{i=1}^N \left(\sum_{j=1}^N w_{ij} \cdot x_j(t) - x_i(t) \right)^2 \quad (1.63)$$

Pentru o memorie auto - asociativă, acest termen de eroare măsoară diferența dintre starea corespunzătoare prototipului dorit și cea obținută prin înmulțirea acestuia cu matricea sinaptică \mathbf{W} . Dacă vom reuși să reducem acest termen de eroare la zero, înseamnă că am memorat prototipul $\mathbf{x}(t)$ prezentat. Va trebui calculat gradientul acestui termen de eroare în funcție de matricea \mathbf{W} .

$$\frac{\partial \varepsilon(t)}{\partial w_{ij}} = \frac{\partial}{\partial w_{ij}} \left(\sum_{i=1}^N \left(\sum_{j=1}^N w_{ij} \cdot x_j(t) - x_i(t) \right)^2 \right) = 2 \sum_{j=1}^N (w_{ij} \cdot x_j(t) - x_i(t)) \cdot x_j(t) \quad (1.64)$$

Dacă vom recurge la o metodă de ajustare identică cu cea propusă de *Widrow* pentru modelul *Adaline*, va fi suficient să ajustăm coeficienții matricii \mathbf{W} după fiecare prezentare a unui vector prototip:

$$w_{ij}(t) = w_{ij}(t-1) - \alpha \frac{\partial \varepsilon(t)}{\partial w_{ij}} \quad (1.65)$$

Parametrul α trebuie astfel ales încât rezultatele să rămână independente de numărul de elemente (puncte) de calcul, atunci când acesta tinde la infinit. Parametrul α va fi un fel de factor de normalizare. Pentru simplificare, vom alege:

$$\alpha = \frac{1}{N} \quad (1.66)$$

Regula de ajustare a ponderilor sinaptice devine:

$$w_{ij}(t) = w_{ij}(t-1) + \frac{1}{N} \left[x_i(t) - \sum_{j=1}^N w_{ij}(t-1) \cdot x_j(t) \right] \cdot x_j(t) \quad (1.67)$$

Această regulă de ajustare este strict locală deoarece nu face apel la cunoașterea matricii \mathbf{W} pentru fiecare element (i, j) . Algoritmul de ajustare a ponderilor va fi:

1. Inițializare:

$$\mathbf{x}(0) = 0 \\ \mathbf{W}(0) = 0$$

2. Aplicarea unui prototip $\mathbf{x}(t)$ și actualizarea matricii \mathbf{W} :

$$w_{ij}(t) = w_{ij}(t-1) + \frac{1}{N} \left[x_i(t) - \sum_{j=1}^N w_{ij}(t-1) \cdot x_j(t) \right] \cdot x_j(t)$$

3. Revenire la pasul 2.

Trebuie să subliniem că, în cadrul acestui algoritm trebuie aplicat de mai multe ori ansamblul de vectori prototip pentru ca eroarea să scadă pe durata acestor prezentări succesive.

1.7.9. Aplicațiile rețelelor recurente

Rețeaua *Hopfield* este larg utilizată pentru implementarea funcției de memorie asociativă. Vom prelua ca exemplu problema recunoașterii caracterelor. Două prototipuri sunt prezentate în figura 1.38. Fiecărui pixel (reprezentat printr-un pătrățel) îi este asociat un neuron. Rețeaua va conține $12 \times 16 = 192$ neuroni.

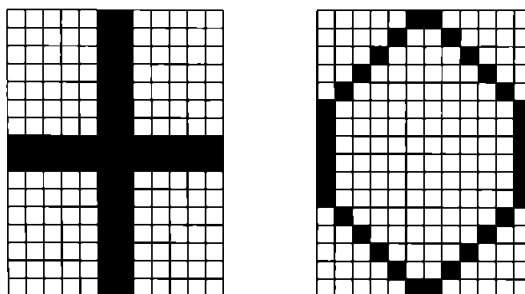


Fig. 1.38. Vectorii prototip

Prin convenție, neuronii activi vor avea ieșirea la valoarea +1 și cei inactivi la valoarea -1. Pentru fiecare vector prototip, se poate reprezenta starea rețelei printr-un grafic 3D, ca în figura 1.39.

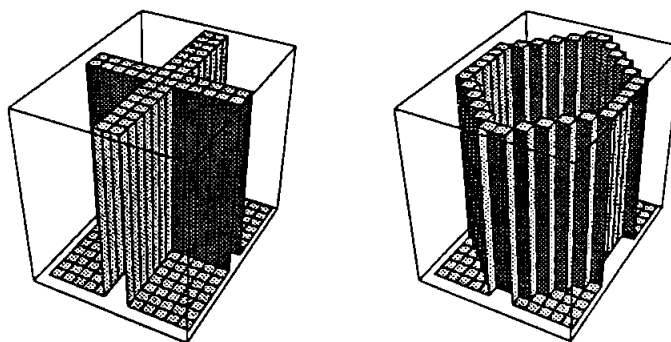


Fig. 1.39. Activitatea rețelei pentru fiecare prototip ce trebuie învățat

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

Aceste 2 prototipuri pot fi învățate utilizând regula lui *Hebb*. Apoi, în timpul fazei normale de lucru (recunoașterea caracterelor), rețeaua va fi exploatată ca o memorie asociativă. Se prezintă rețelei un caracter oarecare (denaturat), ca de exemplu cel reprezentat în figura 1.40. Prezentarea unui caracter constă în forțarea stării neuronilor pe valorile +1 și respectiv -1 impuse de vectorul de intrare (caracterul) respectiv.

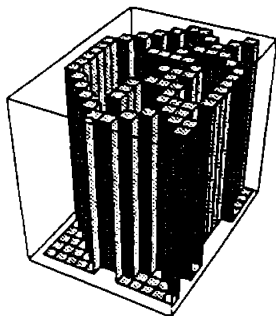


Fig. 1.40. Starea inițială a rețelei pentru caracterul de recunoscut

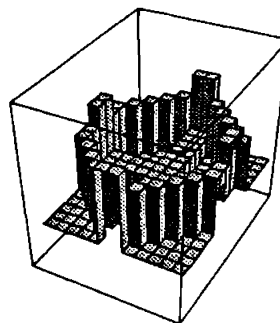


Fig. 1.41. Activitatea rețelei după convergență

Se lasă apoi rețeaua să evolueze liber spre o stare stabilă; pe parcursul evoluției starea neuronilor a fost actualizată utilizând procedura paralelă de actualizare. La convergență se va obține graficul de activitate prezentat în figura 1.41. Să notăm că în acest grafic sunt reprezentate produsul scalar dintre vectorul de stare al neuronilor și matricea sinaptică W obținută în urma fazei de învățare. Activitatea fiecărui neuron va fi reprezentată deci de o valoare reală (analogică).

Se aplică apoi funcția de transfer (de exemplu funcția $\text{sign}(x)$). Astfel, dacă activitatea unui neuron este pozitivă, ieșirea lui va fi forțată la +1, în caz contrar la -1. Nivelul valorii de prag este reprezentat în figura 1.42. În final, după aplicarea funcției $\text{sign}(x)$ asupra potențialelor tuturor neuronilor se va obține noua stare a rețelei, prezentată în figura 1.43.

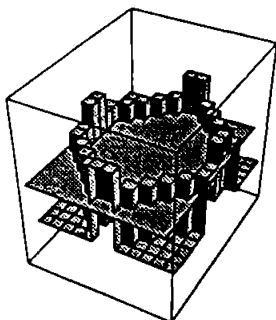


Fig. 1.42. Reprezentarea valorii de prag pe graficul de activitate

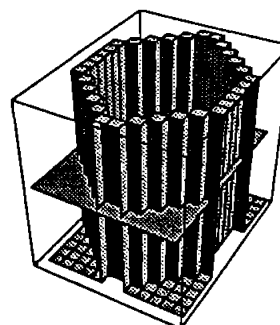


Fig. 1.43. Starea rețelei după aplicarea funcției de transfer

Să notăm că rețeaua a recunoscut prototipul 2: ea a asociat modelului denaturat aplicat la intrare prototipul cel mai asemănător (memorie asociativă). Această structură de rețea complet

interconectată poate fi de asemenea utilizată pentru rezolvarea oricărei probleme de optimizare. *Hopfield* și *Tank* au utilizat-o pentru rezolvarea problemei agentului comercial. Principiul de lucru impune determinarea unei funcții de energie adecvate care va plasa soluțiile problemei de rezolvat în minimele locale aferente acestei funcții.

1.7.10. Concluzii

Modelul *Hopfield* are la bază o rețea recurentă și exploatează proprietățile naturale ale acesteia de convergență spre un punct de stabilitate (stare sau ciclu de stări). Acest model suportă mai multe reguli de învățare cum ar fi: regula lui *Hebb*, regula proiecției și proiecției Delta. Lista nu este exhaustivă și se poate afirma că există atâtea reguli de învățare câte probleme sunt de soluționate. Rețeaua *Hopfield* și dinamica sa are limitări deja cunoscute:

- capacitatea de învățare după regula lui *Hebb* este redusă. După [AMI 85], presupunând că numărul de neuroni N este foarte mare, numărul maxim de prototipuri corect recunoscute este sub $0,138 N$. În aceste condiții, distorsiunea admisă este de ordinul a 3% (estimată prin produsul scalar dintre vectorul prototip normat și vectorul de recunoscut). Dacă numărul de prototipuri depășește limita de $0,138 N$, altfel spus dacă se încearcă "supraincercarea" memoriei în timpul învățării, toate prototipurile învățate anterior se pierd (fenomenul de uitare). Totuși, utilizând regula proiecției, capacitatea de învățare crește și poate atinge circa 40% din numărul de neuroni.
- informațiile nu pot fi stocate sub formă brută în memorie. Trebuie realizată o pre - procesare a datelor în scopul ortogonalizării, dacă se vizează atingerea unei capacități de învățare apropiată de posibilitățile teoretice ale rețelei.
- trebuie evitate erorile, cum ar fi recunoașterea falsă care constă în clasarea greșită a unei configurații de intrare.

1.8. Algoritmul de auto-organizare al lui Kohonen

1.8.1. Observații fiziologice și funcționale

Rețelele biologice reale au o structură organizatorică ordonată care depinde adesea de activitatea pe care sunt specializate. Se constată că, pe cortex, zonele receptoare sunt ordonate într-o manieră identică cu cea a unităților organului senzorial asociat. În plus, se remarcă o corespondență topologică între aceste zone: de exemplu, două zone vecine din cortexul vizual sunt corespondentele a 2 zone vecine de pe retină [HUB 77]. Se poate extinde această observație și asupra cortexului auditiv care prezintă un asemenea aranjament topologic al neuronilor încât, frecvențe vecine sunt detectate de către neuroni vecini.

Mai multe experiențe au pus în evidență că această organizare nu este genetică și că reprezintă rezultatul unui proces de învățare de tip hebbian.

Funcția esențială a algoritmului lui *Kohonen* este de a realiza o cuantificare a unui spațiu de stimuli continuu sub forma unui spațiu de ieșire discret de dimensiune inferioară sau egală cu cea a spațiului de stimuli. Se presupune mai întâi că, în spațiul de stimuli, există o relație de vecinătate, de exemplu cea definită prin distanța euclidiană. Se poate spune în acest caz că spațiul de intrare este structurat. Apoi, în spațiul de ieșire constituit din N neuroni, se fixează o relație de vecinătate de tipul celei descrise prin figura 1.45.

Scopul algoritmului este de a găsi o funcție ϕ care, fiecărui element din spațiul de stimuli face să-i corespundă un neuron, ales după valoarea activității (ieșirii) sale, într-o asemenea manieră încât pentru doi stimuli vecini cele 2 imagini obținute în spațiul de ieșire vor fi de asemenea vecine.

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

Dacă această condiție este îndeplinită pentru toate cuplurile de stimuli, vom spune că funcția ϕ respectă topologia.

1.8.2. Principiul algoritmului

Pe baza simbolismului adoptat, algoritmul de auto-organizare al lui *Kohonen* se reprezintă ca în figura 1.44.

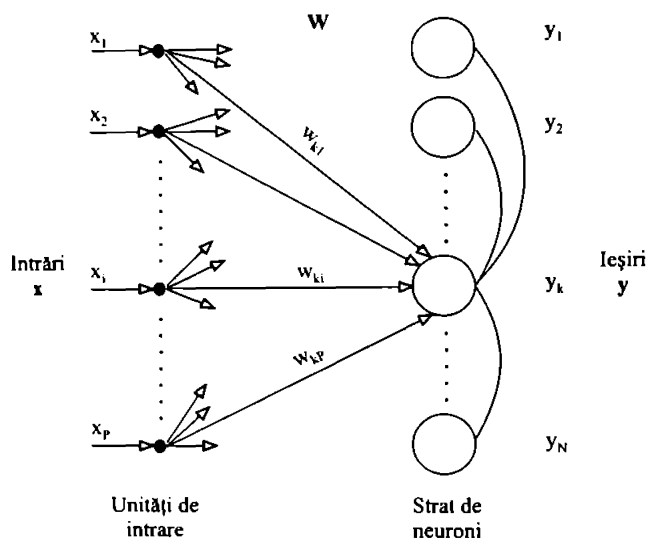


Fig. 1.44. Structura rețelei Kohonen (Self Organizing Feature Map - SOFM)

Rețeaua este constituită din 2 straturi de interconexiuni și un singur strat de neuroni. Primul strat de interconexiuni este așa-numitul strat plastic, deoarece acesta are atribuite ponderile sinaptice w_{ij} și, drept consecință, va suporta ajustarea acestor ponderi. Al doilea strat de interconexiuni este fix și joacă rolul unei rețele competitive.

1.8.3. Funcția primului strat

Primul strat de interconexiuni este utilizat pentru a determina activitatea fiecărui neuron de ieșire, în funcție de vectorul stimul aplicat la intrare. Dacă notăm:

- N - numărul de neuroni din stratul de ieșire
- P - numărul de intrări în fiecare neuron
- \mathbf{x} - vectorul de intrare (de stimuli), $\mathbf{x} = (x_1, x_2, \dots, x_p)^T$
- \mathbf{W} - matricea ponderilor sinaptice de dimensiune $N \times P$

$$\mathbf{W} = \begin{pmatrix} w_{11} & w_{12} & \dots & w_{1p} \\ w_{21} & w_{22} & \dots & w_{2p} \\ \dots & \dots & \dots & \dots \\ w_{N1} & w_{N2} & \dots & w_{Np} \end{pmatrix}$$

Fiecare linie a matricii W formează un vector cu P componente notat w_k ; $k=1, 2, \dots, N$. Vectorul w_k conține ponderile sinaptice aferente conexiunilor de intrare în neuronul k .

y - vectorul de activitate al neuronilor (vectorul de ieșire), $y = (y_1, y_2, \dots, y_N)^T$

atunci:

$$y = W \cdot x$$

Am regăsit deci modelul neuronului *McCulloch - Pitts*.

Considerând că ponderile (vectorii w_k) și intrările (vectorii x) sunt normate, produsul scalar $w_k \cdot x$ va fi maxim când cei doi vectori sunt coliniari. Trebuie să subliniem importanța normei în acest calcul, deoarece principalul scop este de a găsi o măsură a gradului de apropiere sau proximitate dintre vectorii de ponderi sinaptice și vectorul stimul aplicat la intrare (*nearest neighbour classifier*). Dacă vectorii nu sunt corect normați, produsul scalar nu va reprezenta o măsură corectă a gradului de proximitate. Trebuie deci să se utilizeze o normalizare adecvată pentru a compara direcțiile vectorilor fără a pierde norma acestora.

1.8.4. Funcția celui de-al doilea strat

Al doilea strat de interconexiuni este o rețea de inhibiții laterale recurente. Fiecare neuron va fi în relație cu neuronii vecini pe baza unei funcții de interacțiune de tipul celei descrise în figura 1.45.

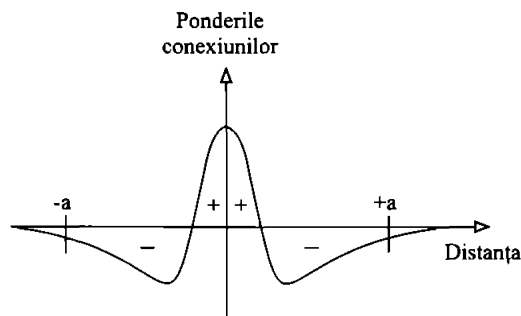


Fig. 1.45. Funcția de interacțiune denumită "pălărie de mexican"

Vom considera inițial o vecinătate monodimensională fixă în timp. Fiecare neuron k este în relație cu vecinii săi ($k - a, \dots, k - 1, k + 1, \dots, k + a$) și fiecare dintre aceștia participă la elaborarea potențialului său prin adăugarea propriului potențial ponderat pe baza funcției "pălărie de mexican". Efectul acestor interacțiuni este de a crea, în ansamblul de neuroni de ieșire, o zonă de activitate stabilă în jurul neuronului care răspunde preferențial la stimulul de intrare aplicat. Această zonă va desemna foarte precis neuronii ale căror ponderi sinaptice vor fi ajustate pe parcursul fazei de învățare. Funcția "pălărie de mexican" poate fi extinsă la cazul bidimensional.

După stabilizarea calculelor, se va determina care este componenta maximală dintre cele N componente ale vectorului de ieșire y .

$$y_k^* = \max_k y_k; \quad k = 1, 2, \dots, N \quad (1.68)$$

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

Se obține astfel vectorul de ponderi sinaptice w_k (asociat neuronului y_k^*) care este cel mai apropiat de vectorul stimul x . Pentru a analiza acest comportament vom reprezenta vectorii x și w_k ($k = 1, 2, \dots, N$) sub formă geometrică. Pentru aceasta vom considera că vectorii de ponderi sunt normalizați ($\|w_k\| = 1$).

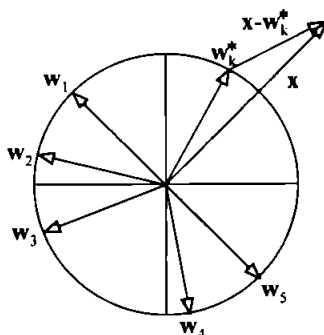


Fig. 1.46. Căutarea proximității relative dintre vectorul x și vectorii w_k

După cum se poate constata din reprezentarea geometrică prezentată în figura 1.46, găsirea produsului scalar maxim dintre vectorii x și w_k se reduce la determinarea distanței euclidiene minime dintre vectorul x și vectorii w_k (dacă vectorii w_k sunt normați). După efectuarea acestei operații și determinarea neuronului câștigător y_k^* , se vor ajusta ponderile vectorului w_k^* (asociat neuronului y_k^*), precum și ponderile care aparțin neuronilor din vecinătatea sa (funcția "pălărie de mexican"), utilizând următoarea regulă de ajustare:

$$w_k(t+1) = w_k(t) + \alpha \cdot y_k(x - w_k(t)) \quad \text{cu } 0 < \alpha < 1 \quad (1.69)$$

Pentru neuronii care nu se află în vecinătatea neuronului y_k^* , ponderile rămân neschimbate:

$$w_k(t+1) = w_k(t) \quad (1.70)$$

Ecuția (1.69) controlează gradul de ajustare a ponderilor sinaptice aferente celui mai activ neuron. Relația menționată arată că se ia o fracțiune (α) din distanța euclidiană dintre x și w_k^* , ($x - w_k^*$), și se "apropie" vectorul w_k^* de x precum și toți ceilalți neuroni din vecinătatea neuronului y_k^* . Această operație este ilustrată în figura 1.47.

Să notăm că, în baza ecuației precedente, nivelul ajustării ponderilor w_k depinde în egală măsură și de valoarea de activare a neuronului y_k . Deci, cu cât răspunsul neuronului este mai puternic, cu atât ajustarea ponderilor va fi mai semnificativă. În sfârșit, să notăm că ponderile care au fost ajustate nu mai sunt normalizate; ele trebuie normalizate pentru a fi pregătite pentru următoarea operație de comparare.

Rezultatul acestui algoritm este că "densitatea spațiului de stimuli converge spre o imagine discretizată astfel încât elementele acestei imagini respectă topologia spațiului de intrare" [COT 89].

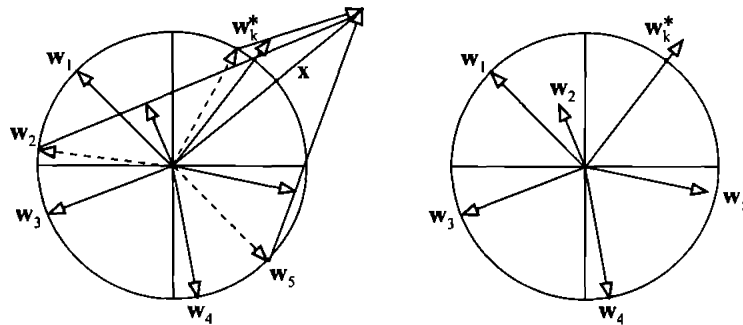


Fig. 1.47. Efectul ajustării ponderilor sinaptice. Vectorii w_1 și w_5 sunt vectorii aferenți neuronilor y_1 și respectiv y_5 din vecinătatea învingătorului y_k^* .

1.8.5. Algoritm SOFM simplificat

Una dintre primele simplificări aduse algoritmului se referă la metoda de evaluare a ieșirii celei mai active. Această măsură a activității se elaborează cu ajutorul produsului scalar dintre vectorul de stimuli și vectorii de ponderi sinaptice. Cum am putut deja constata, este posibilă înlocuirea acestei operații cu una de calcul a distanțelor euclidiene. Deci, determinarea celui mai activ neuron se poate face astfel:

$$\|x - w_k^*\| = \min_k \|x - w_k\| \tag{1.71}$$

Etapa de ajustare a ponderilor sinaptice aferente stratului plastic poate fi și ea modificată. Din moment ce nu se mai calculează produsele scalare pentru determinarea celui mai activ neuron, valoarea y_k poate fi fixată arbitrar. Se poate considera că, pentru neuronul învingător, $y_k^* = 1$ și pentru toți ceilalți neuroni $y_k = 0$. Nu se va selecta astfel decât un singur neuron activ.

În sfârșit, procesul generat de către stratul de inhibiții laterale are ca efect crearea unei "bule" de activitate în jurul neuronului învingător. Acest fenomen este controlat prin forma funcției "pălărie de mexican". Se poate defini o vecinătate în jurul neuronului învingător prin calcul de distanțe euclidiene și se vor ajusta ponderile sinaptice aferente neuronilor care aparțin acestei vecinătăți.

Dacă notăm cu V vecinătatea din jurul neuronului învingător, algoritmul de ajustare a ponderilor se rescrie:

- pentru $|k - k^*| \leq V$

$$w_k(t+1) = w_k(t) + \alpha(x - w_k(t)) \quad 0 < \alpha < 1 \tag{1.72}$$

apoi se efectuează operația de normalizare:

$$w_k = \frac{w_k}{\|w_k\|} \tag{1.73}$$

- pentru $|k - k^*| > V$

$$w_k(t+1) = w_k(t)$$

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

Operația de normalizare poate fi eliminată cu condiția ca x să varieze într-un domeniu convex mărginit în care sunt inițializate ponderile w_k . Într-o astfel de ajustare va fi aplicată norma euclidiană, vectorii ponderilor sinaptice fiind repartizați pe un cerc de rază egală cu o unitate.

Să notăm că vecinătatea din jurul neuronului învingător poate fi o funcție de timp. Această funcție poate lua mai multe forme care depind esențial de topologia planului (stratului) de ieșire. În figura 1.48 se arată evoluția unei vecinătăți, în funcție de timp, potrivit cu o topologie liniară și respectiv planară a stratului de ieșire.

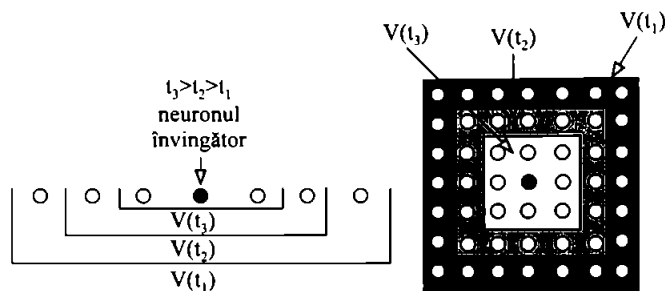


Fig. 1.48. Evoluția unei vecinătăți în jurul neuronului învingător (1D și 2D)

Algoritmul SOFM depinde în faza de învățare de 2 parametri:

V - vecinătatea neuronului învingător, care poate fi o funcție de timp

α - constanta de învățare

Acești doi parametri determină comportamentul rețelei. Trebuie să menționăm că algoritmul prezintă uneori și aberații funcționale care pot fi generate de o alegere incorectă a celor 2 parametri. Aceasta denotă faptul că algoritmi sunt complecși și punerea lor la punct necesită în continuare studii teoretice cu privire la strategiile de optimizare a valorilor parametrilor. Totuși, anumite principii au fost deja clarificate prin simulări numerice și studii teoretice. Algoritmul SOFM poate fi divizat în 2 etape: o organizare grosieră a rețelei urmată de o convergență fină spre soluție. Parametrii α și V vor avea valori mari pe parcursul primei faze (α aproape de 0,5 și V extins la toată suprafața rețelei). Apoi, pe parcursul fazei de convergență, α descrește spre 0 și V spre o vecinătate egală cu 0. După *Kohonen* [KOH 88a], $\alpha(t)$ este o funcție descrescătoare în timp care respectă următoarele condiții:

$$\sum_{s=0}^{\infty} \alpha(s) = \infty \quad \text{și} \quad \sum_{s=0}^{\infty} \alpha^2(s) < \infty \quad (1.74)$$

Alegerea valorilor parametrilor α și V va avea consecințe majore asupra vitezei de auto - organizare a rețelei [COT 89].

1.8.6. Exemple de utilizare a algoritmului SOFM

Considerăm un spațiu de intrare constituit dintr-un plan în care s-au trasat aleator puncte. Coordonatele fiecărui punct (x_1, x_2) vor fi intrările în rețeaua SOFM. În fig. 1.49 se prezintă punctele trasate în plan (spațiul de intrare) precum și dispozitivul de citire a acestor puncte.



Fig. 1.49. a) Distribuția punctelor în plan; b) Dispozitivul de citire

Brațul robotului care citește punctele este constituit dintr-un braț telescopic care generează coordonatele polare aferente fiecărui punct citit. Vom face astfel un anumit număr de citiri pe care le vom aplica pe intrarea rețelei SOFM.

Rețeaua este constituită din 2 intrări și 100 neuroni așa cum indică figura 1.50.

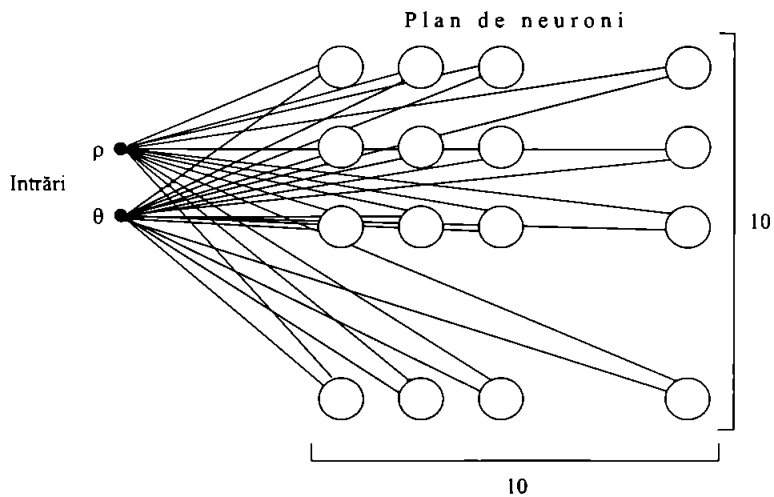


Fig. 1.50. Structura rețelei SOFM

Se inițializează matricea W (ponderile sinaptice aferente stratului de interconexiuni plastice) cu valori aleatoare și se va urmări evoluția acestora pe parcursul procesului de ajustare (învățare).

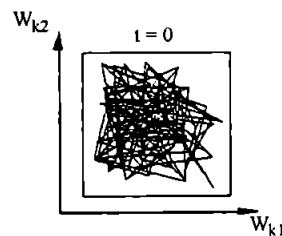


Fig. 1.51. Repartiția inițială a neuronilor în spațiul ponderilor

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

Trebuie să subliniem că termenul de auto - organizare nu înseamnă o migrare a neuronilor în cadrul rețelei. Este vorba de reprezentarea aleasă care indică un aranjament sub formă de plasă sau grilă regulată. Această reprezentare, așa cum arată figura 1.52, descrie în plan "poziția" fiecărui neuron prin cuplul (vectorul) său de ponderi sinaptice. Liniile dintre neuroni rezultă dintr-un etichetaj și indică cei 4 vecini imediați ai fiecărui neuron. Deoarece neuronii sunt aranjați într-un plan, am considerat că fiecare se află în relație cu cei 4 vecini ai săi.

Fenomenul de auto - organizare care se dezvoltă pe parcursul procesului de învățare reflectă 2 componente:

- valorile ponderilor sinaptice tind să acopere planul în care sunt reprezentate urmărind forma și densitatea spațiului de intrare.
- valorile ponderilor se ordonează în planul de ieșire. Această ordine depinde de relația care există între elementele spațiului de intrare.

În figura 1.52 se reprezintă evoluția vectorilor sinaptici w_k în timp (pe parcursul procesului de învățare) unind cu segmente de dreaptă punctele care în spațiul de ieșire corespund la neuroni vecini.

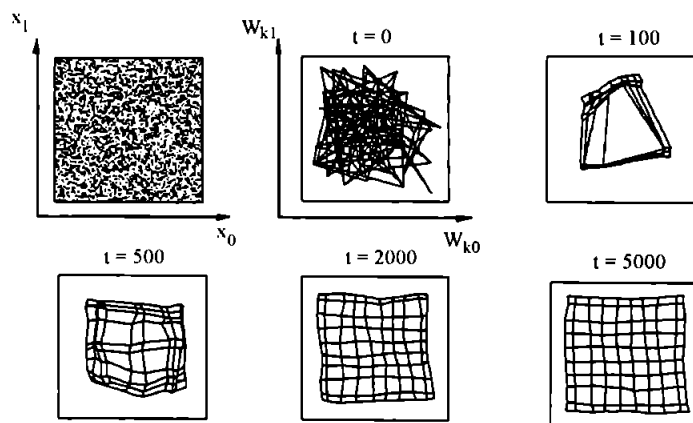


Fig. 1.52. Evoluția organizării în timp

Se observă convergența procesului de auto - organizare care conduce la o imagine discretizată a densității spațiului de stimuli. În exemplul prezentat, densitatea spațiului de stimuli este uniformă (poziționare aleatoare a punctelor în spațiul de intrare). Imaginea discretizată va conține o grilă (plasă) regulată. Să remarcăm numărul mare de iterații necesar pentru derularea procesului de auto - organizare. Este unul dintre principalele inconveniente ale acestui algoritm care are o viteză de convergență redusă.

Reprezentarea aleasă în figura 1.52 plasează neuronii în planul ponderilor sinaptice. Trebuie să înțelegem poziția fiecărui neuron drept o reprezentare a vectorilor w_k , așa cum se arată în figura 1.53.

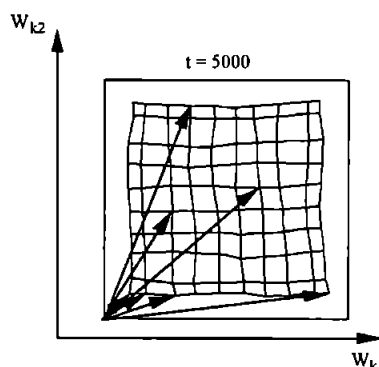
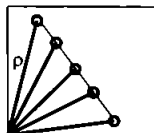


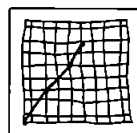
Fig. 1.53 Reprezentarea vectorilor sinaptici în spațiul bidimensional

Această reprezentare presupune că vectorii w_k nu sunt normați, ceea ce înseamnă că, pentru detectarea gradului de proximitate dintre vectorii stimuli și vectorii w_k , s-au utilizat distanțele euclidiene. Dacă se dorește utilizarea metodei produselor scalare, trebuie realizată normalizarea vectorilor w_k într-un spațiu de dimensiune superioară, și apoi calcularea produselor scalare în acest spațiu, pentru a nu modifica norma vectorilor w_k .

După faza de învățare, exploatarea rețelei se poate face în conformitate cu procesul descris în figura 1.54.



a) pozițiile succesive ale brațului robotului



b) răspunsul rețelei SOFM

Fig. 1.54. Exemplu de exploatare a rețelei organizate

Se constată că mișcările executate de brațul robotului în spațiul de intrare sunt reproduse pe planul de neuroni, cu o anumită simetrie, marcată pe grila definită de caroiajul virtual stabilit în faza de învățare. Reproducerea nu este identică ci doar simetrică datorită faptului că, în timpul învățării, neuronii nu au avut repere absolute, ci doar relative constituite prin relația de proximitate dintre stimulii spațiului de intrare.

Dintre reprezentările utilizate pentru a explica funcționarea algoritmului SOFM, amintim "mozaicul lui Voronoi", care definește o zonă de referință în jurul fiecărei ponderi; această zonă corespunde punctelor din spațiul de intrare care excită preferențial același neuron. Este zonă de stimuli de intrare care aparține aceluiași agregat. Figura 1.55 prezintă un exemplu:

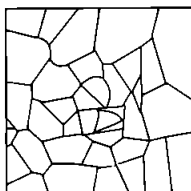


Fig. 1.55. Mozaicul lui Voronoi

1.8.7. Concluzii

După cum am constatat, algoritmul SOFM reprezintă un instrument foarte puternic pentru pre - tratamentul datelor destinate unei viitoare clasificări. Algoritmul elaborează o reprezentare a datelor agregând vectorii vecini care aparțin aceleiași clase și separând agregatele care corespund unor clase diferite. El este de altfel utilizat pe post de strat de pre - tratare în algoritmul *Counter propagation* [HEC 87].

Aplicațiile bazate pe rețeaua SOFM sunt foarte numeroase. Putem cita, spre exemplificare, analiza unei mari cantități de date tehnologice aferente unei linii de fabricație circuite VLSI [TRY 89], crearea unui spectru topografic pentru analiza vorbirii [KOH 88b], controlul unui braț de robot într-un câmp prevăzut cu obstacole [GRA 88], [RUE 90] sau a unui robot mobil într-un labirint [SOR 89].

1.9. Modele avansate

Modelele descrise în capitolele precedente sunt deja modele clasice. O serie de alte modele încep să fie utilizate în diverse aplicații. Dintre aceste modele cităm:

- Modelul *Hérault - Jutten* care prezintă un interes deosebit pentru procesarea de semnal (problema separării surselor) [JUT 87], [JUT 88], [COM 90].
- Modelul ART (*Adaptive Resonance Theory*) elaborat de *Grossberg* [GRO 87] și divizat în 2 clase: ART 1 cu valori de intrare binare și ART 2 cu valori de intrare analogice. Acest model servește la formarea de categorii de date cu ajutorul unui parametru de vigilență care definește finețea categoriilor.
- Modelul RCE (*Restricted Coulomb Energy*) elaborat de *Cooper* (și alții), utilizează învățarea supervizată și este destinat problemelor de clasificare.
- Modelul LVQ (*Learning Vector Quantization*) [KOH 86] care este un clasificator și care se bazează pe un algoritm de învățare semi-supervizată [GUY 90]. Este unul dintre cei mai performanți algoritmi de clasificare. În mod particular se adaptează bine problemelor care revendică un set de antrenament foarte amplu [BEN 90].
- Modelul *Cognitron* al lui *Fukushima* [FUK 75] construit pe o structură multistrat cu o tehnică de auto - organizare pentru învățare. Se aplică în mod esențial la recunoașterea caracterelor.
- Modelul rețelelor selective, propus de *Edelman* [EDE 82] care introduce noțiunea de evoluție a unui sistem (în sensul selecției naturale), noțiune absentă la majoritatea celorlalte modele. Acest model stă la baza seriei de mașini DARWIN.

Mai multe variante au fost de asemenea dezvoltate, pornind de la modelele clasice:

- rețele multistrat recurente cu învățare pe baza algoritmului EBP [ALM 89].
- rețele TDNN (*Time Delay Neural Networks*) [WAI 89] bazate pe perceptroni multistrat și care învață după algoritmul EBP. O caracteristică a acestor rețele este că neuronul recepționează simultan mai multe eșantioane întârziate aferente semnalelor de intrare.
- rețeaua BAM (*Bidirectional Associative Memory*) [KOS 87] construită pe baza a 2 rețele *Hopfield* rebucate.
- rețeaua *Counter propagation* [HEC 87] construită dintr-o rețea *Kohonen* urmată de un perceptron.

Această listă incompletă indică faptul că domeniul este în plină ascensiune. Teoriile sunt încă insuficient elaborate pentru a satisface toate problemele reale. O abordare mai "experimentală" survine adesea pentru a ascunde aceste carențe teoretice, dar reclamă numeroase clarificări pentru aplicațiile concrete.

Algoritmii de învățare pot fi grupați în 4 categorii:

- algoritmi care generează ponderi fixe
- algoritmi de învățare supervizată
- algoritmi de învățare semi - supervizată
- algoritmi de învățare nesupervizată

Se definește învățarea supervizată acea tehnică de învățare care minimizează o funcție de cost, care este în general un criteriu de eroare pătratică. Minimizarea este realizată printr-o metodă de coborâre (pe suprafața de eroare) pe direcția gradientului negativ sau prin metoda gradientului stohastic. Cu această metodă, se cunoaște exact vectorul de eroare, care apare ca diferența între ieșirea dorită și ieșirea reală.

Învățarea semi - supervizată nu operează asupra unui vector de eroare ci asupra unui scalar care indică calitatea răspunsului. Este vorba de un coeficient de similaritate între ieșirea dorită și cea reală care ne oferă o indicație referitoare la mărimea erorii comise.

Învățarea nesupervizată nu face apel la nici o ieșire dorită. Ajustarea ponderilor nu depinde decât de valorile de pe intrarea rețelei și de ponderile sinaptice calculate și actualizate până în momentul respectiv.

În tabelul 1.1. redăm câteva exemple de aplicații posibile în funcție de tipul rețelei; acest tabel nu include evident o tratare exhaustivă.

Aplicație	Tip rețea			
	Fără învățare	Înv. supervizată	Înv. semi-supervizată	Înv. nesupervizată
Procesarea vorbirii sau a imaginilor	Inhibiții laterale recurente	TDNN MLP + <i>Back Propagation</i> MLP rebusată	LVQ	SOFM <i>Hérault - Jutten</i>
Conducere de proces		MLP + <i>Back Propagation</i> MLP rebusată <i>Adaline</i>	<i>Counter propagation</i> <i>Reinforcement learning</i>	ART 3
Optimizare	<i>Hopfield - Tank</i>	<i>Hopfield</i>	<i>Reinforcement learning</i>	SOFM
Analiză de date Clasificare Extragere de caracteristici (compresie date)	Rețele ad - hoc	MLP + <i>Back Propagation</i>	LVQ <i>Neocognitron</i>	RCE SOFM
Procesare de semnal		TDNN <i>Adaline</i>		<i>Hérault - Jutten</i>
Memorie asociativă	<i>Interactive Activation and Competition</i>	<i>Hopfield</i> BAM		

Tabelul 1.1. Aplicații ale RN

După cum se constată din tabelul 1.1, există încă o "latură artistică" în ceea ce privește aplicațiile RN. Care este rețeaua cea mai bine adaptată unei aplicații date, care sunt parametrii optimali

SOLUȚII DE IMPLEMENTARE A REȚELELOR NEURONALE PE ARHITECTURI SISTOLICE

(mărime, constantă de învățare, etc.), care este tipul de învățare, sunt tot atâtea întrebări la care trebuie să răspundă perspicacitatea inginerului. În general, se poate găsi întotdeauna o RN capabilă să rezolve o problemă dată. Dar, rămâne de văzut dacă soluția neuronală este mai bună decât o soluție clasică, dacă soluția neuronală este implementabilă pe calculatoarele disponibile și dacă rezultatele obținute pe cale neuronală au o semnificație corectă.

Cu modelele actuale, inginerul trebuie să integreze o parte din știința sa în rețeaua neuronală. Alegerile sale sunt în mod necesar ghidate de cunoașterea apriorică a problemei tratate.

RN nu trebuie privită ca o cutie neagră care spionează un proces real pentru a-i extrage un model numeric general. Actualmente, conexiunismul nu reprezintă decât un concept relativ nou, care face apel la o reprezentare distribuită, la un paralelism masiv, la noțiunea de învățare și nu de programare și la o tratare fundamental numerică și nu simbolică a informației.

2. IMPLEMENTĂRILE VLSI

Calculatoarele actuale, inclusiv cele bazate pe structuri paralele, sunt construite în jurul unui număr mic de unități de calcul rapide, precise și complexe, fiind perfect adaptate algoritmilor de calcul tradiționali. Aceste mașini nu reprezintă un optim și pentru algoritmi conexioniști care revendică alte aptitudini: precizia este o constrângere minoră, fiecare operație este simplă dar paralelismul execuției este esențial. Pentru domeniul conexionist trebuie dezvoltate mașini care satisfac aceste cerințe.

2.1. Introducere

În figura 2.1. se prezintă modelul neuronului. Acest model stă la baza tuturor rețelelor conexioniste prezentate în această lucrare.

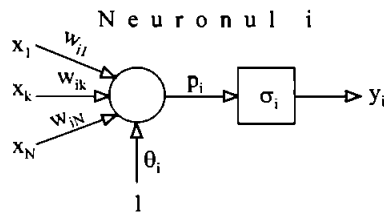


Fig. 2.1. Modelul neuronului

Notățiile utilizate sunt:

- x_k - valoarea de activare a (ieșirea) neuronului k
- w_{ik} - valoarea ponderii sinaptice aferentă conexiunii dintre neuronul k (sursă) și neuronul i (destinație)
- p_i - potențialul neuronului i
- σ_i - funcția de activare (transfer) aferentă neuronului i
- θ_i - valoarea de prag asociată neuronului i
- y_i - valoarea de activare a (ieșirea) neuronului i

Problema construcției unei rețele conexioniste constă în definirea structurii acestei rețele și a parametrilor aferenți. Structura este definită prin conexiunile (locale, totale, parțiale, plastice sau nu etc.) dintre neuroni, iar parametrii sunt reprezentați de funcția de transfer a neuronului (continuă, discontinuă, valorile de activare între 0 și 1 sau între -1 și $+1$ etc.) și valorile de prag.

În funcție de modelul arhitectural ales, se definește o regulă de învățare. Aceasta descrie evoluția coeficienților sinaptici în timp. Ecuația următoare reprezintă o formă generală a acestei reguli:

$$\Delta \mathbf{W} = \alpha(t) \cdot \sigma(\mathbf{u}(t)) \cdot \tau(\mathbf{v}(t))^T \cdot \xi(t) \quad (2.1)$$

- \mathbf{W} - matricea de ponderi sinaptice
- $\alpha(t)$ - constantă (parametru) de învățare, de regulă cuprins între 0 și 1 .
- σ, τ - funcții de transfer
- $\mathbf{u}(t), \mathbf{v}(t)$ - doi vectori.

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

Faza de recunoaștere face apel la un subansamblu de funcții utilizate în faza de învățare. În general, recunoașterea (o fază mai simplă) se derulează astfel:

- elaborarea potențialelor neuronilor

$$p_i(t) = \sum_{k=1}^N w_{ik} \cdot x_k + \theta_i \quad (2.2)$$

- calculul valorilor de activare (ieșirilor):

$$y_i(t) = \sigma(p_i(t)) \quad (2.3)$$

Deci, realizarea unei rețele conexiuniste se reduce la implementarea eficientă a unor funcții ca:

- produsul scalar, produsul matrice - vector și produsul unui vector cu transpusul său;
- aplicarea unor funcții neliniare (funcțiile de activare sau transfer).

2.2. Soluții posibile pentru implementarea rețelelor neuronale

Există trei soluții posibile pentru implementarea rețelelor neuronale:

- Implementări *software* care constau în elaborarea unor simulatoare care rulează pe calculatoare universale sau a unor biblioteci de funcții standard cu ajutorul cărora se pot dezvolta aplicații proprii.
- Implementări *hardware* dezvoltate pe calculatoare universale. Acestea presupun adăugarea unor plăci specializate (acceleratoare) gestionate prin intermediul unor simulatoare adecvate.
- Implementări *hardware* speciale care vizează realizarea unor circuite VLSI specializate pe aplicații particulare.

Simulatoarele *software* reprezintă o soluție ieftină dar care realizează performanțe modeste; performanțele sunt limitate de performanțele mașinii suport. În tabelul 2.1 sunt prezentate câteva dintre cele mai cunoscute simulatoare soft pentru RN, existente pe piață la începutul anului 1998.

Nr. crt.	Simulator	Producător	Caracteristici / Comentarii
1	<i>Brain Maker</i>	<i>California Scientific Software</i>	Rețele <i>feedforward</i> multistrat (maxim 8 straturi), rețele recurente etc.
2	<i>SAS Neural Network Application</i>	<i>SAS Institute</i>	Perceptroni multistrat, counterpropagation, rețele <i>feedforward</i> multistrat etc.
3	<i>Neural Works</i>	<i>Neural Ware</i>	<i>Backpropagation</i> , ART-1, <i>Kohonen</i> , rețele probabilistice, rețele SOFM etc.
4	<i>MATLAB Neural Network Toolbox</i>	<i>Math Works</i>	Rețele <i>feedforward</i> , rețele asociative, rețele SOFM etc.
5	<i>Propagator</i>	<i>ARD Corporation</i>	<i>Backpropagation</i> cu maxim 5 straturi
6	<i>Neuralyst</i>	<i>Cheshire Engineering Corporation</i>	Algoritmul <i>backpropagation</i>
7	<i>Neu Fuz 4</i>	<i>National semiconductor Corporation</i>	Rețele Neuro - fuzzy
8	<i>Neuro Solution v3.0</i>	<i>Neuro Dimension</i>	Rețele de perceptroni multistrat, rețele

			<i>feedforward</i> generalizate, rețele SOFM, rețele recurente etc.
9	<i>Qnet For Windows v2.0</i>	<i>Vesta Services</i>	Sistem de dezvoltare pe 32 biți pentru <i>Windows NT</i> , 95 și 3.1.
10	<i>havBpNet</i> <i>havFmNet</i> <i>havBpNet</i>	<i>hav. Software</i>	Rețele <i>feedforward</i> , recurente, SOFM
11	<i>IBM Neural Network Utility</i>	<i>IBM</i>	Modele multiple de RN, sisteme <i>fuzzy</i> și sisteme combinate: <i>fuzzy</i> - RN
12	<i>Neuro Genetic Optimizer NGO v2.0</i>	<i>BioComp Systems</i>	Utilizează algoritmi genetici pentru optimizarea intrărilor și structurii RN
13	<i>OWL Neural Network Library</i>	<i>HyperLogic Corporation</i>	<i>Backpropagation</i> , DBAM (<i>Discrete Bidirectional Associative Memory</i>), ABAM (<i>Adaptive Bidirectional Associative Memory</i>) etc.
14	<i>Pattern Recognition Workbench Expo PRO</i>	<i>Unica Technologies</i>	<i>Backpropagation</i> și alte arhitecturi de RN utilizate în recunoașterea modelelor

Tabelul 2.1. Simulatoare soft dedicate rețelelor neuronale

Plăcile acceleratoare specializate se implementează într-un sistem gazdă și au la bază procesoare specializate (de exemplu procesoare de semnal) sau structuri multiprocesor special proiectate pentru implementarea modelelor conexioniste.

Implementările VLSI (*Very Large Scale Integration*) presupun realizarea unor circuite specializate dedicate modelelor conexioniste. Aceste circuite pot fi utilizate atât în construcția plăcilor acceleratoare cât și autonom (în aplicații particulare care revendică procesare în timp real).

Vom face în continuare o analiză a implementărilor VLSI, evidențiind problematica pe care proiectantul trebuie să o rezolve.

2.3. Structuri candidate la integrare

Între modelele teoretice distingem RN monostrat și RN multistrat, stratul făcând referință la un grup de neuroni. Această diferențiere nu are influențe profunde în termeni structurali. Este suficient să considerăm RN ca un ansamblu de neuroni între care există sau nu legături.

Diferențele între diferitele arhitecturi de RN apar la nivelul matricii de interconexiuni (matricii de ponderi sinaptice W) care va fi mai mult sau mai puțin populată cu elemente nenule. Implementările VLSI vor putea eventual exploata proprietățile matricii W . Câteva exemple de matrici de interconexiuni sunt prezentate în fig. 2.2.

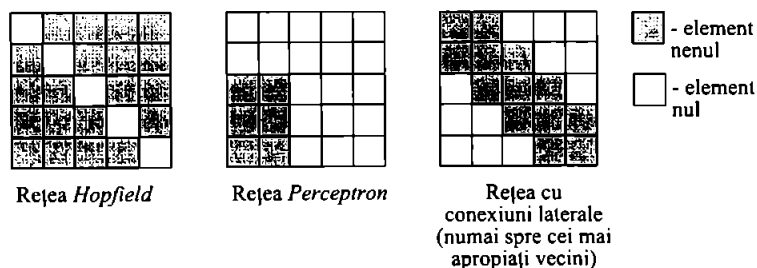


Fig. 2.2. Exemple de matrici de interconexiuni

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

Matricea W descrie structura rețelei de sinapse din cadrul RN. Poziția neuronilor este total determinată prin matricea W și prezența lor nu se justifică decât prin funcția de transfer σ . Această constatare conduce la prima problemă (de decizie) pe care trebuie s-o rezolve proiectantul: va integra sinapse sau va integra neuroni? A integra sinapse înseamnă a integra în structura VLSI o matrice de întreconexiuni care va pune în relație neuronii exteriori circuitului. A integra neuroni înseamnă că blocurile de calcul a funcției de transfer vor fi prezente în circuitul VLSI și, de asemenea, și întreconexiunile dintre aceste blocuri.

2.3.1. Integrarea neuronilor

Dacă întreconexiunile sunt locale, limitate la cei mai apropiați vecini, soluția de integrare a neuronilor este realistă. Dacă însă avem de-a face cu o rețea complet interconectată, cum este cazul rețelei *Hopfield*, integrarea neuronilor devine dificilă din două motive:

- întreconexiunile dintre neuroni vor ocupa cea mai mare parte a suprafeței de implantare;
- extensibilitatea rețelei va fi limitată de numărul de intrări / ieșiri.

Dacă dorim o RN complet interconectată, trebuie permis accesul fiecărui neuron din fiecare circuit spre toți ceilalți neuroni din toate celelalte circuite. Această constrângere nu va putea fi satisfăcută pentru RN cu un număr mare de neuroni decât în cazul aplicării unei tehnici sistolice (de comunicație).

2.3.2. Integrarea sinapselor

Integrarea sinapselor permite extensia rețelei și nu pune probleme de intrare / ieșire (fig. 2.3).

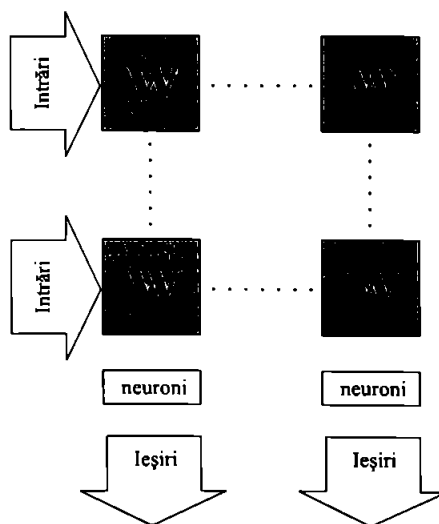


Fig. 2.3. Extensia matricii sinaptice

Dimensiunea RN poate fi mărită printr-un simplu pavaj rectangular de circuite sinaptice. RN multistrat vor putea fi ușor implementate prin conectarea în cascadă a circuitelor sinaptice, așa cum indică figura 2.4.

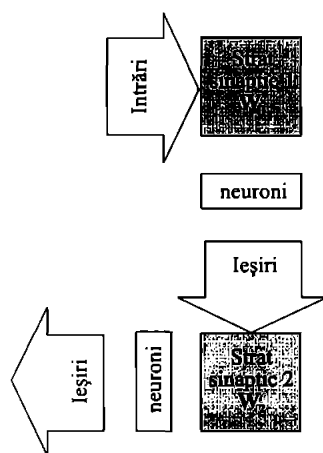


Fig. 2.4. RN multistrat

În capitolele 3, 4 și 5 vom detalia proiectarea unei rețele sistolice (RS) destinată modelelor conexiuniste și care integrează sinapse.

Referitor la implementarea rețelelor neuronale pe structuri VLSI, pe plan mondial se fac cercetări pe mai multe direcții. Implementările VLSI pot fi grupate în patru categorii:

- implementări analogice;
- implementări digitale;
- implementări hibride (digital - analogice);
- implementări opto - electronice.

2.4. Implementările VLSI analogice

Un exemplu de implementare VLSI analogică ar fi retina lui *Mead* și *Mahowald* care mimează funcționarea unei retine biologice [SHE 95]. Implementarea în discuție reprezintă un model analogic al retinei mamiferelor superioare. Intrările în rețea sunt optice. Prototipul realizat integrează o retină de 48 x 48 pixeli ale cărei caracteristici de funcționare se apropie de cele aferente celulelor bipolare din retina vertebratelor.

Este vorba de o rețea hexagonală realizată în tehnologie analogică. Circuitul se descompune în trei părți care corespund unor funcții specifice pe care le regăsim într-o retină biologică:

- fotoreceptorul care calculează logaritmul intensității luminoase pe care o recepționează; este un dispozitiv cu control automat al amplificării;
- evaluatorul mediei temporale și spațiale de la ieșirea fotoreceptorului; acesta corespunde celulelor orizontale din retina biologică;
- evaluatorul diferenței dintre semnalul fotoreceptorului și semnalul furnizat de celulele orizontale; acesta corespunde celulelor bipolare din retina biologică.

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

Rezultatul va fi un circuit de detecție a conturului prin diferențiere spațială, sau de detecție a mișcării prin diferențiere temporală. Vom descompune circuitul în elementele de bază pentru a preciza funcționalitatea și arhitectura acestora.

2.4.1. Fotoreceptorul

Fotoreceptorul transformă semnalul optic recepționat de un fototranzistor într-un potențial proporțional cu logaritmul intensității luminoase. În figura 2.5 se prezintă o implementare posibilă a fotoreceptorului.

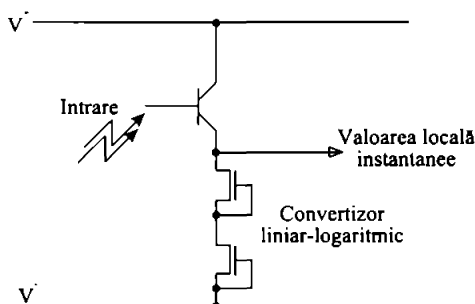


Fig. 2.5. Structura fotoreceptorului cu doi tranzistori

Acest dispozitiv preia curentul generat de elementul fotoreceptor și îl injectează într-un element a cărui caracteristică curent - tensiune este logaritmică. Convertizorul liniar - logaritmic este realizat din două diode construite din tranzistoare MOS, conectate în serie. Să notăm că tranzistorul bipolar utilizat nu este fabricat; este un element parazit generat implicit prin fabricarea tranzistoarelor MOS.

2.4.2. Diferențiatorul (celulele bipolare)

Pornind de la semnalul furnizat de către fotoreceptor, este posibilă percepția unei mișcări într-o imagine. Există o mișcare atunci când se detectează o variație de intensitate luminoasă pe un pixel, adică atunci când se detectează local o derivată temporală nenulă. Pentru această detecție se implementează un dispozitiv simplu, descris în figura 2.6, care elaborează rezultatul comparației dintre valoarea instantanee a unui pixel și valoarea aceluiași pixel convenabil întârziată.

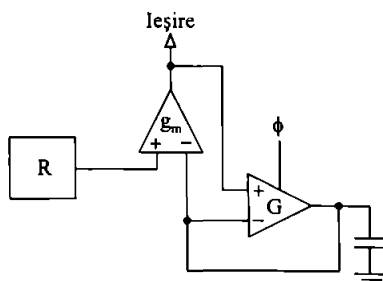


Fig. 2.6. Diferențiatorul

În figura 2.6 R reprezintă fotoreceptorul. Transconductanța G se comportă ca un tranzistor de trecere dar care permite eliminarea sarcinilor parazite generate în momentul comutației. Acest dispozitiv furnizează valoarea locală a derivatei care nu ține cont de valorile celulelor vecine (pixelilor vecini). Cu scopul de a construi un model cât mai apropiat de cel real, autorii au dezvoltat o rețea de rezistențe care pune la dispoziție în fiecare punct (pixel) media ponderată a valorilor vecine (pixelilor vecini).

2.4.3. Rețeaua rezistivă

Structura rețelei rezistive este descrisă în figura 2.7.

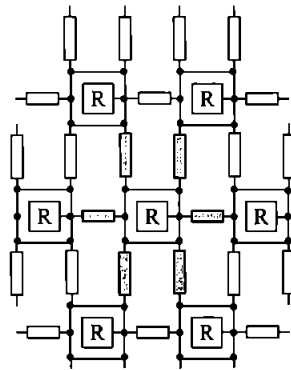


Fig. 2.7. Interconectarea receptorilor

Rețeaua hexagonală de rezistențe furnizează (în fiecare punct) media valorilor de ieșire aferente celor 6 receptori vecini. Acest dispozitiv calculează derivatele spațiale aferente imaginii proiectate pe retină și permite de asemenea detecția conturilor.

2.4.4. Rețeaua completă

Rețeaua completă este constituită din cei trei operatori anterior descriși. În figura 2.8 se prezintă structura globală a unui element din rețea.

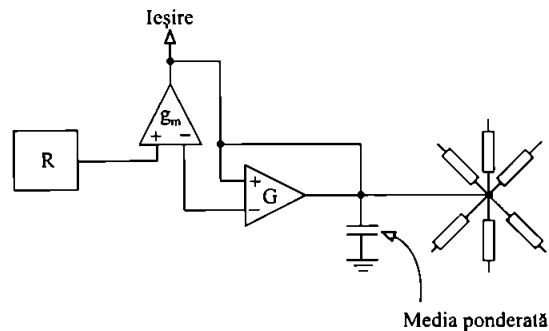


Fig. 2.8. Schema generală a unui receptor

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

În această rețea, amplificatorul calculează diferența dintre curentul generat de către fotoreceptorul R și cel generat de rețeaua de rezistențe. Capacitatea C stochează media spațiu / timp a ieșirilor fotoreceptorilor vecini. Exploatarea rețelei se poate face în două moduri:

- dacă se deconectează rețeaua de rezistențe, dispozitivul calculează derivata temporală pentru detecția de mișcare;
- dacă rețeaua de rezistențe este conectată, dispozitivul calculează diferența relativă dintre intensitatea luminoasă recepționată într-un punct (pixel) și media intensităților luminoase recepționate în cele 6 puncte vecine. Dispozitivul este foarte sensibil în cazul punctelor izolate, mai puțin sensibil în punctele din colțuri și insensibil la un gradient uniform. Această proprietate este de asemenea valabilă dacă obiectul din imagine este în mișcare.

Dispozitivul permite deci detecția unui contur sau a unei mișcări. În scopul integrării acestui dispozitiv într-o structură VLSI, toate elementele sunt realizate cu tranzistoare (nu apar rezistențe și nici capacități integrate).

2.5. Implementările VLSI digitale

Calculatoarele clasice cu procesare serială a instrucțiunilor nu sunt capabile de o emulare eficientă a algoritmilor neuronali, care revendică pentru eficiență ridicată arhitecturi paralele de tip DAP (*Distributed Array Processor*), MIMD (*Multiple Instruction Multiple Data Computer*) sau rețele sistolice RS (*Systolic Array Processor*).

Un neuroprocesor digital (un neurocomputer) este un microsistem VLSI digital ale cărui caracteristici au fost optimizate pentru execuția algoritmilor neuronali. Neuroprocesoarele trebuie să exceleze prin programabilitate, flexibilitate, precizie și *pipeline*-abilitate. Arhitectura VLSI a neuroprocesoarelor digitale trebuie optimizată în acord cu următoarele criterii de performanță:

- Adaptabilitate: *On - chip learning*;
- Flexibilitate: Programabilitate pentru cele mai variate arhitecturi de RN;
- Cost redus: Tehnologii ieftine (de exemplu CMOS);
- Viteză ridicată: Paralelism masiv;
- Integrare la nivelul sistemului: Co-proiectarea *hardware - software*

Un neuroprocesor digital, bazat pe o arhitectură multiprocesor, este prezentat în figura 2.9. În afara ariei de procesoare EP, în structura neuroprocesorului mai apar și alte subsisteme ca: *controller*-ul ariei de procesoare, memoria sistem aferentă ariei de procesoare, interfața cu sistemul gazdă (*host*) etc. Fiecare EP conține în structura sa câteva blocuri funcționale. Un set de registre va fi utilizat pentru memorarea datelor intermediare cum ar fi constanta de învățare, numărul de identificare al elementului EP etc.

O unitate de înmulțire / adunare va fi utilizată în interiorul fiecărui EP pentru calculul sumelor ponderate care intervin în cadrul algoritmilor conexioniști. Numărul de biți necesari pentru reprezentarea datelor procesate se determină pe baza preciziei ce se urmărește a se realiza. Pentru creșterea productivității sistemului se aplică tehnica *pipeline*. Ponderile sinaptice sunt memorate într-o memorie de ponderi a cărei adresare este controlată de un generator de adrese special. Funcția de transfer neliniară aferentă neuronului poate fi implementată sub formă tabelară (*look - up table*), prin memorarea tabelii în memoria neuronului împreună cu valoarea de prag. Opțional, în structura EP pot fi introduse și alte blocuri cum ar fi, de exemplu, o unitate ALU adițională pentru calculul derivatei funcției de transfer pentru algoritmul EBP.

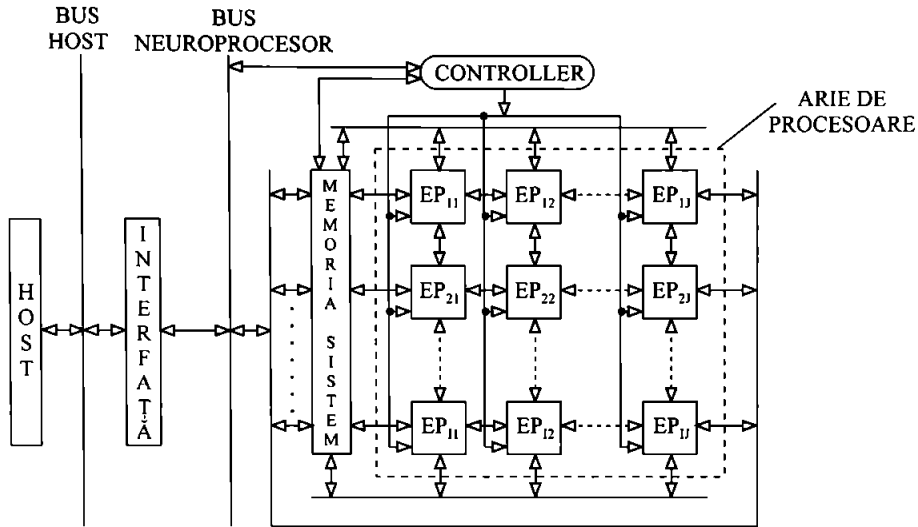


Fig. 2.9. Arhitectura neuroprocesorului digital constituit dintr-o arie de elemente de procesare (EP)

O metodă de implementare care conduce la o eficiență foarte ridicată constă în utilizarea ariilor (rețelelor) sistolice (RS), care permit un grad înalt de paralelism și o *pipeline*-izare completă a algoritmilor conexioniști. RS rezolvă de asemenea și problema comunicațiilor interprocesor.

O RN cu trei straturi este prezentată în figura 2.10. a) iar în figura 2.10. b) se prezintă modul de implementare al acestei RN pe RS bidimensională. Alocarea corespunzătoare a procesoarelor (maparea corespunzătoare a RN pe elementele EP din structura RS) va evita conflictele în transferurile de date dintre procesoarele EP. În cazul schemei de alocare din figura 2.10, sumele ponderate se pot obține de-a lungul fiecărui rând și a fiecărei coloane.

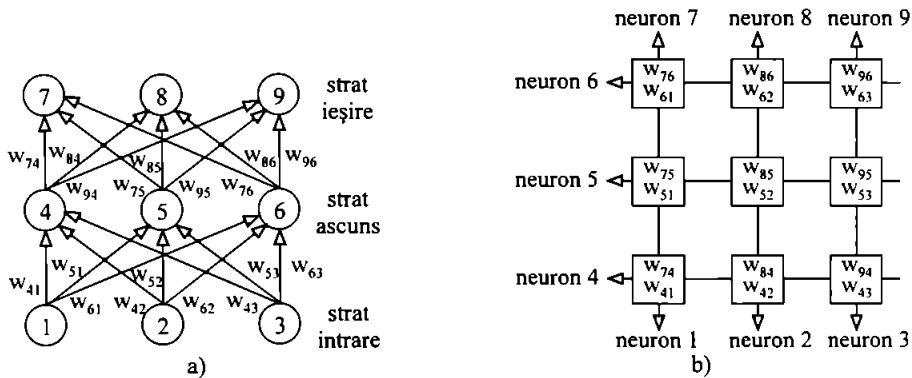


Fig. 2.10. a) Arhitectura RN cu 3 straturi; b) Maparea RN în aria de procesoare

În figura 2.11 este reprezentat fluxul de date în RS în cazul implementării algoritmului EBP (faza *feedforward*).

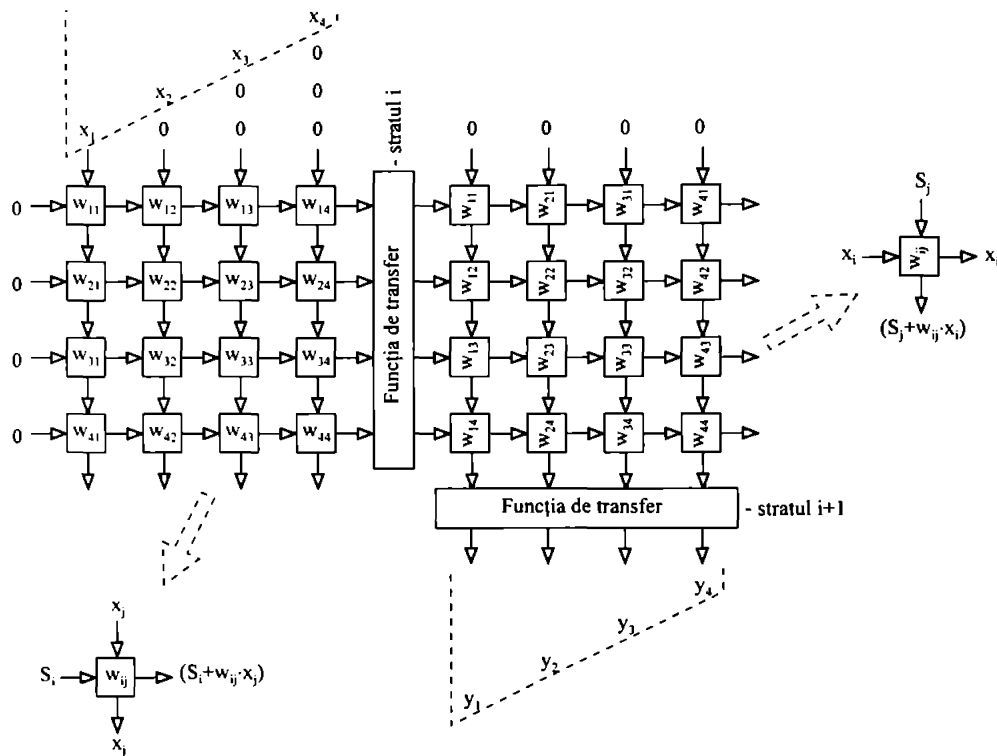


Fig. 2.11. Fluxul de date în RS pentru faza de recunoaștere aferentă algoritmului EBP

Pe parcursul fazei de recunoaștere (*feedforward*) datele se propagă dinspre stratul de intrare spre stratul de ieșire. RS funcționează cu un tact global care asigură sincronizarea operațiilor executate în interiorul celulelor sistolice (CS); fiecare CS execută operațiile aritmetice și de transfer indicate în figura 2.11. În cadrul fazei de învățare fluxul de date își schimbă direcția și va evolua dinspre stratul de ieșire spre stratul de intrare (propagarea erorii înapoi); punctul de start al fluxului de date va fi constituit de semnalele de eroare asociate stratului de ieșire.

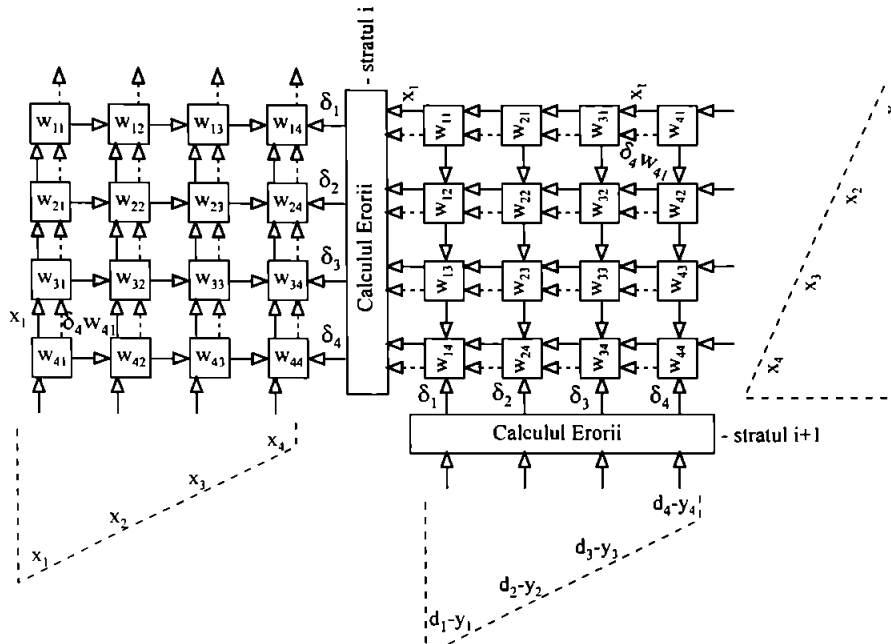


Fig. 2.12. Fluxul de date în RS pentru faza de învățare (*backward*) aferentă algoritmului EBP

Pe baza descrierii algoritmului EBP, prezentată la 1.6.4 putem deduce și operațiile efectuate de CS în faza *backward*. Aceste operații sunt prezentate în figura 2.13.

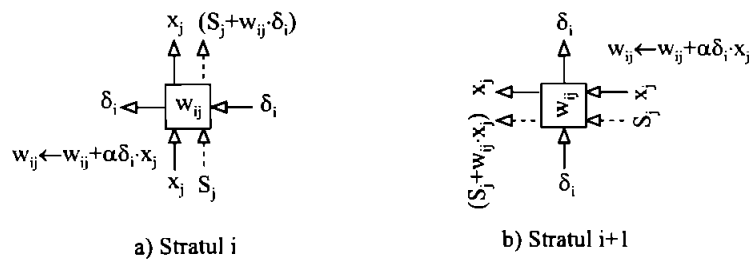


Fig. 2.13. Operațiile efectuate de CS în faza *backward*

Se observă că, indiferent de stratul neuronal din care fac parte, celulele sistolice au structura identică. Să mai notăm că în faza *backward*, în afară de operațiile de transfer evidențiate în figura 2.13, celulele sistolice efectuează și o ajustare a mărimii de stare: $w_{ij} \leftarrow w_{ij} + \alpha \delta_i \cdot x_j$. Această operație realizează ajustarea ponderilor sinaptice (învățarea).

2.5.1. Sistemul CNAPS

Adaptive Solutions Inc. a dezvoltat sistemul CNAPS (*Connected Network of Adaptive Processors*) [MUE 93], [GRI 91]. Comparativ cu mașinile tradiționale bazate pe microprocesoare

sistemul CNAPS este capabil de performanțe mult superioare (câteva ordine de mărime). Aceste performanțe se obțin printr-un paralelism masiv și prin capacitățile de adaptare a RN pentru aplicații de recunoaștere a formelor, a vorbirii, a caracterelor scrise și pentru probleme de control automat. Beneficiind de o coproiectare *hardware - software* eficientă sistemul CNAPS se pretează pentru aplicații flexibile care revendică înaltă programabilitate.

A. Arhitectura CNAPS

În figura 2.14 este prezentată arhitectura sistemului CNAPS care este compus din 2 subsisteme principale: Aria de procesoare CNAPS și Procesorul de Control (CP) [MUE 93]. Aria de procesoare execută algoritmi neuronali pe baza unui program CNAPS. Procesorul de control este un *controller* care controlează operațiile de bază cum ar fi încărcarea programului și interfața de intrare / ieșire cu *hardware*-ul extern. Aceste două subsisteme sunt interconectate prin busul sistemului.

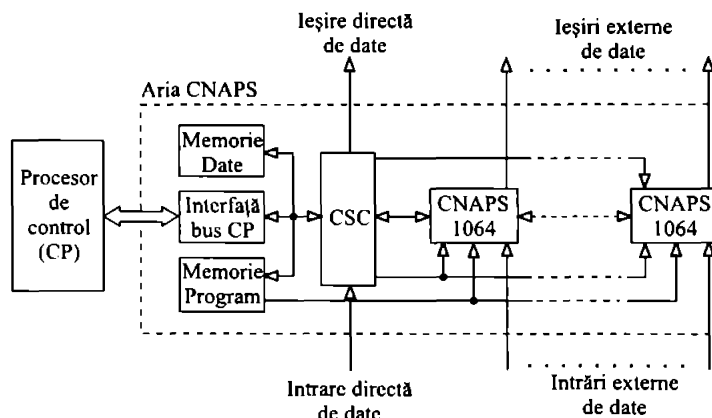


Fig. 2.14. Arhitectura sistemului CNAPS

Aria de procesoare CNAPS conține câteva subsisteme; ea conține mai multe circuite CNAPS - 1064 care procesează concurrent datele numerice. CSC (*CNAPS Sequencer Chip*) controlează operațiile interne sistemului multiprocesor (operațiile dintre *chip*-urile CNAPS - 1064) precum și operațiile aferente celor două module de memorie și interfața cu busul sistemului. Instrucțiunile de 64 biți lungime sunt stocate în memoria program; 32 biți din codul instrucțiunii sunt destinați *chip*-ului CSC. Modulul de interfațare cu busul CP asigură conexiunea între busul sistem și busul intern ariei CNAPS. Memoria de date (memoria DRAM) este utilizată pentru memorarea locală a datelor de intrare și de ieșire. Memoria program (SRAM) este utilizată pentru stocarea programelor care sunt încărcate înainte de startarea execuției procedurilor RN. Datorită modularității și unui control eficient al interfeței aria CNAPS poate fi ușor extinsă dimensional în scopul implementării RN de dimensiuni mari. În figura 2.15 se prezintă configurația unui sistem care conține 4 arii CNAPS. Fiecare arie CNAPS conține 8 *chip*-uri CNAPS - 1064 astfel încât cele 2048 procesoare integrate în acest sistem conduc la o performanță de 40 GCPS (*giga - conexions per second*).

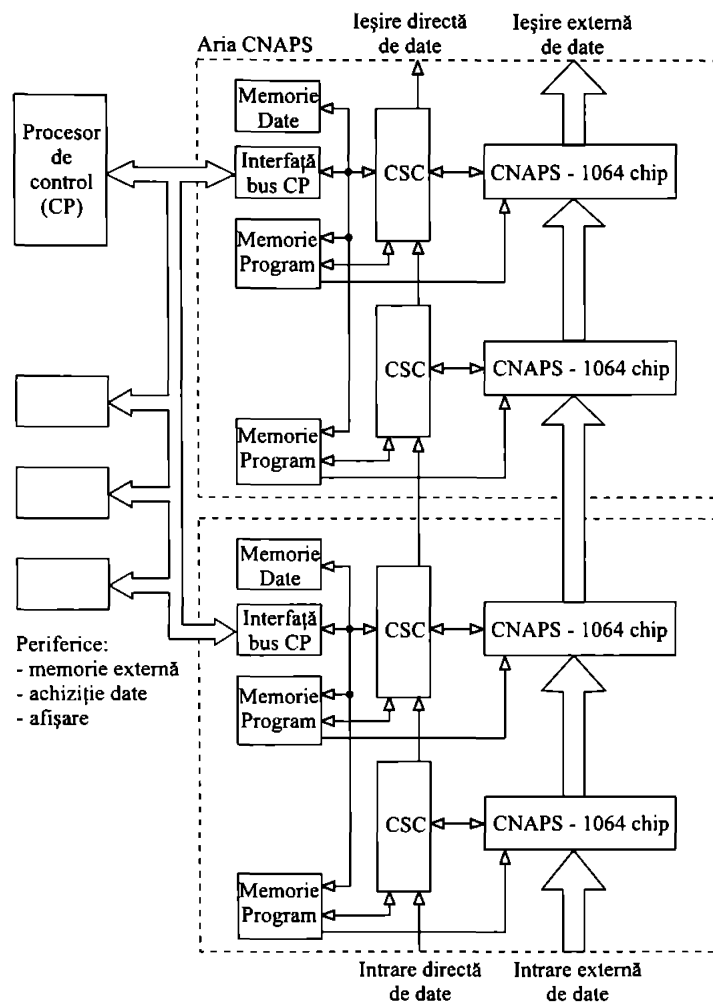


Fig. 2.15. Extensia arhitecturii CNAPS

B. Circuitul CNAPS - 1064

Circuitul CNAPS - 1064 conține o arie de 64 EP. Schema bloc a circuitului precum și structura unui nod EP sunt prezentate în figura 2.16. Procesoarele EP sunt interconectate într-o configurație SIMD. Fiecare EP recepționează date prin intermediul unui bus de 8 biți (IN bus), emite date prin intermediul unui alt bus de 8 biți (OUT bus) și recepționează instrucțiunile prin intermediul unui bus de 32 biți (EPCMD bus). Adicional, procesoarele EP vecine pot comunica direct (date) prin intermediul busului intern de 16 biți (Inter - EP bus).

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

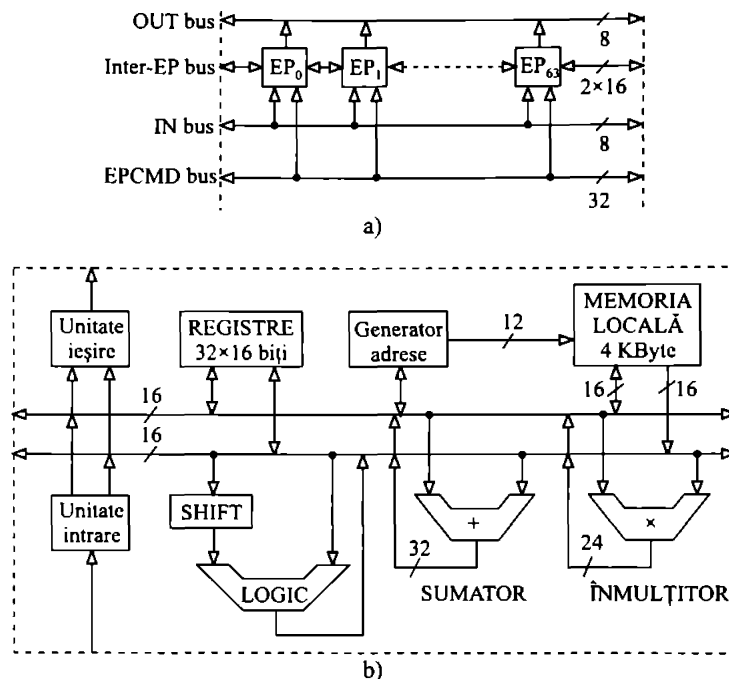


Fig. 2.16. Circuitul CNAPS - 1064. a) Schema bloc; b) Structura EP

Un procesor (EP) conține câteva blocuri. Unitatea de intrare decodifică toate semnalele de control recepționate (pentru elaborarea comenzilor interne) și recepționează date de pe busul IN - bus pe care le memorează temporar. Setul de registre este utilizat pentru memorarea datelor intermediare cum ar fi constanta de învățare și numărul (indexul) procesorului (EP) respectiv. Combinația ÎNMULȚITOR - SUMATOR este utilizată pentru calculul sumelor ponderate proprii algoritmilor conexioniști ($\sum w_{ij} \cdot x_i$). Pentru eficiență maximă se realizează o procesare în tehnică *pipeline*. Ponderile sinaptice sunt memorate în memoria locală. Funcția de transfer neliniară și valoarea de prag pot fi memorate sub formă tabelară (*look-up table*) în memorie. Unitatea de ieșire memorează rezultatele într-un *buffer* de ieșire și conține și circuitul de arbitraj a comunicațiilor globale și / sau inter - EP. Microrutine dedicate controlează atât procesarea la nivelul fiecărui EP cât și configurația întregii rețele de procesoare (EP).

C. Circuitul de secvențiere

Circuitul de secvențiere (CSC) este o unitate de control care controlează operațiile neuronale la nivel *hardware*. CSC controlează 4 stări: IDLE, RUN, SPIN și FLUSH. Diverse semnale de comandă de tip "wait" și "execute" sunt utilizate pentru accesarea memoriilor locale și pentru procesarea condițiilor de întrerupere.

CSC conține 5 module așa cum indică figura 2.17. CPIF (*control processor interface*) controlează interfațarea cu memoria locală.

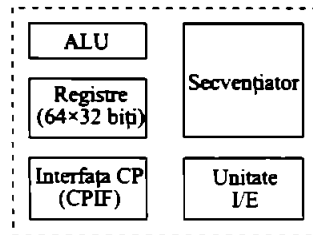


Fig. 2.17. Structura CSC

Semnalele de control emise de CP traversează interfața CPIF cu scopul de a schimba starea CSC. CP poate starta, opri sau înlanțui stările CSC prin intermediul registrelor de stare din CPIF; CPIF are logică proprie de întrerupere. Într-o execuție normală a unui program CNAPS, CPIF supervizează accesul busului local la memoria de date. Secvențiatorul controlează activitatea ariei CNAPS prin execuția instrucțiunilor care specifică secvențierea operațiilor și care controlează operațiile de intrare-ieșire. Conține un "program counter" care poate adresa 64 K instrucțiuni de 64 biți fiecare. Unitatea I/E poate configura diverse fluxuri de date care sunt sincronizate și interacționează cu busurile IN bus și OUT bus.

2.5.2. Sistemul multiprocesor MA 16

Circuitul MA 16 fabricat de Siemens Corporation operează la o frecvență de tact de 50 MHz și conține 4 lanțuri elementare conectate într-o rețea sistolică (RS) [RAM 91a]. MA 16 realizează 800 MCPS, unde o conexiune echivalează cu o operație pe 16 biți. Pentru o RS bidimensională compusă din 16 x 16 circuite MA 16 se estimează un vârf de performanță de 205 GCPS [RAM 91b]. În [RAM 91c] se propun 3 ecuații unificate care acoperă dinamica rețelelor neuronale. Aceste ecuații corespund regulii învățării supervizate:

$$\text{Căutare: } x_i = f_c \left(\sum_{j=1}^N w_{ij} \cdot x_j \right) \cdot \lambda_i \quad (2.4)$$

$$\text{Învățare: } \Delta n = f_c' \cdot \left(\frac{\partial E}{\partial x_n} \cdot o_n + \sum_{i=1}^N w_{in} \cdot \Delta_i \right) \quad (2.5)$$

$$\text{Actualizare: } w_{ij}(t+1) = w_{ij}(t) - \alpha \cdot \left(\sum_{\text{pattern}} \Delta_i \cdot x_j + \frac{\partial E}{\partial w_{ij} \text{ explicit}} \right) \quad (2.6)$$

unde: o_n - funcția caracteristică a neuronului de ieșire

$w_{ij} = 1$

w_{i0} = valoarea de prag

x_i = intrarea neuronului

$x_i = -1$

Δ_i = termen de corecție

În figura 2.18 se prezintă structura unui lanț elementar iar în figura 2.19 se prezintă interconectarea celor 4 lanțuri elementare în cadrul circuitului MA 16. Lanțul elementar din figura 2.18 este capabil să proceseze ecuațiile (2.4), (2.5) și (2.6).

și în consecință vor fi utilizate pentru transferul rezultatelor parțiale. Cele trei faze (*search*, *learning* și *update*) descrise prin ecuațiile (2.4) la (2.6) sunt mapate optimal la nivel *hardware* în operații sistolice așa cum se arată în figura 2.19. Toate operațiile sunt executate în manieră sistolică afectând un bloc de 4 x 4 operanzi.

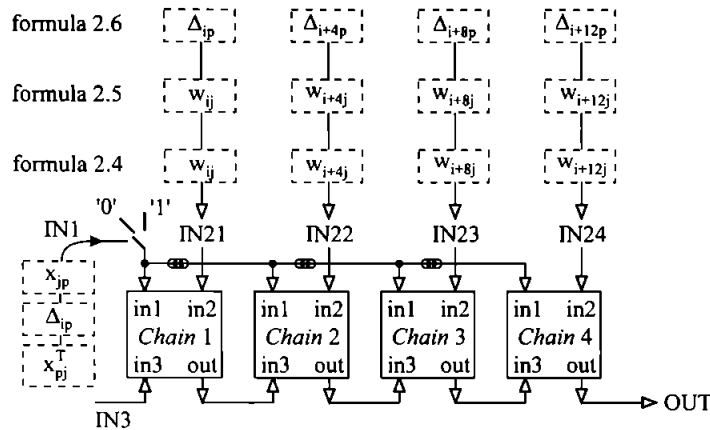


Fig. 2.19. Cele 4 lanțuri elementare utilizate pentru execuția algoritmilor neuronali

Faza de căutare: *Buffer*-ul de intrare (B1) și *buffer*-ul de ponderi sinaptice (B2) recepționează ieșirea neuronului anterior x_j și respectiv ponderile w_{ij} . Deoarece fiecare lanț procesează 4 neuroni ($i = 1, \dots, 4$), vor rezulta 16 sume parțiale care se vor acumula în acumulatorul final. Toate valorile calculate sunt transferate prin lanțul de *buffer*-e B3 așa cum se indică în figura 2.19.

Faza de învățare: *Buffer*-ul de intrare recepționează termenii de eroare Δ_i iar *buffer*-ul B2 recepționează ponderile sinaptice în formă transpusă astfel încât se va calcula suma $\sum_i w_{in} \cdot \Delta_i$

aferețta relației (2.5). Sumarea se poate efectua în formă distribuită conectând acumulatorul final (fig. 2.18) în buclă cu *buffer*-ul B3. Termenul de eroare ($d_n - x_n$) este preîncărcat în B3 iar termenul f'_n este de asemenea încărcat în B4.

Faza de actualizare: *Buffer*-ul de intrare (B1) recepționează intrarea în neuron calculată în faza de căutare iar *buffer*-ul B2 recepționează termenii Δ_i . Constanta de învățare α este preîncărcată în *buffer*-ul B4.

Mai multe circuite MA 16 pot fi configurate fie într-o rețea sistolică liniară fie într-o rețea sistolică bidimensională. MA 16 dispune de circuite de memorie dedicate pentru stocarea ponderilor sinaptice și a valorilor Δ_i .

2.5.3. Circuitul NI 1000

Circuitul digital NI 1000 fabricat de firma Intel este destinat implementării RN de tip RBF (*Radial Basis Function*). Rețelele RBF prezintă anumite avantaje cum ar fi învățarea rapidă și o mai bună estimare, comparativ cu RN multistrat convenționale. O rețea RBF conține 3 straturi neuronale. Primul strat este stratul de intrare iar al doilea strat formează o arie a funcțiilor radiale de bază. Interconexiunile sinaptice dintre primul strat și al doilea strat calculează distanța dintre semnalul (vectorul) aplicat la intrare și semnalul (vectorul) de ponderi memorat. Măsura acestei distanțe D_j poate fi dată de distanța euclidiană dintre cei doi vectori:

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

$$D_j = \sum_i |x_i - w_{ji}| \quad (2.7)$$

unde:

x_i - semnalul de intrare

w_{ji} - ponderea sinaptică memorată

Dacă funcția radială de bază este o funcție simetrică exponențial descrescătoare, ieșirea celui de-al doilea strat va fi:

$$z_j = A_j \cdot e^{-B_j \cdot D_j} \quad (2.8)$$

unde: A_j - factorul de ponderare al funcției exponențiale

B_j - distanța caracteristică a funcției exponențiale

Interconexiunile sinaptice dintre al doilea strat și stratul de ieșire conduc la calculul unei sume ponderate pentru determinarea ieșirii finale:

$$y_k = \sum_j v_{kj} \cdot z_j \quad (2.9)$$

unde v_{kj} reprezintă ponderea sinaptică aferentă conexiunii dintre neuronul j din cel de-al doilea strat și neuronul k din stratul de ieșire.

Circuitul NI 1000 produs de firma INTEL calculează relațiile (2.7) la (2.9) în manieră digitală. Schema bloc a circuitului este prezentată în figura 2.20.

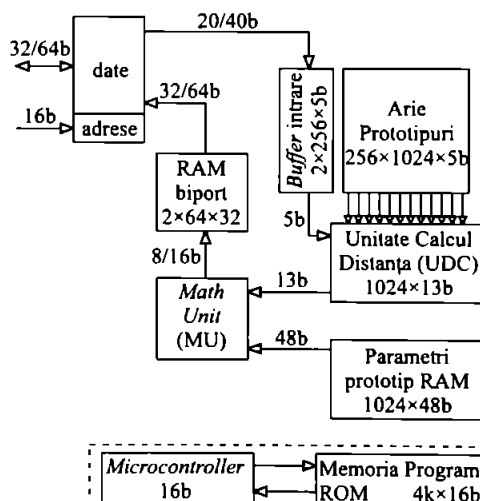


Fig. 2.20. Schema bloc a circuitului NI 1000

Dimensiunile stratului de intrare, stratului RBF și a stratului de ieșire sunt 256, 1024 și respectiv 64 neuroni. Reprezentarea internă a ponderilor sinaptice este realizată pe 5 biți. În unitatea UDC se

calculează distanța dintre semnalul de intrare (memorat în *buffer*-ul de intrare) și semnalul sinaptic (memorat în aria de prototipuri). În unitatea MU se calculează funcția radială de bază și suma ponderată. Pentru performanță procesarea se realizează în tehnică *pipeline* (unitățile UDC și MU formează o cascadă *pipeline* pe 2 nivele). Pentru a mări capabilitățile computaționale mai multe circuite NI 1000 pot fi conectate pe un bus comun. Se poate deci construi un procesor neuronal a cărui structură este prezentată în figura 2.21.

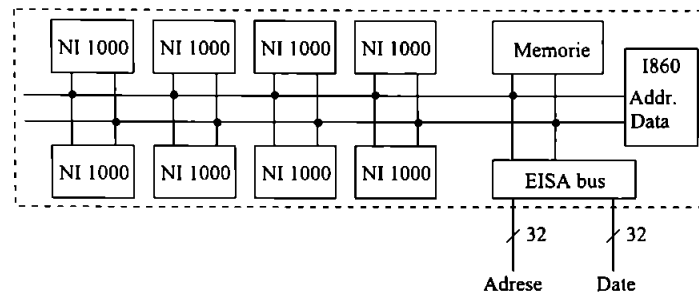


Fig. 2.21. Schema bloc a procesorului neuronal bazat pe circuite NI 1000

2.5.4. Alte neuroprocesoare digitale

Orrey și alții prezintă un neuroprocesor digital performant în [ORR 91]. Circuitul este fabricat în tehnologie CMOS și conține 8 procesoare care operează la 20 MHz. Performanțele estimate sunt de 20 GCPS pentru o RN cu 1000 neuroni. Yasunaga și alții au dezvoltat un neuroprocesor care conține 576 neuroni și 36864 sinapse [YAS 90]. Reprezentarea internă a datelor era realizată pe 8 - 9 biți. Un *chip* neuronal digital de mare performanță a fost implementat de Uchimura și alții [UCH 92]. Circuitul conține 832 sinapse și 13 neuroni având performanțele evaluate la 8 GCPS; datele sunt reprezentate intern pe 8 biți. Watanabe și alții au proiectat un *chip* digital care conține o RN cu 10^6 sinapse [WAT 93]. Celule DRAM realizate cu un singur tranzistor au fost utilizate pentru a memora 10^6 ponderi sinaptice într-o reprezentare întregă pe 8 biți. Un *chip* VLSI digital numit *TinMANN* a fost descris de Melton și alții în [MEL 92]. Circuitul se bazează pe o arhitectură stohastică și un paralelism intern masiv. Un SAU global cu 2048 intrări (neuroni) poate fi evaluat într-un singur ciclu de tact de 67 nsec. Procesorul SPERT care lucrează la 50 MHz a fost proiectat de Wawrzynek și alții [WAW 93]. Problematika învățării "on-chip" în neuroprocesoarele digitale a fost analizată în [LEH 93]. Alte implementări VLSI digitale (neuroprocesoare) sunt descrise în [SHE 95] și [SAN 92].

2.6. Principii de proiectare a neuroprocesoarelor sistolice

Pentru a exploata paralelismul intrinsec caracteristic RN, acestea trebuie implementate pe arhitecturi multiprocesor. Arhitectura care realizează comunicații eficiente între procesoare va furniza și o mare putere de calcul pe unitatea de suprafață. Implementările sistolice digitale dedicate rețelelor neuronale pot opera cu eficiență deosebită datorită procesării în tehnică *pipeline* a datelor. Într-o rețea sistolică (RS) pot coexista mai multe fluxuri de date *pipeline*-izate. Arhitectura multi-EP în inel și respectiv multi-EP matricială [PRZ 91] sunt 2 arhitecturi sistolice reprezentative. Comparativ cu circuitele neuronale analogice, arhitecturile sistolice oferă o mai bună flexibilitate, o mai mare precizie și o completă *pipeline*-izare a fluxurilor de date.

2.6.1. Arhitectura sistemului

Schema bloc a unui sistem bazat pe o arhitectură multi-EP în inel și respectiv multi-EP matricială este prezentată în figura 2.22.

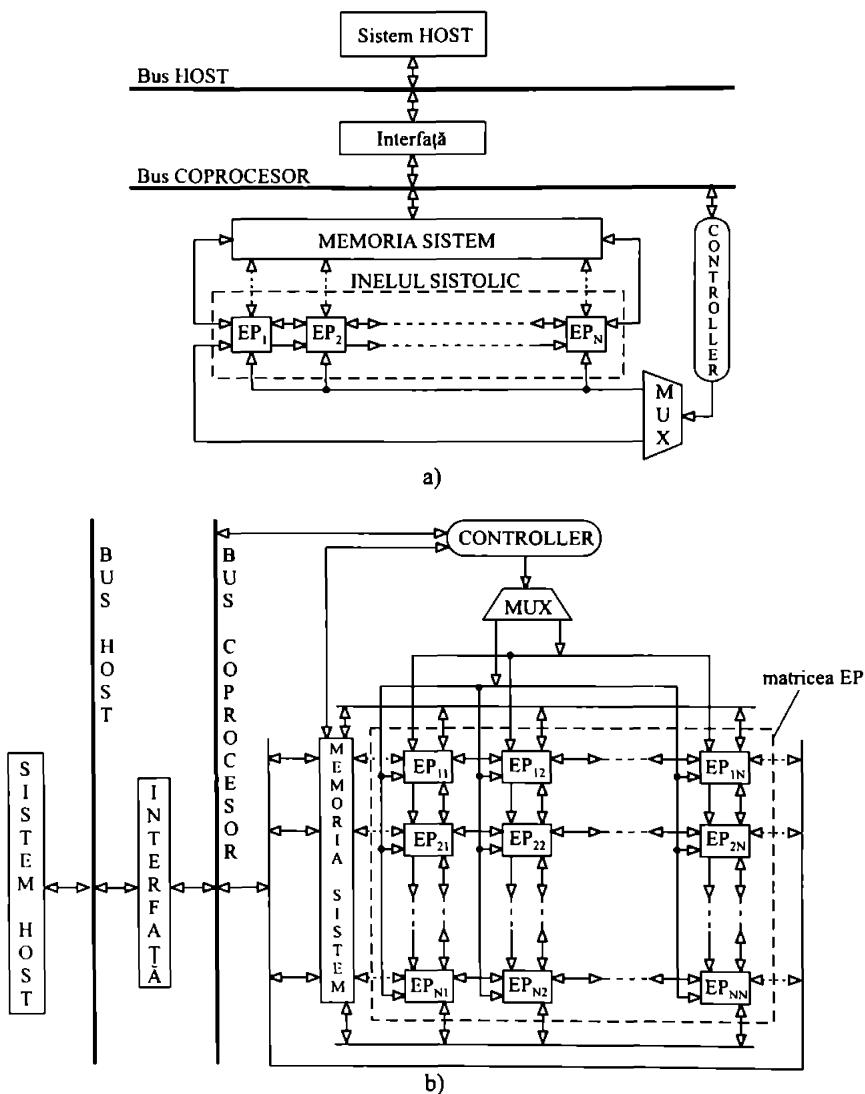


Fig. 2.22. Arhitectura sistemului bazat pe rețele sistolice
 a) RS monodimensională conectată în inel; b) RS bidimensională pătratică

Controller-ul RS poate fi cablat sau microprogramat. În varianta microprogramată acesta va emite o microinstrucțiune pe fiecare ciclu de tact. Microcodul și adresele de acces la memoria sistem pot fi

transmise procesoarelor EP din cadrul RS fie în sistem *broadcasting* fie în sistem *pipeline*; tipul transmisiei este selectat cu ajutorul multiplexorului din figura 2.22. Pentru RS conectată în inel, operațiile I/E dintre RS și memoria sistem vor fi suportate doar de EP-urile extreme (EP_1 și EP_N) și în acest fel memoria sistem biport va fi partajată de procesoarele inelului sistolic. În ciclul de inițializare fiecare EP citește datele locale aferente din memoria sistem prin intermediul busurilor locale indicate în figura 2.22 prin linii punctate. În cazul RS bidimensionale, fiecare EP_{ij} comunică cu cei patru vecini ai săi iar operațiile I / E cu memoria sistem apar pe cele patru laturi ale matricii sistolice. În consecință, memoria sistem va fi o memorie cu 4 porturi de acces partajată de matricea EP. Similar, în ciclul de inițializare, datele locale aferente EP-urilor din matricea sistolică se obțin prin intermediul busurilor locale marcate cu linii punctate și situate pe laturile de est și respectiv de vest aferente matricii sistolice. Un EP din zona de mijloc a unei linii recepționează datele locale proprii printr-o comunicație *pipeline*.

Memoria sistem recepționează date de la sistemul gazdă (*host*) în urma activării unei astfel de cereri; cererea este activată de sistemul gazdă și recepționată de *controller*. Sistemul gazdă servește ca interfață între utilizator și RS. Acesta va furniza informațiile specifice aplicației, cum ar fi vectorii de intrare, ponderile sinaptice inițiale, parametrii de control convergență etc. *Controller*-ul specifică și monitorizează procesarea în cadrul fiecărui EP și guvernează, de asemenea, și comunicațiile de date dintre EP. *Controller*-ul va avea memoria sa dedicată unde va stoca microcodul recepționat de la sistemul gazdă.

2.6.2. Proiectarea elementului de procesare (EP)

Implementările digitale dedicate RN au la bază 4 principii:

- implementări bazate pe microprocesoare de uz general
- implementări bazate pe procesoare de semnal (DSP - *digital signal processor*)
- implementări bazate pe neuroprocesoare de uz general
- implementări bazate pe neuroprocesoare dedicate

În primele 2 cazuri sistemul va avea la bază procesoare existente pe piață. În al treilea caz se utilizează un *hardware* de uz general în ideea realizării unei neuroprocesări flexibile. Acest *hardware* va putea fi utilizat în diverse aplicații de procesare de semnal sau de imagine. În al patrulea caz se construiește un *hardware* dedicat câtorva algoritmi neuronali. În tabelul 2.2 se prezintă o comparație între cele 4 principii de implementare.

Proprietăți	Implementări bazate pe microprocesoare de uz general	Implementări bazate pe DSP	Neuroprocesoare de uz general	Procesoare dedicate
Paralelism	Cel mai slab	Slab	Înalt	Cel mai înalt
Flexibilitate	Cea mai înaltă	Slabă	Înaltă	Cea mai slabă
Precizie	Cea mai înaltă	Înaltă	Slabă	Cea mai slabă
Programabilitate	Cea mai înaltă	Slabă	Înaltă	Cea mai slabă
Integrare / Scalabilitate	Slabă	Înaltă	Cea mai înaltă	Cea mai slabă
Viteză de procesare	Cea mai joasă	Joasă	Înaltă	Cea mai înaltă
Durată ciclu proiectare	Cel mai scurt	Scurt	Lung	Cel mai lung
Costul de producție	Cel mai înalt	Înalt	Scăzut	Cel mai scăzut
Aplicabilitate	De uz general	Procesare de semnal	Diverse	Foarte limitate

Tabelul 2.2. Comparația celor 4 tipuri de implementări digitale dedicate RN

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

Comparând datele prezentate în tabelul 2.2 putem concluziona că neuroprocesoarele de uz general realizează cel mai bun raport performanțe / cost. O propunere de structură internă pentru nodul sistolic EP este prezentată în figura 2.23 [CHA 94]. Pentru memorarea locală a datelor fiecare EP conține un *cache* de date. Multiplicatorul și sumatorul reprezintă principalele unități de procesare din structura EP.

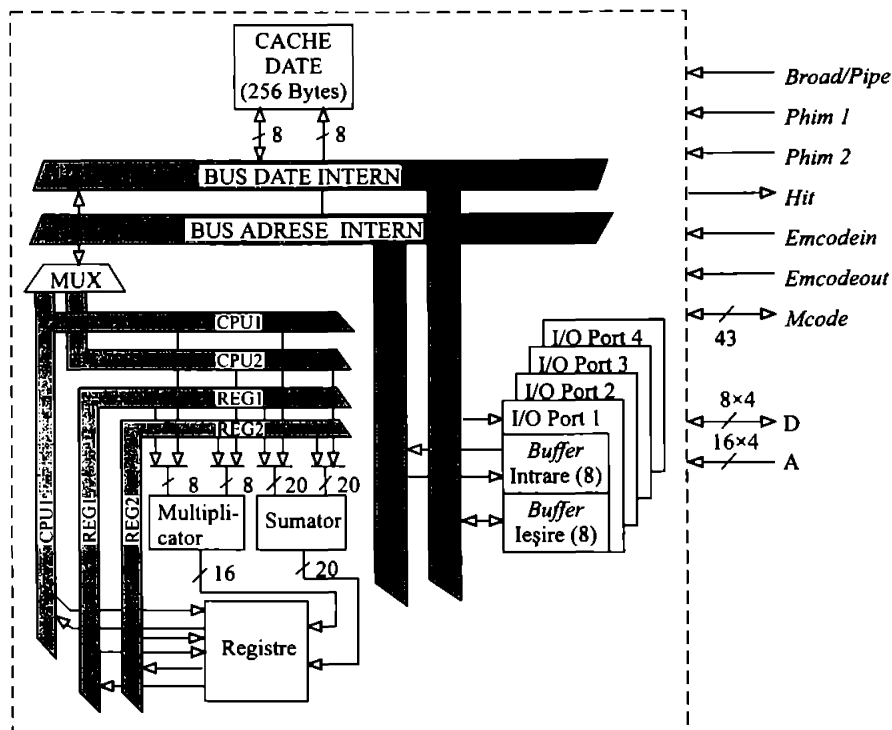


Fig. 2.23. Schema bloc aferentă modulului EP

Structura EP este adaptată RS bidimensionale din figura 2.22. b). Fiecare EP conține 4 porturi I/O pentru a comunica date cu cele 4 EP vecine. O unitate MMU (*memory management unit*) stabilește strategiile de alocare a memoriei precum și strategiile de *fetch* și actualizare a memoriei principale (în raport cu *cache*-urile locale). MMU este locat în *controller* și este partajat de toate unitățile EP din cadrul RS. RS operează sincron, cu un tact global.

O tabelă *look - up* va fi utilizată pentru implementarea funcției de activare aferentă neuronului. Rutinele de nivel înalt executate de sistemul gazdă sunt compilate în macroinstrucțiuni care sunt transmise apoi *controller*-ului. *Controller*-ul decodifică macroinstrucțiunile și le translatează în microinstrucțiuni pe care le emite spre EP-urile din structura RS (*broadcast*) prin intermediul celor 43 linii ale busului *Mcode*. Pentru o *pipeline*-izare cât mai eficientă a instrucțiunilor este recomandată utilizarea unui set de instrucțiuni specific RISC (*reduced - instruction set computer*).

Tot în scopul creșterii vitezei de procesare, registrele EP sunt accesibile prin intermediul a 4 busuri de date: REG1, REG2, CPU1 și CPU2. Drept consecință multe operații interne care implică

registrele EP vor putea fi suprapuse deoarece nu implică ocuparea aceluiași busuri. Un multiplexor este utilizat pentru multiplexarea celor 4 busuri pe busul de date intern. Pentru transferul datelor între *cache*-ul de date și setul de registre este utilizat busul CPU1.

În funcție de modul de adresare utilizat, atât pentru multiplicator cât și pentru sumator, unul dintre operanzi este preluat fie de pe busul CPU1 fie de pe busul REG1, iar celălalt operand este preluat fie de pe busul CPU2 fie de pe REG2. Circuitul a fost proiectat cu ajutorul editorului *Magic* elaborat la Universitatea Berkeley. Performanțele circuitului pot fi îmbunătățite dacă în procesul de proiectare s-ar utiliza pachete CAD comerciale.

3. ANALIZA IMPLEMENTĂRILOR SISTOLICE

În acest capitol vom face o prezentare a metodelor și descrierilor sistolice și vom arăta că aceste metode permit localizarea spațială a algoritmilor conexioniști. Această abordare va permite obținerea unor algoritmi sistolici capabili să execute eficient operațiile aferente modelelor conexioniste actuale. Vom propune o arhitectură sistolică generală din care vor deriva 3 arhitecturi exploatabile. În sfârșit, vom alege arhitectura optimă și vom motiva această alegere.

3.1. Principiul procesării sistolice

Principiul procesării sistolice a fost introdus de *H. T. Kung* și *C. E. Leiserson* în anul 1978 [KUN 78]. Termenul "sistolic" provine dintr-o analogie între circulația fluxului de date în rețeaua sistolică (RS) și cea a sângelui în aparatul circulator aferent organismelor vii. Inima este asimilată elementelor de procesare (EP) din componența RS, care recepționează, procesează și transmit datele. RS reprezintă instrumentul ideal pentru realizarea unei anumite clase de procesoare specializate. Dificil de implementat cu componente discrete din cauza unui cost ridicat, algoritmi sistolici intră în actualitate grație tehnologiilor VLSI (*Very Large Scale Integration*).

Structura în rețea reprezintă suportul tipic pentru algoritmi sistolici. Celulele sistolice (CS) sunt procesoare elementare identice, conectate exclusiv cu celulele vecine. Fluxurile de date traversează rețeaua prin comunicații locale și într-o manieră sincronă. Organizarea în rețea (fig. 3.1) și funcționarea sincronă permit procesarea concurrentă a informațiilor în celulele sistolice conducând la o putere de calcul foarte ridicată. În sfârșit, vom vedea că în anumite condiții, procesarea în tehnică *pipeline* a datelor devine posibilă, conducând la o exploatare maximală a RS.

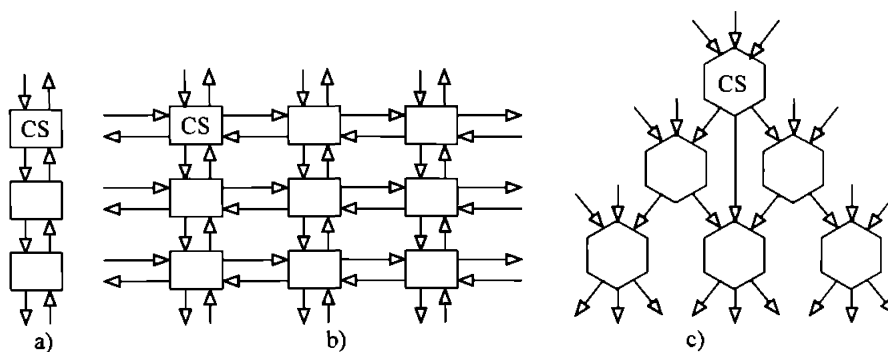


Fig. 3.1. Topologii de RS. a) Liniară; b) Pătrată; c) Hexagonală (CS - celula sistolică referită adesea și cu termenul EP - element de procesare)

Formal, o rețea sistolică este o structură de calcul care posedă următoarele caracteristici:

- Funcționare sincronă;
- Modularitate și regularitate. RS este un ansamblu de unități de calcul (CS) asemănătoare sau identice, conectate prin scheme de interconexiuni omogene, conferind o extensibilitate teoretic infinită;
- Localizare spațială și temporală;
- Posibilitatea procesării datelor în tehnică *pipeline*.

Clasa de aplicații în care RS sunt extrem de eficiente cuprinde problemele în care volumul de calcul de efectuat este mult mai mare decât cel al transferurilor de date de realizat [COS 89].

Fiecărei topologii de RS i se asociază o clasă de probleme [KUN 88], [QUI 89]. Rețelele liniare sunt utilizate pentru operații de convoluție, înmulțiri matrici - vectori, produse scalare etc. Rețelele rectangulare sunt performante în cazul produselor matriciale și, în general, pentru aplicații de filtrare pe imagini (operații cu pixeli). Rețelele triunghiulare sau hexagonale sunt performante în operațiile de decompoziție LU a matricilor (descompunerea unei matrici în 2 matrici triunghiulare, triunghiulară inferioară L și triunghiulară superioară U) și în operațiile de înmulțire a matricilor cu structuri particulare (triunghiulare, bandă diagonală).

Toate operațiile care trebuie implementate pe o RS depind mai întâi de topologia rețelei. Datele tranzitează rețeaua traversând legăturile definite între celulele vecine. Prin urmare, implementarea sistolică a unui algoritm nu este unică ceea ce înseamnă o sarcină relativ complexă pentru proiectant. RS optimă (pentru o anumită aplicație) va fi aleasă în funcție de performanțe (viteza de procesare), complexitate cât mai redusă, posibilități de optimizare (*pipeline*-izare, suplețe în execuția operațiilor de intrare / ieșire etc.). După cum vom vedea, această etapă de definire a topologiei RS este crucială deoarece ea stabilește definitiv capabilitățile viitoarei RS.

Există actualmente mai multe instrumente de concepție a RS care au fost dezvoltate și puse la dispoziția proiectanților. Metodologia lui *Leiserson* și *Saxe* se bazează pe conversia unui sistem sincron reprezentat printr-un graf în algoritm sistolic. Metoda lui *Quinton* permite o sinteză automată a RS, pornind de la o formulare recurentă a calculelor de executat [QUI 84]. Metodologia lui *Li* și *Wah* reprezintă un principiu general pentru optimizarea concepției unei RS planare [LI 85]. *Cosnard* și *Tchunte* au propus, de asemenea, o metodologie bazată pe o analiză descendentă care divizează problema proiecției unui graf de dependențe în 2 etape: construcția unei rețele liniare ale cărei intrări sunt vectori, apoi sistolizarea acestei transformări [COS 89].

În cele ce urmează, vom utiliza metoda lui *Kung* [KUN 88] pentru descrierea RS și vom valida formal funcționarea acesteia pe baza metodei lui *Quinton* [QUI 84]. Metoda lui *Kung* se bazează pe manipulări sistematice de grafuri și ale proiecțiilor acestora; este perfect adaptată pentru calcule simple, cum ar fi produsul matrici - vector, care stă la baza algoritmilor conexioniști.

Această metodologie de descriere a operațiilor permite extragerea soluției de maxim paralelism; este o descriere sub formă de graf de dependență (DG - *Dependence Graph*) [KUN 88] care nu face apel la nici o noțiune arhitecturală, de rebusclaj sau de optimizare temporală. În cadrul DG se iau în considerare numai dependențele dintre nodurile acestuia. Mai aproape de arhitectură, graful fluxului de date (SFG - *Signal Flow Graph*) este o formulare mai completă care pune în evidență relațiile dintre date, întârzierile (*delays*) și eventualele rebusclaje. Acesta definește de asemenea structura suport aferentă algoritmului.

Trecerea de la un DG la un SFG pune următoarele probleme:

- cum vor fi atribuite (repartizate) operațiile pe procesoare;
- în ce ordine vor trebui executate aceste operații.

Răspunsurile la aceste 2 probleme vor defini asignarea procesoarelor și secvențierea operațiilor. Metodologia pe care o vom descrie este cunoscută sub denumirea de mapare canonică [KUN 88].

3.1.1. Graful de dependență (DG)

DG este un graf care descrie dependențele dintre calculele ce compun un algoritm; în nodurile grafului nu vor fi prezentate calculele efectuate.

Pentru a ilustra această definiție, vom construi DG aferent operației de înmulțire a unei matrici cu un vector.

Vom considera operația:

$$y = W \cdot x$$

unde:

- x - vector de dimensiune [N]
- y - vector de dimensiune [P]
- W - matrice [P, N]

Sub formă indicială, ecuația matricială devine:

$$y_i = \sum_{j=1}^N w_{ij} \cdot x_j \quad \text{cu } i = 1, 2, \dots, P \quad (3.1)$$

Indecșii de această formă definesc doar localizarea spațială a elementelor matricii și a celor 2 vectori implicați în operație. Introducând un index temporal k, vom putea rescrie algoritmul (3.1) sub o formă în care fiecare variabilă nu este modificată decât o singură dată (*simple assignment form*).

$$y_i^k = y_i^{k-1} + w_i^k \cdot x^k \quad (3.2)$$

cu: $y_i^0 = 0$, $y_i^N = y_i$ și $w_i^k = w_{ik}$

Această formulare este o ecuație recurentă cu un index spațial (i) și un index temporal (k). Pe baza ecuației (3.2) vom putea construi un graf de dependență, de forma celui prezentat în figura 3.2, unde am considerat vectorul x de dimensiune 4, vectorul y de dimensiune 3 și matricea W de dimensiune [3, 4].

În general, datele de intrare sunt difuzate spre toate nodurile grafului. Această operație revendică un bus de comunicație global care va avea repercusiuni negative asupra realizării fizice.

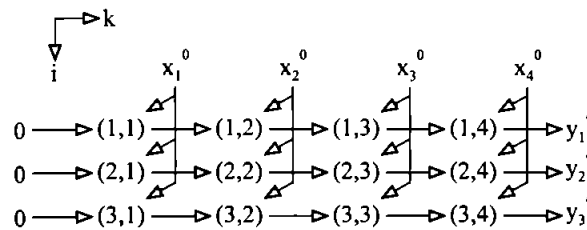


Fig. 3.2. Graful de dependență aferent produsului matrice - vector. În acest graf, datele de intrare sunt distribuite tuturor nodurilor. Arcele nu indică decât o dependență funcțională între noduri și nu induc nici o întârziere (*delay*)

Pentru implementările fizice este interesant de a transforma această difuzie de date spre toate nodurile cu o transmisie de tipul din aproape în aproape. În acest sens, algoritmi trebuie localizați în scopul obținerii unui graf de dependență local (figura 3.3). Există mai multe studii teoretice dezvoltate pe problema localizării algoritmilor pentru a ajunge sub forma unui Sistem de Ecuații Recurente Uniforme (SERU) [QUI 84], dar o metodă generală încă n-a fost adoptată.

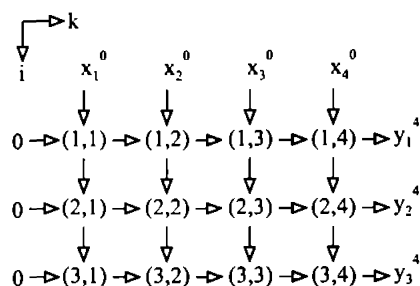


Fig. 3.3. Graful de dependență localizat aferent produsului matrice - vector. Un algoritm este localizat dacă toate variabilele nu depind decât de variabile din noduri vecine

Datele proprii unui nod din graf care nu sunt niciodată deplasate sunt denumite date rezidente. O dată poate fi difuzată (spre mai multe noduri simultan), emisă (de la un nod la altul) și rezidentă.

În DG, operațiile efectuate asupra datelor nu sunt descrise. Sunt descrise doar relațiile dintre noduri care definesc dependențe funcționale. Ele sunt materializate prin arcele de dependență care unesc nodurile grafului.

Totuși, așa cum am subliniat deja, pentru a defini o arhitectură de calcul pornind de la DG, trebuie repartizate (alocate) operațiile pe procesoarele rețelei și trebuie definită o secvențiere. Această operație poate fi realizată pornind de la graful de dependențe localizat. Ea va conduce la un graf de mișcare a datelor (SFG - *Signal Flow Graph*) prin efectuarea unor operații de proiecție spațială și temporală.

3.1.2. Graful de mișcare a datelor (SFG)

Un SFG conține informații funcționale și structurale. Informațiile funcționale descriu operațiile efectuate în noduri. Se consideră că aceste operații nu generează întârzieri. Informațiile structurale descriu legăturile în termenii conexiunilor și întârzierilor. Se separă deci în mod explicit calculul și timpul. Un DG se transformă în SFG prin 2 operații:

- o proiecție spațială a cărei direcție este indicată de un vector de proiecție spațială \mathbf{d} . Aceasta corespunde etapei de asignare a nodurilor la procesoarele rețelei (*processors assignment*);
- o schemă de secvențiere a operațiilor (*scheduling*) reprezentată printr-un vector de secvențiere \mathbf{s} care poartă pe o direcție normală liniilor izocrone. Aceste linii trec prin toate nodurile grafului care efectuează același calcul simultan (dar pe date diferite).

Vectorul \mathbf{s} trebuie să satisfacă 2 condiții:

- toate arcele de dependență trebuie să poarte în aceeași direcție și anume în lungul liniilor izocrone (cauzalitate);
- liniile izocrone nu sunt paralele cu direcția de proiecție spațială \mathbf{d} (identitate de asignare temporală).

Pentru un DG dat căruia îi asociem o direcție de proiecție temporală, există 4 clase de scheme de secvențiere posibile [KUN 88]:

- secvențierea normală: vectorii de proiecție spațială și de secvențiere sunt coliniari;

- secvențierea recursivă: vectorul de secvențiere este coliniar cu vectorul de index spațial al DG;
- secvențierea sistolică: trebuie să existe un *delay* pe fiecare arc al SFG;
- secvențierea optimizată: se obține în urma unei optimizări care ține cont de dependențele specifice, de performanțele procesoarelor și, în general, de tot ceea ce ar putea îmbunătăți concepția.

În cele ce urmează vom insista asupra secvențierii sistolice. Aceasta ne va permite derivarea unor arhitecturi de rețea adaptate implementărilor în circuite VLSI digitale a operațiilor specifice algoritmilor conexioniști. Vom studia mai multe tipuri de proiecție spațială $\mathbf{d} = [i, k]^T$ și temporală $\mathbf{s} = [i, k]^T$ cu scopul de a ilustra definițiile precedente și vom deduce patru SFG cu secvențiere sistolică. Reluând graful DG aferent produsului matrice - vector, putem deriva trei SFG (fig. 3.4) alegând direcții de proiecție coliniare, pointând în cele 3 direcții permise de o secvențiere liniară. Toate grafurile SFG rezultate sunt monodimensionale.

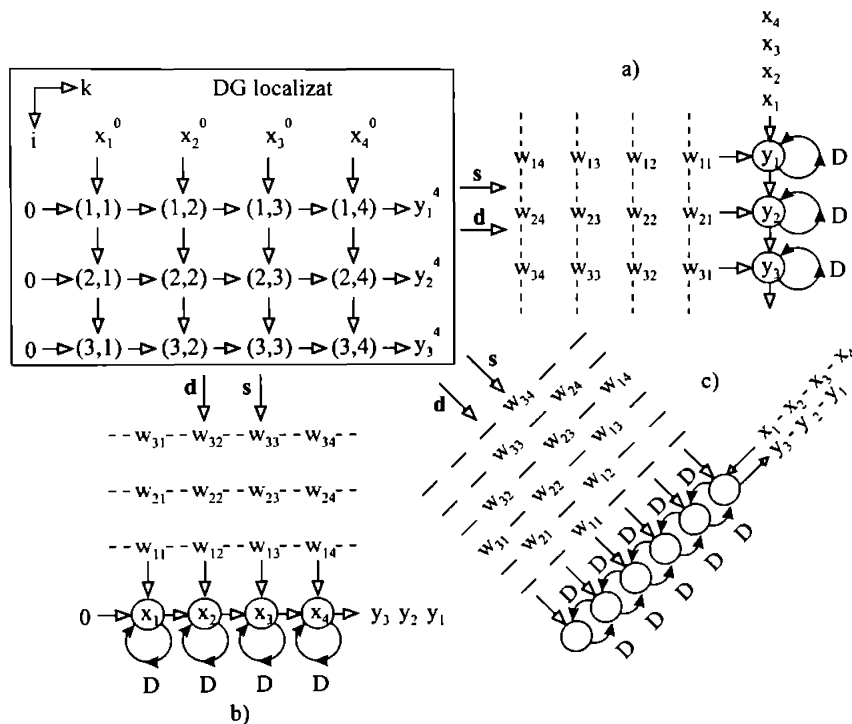


Fig. 3.4. Trei SFG derivate din graful DG localizat aferent produsului matrice - vector

Cele 3 soluții corespund următoarelor proiecții:

- soluția 1 (fig. 3.4. a):

$$\mathbf{d} = \mathbf{s} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

Matricea W este prezentată în linie cu linie grafului SFG. Am obținut secvențierea normală care nu este sistolică deoarece nu există *delay*-uri între noduri. Componentele vectorului x sunt introduse succesiv dar la introducerea fiecăre se distribuie tuturor nodurilor rețelei. Elementele matricii W situate pe cele 4 linii izocrone evidențiate sunt introduse succesiv și sincron cu cele 4 elemente ale vectorului x . Fiecare nod calculează produsul parțial între componenta vectorului x și elementul matricii W corespondent. De asemenea fiecare nod însumează produsele parțiale succesiv elaborate rezultând astfel cele 3 componente aferente vectorului y . Produsul parțial obținut în pasul curent este retrimis în același nod prin arcul de rebusaj cărui îi este asociat *delay*-ul D . În fiecare pas nodul realizează produsul componentelor de pe intrări și adună la acest produs produsul parțial anterior. După 4 pași succesivi componentele vectorului y vor fi acumulate în cele 3 noduri, urmând a fi extrase.

- soluția 2 (fig. 3.4. c))

$$d = s = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

Să remarcăm secvențierea sistolică realizată de acest SFG datorată faptului că pe fiecare arc există un *delay*. Componentele vectorului y circulă în sens opus componentelor vectorului x . Componentele vectorului x circulă fără a fi modificate. În fiecare nod al grafului, elementele matricii W vor intra la intervale constante de 2 impulsuri, așa cum indică cele 6 linii izocrone. Dacă vectorul înmulțitor (x) are N componente, vor fi necesare $N - 1$ impulsuri pentru introducerea primei sale componente (x_1) înainte ca procesul de calcul să înceapă efectiv (în fig. 3.4. c) sunt necesare 3 impulsuri pentru ca x_1 să ajungă în nodul care va efectua primul produs parțial: $w_{11} \cdot x_1$). Pentru o matrice $[P, N]$, vor fi necesare alte $P - 1$ impulsuri din momentul execuției ultimului produs ($w_{34} \cdot x_4$) până la extragerea ultimei componente aferentă vectorului y . În medie, fiecare nod nu va fi activ decât pe durata unui impuls din două.

- soluția 3 (fig. 3.4. b))

$$d = s = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

Este tot o secvențiere standard, nesistolică datorită inexistenței *delay*-urilor pe arcele dintre nodurile grafului. Matricea W este aplicată grafului SFG pe coloane. Componentele vectorului x sunt introduse vertical și simultan în nodurile SFG, înainte de aplicarea elementelor matricii W . Liniile izocrone sunt identice cu liniile matricii W ; elementele unei linii sunt aplicate simultan nodurilor SFG, se elaborează produsele și se obține componenta corespondentă a vectorului rezultat y (nu există *delay*-uri pe arcele dintre noduri). Să notăm că elementele vectorului x sunt rebusate în nodurile SFG (*delay*) pentru a fi "prezente la întâlnirea" cu următoarea linie a matricii W .

Dintre cele 3 SFG prezentate, unul singur are o secvențiere sistolică (soluția 2). Alegând convenabil vectorii d și s , și soluțiile 1 și 3 pot fi transformate în SFG cu secvențiere sistolică.

Vom atribui vectorului s o valoare care implică un *delay* pe fiecare arc din SFG:

$$s = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

și vom trece din nou în revistă cele trei direcții de proiecție spațială prezentată anterior. Vom obține 3 SFG cu secvențiere sistolică descrise în figura 3.5, cu o perioadă *pipeline* $\alpha = 1$ pentru soluțiile 1 și 2 (fig. 3.5. a) și 3.5. b)), și $\alpha = 2$ pentru soluția 3 (fig. 3.5. c)).

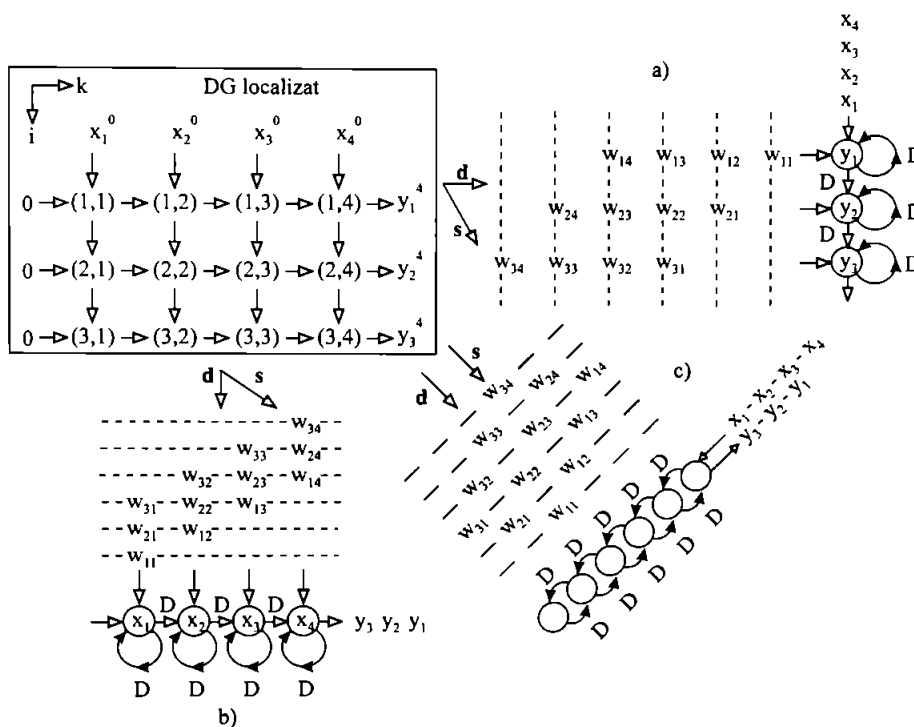


Fig. 3.5. Rețele sistolice (RS) pentru produsul matrice - vector. Liniile izocrone definesc nodurile care procesează simultan și vectorul d este perpendicular pe aceste linii. Această implementare impune un decalaj în aplicarea elementelor matricii W

3.1.3. Prima soluție sistolică

Această soluție (fig. 3.5. a)) este rezultatul unei proiecții după vectorii:

$$d = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad \text{și} \quad s = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

SFG-ul rezultat este compus din P noduri. În termenii transferurilor de date, această soluție presupune emisia vectorului x și a matricii W în timp ce vectorul y este rezident. Un produs parțial

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

y_i^k ($k < N$) se calculează în fiecare nod. Produsul parțial anterior este acumulat (arcu de rebuclaj) pe fiecare impuls produsului curent, pentru a deveni după N impulsuri, componenta corespondentă a vectorului rezultat y . Calculul se descompune în 2 faze:

1. Traversarea grafului SFG de către matricea sinaptică \mathbf{W} și vectorul \mathbf{x} cu acumularea sumelor parțiale y_i^k cu $k < N$.
2. Extragerea vectorului \mathbf{y} în P impulsuri.

Timpul total de calcul va fi deci de $(N + P)$ impulsuri. Această soluție se pretează la o implementare fizică și, după cum vom vedea, a și stat la baza realizării unei rețele conexiuniste.

3.1.4. A doua soluție sistolică

Al doilea SFG monodimensional cu secvențiere sistolică este dualul celui precedent, cu proiecțiile spațială și temporală realizate după vectorii:

$$\mathbf{d} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad \text{și} \quad \mathbf{s} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

Se obține SFG-ul din figura 3.5. b), compus din N noduri. În cadrul acestei variante, componentele vectorului \mathbf{y} și ale matricii \mathbf{W} sunt emise; componentele vectorului \mathbf{x} sunt rezidente în fiecare nod și rămân localizate grație arcelor de rebuclaj. Produsele parțiale y_i^k sunt transmise între noduri și acumulate în modul destinație. Componentele vectorului \mathbf{x} sunt rezidente (grație rebuclajului) și reutilizate pe fiecare impuls pentru elaborarea produsului local (al doilea termen al produsului va fi componenta matricii \mathbf{W} care intră în nod pe impulsul respectiv). Calculul se descompune în 2 faze:

1. Introducerea vectorului \mathbf{x} în N impulsuri;
2. Traversarea grafului SFG de către matricea sinaptică \mathbf{W} (suprapusă cu tranzitul sumelor parțiale y_i^k , cu $k < P$) în P impulsuri.

Timpul total de calcul va fi deci $(N + P)$ impulsuri.

3.1.5. A treia soluție sistolică

Al treilea SFG (fig. 3.5. c)) corespunde unei proiecții după următorii vectori:

$$\mathbf{d} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \text{și} \quad \mathbf{s} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

Particularitatea acestei soluții constă în faptul că nodurile SFG nu conțin date rezidente. Graful SFG este compus din $(N + P - 1)$ noduri. Componentele vectorilor \mathbf{x} și \mathbf{y} circulă în direcții opuse. Componentele matricii \mathbf{W} traversează rețeaua pe o direcție ortogonală acestei circulații. Calculul se derulează în 3 faze:

1. Introducerea vectorului \mathbf{x} pe plasamentul corespunzător inițierii calculului (N impulsuri);
2. Calculul propriu-zis care durează $N + P - 1$ impulsuri pentru fiecare componentă;
3. Extragerea ultimei componente din vectorul \mathbf{y} ($P - 1$ impulsuri).

Timpul total de calcul va fi deci $2(N + P - 1)$ impulsuri. Acest SFG nu este avantajos, nici în spațiu și nici în timp, pentru calcule care implică matrici W complet populate (cu elemente nenule). Soluția devine interesantă pentru matricile cu o structură de tip bandă diagonală, datorită proiecției pe diagonală.

3.1.6. A patra soluție sistolică

Respectând ipoteza unei secvențieri liniare mai rămâne o soluție pe care o vom detalia pe larg în cele ce urmează. Aceasta corespunde unei asocieri directe între DG și SFG, ceea ce înseamnă că vectorii de proiecție vor fi:

$$\mathbf{d} = \mathbf{s} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Vectorul x este introdus din direcția Nord în rețea și va traversa celulele sistolice dispuse într-o grilă de formă rectangulară (fig. 3.6).

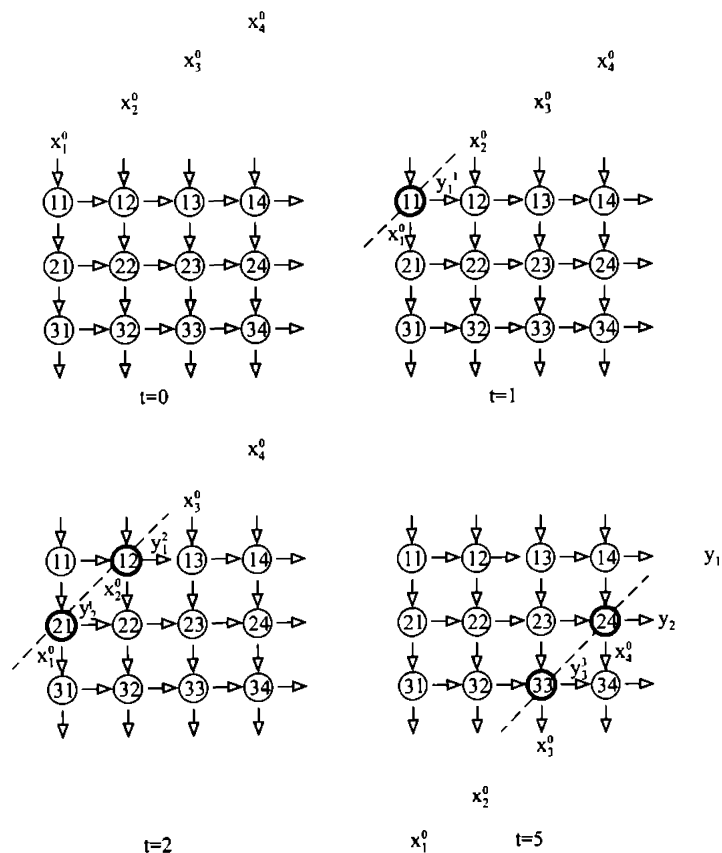


Fig. 3.6. Rețea sistolică (RS) bidimensională

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

Sumele parțiale tranzitează arcele dintre celulele dispuse pe aceeași linie și ies din rețea (sub formă de sume finale care reprezintă componentele vectorului y) pe direcția Est. Considerăm același caz general în care matricea W are dimensiunea $[P, N]$, vectorul x are dimensiunea N și vectorul y dimensiunea P . Timpul de calcul pentru produsul $W \cdot x$ se descompune astfel:

- N unități de timp pentru încărcarea matricii W .
- N unități de timp pentru încărcarea vectorului x .
- P unități de timp pentru ieșirea vectorului y .

Ieșirea vectorului y se face în paralel cu ieșirea vectorului x .

În total, vor fi necesare $2N + P$ unități de timp pentru obținerea vectorului rezultat y .

Timpul total de calcul este deci superior comparativ cu soluțiile anterioare. În ceea ce privește lărgimea busurilor de date, regăsim inconvenientul asociat primei soluții sistolice, și anume acela că lărgimea busurilor pe care circulă sumele parțiale depinde de numărul de celule sistolice de pe traseu. La o primă privire, această soluție pare costisitoare în ceea ce privește suprafața ocupată deoarece ridică numărul de CS la o valoare egală cu numărul de ponderi sinaptice din matricea W . Totuși, după cum vom vedea în cele ce urmează, performanțele cresc spectaculos în cazul execuției unor produse succesive între matricea sinaptică W și un set de vectori de intrare succesivi $x^{[k]}$, datorită *pipeline*-izării acestor operații. Acest aspect este extrem de interesant pentru algoritmi conexioniști în care se execută produse succesive între matricea sinaptică W și vectorii de intrare $x^{[k]}$. Aceasta va conduce și la distribuirea costului încărcării matricii W pe un număr mare de produse efectuate; la nivel global, costul încărcării matricii W devine neglijabil.

Pentru toate tipurile de SFG cu secvențiere sistolică (anterior descrise), se pot deriva direct rețele sistolice (RS) ale căror performanțe și constrângeri legate de integrare variază. Aceasta presupune că proiectantul, care vizează realizarea unui circuit VLSI, va selecta arhitectura pe care o va implementa în funcție de limitările tehnologice pe care le cunoaște. Vom trece în revistă câteva dintre acestea și vom analiza incidența lor asupra fazei de alegere a RS adecvate implementării VLSI.

3.1.7. Implantabilitatea rețelelor sistolice

O rețea sistolică este formată dintr-o structură regulată de celule sistolice identice. Realizarea RS pretinde deci un efort care se concentrează asupra definirii și optimizării unei singure celule. Rețeaua va fi apoi formată prin multiplicarea liniară sau bidimensională a acestui element de bază. Schema de interconexiuni, dictată de poziția punctelor de intrare / ieșire în CS, este locală exceptând interconexiunile aferente semnalelor de sincronizare. Aceasta conduce la a doua proprietate a RS: conexiunile locale facilitează implantabilitatea. Până în prezent, comunicațiile au fost considerate ca un element de cost mic în concepția procesoarelor clasice. Realizarea unor structuri masiv paralele infirmă acest punct de vedere deoarece costul interconexiunilor din structura VLSI devine superior celui aferent elementelor active. Trebuie deci localizate schimburile de date dintre elementele active, ceea ce este în mod natural realizat în RS.

Granularitatea, adică numărul de elemente de procesare (EP) utilizate pentru calcul reprezintă de asemenea un factor esențial pentru implementarea rețelelor sistolice. Acesta rezultă în general dintr-un compromis între constrângerile de intrare / ieșire, contextul procesării în interiorul RS și tehnologia utilizată.

Robustețea RS reprezintă de asemenea un factor care trebuie luat în considerare. Pentru o rețea compusă din N elemente EP de robustețe r , robustețea rețelei va fi r^N . Cum r este întotdeauna mai mic decât 1, robustețea de ansamblu aferentă unui sistem sistolic este foarte slabă. Totuși,

utilizată în contextul algoritmilor conexioniști, RS poate fi perfect adaptată pentru a beneficia de robustețea intrinsecă a modelelor conexioniste.

Puterea de calcul nu este singurul factor care definește eficiența unei arhitecturi sistolice. Această putere trebuie considerată în contextul unui sistem (sistolic) care schimbă date cu un sistem gazdă. Lărgimea de bandă aferentă canalului de comunicație cu sistemul gazdă reprezintă un parametru esențial care stabilește o limită superioară pentru puterea de calcul a sistemului luat în ansamblu. Trebuie deci aleasă o arhitectură care să permită echilibrarea timpului de calcul cu cel de comunicație (intrare / ieșire din RS).

În sfârșit, partiționabilitatea unei aplicații mari pe RS trebuie luată în considerare pentru a realiza RS de mărime rezonabilă (adaptată necesităților).

În cele ce urmează vom descrie 4 RS capabile să calculeze produsul matrice - vector. Vom alege o structură care nu impune partiționarea acestui produs, adică numărul de EP va fi egal cu pătratul numărului de componente aferente vectorului înmulțitor. Această constrângere va sta la baza tuturor dezvoltărilor pe care le vom lua în considerare în cele ce urmează.

3.2. Arhitecturi de RS destinate algoritmilor conexioniști

După cum am menționat deja, produsul matrice - vector este operația de bază, comună tuturor modelelor conexioniste. În faza de recunoaștere acest produs apare fie individual, fie ca operație intermediară conjugată cu alte operații matriciale sau vectoriale. În ipoteza unei secvențieri sistolice, pornind de la SFG-urile descrise anterior, vom deduce 4 tipuri de RS. Pentru fiecare tip vom analiza implementabilitatea, extensibilitatea și raportul performanțe / complexitate.

Reluând notațiile utilizate pentru descrierea modelelor conexioniste, vom defini câțiva termeni adecvați descrierii datelor care intră într-un calcul sistolic:

x - vectorul de intrare în rețea, de dimensiune $[N]$

y - vectorul de ieșire, de dimensiune $[P]$

W - matricea sinaptică, de dimensiune $[P, N]$

PS - vectorul de sume parțiale care tranzitează conexiunile dintre celulele sistolice

Să notăm că, în cazul rețelei *Hopfield*, nu vom diferenția între intrare și ieșire și ne vom referi doar la starea celulelor. Aplicarea unui vector de intrare x revine la a forța starea celulelor la valoarea x , iar ieșirea vectorului y revine la citirea acestei stări.

Arhitecturile sistolice dedicate produsului matrice - vector pe care urmează să le prezentăm nu vor fi toate foarte bine adaptate pentru implementarea algoritmilor conexioniști. Ele nu iau în considerare secvențierea specifică a calculelor executate în cadrul algoritmilor conexioniști. Cascadarea (înlănțuirea) produselor matrice - vector este specifică, de exemplu, fazei de recunoaștere aferentă rețelei *Hopfield*. Accesibilitatea la elementele matricii W (ponderi sinaptice) trebuie asigurată dacă dorim să implementăm și faza de învățare. Vectorul rezultat din produsul $W \cdot x$ trebuie transformat apoi printr-o funcție neliniară (funcția de activare). Operațiile de intrare / ieșire în / din RS trebuie studiate pentru a obține o interfață a cărei lărgime de bandă să reprezinte optimul dintre viteza de procesare a RS și cea a procesorului gazdă. Toate aceste probleme reprezintă constrângeri conjugate, proprii fazei de concepție a unui circuit specializat dedicat unei clase de aplicații, sau a unui coprocesor studiat în mediul său de lucru.

Vom relua sistematic soluțiile sistolice prezentate anterior și vom studia arhitecturile care vor putea fi exploatate pentru cele două operații esențiale ale calculului conexionist: faza de recunoaștere și faza de învățare. Recunoașterea, după cum am subliniat și în capitolele precedente, este identică pentru toate modelele conexioniste. Este un produs matrice - vector unic (*Adaline*, *Perceptroni mono și multistrat etc.*) sau iterativ (*Hopfield*, *Boltzmann etc.*). Învățarea va fi evaluată

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

din punctul de vedere al accesibilității la coeficienții sinaptici și la operațiile necesare pentru calculul valorilor de ajustare aferente acestor coeficienți.

3.2.1. Arhitectura generală

Arhitectura generală pe care o vom studia realizează faza de recunoaștere aferentă rețelei *Hopfield*. Reamintim că acest model reclamă un calcul iterativ al produselor dintre un vector de stare și matricea de ponderi sinaptice. După fiecare înmulțire, vectorul rezultat suferă o transformare neliniară (funcția de transfer) generând astfel noul vector de stare al rețelei *Hopfield*. Acesta este din nou înmulțit cu matricea sinaptică și operațiile se repetă până la convergență (egalitate între 2 vectori de stare succesivi).

Formal, arhitectura generală trebuie să permită execuția celor 3 operații expuse:

- produse matrice - vector iterative

$$\mathbf{p}(t + 1) = \mathbf{W} \cdot \mathbf{x}(t) \quad (3.3)$$

- aplicarea unei funcții neliniare (funcția de activare)

$$\mathbf{x}(t + 1) = \sigma(\mathbf{p}(t + 1)) \quad (3.4)$$

- detecția convergenței

$$\mathbf{x}(t + 1) = \mathbf{x}(t) \quad (3.5)$$

O schemă generală a acestei arhitecturi este prezentată în figura 3.7.

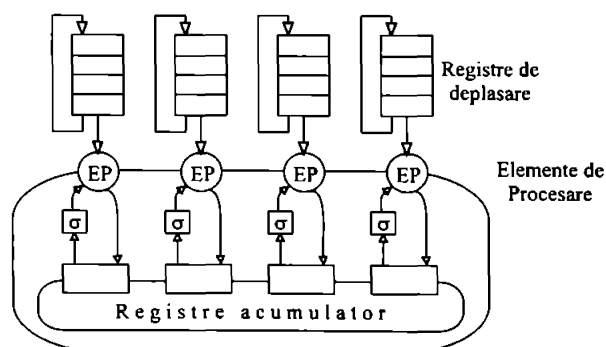


Fig. 3.7. Arhitectura generală în inel

Am definit această arhitectură generală pornind de la analiza RS liniare descrise anterior. Am completat legăturile de date aferente rețelelor elaborate cu rebuslăje care vor permite cascada produselor $\mathbf{W} \cdot \mathbf{x}$ succesive. În figura 3.7 sunt reprezentate toate căile de date posibile: comunicațiile între EP, între registrele acumulator și între registrele de deplasare care memorează ponderile sinaptice.

Arhitecturile care vor deriva din această formă generală se vor caracteriza prin amplasări specifice ale acestor căi de comunicație. Le vom referi prin menționarea operanzilor ficți (rezidenți).

Trei arhitecturi implementabile vor fi propuse: cu sume parțiale fixe (fără căi de comunicație între registrele acumulator), cu componentele vectorului fixe (fără căi de comunicație între procesoarele EP) și cu coeficienți sinaptici ficși (fără căi de date între coeficienții matricii W). Să notăm, cu titlu de generalitate, că este posibilă de asemenea o soluție fără operanzi ficși. Aceasta n-o vom descrie datorită faptului că presupune un volum de transferuri de date comparabil cu cel al calculelor efectuate. În plus, soluția în discuție nu prezintă eficiență în cazul algoritmilor conexioniști.

Operatorul funcției neliniare σ este plasat pe calea de comunicație dintre fiecare EP și acumulatorul său. Acest amplasament rezultă dintr-o transcriere directă a ecuației (3.2) care specifică faptul că funcția neliniară se aplică vectorului rezultat din produsul $W \cdot x$; fiecare registru acumulator va conține o componentă a acestui vector rezultat.

Resursele nodurilor rețelei sistolice (celulelor sistolice) sunt specificate în figura 3.8. Această structură respectă de asemenea secvențierea sistolică definită. Singurul element suplimentar care specializează această structură pentru algoritmii conexioniști este funcția neliniară σ .

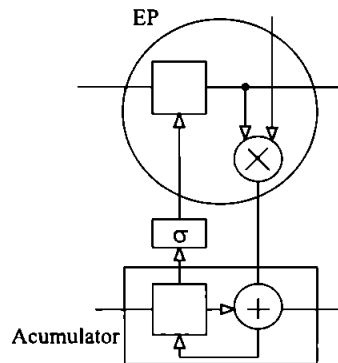


Fig. 3.8. Structura celulei sistolice

În sfârșit, pentru transferurile de date sunt posibile 2 tehnici: transmisia serială sau paralelă. Soluția serială reduce numărul de linii de comunicație dintre celulele sistolice dar conduce la creșterea timpilor de încărcare și vidare a inelului sistolic. Soluția paralelă reduce timpii de comunicație dar crește numărul de conexiuni. Vom detalia în capitolele următoare avantajele și inconvenientele fiecărei soluții.

Mai rămâne de stabilit repartitia coeficienților sinaptici în registrele de deplasare. Această repartitie depinde exclusiv de datele pe care vrem să le tranzitam între nodurile RS. În cele ce urmează, vom defini diferitele repartiții posibile și vom deriva arhitecturile rezultate.

3.2.2. Inel sistolic cu componentele vectorului fixate

Datele mobile vor fi sumele parțiale și componentele matricii W . Componentele matricii W trebuie repartizate în registrele de deplasare, astfel încât să obținem secvențierea sistolică definită. Componentele vectorului x sunt rezidente în fiecare nod.

Disponerea coeficienților sinaptici depinde de disponerea componentelor vectorului (în timpul procesului de calcul). În faza de inițializare, componentele vectorului x vor fi introduse secvențial în inelul sistolic iar coeficienții matricii W în registrele de deplasare. În faza de calcul celulele sistolice vor recepționa componentele matricii simultan. Decalajul de prezentare a

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

coeficienților matricii W impus de algoritmul sistolic trebuie implementat printr-o repartitie corespunzătoare a acestor coeficienți în registrele de deplasare, așa cum indică figura 3.9.

Pe fiecare impuls, cele N produse parțiale sunt calculate și acumulate în memoriile locale (registrele acumulator). Sumele parțiale PS sunt apoi transmise între nodurile rețelei pentru o nouă acumulare.

Dacă luăm în considerare numai faza de calcul, vor fi necesare N operații de înmulțire și de transmisie a sumelor parțiale, urmate de conversia rezultatelor obținute cu ajutorul operatorilor care implementează funcția neliniară σ .

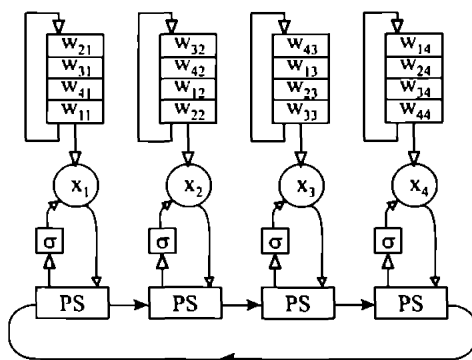


Fig. 3.9. Inel sistolic cu componentele vectorului fixe

Coeficienții matricii W sunt deplasați în registrele de deplasare și, după N impulsuri, se vor regăsi în poziția inițială. De asemenea, componentele vectorului rezultat se regăsesc în poziția corectă, necesară declanșării unui nou produs.

În cadrul acestei soluții, componentele vectorului de multiplicat sunt rezidente în fiecare CS. Ele nu vor fi afectate în timpul înmulțirilor și deci reprezentarea lor (numărul de biți necesari pentru reprezentare) nu va varia. Sumele parțiale PS care tranzitează între celulele sistolice vor avea o reprezentare variabilă, care crește în funcție de numărul de celule traversate. Pentru o rețea cu N celule, fiecare ieșire y_i se va obține prin însumarea a N produse $w_{ij} \cdot x_j$. Presupunem că există o valoare maximă stabilită pentru componentele x_j și respectiv w_{ij} . Numărul de biți necesari pentru reprezentarea acestora va fi:

$$\begin{aligned} T_{\max}(x_j) &= \text{Log}_2(\text{Max}[x_j]) \\ T_{\max}(w_{ij}) &= \text{Log}_2(\text{Max}[w_{ij}]) \end{aligned}$$

Numărul maxim de biți necesari reprezentării componentelor vectorului rezultat va fi:

$$T_{\max}(y_i) = \text{Log}_2(N) + T_{\max}(x_j) + T_{\max}(w_{ij}) \quad (3.6)$$

Transferul sumelor parțiale între celule trebuie realizată printr-o comunicație care trebuie să țină cont de relația de mai sus. Avem la dispoziție mai multe soluții:

- calculul lățimii căii de comunicație dintre celule (numărul de biți necesari) pe baza relației de mai sus și construcția unei arhitecturi fixe

- fixarea unei căi de comunicație a cărei lățime (număr de biți) care depinde de valoarea maximă admisibilă pentru PS. Această dimensiune fixă limitează, de asemenea, numărul maxim de celule din RS, dacă dorim evitarea depășirilor de capacitate de reprezentare
- implementarea unei transmisii seriale (bit după bit) a sumelor PS, tehnică ce introduce noțiunea de *bit level systolic array* [KNO 88]. Vom distinge în acest caz două nivele de *pipelining*, primul situat la nivelul celulelor sistolice și al doilea la nivelul comunicației dintre celule. Să notăm că dimensiunea registrului destinat memorării componente PS în celula sistolică va reprezenta de asemenea o limitare dar această dimensiune poate fi fixată la o valoare acoperitoare.

Inelul sistolic cu componentele vectorului fixate este teoretic extensibil. Lungimea rețelei poate fi mărită prin simpla adăugare a unor noduri suplimentare. Este, de asemenea, posibilă implementarea unor produse ample, în care numărul de CS este inferior numărului de componente vectoriale; este suficientă completarea rețelei cu un dispozitiv destinat acumulării produselor parțiale.

Timpul de încărcare a matricii sinaptice se evaluează în funcție de tipul comunicației implementate între RS și sistemul gazdă. Dacă accesul la RS este secvențial, vor fi necesare N^2 cicluri de încărcare și un singur port de acces (matricea W este o matrice pătrată de dimensiune $[N \times N]$ în cazul rețelei *Hopfield*). Dacă accesul este paralel vor fi necesare N cicluri de încărcare și N porturi de acces (această remarcă va fi valabilă și pentru inelul sistolic cu sume parțiale fixe care va fi prezentat în secțiunea 2.3.2).

Componentele vectorului înmulțitor sunt localizate în fiecare celulă sistolică. Încărcarea sau citirea acestui vector se poate face fie secvențial, fie paralel. În cazul paralel sunt necesare N porturi de intrare / ieșire, dar timpul de acces se reduce la un ciclu.

Dacă luăm în considerare comunicația paralelă, timpul total de procesare va fi:

N cicluri pentru încărcarea matricii W
1 ciclu pentru încărcarea vectorului x
 $k \cdot N$ cicluri de calcul, k reprezentând numărul de iterații necesare atingerii unei stări stabile (convergenței)
 k cicluri pentru aplicarea funcției neliniare σ
1 ciclu pentru citirea vectorului rezultat

Rezultă un total de $N + k(N + 1) + 2$ cicluri de procesare, dacă presupunem că toate operațiile revendică cicluri de durate identice.

Acest tip de arhitecturi sistolice este destinat implementării rețelelor de dimensiuni mari datorită proprietății de extensibilitate. Performanțele canalului de comunicație dintre RS și calculatorul gazdă sunt esențiale din punctul de vedere al performanțelor globale deoarece N componente trebuie transferate simultan.

În cazul comunicației seriale, timpii de încărcare cresc considerabil, dar teoretic sunt suficiente 2 porturi de intrare / ieșire pentru orice dimensiune de RS. Timpul total de procesare se evaluează astfel:

N^2 cicluri pentru încărcarea matricii W
 N cicluri pentru încărcarea vectorului x în celulele sistolice
 $k \cdot N$ cicluri de calcul, k reprezentând numărul de iterații necesare atingerii unei stări stabile
 k cicluri pentru aplicarea funcției neliniare σ
 N cicluri pentru citirea vectorului rezultat

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

Rezultă un total de $N^2 + (k + 2)N + k$ cicluri de procesare, presupunând că toate operațiile necesită cicluri identice (ca durată).

Soluția secvențială este avantajoasă deoarece decongestionează comunicația cu sistemul gazdă. Ea revendică însă o cale de comunicație suplimentară între elementele EP (pentru încărcarea vectorului x) care vine să dubleze pe cea existentă între acumulatori; aceasta conduce la un cost suplimentar în ceea ce privește suprafața de integrare ocupată. Este preferabilă alegerea unei arhitecturi de RS care conține deja o cale de comunicație între elementele EP (necesară pentru calcule), cale care va fi utilizată și pentru operațiile de intrare / ieșire.

3.2.3. Inel sistolic cu sume parțiale fixe

Această arhitectură presupune circulația componentelor vectorului x între elementele EP; așa cum am subliniat mai sus, va fi o arhitectură bine adaptată operațiilor de comunicație cu sistemul gazdă. Din aceleași motive prezentate la 3.2.2, repartiția elementelor matricii W în registrele de deplasare se va realiza ca în figura 3.10.

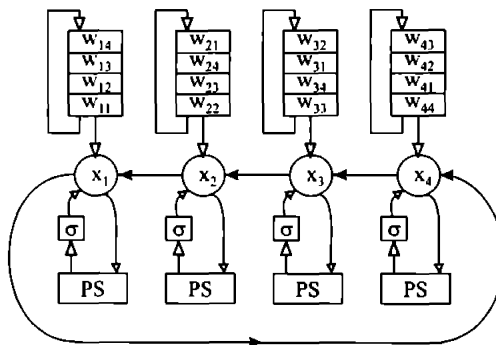


Fig. 3.10. Inel sistolic cu vector mobil și sume parțiale fixe

În această arhitectură, componentele vectorului x se deplasează între elementele EP și, în interiorul fiecărui EP, sunt înmulțite cu coeficienții matricii W prezenți pe intrări. Sumele parțiale sunt acumulate în interiorul EP. După N cicluri de calcul / transmisie, sumele parțiale acumulate sunt transformate prin funcția neliniară σ (implementată în fiecare EP). Vectorul rezultat este astfel gata pentru execuția următorului produs $W \cdot x$. La finele operațiilor (convergență) componentele vectorului de stare sunt extrase secvențial din RS și încărcate în sistemul gazdă.

Această arhitectură, în termeni de implantabilitate, este mai bună decât precedenta. Datele mobile sunt aici elementele matricii W și componentele vectorului x . Aceste componente nu sunt modificate pe parcursul calculului și prin urmare dimensiunea (numărul de biți alocați) nu variază. Să notăm că elementele vectorului x se modifică periodic dar în urma aplicării funcției σ ; această funcție neliniară restrânge valoarea componentelor la domeniul de reprezentabilitate, care este domeniul în care ia valori σ $([-1, +1])$. În consecință, căile de comunicație între EP au o lățime fixă indiferent de lungimea inelului sistolic. Registrele de acumulare pot fi dimensionate pornind de la dimensiunea operanzilor, așa cum s-a arătat în paragraful precedent.

Dacă rețeaua este constituită din N celule, fiecare registru PS va conține componenta finală a vectorului rezultat după traversarea a N componente x_j prin EP-ul respectiv. Dimensiunea registrului PS va fi calculată astfel încât acesta să poată memora cel mai mare rezultat posibil:

$$\text{Dimensiune(PS)} = \text{Log}_2(y_{\text{imax}}) = \text{Log}_2(N) + T_{\text{max}}(x_j) + T_{\text{max}}(w_{ij}) \quad (3.7)$$

Pentru această arhitectură, vom considera că încărcarea și citirea datelor se face secvențial. Timpul total de procesare va conține:

- N^2 cicluri pentru încărcarea matricii W
- N cicluri pentru încărcarea vectorului x
- $k \cdot N$ cicluri de calcul, k reprezentând numărul de iterații până la convergență
- k cicluri de aplicare a funcției neliniare σ
- N cicluri pentru citirea vectorului rezultat

Rezultă un total de $N^2 + N(k + 2) + k$ cicluri, dacă presupunem că toate operațiile revendică cicluri identice (ca durată) pentru execuție. Timpul de procesare este identic cu cel obținut pentru arhitectura anterioară (cazul comunicației seriale). Această arhitectură este însă mai bună din 2 motive:

- schimburile de date cu sistemul gazdă sunt favorizate de comunicația dintre EP-uri
- căile de comunicație dintre EP-uri au lățimea fixă deoarece datele vehiculate nu depășesc domeniul de reprezentare stabilit. Funcția neliniară σ este o funcție de compresie care comprimă valorile într-un interval fixat. Se poate deci determina în avans lățimea căilor de comunicație; aceasta rămâne fixă indiferent de lungimea inelului sistolic.

Punctele slabe care revin acestei arhitecturi sunt:

- extensibilitate dificilă datorată registrelor de deplasare. Elementele matricii W sunt stocate în aceste registre. Lungimea acestor registre (numărul de locații disponibile pentru elementele w_{ij}) trebuie să fie egală cu numărul de noduri din RS. Ori, dacă dorim o rețea extensibilă, numărul de noduri nu va fi cunoscut a priori.
- dimensiunea acumulatorilor (numărul de biți) depinde de dimensiunea RS (numărul de noduri). Dacă dorim o precizie ridicată de reprezentare a sumelor parțiale, trebuie prevăzută o dimensiune corespunzătoare pentru acumulator. Această dimensiune trebuie definită la construcția RS și va limita ulterior extinderea inelului sistolic.

Această soluție a fost propusă și de *Weinfeld* [WEI 89] pentru integrarea unei rețele sistolice liniare dedicată fazei de recunoaștere aferentă modelului *Hopfield*.

Cea mai importantă limitare a sistemelor sistolice liniare prezentate constă în extensibilitatea limitată datorită atribuirii unui număr relativ mare de coeficienți sinaptici fiecărui EP. O rețea sistolică bidimensională poate depăși această limitare deoarece coeficienții matricii W devin locali (câte unul alocat pe fiecare EP). Suprafața totală va crește desigur deoarece elementele EP vor fi multiplicat. Acest cost suplimentar va fi fără îndoială compensat de posibilitățile de extindere precum și de puterea de calcul mult crescută grație unui *pipelining* vectorial. Această soluție urmează a fi prezentată.

3.2.4. Arhitectura cu coeficienți sinaptici fiși

În această rețea bidimensională, numărul de elemente EP va fi egal cu numărul de coeficienți sinaptici din matricea W . Sumele parțiale și componentele vectorului x vor fi mobile pentru a realiza două funcții:

- punerea la dispoziție a componentelor x , necesare în fiecare linie a rețelei sistolice

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

- transmiterea rezultatelor parțiale de-a lungul liniilor RS.

Transmisia datelor este descrisă în figura 3.11.

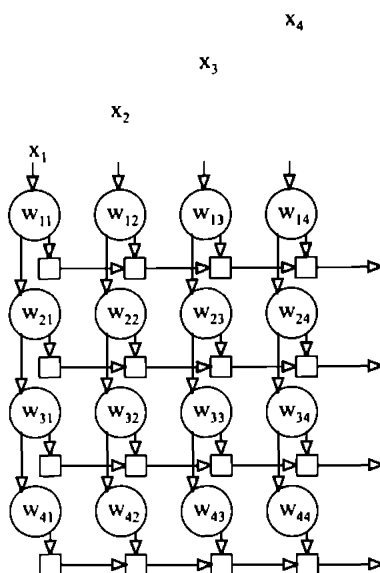


Fig. 3.11. Arhitectura rețelei sistolice bidimensionale

Componentele vectorului x sunt introduse secvențial în rețea. Produsul local este calculat imediat ce o componentă x_i este introdusă într-un EP. Produsele locale sunt apoi transmise celulelor vecine pe direcția est și acumulate la nivelul fiecărei celule. În momentul în care ultima componentă a vectorului x părăsește rețeaua pe direcția sud, ultima componentă a vectorului rezultat ($W \cdot x$) va părăsi rețeaua pe direcția est. Apare deci o dublă mișcare a datelor:

- componentele vectorului x tranzitează rețeaua dinspre nord spre sud între EP-urile aceleiași coloane
- sumele parțiale tranzitează rețeaua de la vest spre est între EP-urile aceleiași linii.

Pentru a avea celule identice indiferent de amplasamentul lor în cadrul RS, dimensiunea registrelor PS trebuie să permită memorarea celui mai mare rezultat posibil:

$$\text{Dimensiune PS} = \text{Log}_2(N) + \text{Tmax}(x_j) + \text{Tmax}(w_{ij}) \quad (3.8)$$

Pentru această arhitectură vom considera că încărcarea și citirea datelor în / din RS se face secvențial, dar încărcarea matricii W se face în paralel. Timpul de calcul va include:

- N cicluri pentru încărcarea matricii W
- N cicluri pentru încărcarea vectorului x
- $k \cdot N$ cicluri de calcul, k reprezentând numărul de iterații până la convergență
- k cicluri de aplicare a funcției neliniare
- N cicluri pentru citirea vectorului rezultat

Rezultă un timp de procesare total de $k + N(k + 3)$ cicluri, dacă presupunem că toate operațiile revendică cicluri identice (ca durată) pentru execuție.

Această arhitectură prezintă 2 avantaje:

- viteză ridicată în raport cu rețelele sistolice liniare deoarece permite procesarea vectorilor de intrare în tehnică *pipeline*
- extensibilitate. Fiecare EP conține un singur coeficient sinaptic. Este suficientă extinderea pavajului bidimensional pentru extinderea rețelei.

Constrângerile referitoare la lățimea căilor de date și cele referitoare la dimensiunea acumulatorilor care servesc la tranzitarea sumelor parțiale rămân în vigoare și sunt identice cu cele din cadrul arhitecturilor precedente.

Inconvenientul major constă în suprafața ocupată de rețeaua bidimensională. În plus, gradul de utilizare a celulelor sistolice care este 100 % în cazul liniar, în cazul bidimensional și pentru această soluție elementară, depinde de momentul evaluării și este maxim 50 %. Această limitare este compensată de posibilitatea *pipeline*-izării fluxului de vectori de intrare, ceea ce conduce la un grad de utilizare de 100 % în toate momentele de timp. Vom detalia această *pipeline*-izare în capitolele următoare. În sfârșit, vectorul x este introdus de pe direcția nord și noua sa valoare se obține pe latura de est a RS. Dacă dorim înlănțuirea mai multor produse matrice - vector pe RS diferite (succesive), va trebui studiată topologia de implementare în VLSI pentru a liniariza transferurile de date.

3.3. Concluzii

Dintre cele 3 soluții prezentate vom reține rețeaua bidimensională rectangulară din mai multe motive:

În primul rând, datorită posibilităților de extensie care permit construirea unor rețele mari. Mărirea rețelei se face printr-o extensie a pavajului bidimensional (de celule) și această extensie nu afectează resursele celulei sistolice. Fiecare celulă memorează o singură pondere sinaptică. Creșterea numărului de sinapse va fi egală cu pătratul numărului de neuroni adăugați. Creșterea suprafeței ocupate de o arhitectură bidimensională este de asemenea o funcție pătratică. Raportul (număr sinapse) / (suprafața ocupată) rămâne constant ceea ce indică faptul că extensibilitatea nu va fi limitată de un factor de suprafață disponibilă.

În al doilea rând, principiul distribuției calculului pe un număr mare de procesoare corespunde mult mai bine principiului conexioniștilor. Robustețea furnizată de modelele conexioniștilor poate fi mai bine exploatată printr-o amenajare adecvată a unor linii de reconfigurare care vor permite izolarea și înlocuirea celulelor defecte.

În al treilea rând, soluția sistolică permite efectuarea calculelor cu o precizie suficientă pentru majoritatea algoritmilor conexioniști. Tehnologia actuală permite integrarea unui număr suficient de celule sistolice cu un randament ridicat.

În sfârșit, după cum vom constata în cele ce urmează, această arhitectură bidimensională permite implementarea unor mecanisme mai complexe, proprii algoritmilor conexioniști: detecție de convergență, gestiune eficientă a operațiilor de intrare / ieșire, învățare după reguli multiple, interconectarea subrețelelor etc.

4. PROIECTAREA RS DEDICATE ALGORITMILOR CONEXIONIȘTI

În acest capitol vom detalia funcționarea rețelei sistolice bidimensionale obținută în capitolul anterior pentru produsul matrice - vector. Vom arăta că există un flux de date și o secvențiere care permit înlănțuirea unor astfel de produse matrice - vector; funcția neliniară σ se va aplica tuturor componentelor vectorului rezultat. Aceste operații reprezintă esența majorității algoritmilor conexioniști.

4.1. Soluția cu valori de activare și sume parțiale mobile

Rețeaua bidimensională obținută în cap. 3 alocă un EP per sinapsă. Pentru un vector cu N componente, rețeaua va conține N^2 sinapse cărora li se vor alocă N^2 procesoare EP. Fiecare EP va conține un singur coeficient sinaptic și va acumula rezultatele parțiale pe care le-am denumit sume parțiale (PS). Pentru RS liniare, am arătat că este posibilă înlănțuirea produselor matrice - vector și aplicarea funcției neliniare fiecărei componente din vectorul rezultat, dacă se repartizează corect coeficienții matricii W în interiorul fiecărei celule și dacă se realizează un rebucraj al rețelei sistolice liniare.

Circulația datelor în interiorul RS bidimensionale este descrisă în figura 4.1. Se poate constata că, inițial, vectorul înmulțitor este introdus de pe direcția Nord în rețea și vectorul rezultat este emis pe direcția Est din rețea.

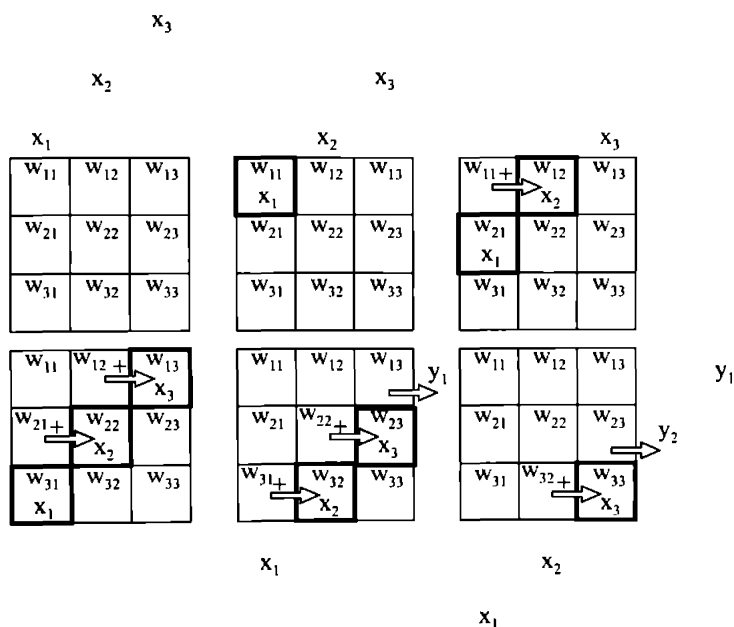


Fig. 4.1. Circulația datelor în RS bidimensională

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

Datele mobile sunt:

- Componentele vectorului x care se deplasează dinspre Nord spre Sud;
- Sumele Parțiale (PS) care se deplasează dinspre Vest spre Est.

Fiecare celulă sistolică (i, j) recepționează componenta x_j a vectorului x în momentul (t) , calculează produsul local $w_{ij} \cdot x_j$ și îl adaugă sumei parțiale PS_i recepționată din direcția vest. Rezultatul acestei acumulări este transmis celulei vecine de pe direcția Est. Dacă vectorul x are N componente, pentru a obține fiecare componentă a vectorului y sunt necesare N operații de tipul celei descrise anterior, și fiecare componentă se calculează cu un decalaj de o unitate de timp față de cea precedentă. Vectorul y va fi obținut deci după N unități de timp. Vom denumi pas de recurență etapa care constă în calculul unui produs $W \cdot x$.

În momentul în care o componentă a vectorului rezultat (y_i) apare pe latura de Est a RS, se poate aplica funcția neliniară σ . Funcția de activare σ nu va fi integrată în interiorul RS; ea se va calcula cu ajutorul unor operatori specializați plasați pe latura de est a RS. Apoi, vectorul rezultat y va trebui reintrodus în RS ca vector de intrare x , pentru calculul următorului produs $W \cdot x$. Acest reuclaj se poate realiza cu ajutorul unor căi de date care leagă latura de Est cu latura de Nord, așa cum se arată în fig. 4.2.

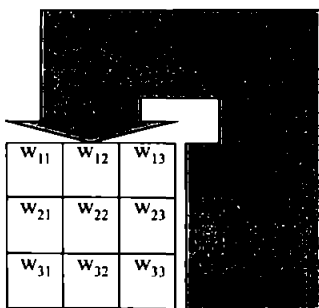


Fig. 4.2. Reuclajul datelor pe intrarea RS

Nemaifiind o transmisie locală, această soluție distruge omogenitatea transmisiilor de date și poate genera timpi suplimentari de propagare care vor afecta performanțele globale ale RS. Este preferabil să se găsească o soluție care respectă regula de transmisie la cel mai apropiat vecin, fără legături la distanțe mari.

4.2. O soluție care implică transmisiile locale de date

Vom discuta această soluție asociată cu faza de recunoaștere pe o RN de tip *Hopfield*. Soluția se bazează pe o circulație locală a datelor care respectă 3 reguli:

- (r1): Secvența operațiilor de calcul este identică cu cea prezentată în paragraful 4.1.
- (r2): Fiecare celulă sistolică (CS) nu cunoaște decât propria stare și stările celor 4 celule vecine.
- (r3): La sfârșitul fiecărui pas de recurență, vectorul înmulțitor trebuie să se afle în aceeași poziție inițială în care se afla și vectorul înmulțitor precedent (cel care l-a generat).

Ideea de bază se referă la poziția inițială a vectorului înmulțitor: acesta trebuie plasat pe diagonala principală a RS, așa cum indică figura 4.3.

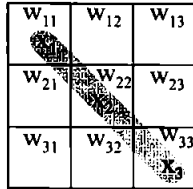


Fig. 4.3. Poziția inițială a vectorului x

Toate celulele sistolice sunt identice și nu vor schimba informații decât cu celulele vecine (cei 4 vecini cei mai apropiați). Structura CS este prezentată în figura 4.4.

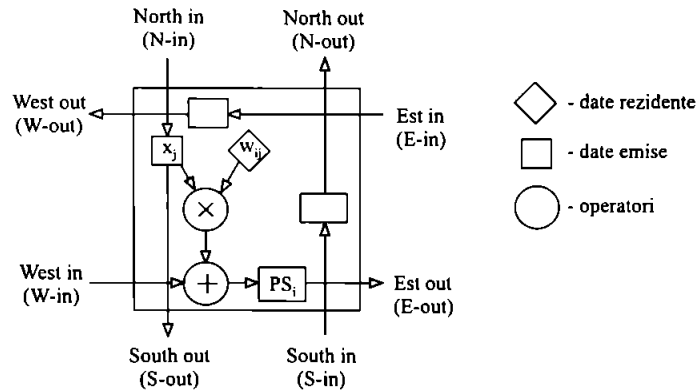


Fig. 4.4. Schema logică a unei CS

4.2.1. Derularea algoritmului

Algoritmul de calcul al unui vector rezultat (un pas de recurență) poate fi descompus în 3 faze succesive:

- elaborarea progresivă a sumelor parțiale
- aplicarea funcției neliniare
- re poziționarea vectorului rezultat în poziția inițială pentru un nou pas de recurență

Vom prezenta derularea completă a unui pas de recurență pentru o matrice W de dimensiune 3×3 .

A. Prima etapă

După cum am precizat deja, vectorul înmulțitor x va fi plasat inițial pe diagonala RS. Respectând regula r1 nu vom putea efectua decât un singur produs, așa cum se indică în figura 4.5.

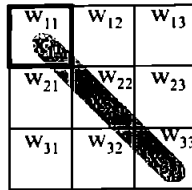


Fig. 4.5. Prima etapă (r - pasul de recurență)

Dacă notăm:

$PS_{i,j}$ - suma parțială prezentă în celula (i, j)

$x_{i,r}$ - componenta i a vectorului x în pasul de recurență r

vom calcula prima sumă parțială:

$$PS_{1,1} = w_{11} \cdot x_{1,1} \quad (4.1)$$

B. A doua etapă

În a 2-a etapă vor intra în joc celulele w_{21} și w_{12} . Vom propaga datele în RS astfel încât să le punem la dispoziția acestor celule:

- coeficienții matricii W sunt date rezidente; nu vor fi deplasați
- componentele vectorului x vor fi date emise și se vor propaga astfel:
 - ◆ componenta $x_{1,1}$ va fi propagată spre sud; ea va fi astfel prezentă pentru calculul produsului și sumei parțiale în celula w_{12}
 - ◆ componenta $x_{2,1}$ va fi propagată spre nord pentru a fi pusă la dispoziția celulei w_{21}
 - ◆ componenta $x_{3,1}$ va fi propagată de asemenea spre nord pentru a o apropia de celula w_{13} (anticipare)
- suma parțială $PS_{1,1}$ va fi transmisă spre est.

După aceste transmisii locale de date vor fi calculate produsele locale, care vor fi adăugate sumelor parțiale care intră în celule din direcția vest:

$$PS_{1,2} = PS_{1,1} + w_{12} \cdot x_{2,1} \quad (4.2)$$

$$PS_{2,1} = w_{21} \cdot x_{1,1} \quad (4.3)$$

Figura 4.6. explicitează deplasarea datelor și calculele efectuate.

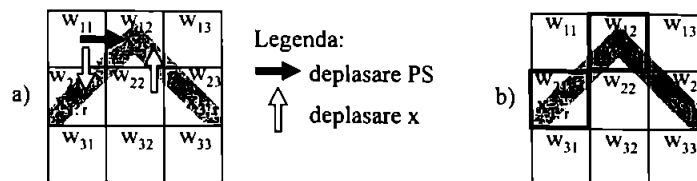


Fig. 4.6. A doua etapă. a) comunicații; b) calcul

La sfârșitul acestei etape, sumele parțiale $PS_{1,2}$ și $PS_{2,1}$ vor fi prezente în celulele w_{12} și respectiv w_{21} .

C. A treia etapă

Datele mobile vor fi transmise ca în etapa precedentă. Calculele se vor efectua în celulele w_{31} , w_{22} și w_{13} .

- $x_{1,1}$ și $x_{2,1}$ sunt propagate spre sud
- $x_{3,1}$ este propagat spre nord
- sumele parțiale $PS_{1,2}$ și $PS_{2,1}$ sunt transmise spre est.

Vor fi calculate produsele locale și vor fi acumulate sumelor parțiale care intră dinspre vest:

$$PS_{1,3} = PS_{1,2} + w_{13} \cdot x_{3,1} \tag{4.4}$$

$$PS_{2,2} = PS_{2,1} + w_{22} \cdot x_{2,1} \tag{4.5}$$

$$PS_{3,1} = w_{31} \cdot x_{1,1} \tag{4.6}$$

La sfârșitul acestei etape, prima componentă a vectorului rezultat (prima sumă finală) va fi prezentă în celula w_{13} (de pe marginea de est a RS) așa cum indică și figura 4.7.

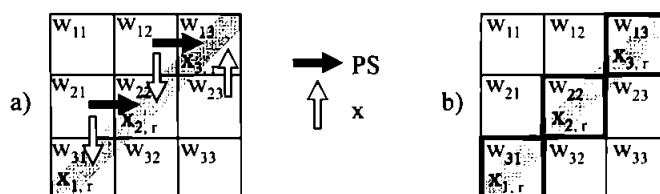


Fig. 4.7. A treia etapă. a) comunicații; b) calcul

Următoarele sume parțiale sunt încă în curs de calcul și, pentru a înlănțui recurențele, va trebui să reamplasăm vectorul rezultat în diagonala rețelei (vectorul rezultat devine vectorul înmulțitor pentru următorul pas de recurență).

D. A patra etapă

Cu scopul de a respecta regula r3, vom încerca să exploatăm decalajul temporal care apare pe latura de est a RS (componentele vectorului rezultat se obțin în etape succesive în celulele de pe latura de est a RS). Vom profita de acest fapt și vom retransmite componentele deja calculate pe o cale de întoarcere spre diagonala principală. Prima componentă a vectorului rezultat este deja disponibilă pe latura de est a RS. Ea va fi transformată prin aplicarea funcției de transfer σ în componenta corespondentă a noului vector x , care va fi vectorul înmulțitor pentru următorul pas de recurență:

$$x_{1,2} = \sigma(PS_{1,3}) \tag{4.7}$$

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

Datele mobile sunt transmise:

- $x_{1,2}$ se va propaga dinspre est spre vest (Să remarcăm trecerea la următorul pas de recurență!);
- $x_{2,1}$ și $x_{3,1}$ sunt propagate spre sud;
- componenta $x_{1,1}$, prezentă pe latura de sud, va fi retrimisă spre nord. Această retransmisie spre nord va servi mai târziu la detecția convergenței pentru faza de recunoaștere aferentă rețelei *Hopfield* (deocamdată nu detaliem);
- sumele parțiale $PS_{3,1}$ și $PS_{2,2}$ sunt transmise spre est.

Calculul efectuat sunt:

$$PS_{2,3} = PS_{2,2} + w_{23} \cdot x_{3,1} \quad (4.8)$$

$$PS_{3,2} = PS_{3,1} + w_{32} \cdot x_{2,1} \quad (4.9)$$

Noutatea adusă de această etapă constă în aplicarea funcției neliniare σ componentei $PS_{1,3}$ prezente deja pe latura de est a RS (ec. (4.7)) și retransmisia rezultatului $x_{1,2}$ spre diagonală.

Circulația datelor și calculele efectuate sunt evidențiate în figura 4.8. Lungimea căii de revenire pe diagonală (de pe latura de est) este evident egală cu cea parcursă dinspre diagonală spre latura de est. Această lungime se diminuează cu o unitate pe măsură ce avansăm de la o linie la alta dinspre nord spre sud.

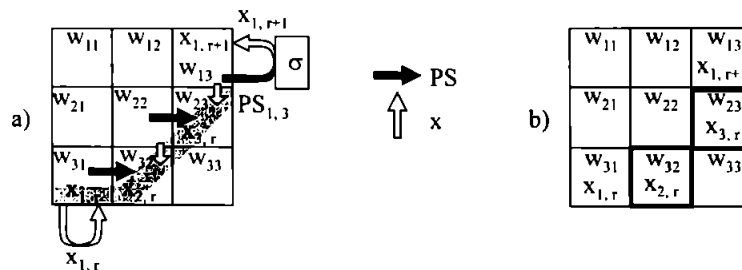


Fig. 4.8. Etapa a patra: aplicarea funcției neliniare σ și retransmiterea spre diagonală
a) comunicații; b) calcul

Decalajul temporal între sosirile succesive pe latura de est ale componentelor vectorului rezultat este de asemenea de o unitate de timp. Rezultă că toate componentele vectorului rezultat se vor găsi în același moment pe diagonală și vor fi gata pentru declanșarea următorului pas de recurență.

E. A cincea etapă

Etapa a cincea este similară celei anterioare; se va retransmite spre diagonală a 2-a componentă a vectorului rezultat (după aplicarea prealabilă a funcției neliniare σ).

$$x_{2,2} = \sigma(PS_{2,3}) \quad (4.10)$$

Vom avea de asemenea următoarele transferuri de date:

- $x_{1,2}$ și $x_{2,2}$ sunt propagate dinspre est spre vest;

Proiectarea Rețelelor Sistolice dedicate algoritmilor conexiști

- $x_{3,1}$ este propagat spre sud;
- $x_{1,1}$ și $x_{2,1}$ (prezentă pe latura de sud) sunt retrimise spre nord;
- suma parțială $PS_{3,2}$ este transmisă spre est.

$$PS_{3,3} = PS_{3,2} + w_{33} \cdot x_{3,1} \quad (4.11)$$

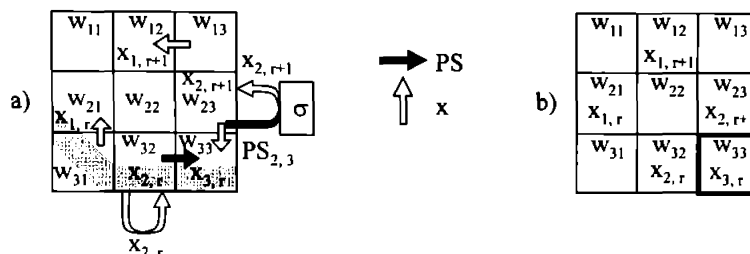


Fig. 4.9. Etapa a cincea. a) comunicații; b) calcul

În etapa a șasea va fi obținut rezultatul așteptat: produsul $W \cdot x$ va fi finalizat, funcția neliniară σ va fi aplicată pentru ultima componentă și vectorul rezultat va fi re poziționat pe diagonala principală a RS.

F. A șasea etapă

Suma parțială $PS_{3,3}$ este prezentă pe latura de est a RS. Se va aplica funcția de transfer σ și se va obține:

$$x_{3,2} = \sigma(PS_{3,3}) \quad (4.12)$$

Se vor executa de asemenea următoarele transferuri de date:

- $x_{1,2}$, $x_{2,2}$ și $x_{3,2}$ sunt propagate dinspre est spre vest;
- $x_{1,1}$, $x_{2,1}$ și $x_{3,1}$ sunt retrimise spre nord.

În această etapă nu se execută calcule. Figura 4.10. prezintă situația RS la finele etapei a șasea (care coincide cu sfârșitul primului pas de recurență).

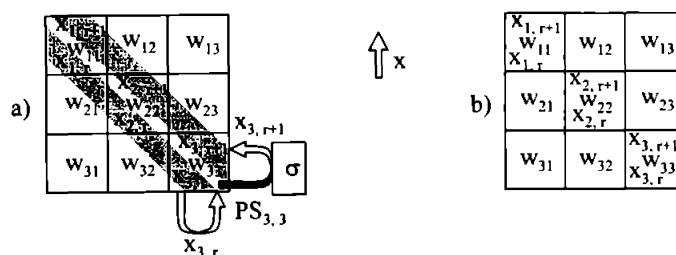


Fig. 4.10. Etapa a șasea. a) comunicații; b) calcul

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

La finele celei de-a șasea etape, atât componentele noului vector \mathbf{x}_{r+1} cât și cele ale vechiului vector \mathbf{x}_r vor fi prezente în celulele sistolice de pe diagonala RS. Regulile r1, r2 și r3 au fost respectate și rețeaua este pregătită pentru lansarea următorului pas de recurență.

Prezența simultană a vechiului vector și a noului vector în celulele de pe diagonală ne permite să comparăm cei doi vectori termen cu termen. Vom dispune astfel de N semnale de comparare locală în diagonala RS care vor trebui compuse pentru a genera semnalul de comparare globală (rezultatul comparării). Această funcție este esențială pentru faza de recunoaștere aferentă algoritmului *Hopfield*, pe parcursul căreia trebuie detectată convergența spre o stare stabilă. Algoritmul presupune calculul:

$$\mathbf{x}_r = \sigma(\mathbf{W} \cdot \mathbf{x}_{r-1}) \quad (4.13)$$

până la atingerea convergenței (unei stări stabile):

$$\mathbf{x}_r = \mathbf{x}_{r-1} \quad (4.14)$$

4.2.2. Detecția convergenței

După un pas de recurență fiecare celulă diagonală va conține $x_{i,r}$ și $x_{i,r+1}$. Compararea acestor 2 valori ne va furniza un rezultat local în legătură cu convergența vectorului \mathbf{x} . Semnalul de convergență global, pentru o RS de dimensiune $N \times N$, se va calcula:

$$\text{Conv}(N) = \prod_{i=1}^N (x_{i,r} = x_{i,r+1}) = \prod_{i=1}^N CV_i \quad (4.15)$$

Ecuția (4.15) se poate scrie și sub formă recursivă:

$$\begin{aligned} \text{Conv}(0) &= 1 \\ \text{Conv}(N) &= \text{Conv}(N-1) \cdot \text{Conv}(N) \end{aligned} \quad (4.16)$$

Fiecare celulă diagonală va conține un dispozitiv de comparare care va genera valoarea booleană CV_i . Acest semnal va fi transmis local (spre vest) și va fi sintetizat pe latura de vest a RS printr-un dispozitiv care calculează sistolic conjuncția prezentă în ecuația (4.16). Vom abandona pentru moment propagarea componentelor vectorilor \mathbf{x} și \mathbf{PS} . Poziția inițială a vectorului \mathbf{CV} este prezentată în figura 4.11.

CV_1		
w_{11}	w_{12}	w_{13}
	CV_2	
w_{21}	w_{22}	w_{23}
		CV_3
w_{31}	w_{32}	w_{33}

Fig. 4.11. Poziția inițială a semnalelor de convergență locale

Toate aceste semnale vor fi transmise spre latura de vest. Dispozitivul sistolic care implementează ecuația 4.16 combină apoi aceste semnale locale pentru a genera semnalul de convergență global (fig. 4.12).

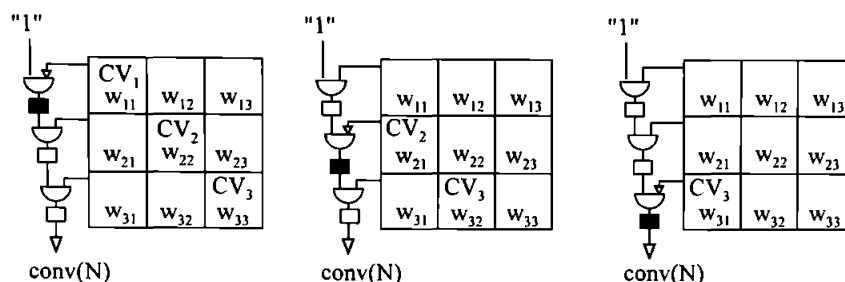


Fig. 4.12. Propagarea sistolică a semnalului de convergență

Decalajul temporal implicat de transmisia diagonalei spre latura de vest a RS este exploatat pentru sincronizarea calculului semnalului de convergență. Această subrețea poate fi descrisă de asemenea sub forma unui sistem de ecuații recurente uniforme și validitatea sa poate fi demonstrată.

4.2.3. Definierea unei arhitecturi logice

Vom propune o arhitectură logică pentru RS bidimensională care implementează funcțiile anterior descrise (fig. 4.13). Rebuclajele prevăzute pe laturile de Nord, Sud și Est vor permite re poziționarea noului vector x pe diagonala RS în scopul înlănțuirii pașilor de recurență succesivi. În ceea ce privește calculele efectuate, toate celulele sistolice sunt identice. În ceea ce privește comunicațiile, celulele diagonale vor trebui specializate. După cum a fost deja descris, în timpul primei faze de calcul, componentele vectorului x , prezente inițial în diagonală, sunt propagate mai întâi spre nord, pentru ca apoi să fie rebucate spre sud. Componentele noului vector x sunt propagate de pe latura de est și sunt poziționate pe diagonală pentru declanșarea următorului pas de recurență.

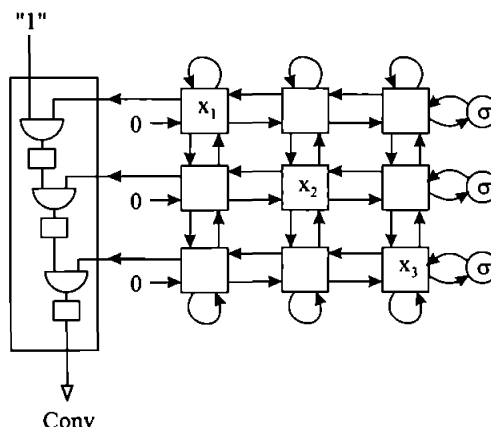


Fig. 4.13. Arhitectura logică a RS

În diagonală componentele noului vector vor fi reflectate spre nord. Acest moment al schimbării direcției spre nord coincide cu startul următorului ciclu de recurență. Figura 4.14 ilustrează reflexia în diagonală a componentelor noului vector.

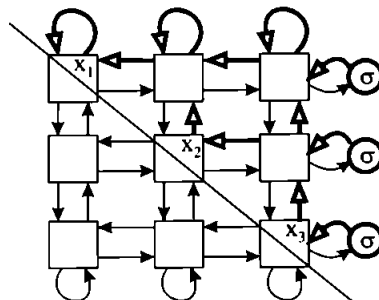
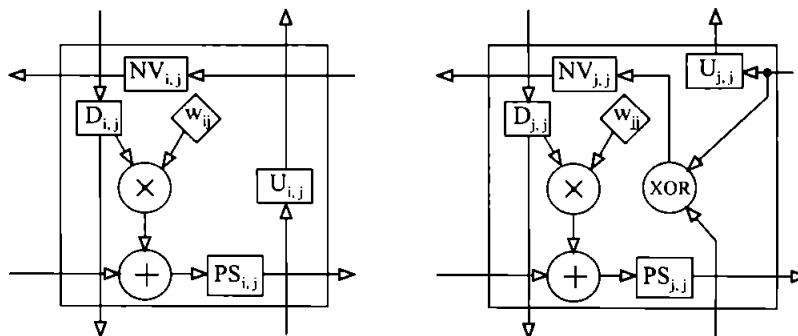


Fig. 4.14. Propagarea vectorului x cu schimbarea direcției pe diagonală

Vom defini în continuare arhitectura logică a CS diferențiind celulele diagonale de cele non-diagonale. Pe baza acestor definiții arhitecturale, vom detalia toate calculele realizate într-o celulă.

4.2.4. Specificațiile celulei sistolice

Arhitectura generală a RS a fost prezentată în fig. 4.14. Vom preciza în detaliu operațiile efectuate în fiecare celulă, precum și căile de comunicație asociate fiecărei celule. Figura 4.15 pune în evidență aceste aspecte.



Transmisii date:

$$\begin{aligned} D_{i,j} &= D_{i-1,j} \\ U_{i,j} &= U_{i+1,j} \\ NV_{i,j} &= NV_{i,j+1} \end{aligned}$$

Calcul:

$$PS_{i,j} = PS_{i,j-1} + w_{ij} \cdot D_{i,j}$$

a) non - diagonale

Transmisii date:

$$\begin{aligned} D_{j,j} &= D_{j-1,j} \\ U_{j,j} &= N_{j,j+1} \end{aligned}$$

Calcul:

$$\begin{aligned} NV_{j,j} &= U_{j+1,j} \oplus NV_{j,j+1} \\ PS_{j,j} &= PS_{j,j-1} + w_{jj} \cdot D_{j,j} \end{aligned}$$

b) diagonale

Fig. 4.15. Specificațiile celulelor sistolice

Celula diagonală realizează virtual trecerea de la un pas de recurență la pasul următor, prin încărcarea noii valori NV (noul vector x) în registrul U (Up). Din acest moment, noua valoare devine valoarea curentă (actuală) a vectorului x deoarece s-a declanșat pasul de calcul al unui nou vector x .

Pentru completarea specificațiilor din figura 4.15 mai trebuie precizate două condiții de limită. Condiția aferentă laturii de nord care realizează rebuclajul se exprimă prin:

$$D_{1,i} = U_{1,i}$$

iar condiția de inițializare aferentă laturii de vest va fi:

$$PS_{i,1} = 0$$

ceea ce înseamnă că se demarează pasul recurențial cu o sumă parțială nulă. Vom vedea că această condiție de inițializare cu zero nu va mai fi respectată în faza de învățare.

Complexitatea circulației datelor în cadrul RS reclamă o etapă de sinteză, care ne va permite introducerea unui simbolism de care vom face uz în cele ce urmează. Acesta ne va permite să evidențiem mișcarea conjugată a datelor care circulă în RS.

Prima etapă constă în introducerea vectorului înmulțitor (x_r) în rețea. Această etapă va consuma N unități de timp.

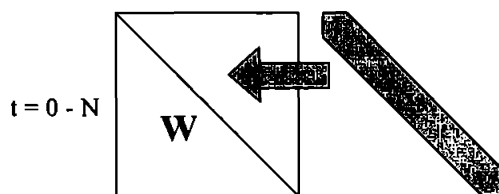


Fig. 4.16. Încărcarea vectorului x_r

În momentul în care vectorul x_r atinge diagonala, va continua să se propage spre nord (reflectat spre nord), până pe latura de nord a rețelei. Aici este rebuclat (reflectat) pe direcția sud. Din momentul rebuclării, sumele parțiale (PS) vor fi elaborate progresiv și propagate dinspre vest spre est (fig. 4.17).

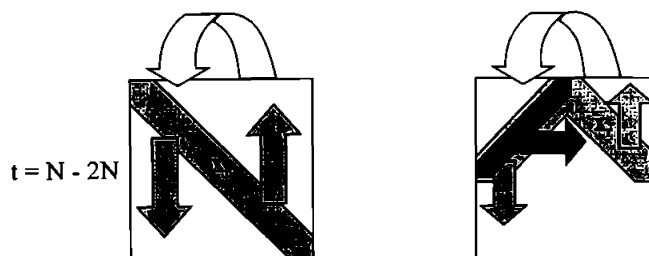


Fig. 4.17. Rebuclarea vectorului x_r pe latura de nord și generarea sumelor parțiale PS

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

În momentul în care prima sumă parțială (devenită finală) atinge latura de est, prima componentă a vectorului x_t participă la elaborarea PS atinge latura de sud (figura 4.18).

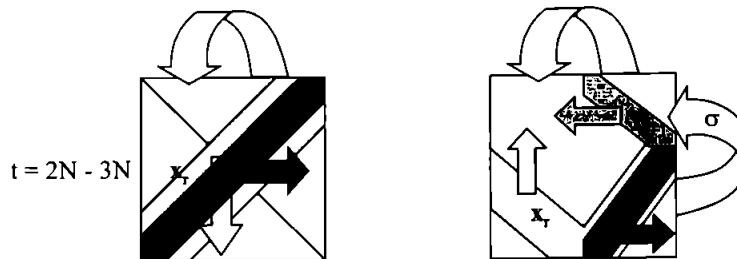


Fig. 4.18. Elaborarea sumelor parțiale PS și aplicarea funcției de activare σ

Din momentul în care prima sumă parțială (finală) atinge latura de est, începe generarea componentelor noului vector x_{t+1} prin aplicarea funcției neliniare σ . Componentele noului vector vor fi repropagate spre vest pentru a se întâlni pe diagonală cu componentele vechiului vector x_t (care revin din direcția sud). Compararea celor 2 vectori (pentru detecția convergenței) se va face deci în celulele diagonale (fig. 4.19).

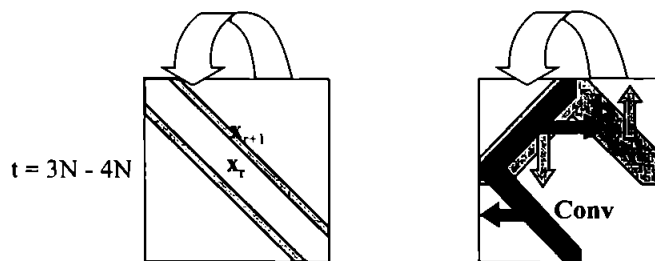


Fig. 4.19. Revenirea pe diagonală și calculul convergenței

Elaborarea semnalului de convergență global se realizează în paralel cu evoluția calculelor pentru al doilea pas de recurență. Acest semnal de convergență este procesat sistolic dinspre nord spre sud (lanțul de porți și elemente de memorare din figura 4.12) și procesarea sa va dura N cicluri (timpul în care prima sumă parțială parcurge distanța de la diagonală până la marginea de est a RS). În cazul algoritmului Hopfield, dacă se atinge convergența, starea nu mai variază (sumele parțiale nu se mai modifică de la un pas de convergență la altul). Produsele suplimentar calculate nu vor modifica rezultatele. Vectorul care sosește pe latura de est poate deci fi extras, componentă cu componentă imediat după detecția convergenței. În general, această proprietate va fi exploatată pentru gestionarea operațiilor de intrare / ieșire ale RS, pentru toți algoritmi care fac apel la o recursivitate a cărei profunzime este limitată de detecția unei stări stabile.

Această prezentare geometrică a circulației datelor în diferitele faze ale algoritmului nu constituie o demonstrație a validității acestuia. Pentru a demonstra validitatea algoritmului vom relua ecuațiile elaborate în fig. 4.15 și le vom exprima sub forma unui Sistem de Ecuații Recurente Uniforme (SERU). Vom demonstra apoi că operațiile asociate transferurilor de date realizează corect calculul dorit. Acest principiu va fi reluat pentru a valida, în secțiunile următoare, algoritmi destinați fazei de învățare.

4.2.5. Descrierea Sistemului de Ecuații Recurente Uniforme

Pentru rețeaua *Hopfield*, funcția de recunoaștere se exprimă prin ecuația recurentă:

$$x_{i,r} = \sigma \left(\sum_{j=1}^N w_{ij} \cdot x_{j,r-1} \right)$$

unde:

- i și j sunt indicii spațiali
- r este pasul de recurență
- t este indice temporal (va fi utilizat în continuare)
- N ordinul matricii pătrate **W** (matricea de coeficienți sinaptici)
- σ este funcția de transfer neliniară

Vom scrie ecuațiile SERU care vor exprima funcția fiecărui element component al celulei sistolice precum și comunicațiile acestuia. Această descriere nu este posibilă decât pentru o RS localizată.

Sumele parțiale evoluează în rețea dinspre latura de vest spre latura de est (ec. 4.16); aceste sume parțiale sunt inițializate pe latura de vest prin injectarea unor valori nule (ec. 4.17).

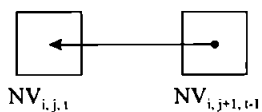
$$PS_{i,j,t} = PS_{i,j-1,t-1} + w_{i,j} \cdot D_{i,j,t-1} \quad (4.16)$$

$$PS_{i,1,t} = 0 \quad (4.17)$$

Noile valori (componentele noului vector **x**) sunt introduse din direcția est și tranzitează RS spre vest. Aceste valori se obțin pe latura de est (în exteriorul RS) prin aplicarea funcției neliniare σ sumelor parțiale obținute în celulele de pe latura de est a RS.

$$NV_{i,N,t} = \sigma(PS_{i,N,t-1}) \quad (4.18)$$

În interiorul rețelei noile valori tranzitează spre vest (fără a fi modificate) până la diagonală:



$$NV_{i,j,t} = NV_{i,j+1,t-1} \quad (4.19)$$

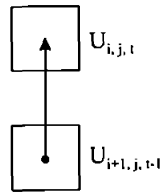
Începând de la diagonală și până pe latura de vest, noile valori vor reprezenta de fapt rezultatul comparațiilor locale efectuate în celulele diagonale (convergența):

$$NV_{j,j,t} = U_{j+1,j,t-1} \oplus NV_{j,j+1,t-1} \quad (4.20)$$

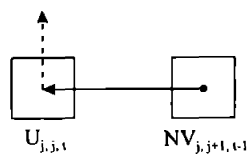
Valorile care urcă traversând rețeaua dinspre sud spre nord sunt descrise de ecuația (4.21). Acestea sunt noile valori recepționate în diagonală dinspre est și reflectate în diagonală spre nord (ec. 4.22).

$$U_{i,j,t} = U_{i+1,j,t-1} \quad (4.21)$$

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

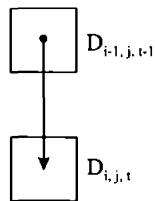


$$U_{j,j,t} = NV_{j,j+1,t-1} \quad (4.22)$$

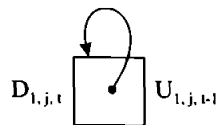


Valorile descendente care traversează rețeaua dinspre nord spre sud (ec. 4.23) provin din valorile care au urcat spre nord și care au fost rebuflate pe latura de nord a RS (ec. 4.24):

$$D_{i,j,t} = D_{i-1,j,t-1} \quad (4.23)$$



$$D_{1,j,t} = U_{1,j,t-1} \quad (4.24)$$



Vom urmări căile de acces la o celulă sistolică cu scopul de a ajunge la marginile RS; cunoaștem de asemenea valorile de intrare și momentele de timp la care acestea sunt aplicate la marginile RS. Vom analiza mai întâi traseul notat cu 2 în figura 4.20. Vom porni de la celula (i, j) și vom urca până la latura de nord a RS aplicând succesiv ecuația (4.23):

$$D_{i,j,t} = D_{1,j,t+i} \quad (4.25)$$

Pe latura de nord, valorile care urcă dinspre sud (U), sunt rebuflate din nou spre sud, pe calea descendentă (D). Acest rebuflaj descris de ecuația (4.24) ne permite să scriem:

$$D_{1,j,t+i} = U_{1,j,t+i} \quad (4.26)$$

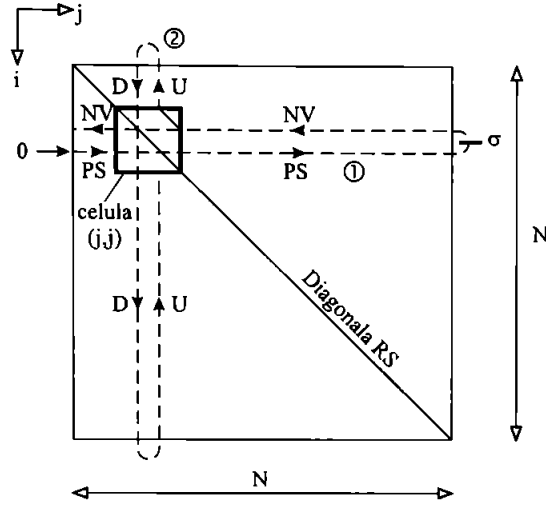


Fig. 4.20. Căile de acces la o celulă sistolică

Calea ascendentă (U) pornește din diagonală având ca sursă calea NV de pe traseul notat cu 1. Vom aplica succesiv ecuația (4.21) până la diagonală:

$$U_{i,j,t-i} = U_{i+j-1,j,t-i-(j-1)}$$

$$U_{i,j,t-i} = U_{j,j,t-i+j+1} \quad (4.27)$$

Valorile ascendente (U) se obțin din valorile NV calculate la marginea de est a RS (prin aplicarea funcției neliniare σ) și care sunt emise spre vest până la diagonală (ec. 4.22).

$$U_{j,j,t-i+j+1} = NV_{j,j+1,t-i-j} \quad (4.28)$$

În sfârșit, dacă aplicăm succesiv ecuația (4.19) pe calea NV (traseul 1) până la marginea de est a RS, obținem:

$$NV_{j,j+1,t-i-j} = NV_{j,j+1+(N-j-1),t-i-j-(N-j-1)}$$

$$NV_{j,j+1,t-i-j} = NV_{j,N,t-i-N+1} \quad (4.29)$$

Combinând ecuațiile (4.25) și (4.29) obținem:

$$D_{i,j,t} = NV_{j,N,t-i-N+1} \quad (4.30)$$

Vom substitui în ecuația (4.16) expresia lui $D_{i,j,t}$ obținută în (4.30):

$$PS_{i,j,t} = PS_{i,j-1,t-1} + w_{i,j} \cdot NV_{j,N,t-i-N} \quad (4.31)$$

Pe latura de est a RS, conform (4.16), vom avea valorile:

$$PS_{i,N,t} = \sum_{j=1}^N w_{i,j} \cdot D_{i,j,t-N+j-1}$$

care în baza relației (4.30) devin:

$$PS_{i,N,t} = \sum_{j=1}^N w_{i,j} \cdot NV_{j,N,t-2N+i+j} \quad (4.32)$$

Aplicând (4.18) și (4.19) obținem:

$$NV_{i,N,t} = \sigma \left(\sum_{j=1}^N w_{i,j} \cdot NV_{j,N,t-2N+i+j-1} \right) = \sigma \left(\sum_{j=1}^N w_{i,j} \cdot NV_{j,N-1,t-2N+i+j} \right) \quad (4.33)$$

Considerând originea timpului momentul introducerii primei componente din vectorul \mathbf{x} în RS (starea inițială prezentată în fig. 4.16), și observând că un pas de recurență are durată de $2N$ unități de timp, atunci a i -a componentă a vectorului \mathbf{x} va atinge latura de est a rețelei (venind dinspre vest) la momentele $2rN + i$. Deci:

$$x_{i,r} = NV_{i,N,2rN+i} \quad (4.34)$$

Utilizând relația (4.34) vom putea rescrie ecuația (4.33) astfel:

$$NV_{i,N,2rN+i} = \sigma \left(\sum_{j=1}^N w_{i,j} \cdot NV_{j,N-1,2N(r-1)+j} \right)$$

Rezultă astfel:

$$x_{i,r} = \sigma \left(\sum_{j=1}^N w_{i,j} \cdot x_{j,(r-1)} \right) \quad (4.35)$$

Am arătat deci că SERU definit realizează corect produsele succesive $\mathbf{W} \cdot \mathbf{x}$ (pașii de recurență sunt corect executați). Operațiile descrise sunt implementate în celulele sistolice și, prin urmare RS va efectua corect calculul iterativ. Demonstrarea validității mecanismului de detecție a convergenței se face după același principiu.

Demonstrația prezentată nu se referă decât la faza de recunoaștere aferentă unui singur vector de intrare. Gradul de utilizare a celulelor sistolice este dat de raportul dintre numărul de celule active și numărul total de celule disponibile (pe durata unui pas de recurență). Gradul mediu de utilizare se poate calcula la nivelul unui pas de recurență, contorizând celulele utilizate pe parcursul celor 6 etape conținute într-un ciclu de recurență; gradul mediu va fi foarte redus: $1/2 N$. Totuși, performanțele rețelei pot fi mult îmbunătățite și se poate ajunge chiar la un grad de utilizare de 100 %, dacă se ia în considerare o procesare în tehnică *pipeline* a vectorilor de intrare. În acest scop, vom considera un șir de vectori aplicați succesiv la intrarea rețelei. Fiecare vector este notat cu un indice p care indică numărul său de ordine în cadrul șirului. SERU definit anterior nu se va modifica. Știind că un pas de recurență are durată $2N$ unități de timp, a i -a componentă a vectorului de indice p va ieși din RS la momentul $2rN + i + p$ (cu $p \leq 2N$).

Ecuția (4.34) va deveni:

$$x_{i,r,p} = NV_{i,N,2rN+i+p} \quad p \leq 2N \quad (4.36)$$

Dacă substituim (4.36) în (4.33) obținem:

$$x_{i,r,p} = \sigma \left(\sum_{j=1}^N w_{ij} \cdot NV_{j,N-1,2N(r-1)+j+p} \right) \quad (4.37)$$

și, în final:

$$x_{i,r,p} = \sigma \left(\sum_{j=1}^N w_{ij} \cdot x_{j,r-1,p} \right) \quad p \leq 2N \quad (4.38)$$

După cum se poate constata din (4.38), indicele de aplicare a vectorilor p rămâne neschimbat pentru fiecare pas de recurență. Aceasta înseamnă că nu poate să apară confuzie între vectori și procesarea în tehnică *pipeline* se va efectua corect. În aceeași manieră se poate arăta că *pipelining*-ul nu va afecta funcționarea corectă a mecanismului de detecție de convergență.

4.3. Extensia arhitecturii pentru învățare

Principalii algoritmi conexioniști sunt constituiți din operații matriceale și vectoriale simple și din transformări neliniare (funcția de activare). În cele ce urmează, vom arăta că regula Perceptronului, regula lui *Hebb* și cea a lui *Kohonen* pot fi implementate sub forma unui algoritm sistolic. Pentru toate modelele, faza de recunoaștere este identică și se rezumă la un simplu produs între matricea sinaptică W și vectorul de intrare x , urmat de aplicarea funcției neliniare σ . Diferențele apar la faza de învățare.

4.3.1. Resursele comune diferitelor modele

Regulile de învățare se bazează pe câteva operații matriceale și vectoriale rezumate în tabelul 4.1.

Denumire Operație	Expresie matematică	Identificator
Produsul unui vector cu transpusul său = matrice	$v \cdot v^T$	(1)
Sumă / Diferență matriceală	$W_1 \pm W_2$	(2)
Funcția neliniară	$\sigma(v)$	(3)
Produsul matrice - vector	$W \cdot v$	(4)

Tabelul 4.1. Operațiile de bază utilizate în faza de învățare

unde:

- v - vector coloană
- W - matrice
- σ - funcție neliniară

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

Toate regulile de învățare vor putea fi construite cu ajutorul acestor operații de bază (set de primitive). Tabelul 4.2 prezintă câteva reguli de învățare precum și operațiile de bază utilizate pentru implementarea acestora.

Regula de învățare	Faza de recunoaștere	Faza de învățare	Operații de bază
<i>Hebb</i>	$y = \sigma(\mathbf{W} \cdot \mathbf{x})$	$\mathbf{W} := \mathbf{W} + \mathbf{x} \cdot \mathbf{x}^T$	(1) (2) (3) (4)
<i>Perceptron, Adaline</i>	$y = \sigma(\mathbf{W} \cdot \mathbf{x})$	$\mathbf{W} := \mathbf{W} + \alpha \cdot (\mathbf{d} - \mathbf{W} \cdot \mathbf{x}) \cdot \mathbf{x}^T$	(1) (2) (3) (4)
<i>Kohonen</i>	$y = \sigma(\mathbf{W} \cdot \mathbf{x})$	$\mathbf{W}_i := \mathbf{W}_i + \alpha \cdot \lambda(\Delta_{\text{topo}}(i, I))(\mathbf{x}^T - \mathbf{W}_i)$	(2) (3) (4)
<i>Back Propagation</i>	$y = \sigma(\mathbf{W} \cdot \mathbf{x})$	$\delta^{[s]} = \sigma'(\mathbf{p}^{[s]}) \otimes (\mathbf{d} - \mathbf{y})$ - strat de ieșire $\delta^{[q]} = \sigma'(\mathbf{p}^{[q]}) \otimes (\mathbf{W}^{[q+1]})^T \cdot \delta^{[q+1]}$ - straturi ascunse $\mathbf{W} := \mathbf{W} + \alpha \cdot \delta^{[q]} \cdot (\mathbf{x}^{[q]})^T$	(1) (2) (3) (4)

Tabelul 4.2. Câteva reguli de învățare

unde:

- \mathbf{x} - vectorul aplicat la intrare
- \mathbf{d} / \mathbf{y} - vectorul dorit / real de ieșire
- \mathbf{W} - matricea de coeficienți sinaptici
- \mathbf{W}_i - linia i a matricii \mathbf{W}
- α - constanta de învățare
- λ - funcția de vecinătate (restricționează actualizarea ponderilor doar pentru neuronii i aflați în vecinătatea învingătorului I)
- $\Delta_{\text{topo}}(i, I)$ - distanța între neuronul i și învingătorul I în spațiul topologic definit
- σ - funcția de transfer neliniară (funcția de activare)
- $\delta^{[s]}$ - vectorul de eroare asociat stratului de ieșire
- $\delta^{[q]}$ - vectorii de eroare asociați straturilor interne (ascunse)
- $\sigma'(\mathbf{p}^{[q]})$ - vectorul derivatelor funcției de transfer în raport cu potențialele neuronilor din stratul $[q]$

Am arătat deja în paragraful alocat funcției de recunoaștere că operația (4) poate fi implementată pe rețeaua sistolică definită. Operația (3) a fost rejectată în exteriorul RS dar poate fi aplicată facil datorită accesibilității la datele care circulă în interiorul RS (în momentul rebucării acestora pe latura de est). În cele ce urmează vom arăta că și învățarea (după diferitele reguli) este implementabilă pe RS definită.

4.3.2. Regula lui Hebb

Vom considera deocamdată regula lui *Hebb* ca un mecanism de calcul și nu ca unul de învățare (pentru simplificare). Performanțele acestei reguli sunt insuficiente pentru a permite aplicații concrete. Totuși, mecanismul definit de regula lui *Hebb* este comun și altor reguli mai evaluate (de exemplu regula "*Delta projection*") capabile de performanțe superioare.

Pentru actualizarea ponderilor w_{ij} trebuie să dispunem, în fiecare celulă (i, j) , de valorile x_i, x_j . Aceasta se va realiza prin introducerea vectorului \mathbf{x} în diagonala RS și obligându-l să circule simultan spre Nord și spre Vest. Baleiajul realizat de cele două direcții perpendiculare va pune la dispoziția fiecărei celule sistolice (i, j) , dubletele (x_i, x_j) . Va fi suficient să se calculeze produsul $x_i \cdot x_j$ și să se adauge valoarea obținută la valoarea coeficientului sinaptic local. Circulația datelor este prezentată în figura 4.21.

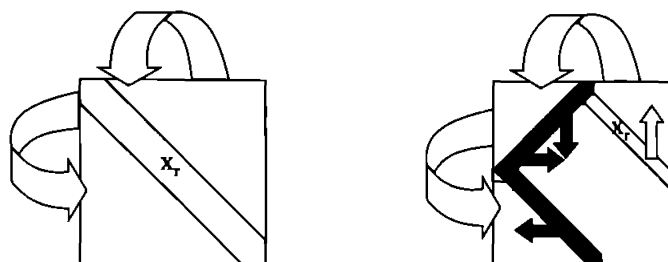


Fig. 4.21. Circulația datelor în RS pentru regula lui Hebb

Să notăm că, pe durata fazei de învățare, ponderile sinaptice W nu trebuie să interacționeze cu datele care traversează RS de la Vest spre Est. Dispozitivul de înmulțire din CS este deci liber și poate fi utilizat pentru calculul produselor locale $x_i \cdot x_j$. Fiecare CS trebuie să memoreze activitatea neuronilor i și j , să calculeze produsul $x_i \cdot x_j$ și să adauge rezultatul la valoarea ponderii sinaptice locale. Aceste operații sunt realizabile cu o CS a cărei structură internă este descrisă în figura 4.22. a).

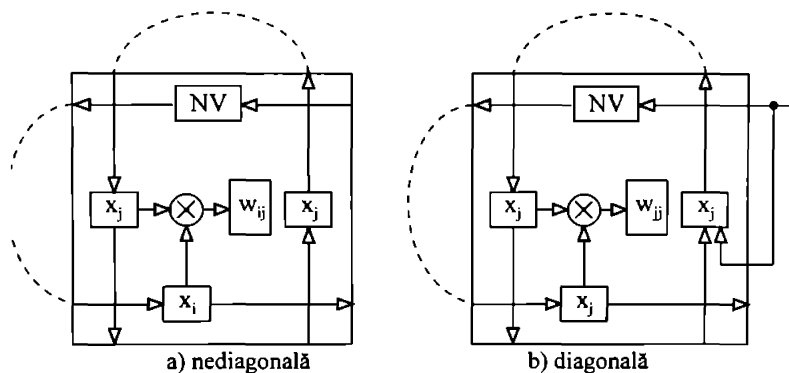


Fig. 4.22. Resursele interne ale CS

Circuitele de înmulțire și împărțire au fost prevăzute și pentru faza de recunoaștere. Principalele diferențe se vor referi la secvențiere și la liniile de date prin care se accesează resursele interne ale CS. Ca și în cazul fazei de recunoaștere, celulele diagonale au o structură de comunicație ușor diferită de celelalte celule. Celulele diagonale trebuie să asigure retransmiterea componentelor vectorului de stare spre latura de nord a RS. Dar, pentru faza de învățare, ea trebuie să rămână de asemenea transparentă pentru a transmite vectorul de stare spre latura de vest. Aceasta permite o retransmitere care sincronizează baleiajul planului de ponderi sinaptice de către cei doi vectori de stare identici (dar pe direcții perpendiculare). Aceste modificări sunt redată în figura 4.22. b), unde se poate observa legătura dintre registrul NV și registrul care emite componenta x_j pe direcția nord. Această legătură se va putea realiza cu ajutorul unui multiplexor suplimentar în structura celulei.

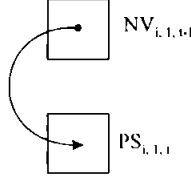
În cele ce urmează vom arăta că această structură definită poate calcula produsul unui vector cu transpusul său și că matricea rezultat poate fi adăugată matricii curente de ponderi sinaptice.

4.3.3. Funcționarea RS pentru faza de învățare Hebb

Vom defini ecuațiile SERU aferente acestei faze. Sumele parțiale vor evolua în RS dinspre marginea de vest spre cea de est (ec. 4.39). Sumele parțiale injectate pe latura de vest a rețelei se obțin din valorile NV care tranzitează RS dinspre est spre vest (ec. 4.40).

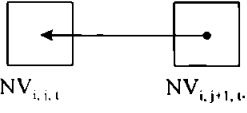


$$PS_{i,j,t} = PS_{i,j-1,t} \quad (4.39)$$



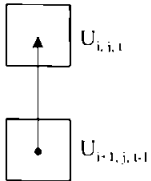
$$PS_{i,j,t} = NV_{i,j,t} \quad (4.40)$$

Noile valori (NV) sunt introduse prin latura de est a rețelei și tranzitează RS până la marginea de vest.

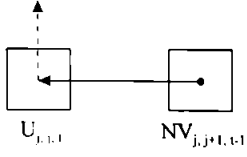


$$NV_{i,j,t} = NV_{i,j+1,t} \quad (4.41)$$

Valorile ascendente traversează rețeaua dinspre sud spre nord (ec. 4.42) și sunt obținute din noile valori (NV) reflectate pe direcția nord de către celulele diagonale (ec. 4.43).

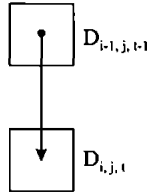


$$U_{i,j,t} = U_{i-1,j,t} \quad (4.42)$$

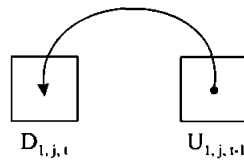


$$U_{j,j,t} = NV_{j,j-1,t} \quad (4.43)$$

Valorile descendente tranzitează din nord spre sud (ec. 4.44), și provin din valorile ascendente care sunt rebucate pe latura de nord (ec. 4.45).

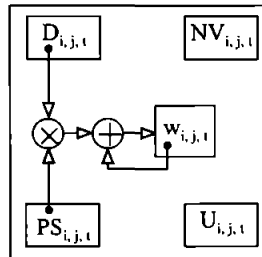


$$D_{i,j,t} = D_{i-1,j,t-1} \quad (4.44)$$



$$D_{i,j,t} = U_{i,j,t-1} \quad (4.45)$$

Fiecare pondere sinaptică este actualizată local cu produsul dintre valoarea descendentă și cea tranzitată prin registrele PS.



$$W_{i,j,t} = W_{i,j,t-1} + PS_{i,j,t-1} \cdot D_{i,j,t-1} \quad (4.46)$$

Pentru o celulă (i, j) , vom relua traseul în spațiu și timp până la atingerea unei laturi a RS. Vom prezenta mai întâi traseul notat cu 1 în figura 4.23. Pornind dintr-o celulă oarecare (i, j) , vom urmări traseul spre direcția vest până la rebucraj. Această operație va fi efectuată prin aplicarea ecuației (4.39) pentru $j < 1$.

$$PS_{i,j,t} = PS_{i,j-(j-1),t-(j-1)} = PS_{i,1,t-j+1}$$

Am ajuns pe latura de vest a rețelei și aici vom aplica ecuația (4.40):

$$PS_{i,j,t} = NV_{i,1,t-j} \quad (4.47)$$

Refacem apoi traseul până pe latura de vest a RS, pe traseul NV, prin aplicarea succesivă a ecuației (4.41).

$$NV_{i,1,t+j} = NV_{i,1+(N-1),t+j-(N-1)} = NV_{i,N,t+j-N+1} \quad (4.48)$$

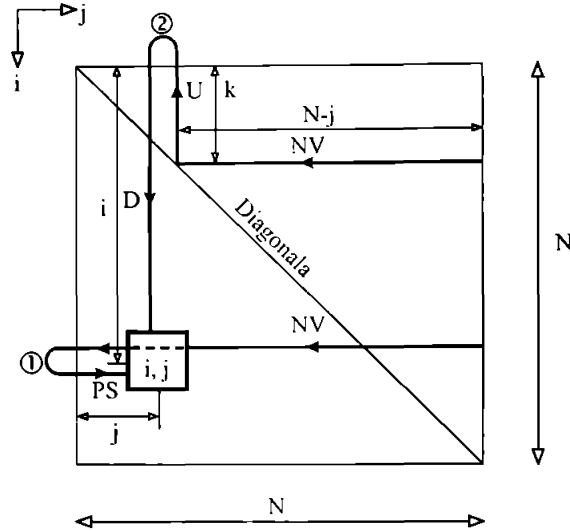


Fig. 4.23. Căile de acces la celula (i, j)

Deci, pentru traseul 1 am obținut:

$$PS_{i,j,t} = NV_{i,N,t+j-N+1} \quad (4.49)$$

Vom urmări acum traseul 2, așa cum este indicat în figura 4.23. Aplicând ecuația (4.44) până la latura de nord a RS ($i < 1$), obținem:

$$D_{i,j,t} = D_{i-(i-1),j,t-(i-1)} = D_{1,j,t+i-1}$$

Rebuclajul pe latura de nord se va face conform ecuației (4.45):

$$D_{1,j,t+i-1} = U_{1,j,t+i}$$

Coborâm apoi până la diagonala RS, aplicând succesiv ecuația (4.44):

$$U_{1,j,t+i} = U_{1-(j-1),j,t+i-(j-1)} = U_{j,j,t+i-j+1}$$

Reflexia direcției pe diagonală se face după ecuația (4.43):

$$U_{j,j,t+i-j+1} = NV_{j,j+1,t+i-j}$$

Pe baza ecuației (4.41), pe care o aplicăm succesiv până la $j = N$, vom reface traseul până la latura de est a RS:

$$NV_{j,j+1,t+i-j} = NV_{j,j+1+(N-j-1),t+i-j-(N-j-1)} = NV_{j,N,t+i-N+1}$$

Rezumând, pentru traseul 2 obținem:

$$D_{i,j,t} = NV_{j,N,t-i-N+1} \quad (4.50)$$

Aplicând ecuația (4.46), obținem ecuația de actualizare a ponderilor sinaptice:

$$w_{i,j,t} = w_{i,j,t-1} + NV_{i,N,t-j-N} \cdot NV_{j,N,t-i-N} \quad (4.51)$$

Din expresia (4.51) deducem că indicii temporali aferenți variabilelor NV depind de localizarea spațială a celulei. Dacă vom considera că toate componentele vectorului sunt aplicate simultan pe latura de est a rețelei, celula nu va recepționa corect componentele (x_i, x_j) ale vectorului din cauza diferenței de lungime între cele 2 trasee (fig. 4.24).

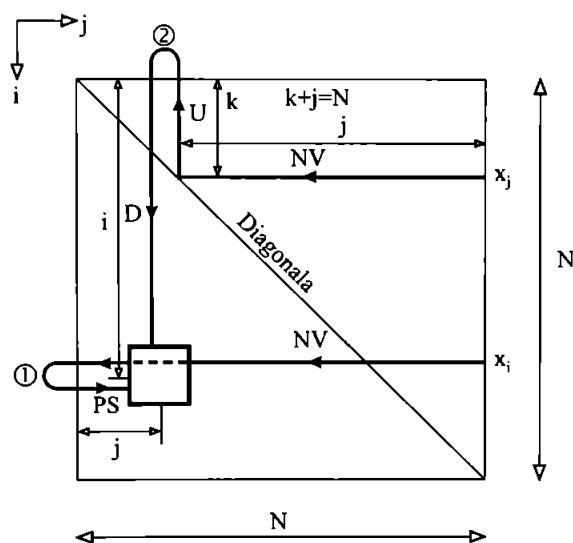


Fig. 4.24. Lungimea celor două căi de acces la celula i, j

Traseul 1 are lungimea $N + j$ iar traseul 2 are lungimea $N + i$. Va trebui să introducem un decalaj temporal în aplicarea componentelor vectorului x care va compensa diferența de lungime. Aplicarea vectorului x se va face astfel:

$$NV_{k,N,t} = x_{k,t-k} \quad (4.52)$$

În concluzie, vom respecta constrângerea fixată deja de a aplica componentele vectorului în paralel și simultan celulelor de pe diagonala RS, așa cum indică figura 4.25.

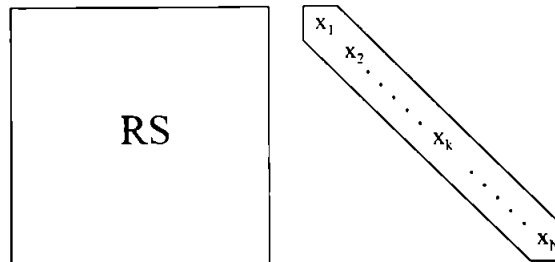


Fig. 4.25. Aplicarea vectorului x

Relația (4.51) poate fi rescrisă utilizând componentele vectorului x :

$$w_{i,j,t} = w_{i,j,t-1} + x_{i,N-t+j-N} \cdot NV_{j,N,t-i-N+j}$$

$$w_{i,j,t} = w_{i,j,t-1} + x_{i,t-i-j-N} \cdot x_{j,t-i-j-N}$$

Dacă eliminăm indicele temporal care este identic pentru cele două componente ale vectorului x , obținem:

$$w_{i,j,t} = w_{i,j,t-1} + x_i \cdot x_j \quad (4.53)$$

Circulația datelor și operațiile locale descrise, asociate cu decalajul temporal la introducerea componentelor vectorului x , permit calculul matricii W conform cu regula lui *Hebb*.

Ansamblul de ecuații (4.39) până la (4.48) constituie un SERU deoarece nu conține decât translatari independente care operează asupra indicilor. Pentru a finaliza demonstrația, mai trebuie definită o funcție de timp T compatibilă cu dependențele definite prin ecuațiile elaborate, și aleasă o funcție de alocare A care va atribui operații concurente diverselor procesoare. Pentru a obține rețeaua descrisă, funcțiile T și A vor fi:

$$T(i, j, t) = t$$

$$A(i, j, t) = (i, j)$$

Să notăm că sunt posibile și alte funcții T și A dar vor conduce la implementări diferite. Aceste funcții corespund vectorilor de proiecție d și s definiți în capitolul 3.

Am arătat că operațiile (1) și (2) din tabelul 4.1 pot fi implementate pe RS bidimensională definită, asigurând funcționarea RS pentru regula lui *Hebb*. Vom arăta că aceste operații vor permite deopotrivă implementarea regulii Perceptronului.

4.3.4. Regula Perceptronului

Algoritmul Perceptronului, prezentat sintetic la 1.5, poate fi simplificat pentru a pune în evidență produsul scalar $W \cdot x$, ca indicator în procesul de ajustare a ponderilor. Regula poate fi rescrisă:

pas 1: Dacă ieșirea dorită d asociată vectorului de intrare x este negativă, schimbă semnul lui x :

- dacă $d = 1$ atunci,
 $z = x$
- dacă $d \neq 1$ atunci,
 $z = -x$

pas 2: Dacă $\mathbf{W} \cdot \mathbf{z} < 0$ atunci,
 $\mathbf{W} = \mathbf{W} + \alpha \cdot \mathbf{z}$
 unde: α - constanta de învățare
 \mathbf{W} - vectorul de ponderi sinaptice de dimensiune P

Primul pas al algoritmului va fi implementat în exteriorul RS. El constă în pregătirea datelor de intrare; semnul vectorului de intrare \mathbf{x} va fi setat la valoarea celui al ieșirii dorite d (se obține vectorul \mathbf{z}). Pasul 2 va fi executat pe RS care va fi configurată ca pentru faza de recunoaștere (paragraful 4.2.1). Rețeaua sistolică va fi identică cu cea prezentată la 4.2.1 dar matricea \mathbf{W} se reduce la un simplu vector de dimensiune P (numărul de intrări în RN). Va exista apoi o singură ieșire (cea care corespunde singurului neuron din structura Perceptronului).

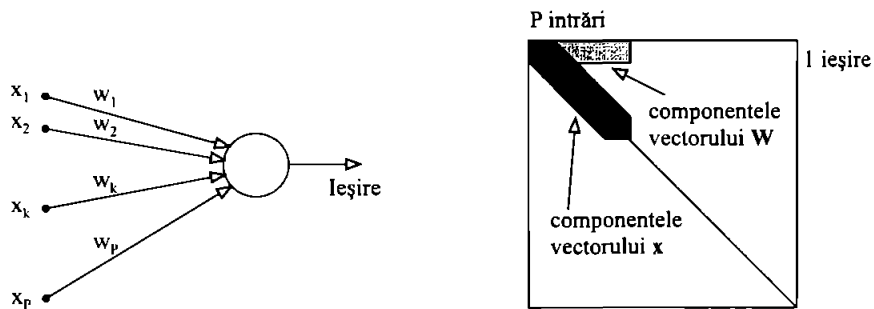


Fig. 4.26. Structura RS pentru regula Perceptronului

Pentru execuția pasului 2 trebuie examinat semnul produsului scalar $\mathbf{W} \cdot \mathbf{z}$ care ne va indica dacă ponderile sinaptice trebuie ajustate sau nu. Această operație va fi realizată simplu printr-o funcție de tip treaptă (fig. 4.27).

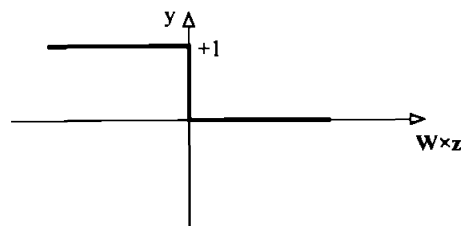


Fig. 4.27. Funcția σ de tip treaptă

Regula de ajustare a ponderilor sinaptice va fi:

$$\mathbf{W} := \mathbf{W} + \alpha \cdot \mathbf{z} \cdot \sigma(\mathbf{W}^T \cdot \mathbf{z})$$

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

unde:

$$\mathbf{W} = \begin{pmatrix} w_1 \\ w_2 \\ \dots \\ w_p \end{pmatrix} \quad \text{- vectorul ponderilor sinaptice}$$

$$\mathbf{z} = \begin{pmatrix} z_1 \\ z_2 \\ \dots \\ z_p \end{pmatrix} \quad \text{- vectorul de intrare (cu semnul ajustat conform pasului 1)}$$

Regula de ajustare a ponderilor revendică operația (1) din tabelul 4.1; corectitudinea execuției acestei operații pe RS definită a fost demonstrată la 4.3.3.

Pe parcursul calculului valorii de ieșire (y), vectorul de intrare (\mathbf{z}) va continua să circule în rețea pe direcția verticală. El va fi disponibil, pentru fiecare CS, exact în momentul în care și indicatorul de modificare $\sigma(\mathbf{W}^T \cdot \mathbf{z})$ va fi prezent în celula respectivă. Pentru ajustarea ponderilor sinaptice va fi suficientă adăugarea valorii $\alpha \cdot z_i$ la valoarea curentă a ponderii w_i . Circulația datelor în RS este prezentată în figura 4.28.

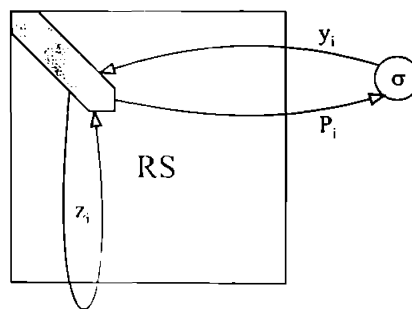


Fig. 4.28. Circulația datelor în RS pentru regula Perceptronului

Arhitectura fiecărei CS va fi identică cu cea dezvoltată pentru regula lui *Hebb*. Sumatorul rămâne nemodificat, iar multiplicarea cu constanta de învățare α (subunitară) se poate realiza printr-o simplă deplasare la dreapta (împărțire cu o constantă supraunitară).

În cazul implementării unui singur Perceptron, RS va fi subutilizată. Am putea exploata paralelismul și independența operațiilor din RS pentru a implementa atâția Perceptroni câte linii de celule există în RS, realizând astfel un sistem asociativ de tip memorie asociativă sau clasificator liniar. Numărul de intrări, pentru fiecare Perceptron, va fi identic și este limitat superior de numărul de coloane din RS (figura 4.29).

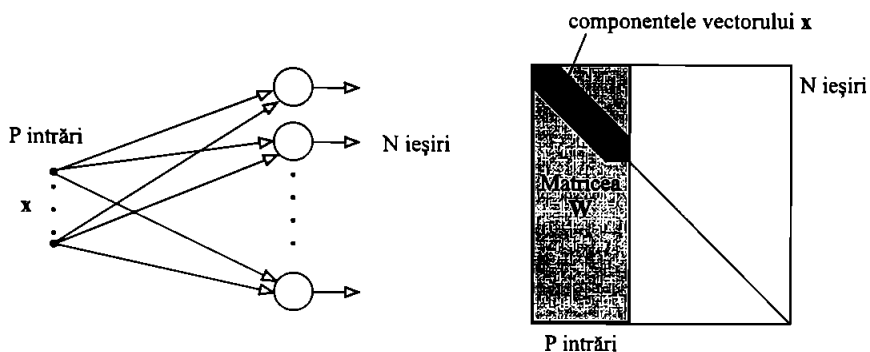


Fig. 4.29. Rețea formată din mai mulți Perceptroni

Această formă de rețea o regăsim în cadrul RN *feedforward* multistrat utilizate frecvent în aplicații. Analiza efectuată până în acest moment ne permite să extindem posibilitățile RS și la o altă regulă de învățare: auto-organizarea elaborată de *Kohonen*.

4.3.5. Regula lui Kohonen

Regula lui *Kohonen* (în varianta inițială), se aplică pe o structură cu două straturi de interconexiuni. Primul strat este stratul de interconexiuni plastice asupra căruia va acționa algoritmul de învățare (paragraful 1.8). Al doilea strat este un strat de interconexiuni laterale, care cablează explicit un operator Laplacian destinat ortogonalizării vectorilor de ieșire. Vom implementa fizic ambele straturi pe aceeași RS bidimensională. Figura 4.30 indică dispunerea datelor în cadrul RS.

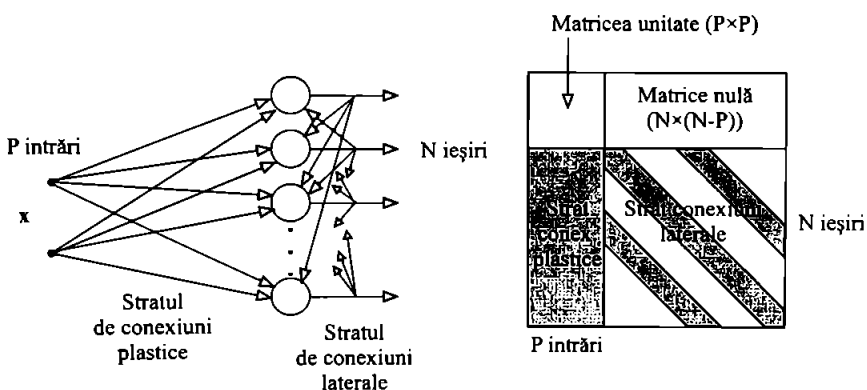


Fig. 4.30. Dispunerea matricilor de interconexiuni în cadrul RS

RS va conține patru matrici de coeficienți. Matricea unitate de dimensiune $(P \times P)$ va servi la recopierea valorilor vectorului de intrare cu scopul de a facilita circulația acestuia în cadrul RS. Matricea nulă de dimensiune $(N \times P)$ va interzice toate interacțiunile dintre valorile de intrare și potențialele neuronilor în timpul fazei de formare a zonei din jurul neuronului învingător (zona de activitate maximă). Matricea stratului plastic de dimensiune $(N \times P)$ va conține toate ponderile

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

sinaptice aferente stratului de conexiuni plastice. Matricea de conexiuni laterale este o matrice formată din benzi diagonale, de dimensiune $(N \times N)$ și care definește vecinătatea de interacțiune dintre neuroni (funcția de vecinătate). Ea va fi utilizată în faza de formare a zonei de activitate maximă. După definirea structurii RS, vom descrie traseele de comunicație prin care se vor pune la dispoziția fiecărei CS informațiile necesare ajustării ponderilor sinaptice.

Regula de învățare originală formulată de Kohonen se descompune în două faze de calcul:

- Evaluarea activității neuronilor în urma aplicării unui vector de intrare
- Actualizarea ponderilor sinaptice în funcție de activitatea calculată și vectorul de intrare aplicat.

Prima fază se descompune, la rândul ei, în 2 sub-faze:

- Calculul activității inițiale (potențialului) fiecărui neuron în urma aplicării unui vector de intrare. Aceasta înseamnă efectuarea produsului dintre matricea sinaptică W ($N \times P$) și vectorul de intrare $x(P)$; rezultă un vector de lungime N ale cărui componente reprezintă activitatea inițială a fiecărui neuron;
- Evoluția și stabilirea activității neuronilor prin interconexiunile laterale ($N \times N$). Aceasta înseamnă calculul vectorului de activitate (de potențiale) final; această etapă necesită un anumit număr de iterații pentru a asigura instalarea unei activități (potențiale) stabile. Kohonen denumea această fază: formarea "bulei" de activitate.

Cele două sub-faze sunt descrise în figura 4.31.

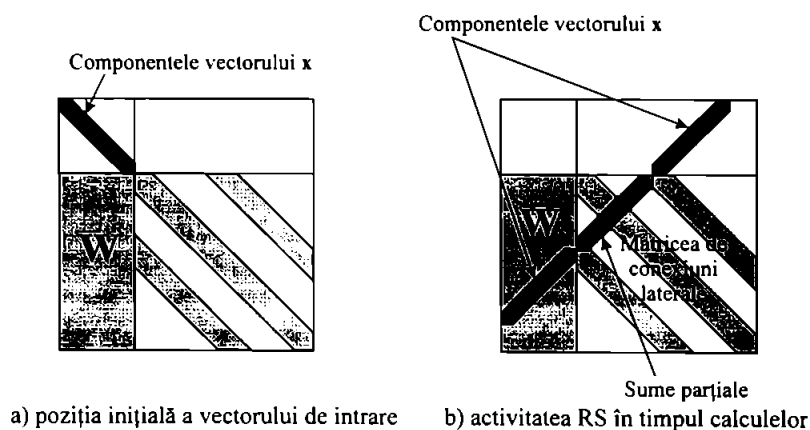


Fig. 4.31. Calculul vectorului de activitate

Procesele ilustrate în fig. 4.31 sunt identice cu cele descrise anterior. Vectorul x este plasat în diagonala RS, tranzitează spre nord unde va fi rebusat spre sud. Când atinge zona de coeficienți sinaptici W , va fi generat vectorul de sume parțiale PS care va evolua până în momentul în care ultima componentă a vectorului x atinge latura de sud a RS.

Procesul de evoluție și stabilizare a "bulei" de activitate are o dinamică identică cu cea a fazei de recunoaștere aferentă rețelei *Hopfield*. Acesta constă în calculul iterativ al produselor dintre matricea de interconexiuni laterale și vectorul de activitate obținut anterior, până la stabilizarea

vectorului de activitate. Am arătat că această operație poate fi implementată pe RS definită (paragraful 4.2). În acest caz funcția neliniară σ nu se va aplica după calculul fiecărui produs.

Pe durata acestor iterații, vectorul de intrare x trebuie conservat și sincronizat cu mișcarea vectorilor de activitate. Structura pătrată a RS, cu căi de comunicație de lungimi egale, asigură această sincronizare, după cum rezultă din figura 4.32.

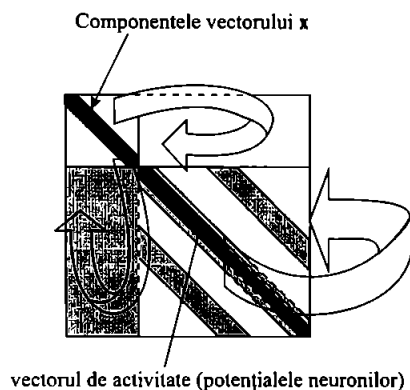


Fig. 4.32. Sincronizarea circulației vectorilor de intrare și de activitate în interiorul RS

După faza de stabilizare a “bulei” de activitate, vom observa că există neuroni cu activitate intensă (potențial ridicat) și neuroni cu activitate slabă. Potențialul neuronilor este comparat cu o valoare de prag, convertit apoi de o funcție treaptă în valorile +1 și 0. Funcția prag transformă potențialul într-un indicator de stare care va indica neuroni ale căror ponderi sinaptice trebuie ajustate. Ajustarea va fi calculată pe baza componentelor vectorului de intrare x , prezent în primele P CS de pe diagonală (fig. 4.33).

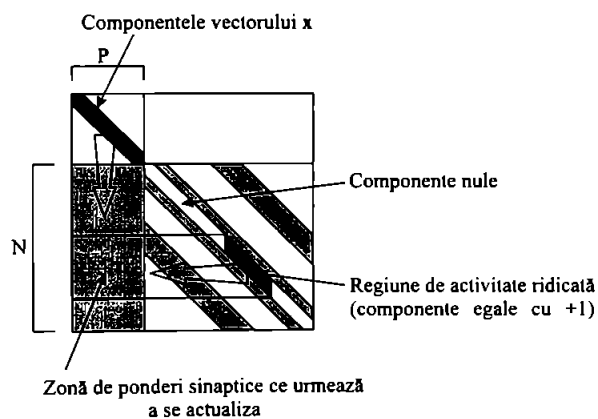


Fig. 4.33. Circulația datelor în RS pentru actualizarea ponderilor sinaptice

Se aplică sistemul de circulație a datelor prezentat anterior; vectorul de activitate (transformat în vector indicator de actualizare în momentul aplicării funcției treaptă) va traversa planul de conexiuni plastice. Zona de intersecție dintre coloanele vectorului de intrare și liniile vectorului

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

indicator de actualizare determină ponderile sinaptice ce vor fi actualizate. Ajustarea se va face după următoarea regulă:

- pentru neuronii i a căror potențial este $+1$:

$$w_{ij} := w_{ij} + \alpha \cdot (x_i - w_{ij}) \quad (4.54)$$

- pentru ceilalți neuroni (care au potențialul 0), ponderile rămân nemodificate.

Considerând vectorul de activitate $y \in [0,1]^N$, formula de actualizare a ponderilor sinaptice devine:

$$w_{ij} := w_{ij} + \alpha \cdot (x_i - w_{ij}) \cdot y_j \quad (4.55)$$

Aceasta înseamnă că, pentru fiecare coeficient w_{ij} , se evaluează diferența dintre componenta x_i a vectorului de intrare și ponderea w_{ij} aferentă. Se calculează apoi produsul dintre diferența obținută și componenta y_j a vectorului de activitate. Dacă $y_j = 1$ atunci ponderea w_{ij} va fi ajustată, dacă $y_j = 0$ atunci ponderea w_{ij} rămâne nemodificată.

Învățarea este deci integrată în sistemul de circulație a datelor. Avantajul metodei rezidă în localizarea fizică a operațiilor care, printre altele, pot fi realizate în paralel. Ca și în cazul funcției de recunoaștere, rețeaua sistolică va admite un flux de vectori x aplicați în sistem *pipeline* pe intrare. RS, compusă din $(N + P)^2$ celule sistolice poate conține simultan $2(N + P)$ vectori de intrare procesați în tehnică *pipeline*. Suplimentar, structura care se bazează numai pe conexiuni locale permite o extensie simplă a RS prin simpla interconectare a mai multor subrețele sistolice.

Vom prezenta în cele ce urmează forma matricii de interconexiuni laterale și vom discuta câteva probleme legate de normalizarea datelor.

A. Construcția matricii de conexiuni laterale

Faza de formare a bulei se bazează pe un mecanism de competiție între neuroni. Competiția se realizează printr-o interacțiune între neuroni. Interacțiunea se definește prin valoare și prin rază (în general o funcție de tipul "Pălărie de mexican"). În cazul implementării noastre această interacțiune poate fi modelată printr-o matrice de interconexiuni care realizează o formă de convoluție între celule. Pentru o rețea de N neuroni, această matrice are N^2 elemente. Ea este definită de către utilizator în funcție de aplicația aleasă. Prin această matrice, se pot defini vecinătăți de forme și dimensiuni diferite.

Să considerăm o rețea SOFM formată din 9 neuroni; această rețea poate fi organizată în linie (rețea monodimensională), în plan (rețea bidimensională) sau sub formă de cub (rețea tridimensională). Rețeaua SOFM realizează o proiecție a unui spațiu de dimensiune mare (spațiu de intrare) pe un spațiu de dimensiune redusă (spațiu de ieșire). În spațiul de reprezentare ales, mărimea zonei de interacțiune poate varia, adică fiecare neuron poate interacționa cu vecinii săi până la o anumită distanță.

Matricea de interconexiuni va fi deci o matrice formată din benzi, diagonale în cazul unui spațiu monodimensional și a unei vecinătăți directe. Va fi o matrice pătrată construită astfel:

- se etichetează (numerotează în ordine) fiecare neuron din rețea;
- se completează elementele matricii de interconexiuni cu valorile de interacțiune dintre fiecare cuplu de neuroni.

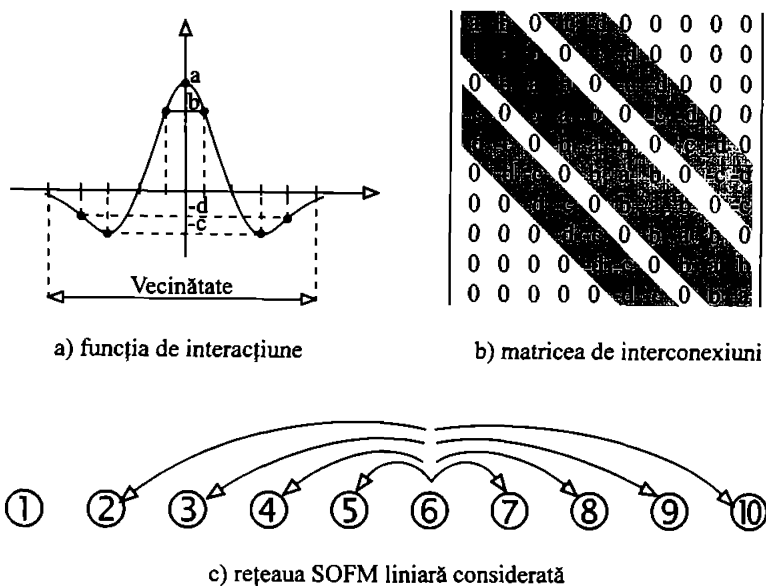


Fig. 4.34. Rețele SOFM monodimensionale; exemplu de vecinătate și matricea de interconexiuni asociată

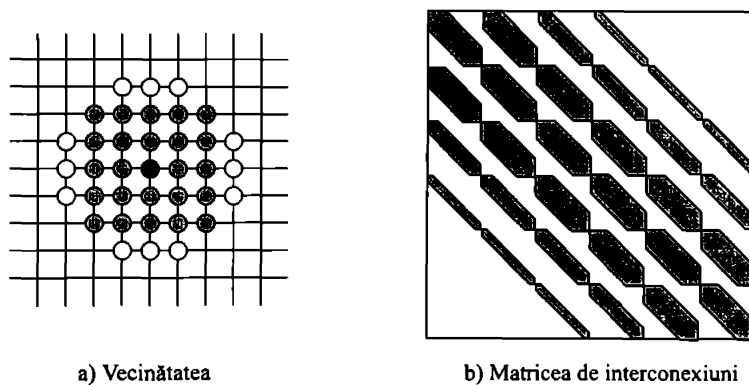


Fig. 4.35. Rețele SOFM bidimensionale; exemplu de vecinătate și matricea de interconexiuni asociată

Soluția utilizării unei matrici de interconexiuni prezintă mai multe avantaje:

- mărimea vecinătății nu este fixă; ea depinde de elementele matricii;
- utilizatorul poate optimiza funcția de interacțiune în funcție de rezultatele obținute;
- RN poate fi extinsă fără nici o dificultate.

Dezavantajul major al acestei soluții rezidă în dimensiunea matricii de interconexiuni; ea are dimensiunea $N \times N$ pentru o rețea SOFM compusă din N neuroni.

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

Validitatea acestui principiu a fost experimentată prin simulări pe o rețea SOFM de 15 x 15 neuroni dispuși în plan. După prezentarea vectorului de intrare toți neuronii au prezentat un nivel de activitate (potențial) dependent de gradul de apropiere (proximitate) dintre vectorul de intrare aplicat și vectorul ponderilor sinaptice aferent neuronului în cauză. Pe durata formării "bulei de activitate" (procesul iterativ care implică matricea de interconexiuni), potențialele au evoluat spre o stare stabilă de tipul celei prezentate în figura 4.36.

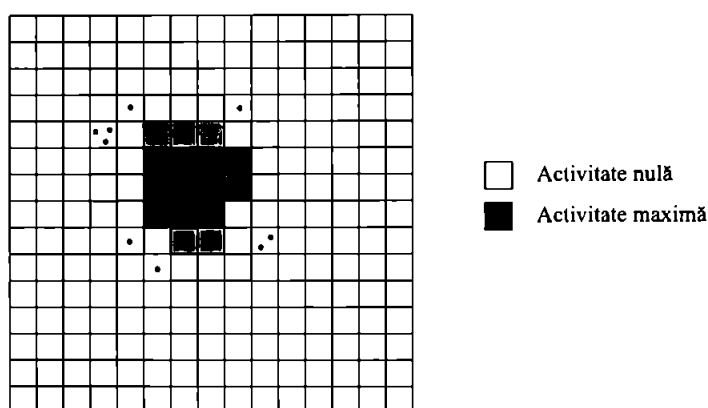


Fig. 4.36. Formarea unei zone de activitate stabilă

Toți neuronii din zona (zonele) de activitate intensă vor suporta procesul de ajustare a ponderilor sinaptice asociate; restul neuronilor nu vor beneficia de ajustare.

B. Normalizarea datelor

Algoritmul lui *Kohonen* se bazează pe măsurarea distanțelor dintre vectorul de intrare și vectorii de ponderi sinaptice. Această măsură se poate realiza prin produs scalar, calcul de distanțe euclidiene, Manhattan etc. Rețeaua sistolică propusă favorizează măsurarea distanțelor cu ajutorul produselor scalare. Vectorii de intrare și de ponderi vor trebui să aibă norme comparabile astfel încât potențialele neuronilor să reprezinte cât mai fidel diferențele unghiulare dintre vectori.

Trebuie deci, teoretic, să se parcurgă un proces de normalizare a vectorilor de ponderi, după fiecare proces de ajustare. Normalizarea impune un calcul global care este incompatibil cu arhitectura distribuită a RS care lucrează exclusiv local. Această constrângere poate fi eliminată ușor dacă se operează cu vectori de intrare normalizați. În acest caz, a fost demonstrat [OJA 82] că norma vectorilor de ponderi sinaptice tinde statistic spre norma vectorilor de intrare. În baza acestui fapt, produsul scalar va reprezenta o măsură corectă a unghiului dintre vectori și deci măsoară corect proximitatea acestora. Această proprietate a fost demonstrată prin simulări.

4.4. Evaluarea Rețelei Sistolice propuse

Calitatea unei implementări depinde de criteriile de evaluare utilizate. Vom utiliza un set de criterii cu scopul de a pune în evidență principalele calități precum și deficiențele soluției sistolice propuse.

RS propusă reprezintă o implementare care integrează sinapse. În contextul unei implementări VLSI, paralelismul masiv reprezintă un factor esențial. Paralelismul trebuie să fie cât mai fin posibil pentru a reproduce structura distribuită aferentă RN. Orice multiplexare a calculelor va conduce la reducerea calităților intrinseci aferente rețelelor conexioniște: reprezentare distribuită și toleranță la defect. RS propusă satisface aceste condiții deoarece fiecare sinapsă este integrată într-o celulă sistolică proprie (cel mai fin paralelism posibil).

Viteza de procesare reprezintă de asemenea un factor determinant care trebuie să justifice soluția propusă. Vom vedea că RS propusă este capabilă de performanțe deosebite (comparativ cu alte implementări) deoarece permite procesarea în tehnică *pipeline* a vectorilor de intrare.

Observabilitatea funcționării este facilitată de circulația datelor în interiorul RS. Potențialele neuronilor pot fi accesate (observate) pe latura de est a RS, înainte și după aplicarea funcției de transfer. Controlul algoritmilor și vizualizarea stării RN sunt de asemenea asigurate. Evoluția ponderilor sinaptice nu poate fi observată "on - line" datorită căilor de date comune, care permit atât încărcarea matricii W cât și tranzitarea componentelor vectorilor de intrare. Această situație conduce la ideea implementării unor căi de date separate pentru încărcarea matricii W . Aceasta va permite implementarea eficientă a RN multistrat deoarece în timpul calculelor aferente stratului curent, în paralel se poate încărca matricea de ponderi aferentă stratului următor. Vom analiza această soluție în capitolul dedicat perfecționării RS propuse.

Integrarea fazei de învățare este posibilă prin compunerea operațiilor de bază executate de RS. RS poate fi configurată după necesități, în ideea implementării celor mai diverse reguli de învățare.

Extensibilitatea rețelei este posibilă datorită interconexiunilor exclusiv locale. Comunicațiile se reduc la legături cu cei mai apropiați vecini ceea ce permite o creștere virtual nelimitată a dimensiunii RS. Vom vedea în cele ce urmează că formatul de reprezentare a datelor reprezintă un factor de limitare fizică, care poate fi în parte compensat prin utilizarea unei aritmetici care ține cont de depășirile de capacitate de reprezentare. Dimensiunea RS propuse (care integrează sinapse) reprezintă un factor limitativ care pare prohibitiv dacă îl comparăm cu soluțiile care integrează neuroni. Numărul de sinapse crește proporțional cu pătratul numărului de neuroni dar această creștere este indispensabilă dacă dorim să asigurăm caracterul de generalitate al RS. În cazul circuitelor neuronale dedicate, va fi posibilă optimizarea numărului de CS în funcție de forma matricii W (diagonală, simetrică etc.).

Simplitatea și identitatea structurii celulelor sistolice reprezintă un mare avantaj în faza de proiectare. Odată construită o CS, ea va fi multiplicată pentru realizarea RS în ansamblu. După amplasament (pe diagonală sau în afara diagonalei RS), comunicațiile CS vor fi specializate și se va putea construi orice formă de rețea.

4.5. Arhitectura RS

4.5.1. Specificațiile celulei sistolice

Fiecare CS execută următoarele operații:

- emisie / recepție date la / de la celulele vecine: Nord, Sud, Est, Vest;
- memorarea coeficientului sinaptic w_{ij} , a rezultatului parțial PS_{ij} , a elementelor vectorilor care se transferă în sens ascendent și descendent în interiorul RS (Up_{ij} , $Down_{ij}$), și a unui element din vectorul rezultat NV_{ij} ;
- calculul unei sume ponderate (PS_{ij}) cu detecția depășirii capacității de reprezentare;
- compararea a două valori pentru detecția convergenței;
- modificarea ponderii sinaptice (w_{ij}) în procesul de învățare.

4.5.2. Formatul operanzilor

Oricare ar fi faza executată în RS (recunoaștere sau învățare), toate CS execută aceleași operații de calcul:

- în cadrul fazei de recunoaștere:

$$PS_{i,j} = PS_{i,j-1} + D_{i-1,j} \cdot w_{i,j} \quad (4.56)$$

- în cadrul fazei de învățare:

$$w_{i,j} = w_{i,j} + \alpha \cdot D_{i-1,j} \cdot PS_{i,j-1} \quad (4.57)$$

$$w_{i,j} = w_{i,j} + \alpha \cdot (D_{i-1,j} - w_{i,j}) \quad (4.58)$$

Variabilele care intervin în cadrul algoritmilor neuronali sunt booleene (*Hopfield, Perceptron*) și respectiv reale (*Adaline, Kohonen* etc.). Reprezentarea datelor se poate face în virgulă fixă sau în virgulă mobilă. Reprezentarea în virgulă fixă este caracterizată de un domeniu de reprezentare și o precizie a reprezentării mai reduse; ea revendică însă dispozitive de adunare / înmulțire mai simple din punct de vedere structural. Reprezentarea în virgulă flotantă extinde domeniul de reprezentare și precizia dar revendică dispozitive de adunare / înmulțire mai complexe. Evaluările făcute pe cale experimentală au relevat că o reprezentare în virgulă fixă pe 8 biți oferă suficientă precizie (pentru algoritmi conexioniști studiați). Deci, registrele NV, Down și Up vor avea lungimea de 8 biți (într-o primă variantă care conduce la o structură mai simplă a CS).

Registrul coeficientului sinaptic (*W*) va conține valori care rezultă din două surse posibile:

- învățare externă;
- învățare internă.

În cazul învățării externe, calculul coeficienților sinaptici va fi realizat de calculatorul gazdă, în virgulă flotantă. În acest caz, reprezentarea în virgulă flotantă trebuie redusă la o reprezentare în virgulă fixă pe *p* biți, după ce în prealabil coeficienții fiecărei linii au fost normalizați prin procedura următoare:

- se va utiliza o funcție *MaxAbs(.)* care generează cea mai mare valoare (absolută) dintre toate valorile coeficienților ce aparțin unei linii a matricii *W*;
- pentru *i = 1, ..., N*
 pentru *j = 1, ..., N*

$$w_{ij} = \frac{w_{ij}}{\text{MaxAbs}(W_i)}$$

unde *W_i* - linia *i* a matricii *W*

Această operație de normalizare are 2 consecințe asupra matricii *W*:

- pierderea simetriei;
- scăderea preciziei.

Proiectarea Rețelelor Sistolice dedicate algoritmilor conexioniști

În cazul rețelei *Hopfield*, pierderea de simetrie este fără consecințe pentru faza de recunoaștere. Dacă funcția de transfer este $\text{sign}(\cdot)$, avem următoarea proprietate:

$$\sigma(a \cdot \mathbf{x}) = \sigma(\mathbf{x}) \quad \text{dacă } a > 0 \quad (4.59)$$

dar:
$$\frac{1}{\text{MaxAbs}(\mathbf{W}_i)} > 0 \quad (4.60)$$

ceea ce înseamnă că:

$$\sigma\left(\frac{1}{\text{MaxAbs}(\mathbf{W}_i)} \cdot \mathbf{x}\right) = \sigma(\mathbf{x}) \quad (4.61)$$

În ceea ce privește precizia, mai multe simulări au arătat că pierderea de precizie este fără consecințe asupra comportamentului RN în faza de recunoaștere, atâta timp cât:

$$p \geq \log_2(N) + 1 \quad (4.62)$$

Optând pentru $p = 8$, rezultă că vom putea construi o RN de 128 neuroni fără a suporta o degradare a performanțelor din cauza preciziei de reprezentare.

Reprezentarea întreagă pe 8 biți va conduce la o structură foarte simplă a celei sistolice. În cele ce urmează vom prezenta etapele de proiectare a CS pentru $p = 8$. Pentru modelele conexioniste mai complexe, valoarea $p = 8$ va fi revizuită și CS reproiectată.

În cazul învățării interne, calculele se vor executa în interiorul RS în virgulă fixă iar rezultatele vor fi approximate prin trunchiere.

Registrul PS va avea o lungime care depinde de lungimea registrului W și de numărul de neuroni simulați de rețeaua sistolică. Dacă notăm:

- k - numărul de biți alocați pentru reprezentarea componentelor vectorului \mathbf{x} (lățimea registrelor Down și respectiv Up).
- p - lățimea registrului W (numărul de biți)
- m - lățimea registrului PS (numărul de biți)
- N - dimensiunea rețelei sistolice ($N \times N$ celule)

în baza relației (4.55) putem scrie:

$$m = k + p + \log_2 N \quad (4.63)$$

Pentru prima variantă de CS, vom avea $m = 8 + 8 + 7 = 23$ biți. Vom aloca registrului PS 24 biți (multiplu de 8).

Cei 24 biți ai registrului PS trebuie transmiși de la o celulă la alta din direcția vest spre est în timpul elaborării potențialelor neuronilor. Transmisia se poate realiza în serie sau în paralel. În cazul paralel, transmisia se face într-o singură perioadă de tact dar revendică:

- un dispozitiv de înmulțire paralel;
- un bus de comunicație cu lățimea de 24 biți;
- 24 pini de intrare / ieșire pentru fiecare linie a matricii W.

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

În cazul serial, transmisia este de 24 ori mai lentă dar comunicațiile între celule se vor realiza pe un singur fir. Sumatorul și circuitul de înmulțire din interiorul CS vor lucra secvențial și vor avea o structură mult simplificată.

Vom reține (pentru început) transmisia serială, pe care o justificăm pe baza a 2 considerente:

- structura internă a celulei sistolice trebuie să fie cât mai simplă pentru a putea implementa un număr cât mai mare de celule (sinapse);
- viteza de calcul a RS trebuie corelată cu lărgimea de bandă a canalului de comunicație dintre RS și calculatorul gazdă. Dacă lărgimea de bandă a canalului de comunicație este inferioară vitezei de calcul a RS, RS va recurge la stări "wait".

Patru comenzi determină funcționarea CS în 4 moduri de lucru:

- RESET: șterge toate registrele interne
CLEARX: șterge toate registrele cu excepția registrului W (coeficienții sinaptici).
SHIFTW: conținutul registrului W se deplasează cu o poziție pe fiecare impuls de tact. Comanda va fi utilizată pentru încărcarea matricii sinaptice **W** în RS, din direcția nord spre sud.
COMPUTE: conținutul tuturor registrelor se deplasează cu un bit pe fiecare impuls de tact și se elaborează suma parțială PS care se transmite, bit cu bit, spre est (serial).

Celulele diagonale se deosebesc de cele nondiagonale prin modul de comunicație. Pentru diagonale, intrarea din est este "reflectată" spre ieșirea de nord pentru a permite transmisia noii valori a vectorului x în RS, pentru iterația următoare. Pe ieșirea spre vest celula diagonală transmite rezultatul comparației locale dintre vechea valoare x_{ij} și noua valoare x_{ij} (evaluarea convergenței algoritmului în faza de recunoaștere - rețeaua *Hopfield*). În faza de învățare, noua valoare tranzitează de la vest spre est atât celulele diagonale cât și cele nondiagonale.

Interfațarea CS se face prin 4 intrări și 4 ieșiri seriale, și câteva semnale de comandă:

- CLK - semnalul de tact principal
RESET - inițializează la zero registrele UP, DOWN, NV, PS și W
CLEARX - inițializează la zero registrele UP, DOWN, NV și PS
SHIFTX - deplasarea registrelor UP, DOWN și NV
SHIFT PS - deplasarea registrului PS
SHIFT W - deplasarea registrului W (încărcarea matricii sinaptice)
DIAG - celulă diagonală (DIAG = 1) și respectiv nondiagonală (DIAG = 0).

Structura celulei sistolice este redată în figura 4.37. Structura celulei este simplă și poate fi multiplicată într-un număr mare de exemplare. Comunicațiile se realizează numai cu celulele vecine. Celula poate fi programată diagonală (DIAG = 1) sau nondiagonală (DIAG = 0); acest lucru va permite realizarea RS reconfigurabile.

Vom expune în continuare funcțiile fiecărei unități din componența CS.

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

curent al înmulțitorului (x). Cei 16 biți ai produsului parțial obținut pe tactul curent vor fi decalajați cu o poziție și apoi adăugați produsului parțial obținut pe tactul următor.

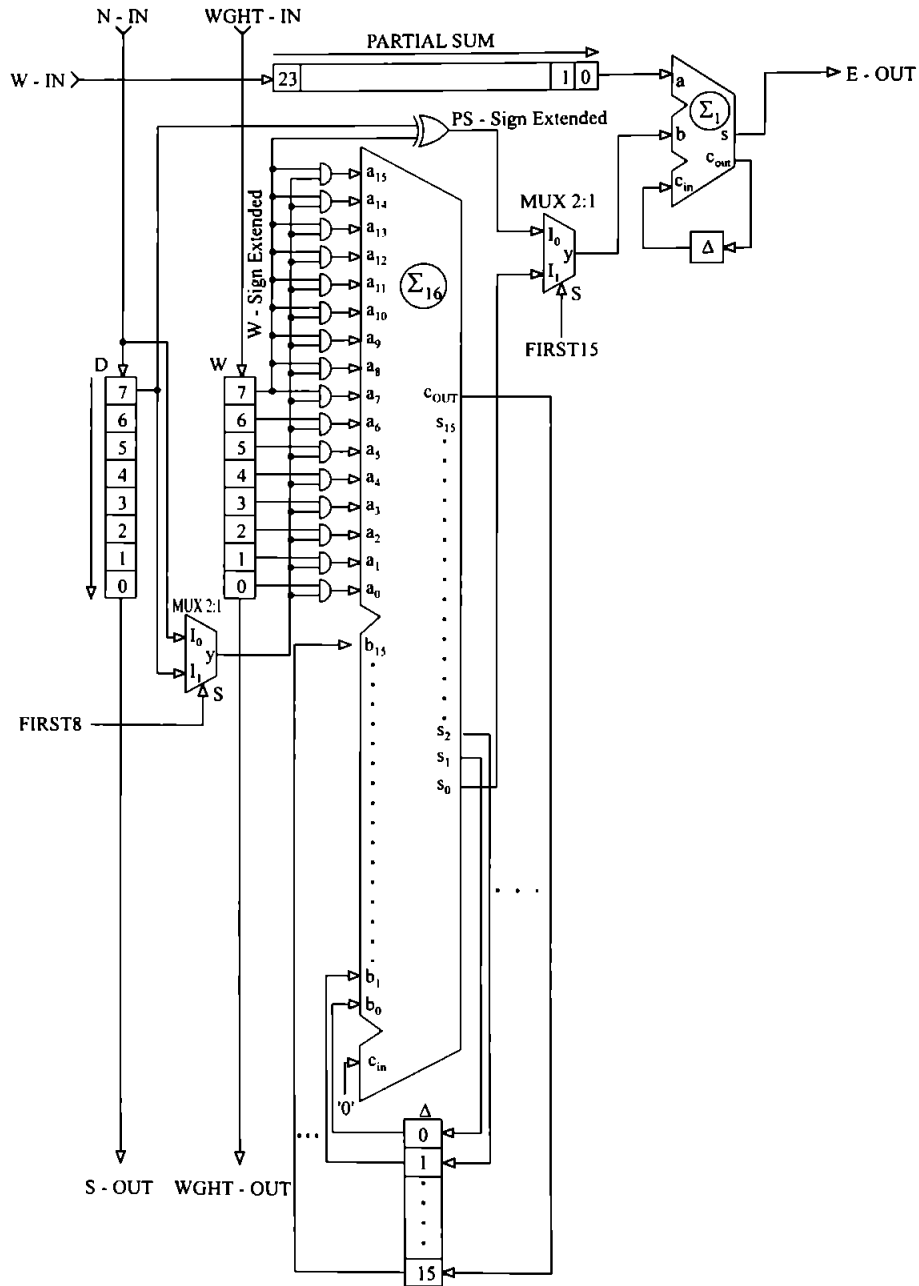


Fig. 4.38. Unitatea de calcul

Unitatea de calcul a CS va fi deci compusă din două sumatoare și registrele de memorare / deplasare aferente.

4.5.4. Unitatea de comunicație

Unitatea de comunicație se ocupă de transmisiile de date între celule. Am stabilit două tipuri de comunicații asociate celor 2 tipuri de celule: diagonale și respectiv nondiagonale (fig. 4.39).

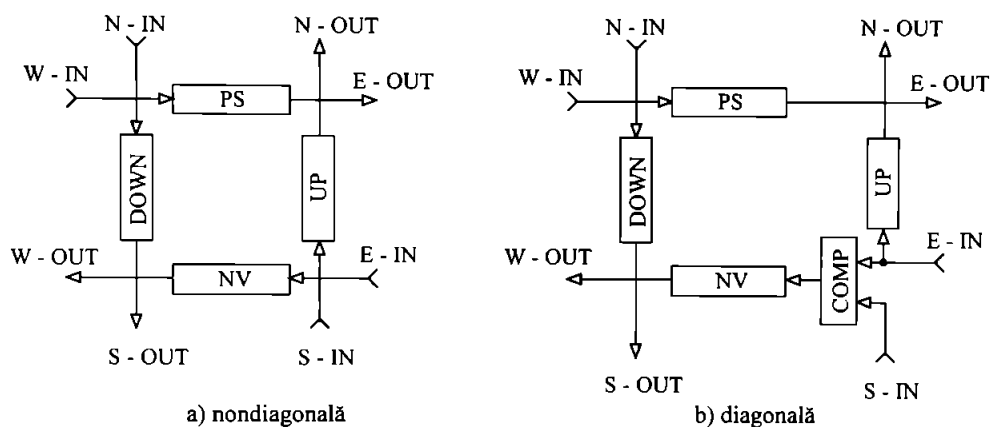


Fig. 4.39. Comunicațiile celulei sistolice

Gestiunea comunicațiilor reprezintă un mecanism esențial în RS. Toate celulele sistolice au resurse de comunicație identice; fiecare celulă poate fi programată diagonală sau nondiagonală. Aceasta permite extensia RS printr-un simplu pavaj de circuite VLSI așa cum se arată în figura 4.40.

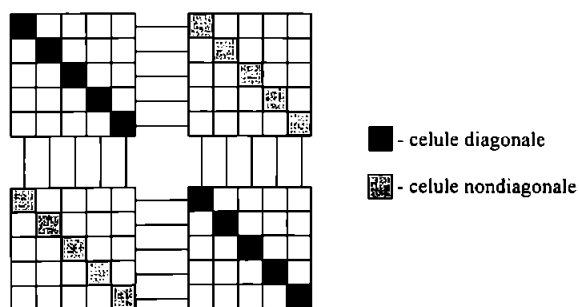


Fig. 4.40. Extensia RS prin conectarea mai multor circuite VLSI și programarea corespunzătoare a celulelor diagonale și nondiagonale

Mecanismul de comparare, activ în celulele diagonale, servește în faza de recunoaștere (rețeaua Hopfield) la detecția identității dintre vectorii de stare care corespund la două iterații succesive (detecția convergenței). Procesul de comparare se va realiza printr-un comparator serie.

Concepția părții de comunicație se reduce la concepția unor registre de deplasare și a unui comparator.

4.5.5. Adaptarea arhitecturii pentru învățare

Pentru regulile de învățare derivate din regula lui *Hebb* (*Delta Projection*, *Adaline*, *Perceptron*) ajustarea arhitecturală necesară este minoră. Trebuie adăugate anumite căi de date care să faciliteze implementarea ecuațiilor impuse de funcția de învățare. Arhitectura prezentată în fig. 4.41 și care conține 4 multiplexoare suplimentare satisface aceste cerințe.

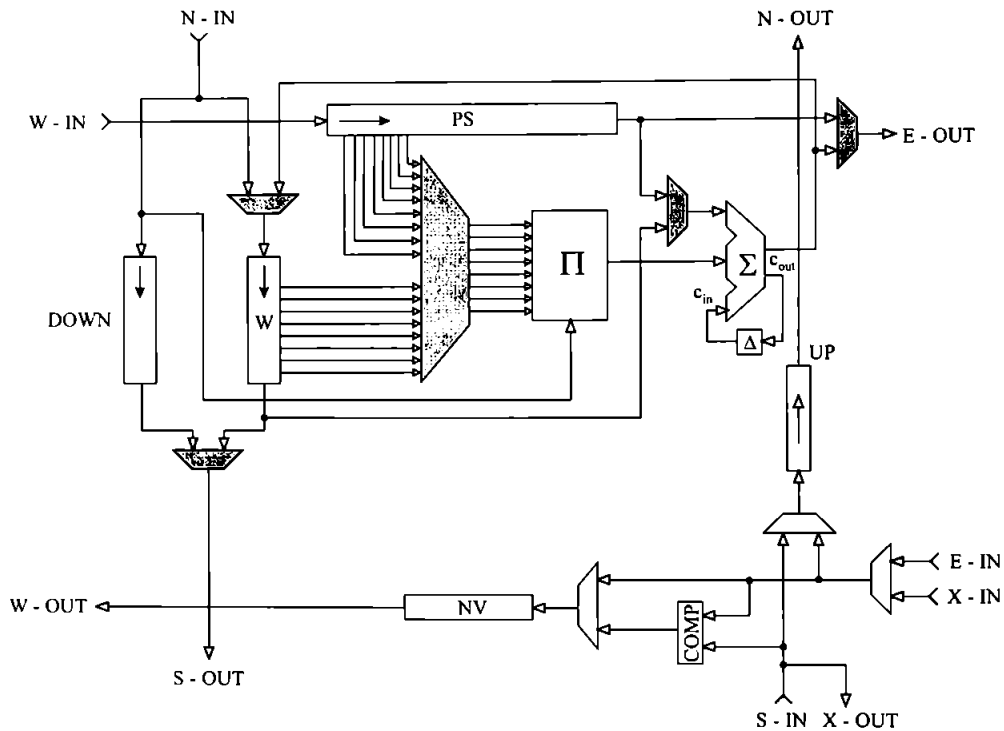


Fig. 4.41. Structura CS dotată cu funcția de învățare

Multiplexoarele evidențiate nu provoacă o creștere spectaculoasă a complexității CS dar vor permite ajustarea ponderilor sinaptice pe baza regulii de învățare dorite.

4.6. Evaluarea performanțelor

Pentru evaluarea performanțelor diferitelor implementări (dedicate algoritmilor conexioniști) se utilizează două măsuri:

- Numărul de conexiuni pe secundă (CPS - *Connections Per Second*). Trebuie să menționăm că CPS are o semnificație dependentă de model (arhitectură, precizie etc.). De exemplu, pe o rețea de tip *Adaline* (directă și fără rebusclaj) care operează în faza de recunoaștere, CPS măsoară numărul de operații de înmulțire + acumulare. În cazul unei rețele recurente care operează în faza de învățare CPS nu mai reprezintă același lucru; va

Proiectarea Rețelelor Sistolice dedicate algoritmilor conexioniști

reprezenta numărul de conexiuni ajustate pe secundă. De asemenea, CPS nu ține cont de timpii de intrare / ieșire;

- Numărul de actualizări pe secundă (UPS - *Updates Per Second*) reprezintă o măsură mai globală care ține cont de timpii de intrare / ieșire.

În RS bidimensională propusă, durata unei iterații se definește ca timpul scurs între momentul plecării și momentul revenirii unui vector din / în diagonală. Dacă N reprezintă numărul de neuroni din rețea (RS va avea dimensiunea $N \times N$), o iterație va fi executată în $2 \cdot 24 \cdot N$ perioade de tact (fiecare sumă parțială PS este codificată pe 24 biți). Pentru RS sistolică bidimensională cele două măsuri se vor calcula astfel:

- Numărul de UPS (notat cu δ) va fi:

$$\delta = 2NL / T_u \quad (4.64)$$

$$\text{unde: } T_u = 2Nm / f \quad (4.65)$$

f - frecvența semnalului de tact

T_u - timp de *update* (timp de actualizare a matricii W). Pentru actualizarea tuturor elementelor w_{ij} se va parcurge un ciclu complet; pentru un ciclu complet vectorul x pornește din diagonală și revine în diagonală, după timpul T_u .

m - dimensiunea PS (număr de biți).

L - rata *pipeline*-ului.

Dacă $L = 1$ (100 %), în timpul T_u , funcția neliniară σ (vezi cei 6 pași din cadrul ciclului T_u prezentați la 4.2) se va aplica la 2NL sume parțiale PS, conducând la 2NL actualizări.

Din (4.64) și (4.65) rezultă că:

$$\delta = L \cdot f / m \quad (4.66)$$

ceea ce înseamnă că numărul de actualizări pe secundă este independent de numărul de neuroni (N).

- Numărul de CPS (notat cu γ) va fi:

$$\gamma = LN^2 / T_s \quad (4.67)$$

$$\text{unde: } T_s = m / f \quad (4.68)$$

f - frecvența tactului (tactul de bază care va realiza și deplasarea biților în PS)

m - dimensiunea PS (număr de biți).

T_s - timp de calcul sinapsă (timpul de calcul al unui set de sume parțiale).

Dacă rata *pipeline*-ului ar fi de 100 % ($L = 1$) atunci se vor calcula N^2 sume PS în timpul T_s , ceea ce confirmă relația (4.67).

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

Pentru o frecvență de tact $f = 100$ MHz (valoare permisă de tehnologiile actuale), $m = 24$ biți, $L = 1$ (100 %) - permis de RS propusă, și pentru un număr de neuroni rezonabil $N = 32$, obținem:

$$\delta = 1 \cdot 100 \cdot 10^6 / 24 = 4,1(6) \text{ MUPS}$$
$$\gamma = 1 \cdot 32^2 \cdot 100 \cdot 10^6 / 24 = 4,2(6) \text{ GCPS}$$

Aceste performanțe estimate ne îndreptătesc să considerăm că RS analizată conduce la implementări VLSI utilizabile în aplicații reale. Conform [JUT 90], identificarea surselor radar reclamă aproximativ 625 MCPS, controlul unui robot aproximativ 350 KCPS, recunoașterea cuvintelor aproximativ 24 MCPS. Se observă că măsurile de performanță estimate depășesc confortabil aceste valori. Cu titlu de comparație, în tabelul 4.3 prezentăm performanțele raportate pentru câteva implementări comerciale:

INTEL 80170NW - ETANN (*)	Circuit VLSI	2 GCPS
Fujitsu Neurocomputer	Placă DSP	500 MCPS
Cray X - MP	Parallel computer	50 MCPS

(*) ETANN - Electronically Trainable Analog Neural Network

Tabelul 4.3. Performanțele câtorva implementări comerciale

Datele prezentate în tabelul 4.3 trebuie interpretate cu anumite precauții. Parametrul CPS dă o indicație aproximativă asupra performanțelor mașinilor citate. El nu reprezintă o măsură precisă deoarece el nu înglobează timpii de intrare / ieșire, timpii de calcul ai unor funcții cum ar fi funcția de transfer etc. În ceea ce privește RS propusă se impun trei precizări:

- Pentru RS, parametrul CPS este o funcție care depinde de pătratul numărului de neuroni. Pentru N mare se poate ajunge teoretic la valori mult mai mari (de exemplu 68,2 GCPS pentru $N = 128$ neuroni);
- RS este bine adaptată execuției algoritmilor conexioniști dar funcționarea sa depinde de lățimea de bandă a canalului de intrare / ieșire (viteza de transfer a interfeței cu sistemul gazdă);
- Numărul de circuite VLSI necesar pentru implementarea RS aferentă, de exemplu, unei RN compusă din 128 neuroni reprezintă un alt element de care trebuie să se țină cont. Numărul de circuite depinde de numărul de sinapse integrate (crește cu pătratul lui N). Totuși, CS propusă implică comunicații seriale locale, are o structură foarte simplă și RS (care se bazează pe comunicații strict locale) se pretează excelent pentru implementările VLSI.

4.7. Implementarea RS în structuri FPGA

Proiectul ce urmează a fi descris în acest paragraf a fost realizat în laboratorul de cercetare "*Systèmes logiques et numériques*" din cadrul *Universității Libere Bruxelles*. Utilizând Sistemul de Dezvoltare al firmei XILINX, RS descrisă a fost proiectată, implementată în structuri FPGA, testată funcțional și evaluată. Sistemul de dezvoltare al firmei XILINX este constituit dintr-un pachet de programe care permit proiectarea structurilor logice, testarea funcțională, evaluarea performanțelor structurilor proiectate și implementarea acestora în circuite FPGA (*Field Programmable Gate Array*). FPGA este un circuit VLSI programabil; fiecare FPGA conține un anumit număr de blocuri

Proiectarea Rețelelor Sistolice dedicate algoritmilor conexioniști

logice (CLB - *Configurable Logic Block*). De exemplu, circuitul XC 4025 E produs de firma XILINX conține 2^{10} CLB; structura unui CLB este prezentată în figura 4.42.

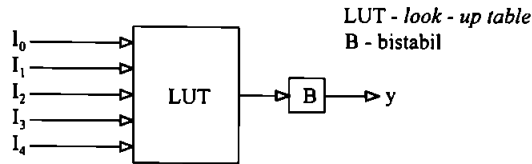


Fig. 4.42. Structura CLB din cadrul circuitului FPGA XC 4025 E

Fiecare CLB conține o tabelă (LUT - *look - up table*) cu 5 intrări și un element de memorare (B - circuit bistabil). Programarea circuitelor FPGA constă în:

- stabilirea conținutului tabelelor LUT din CLB-urile circuitului FPGA;
- stabilirea interconexiunilor dintre intrările și ieșirile diverselor CLB din structura FPGA.

RS proiectată în această lucrare este constituită dintr-un anumit număr de CS identice. Cu ajutorul sistemului de dezvoltare XILINX a fost proiectată și testată funcțional CS din fig. 4.37. Schema bloc (afărentă CS) rezultată în urma acestei etape de proiectare este prezentată în fig. 4.43.

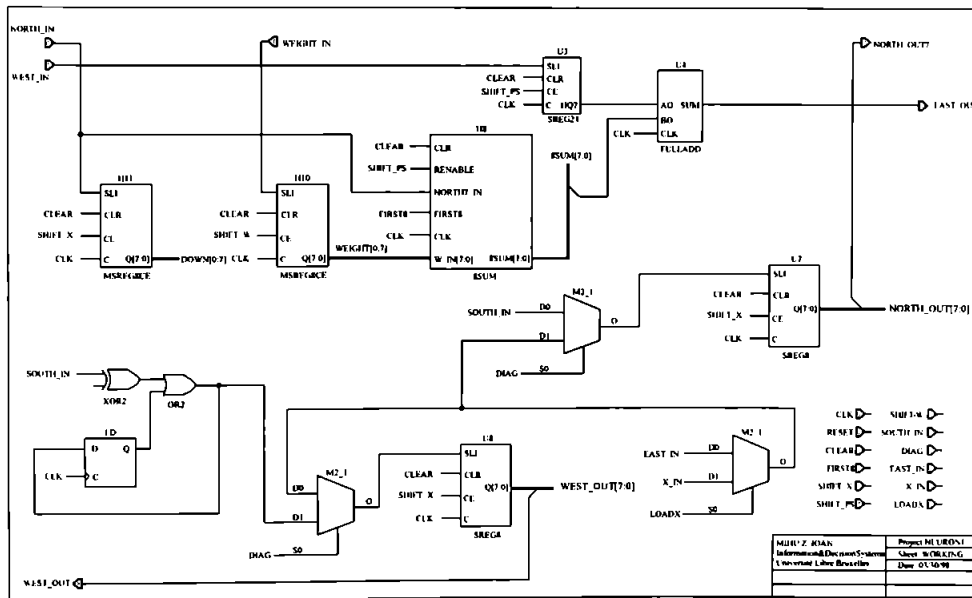


Fig. 4.43. Schema CS proiectată în sistemul de dezvoltare XILINX

4.7.1. Testarea funcțională a CS

Celula sistolică a fost în întregime simulată și testată funcțional în sistemul de dezvoltare XILINX. Pentru testarea funcțională s-au utilizat seturi adecvate de vectori de test. Cu titlu de

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

exemplu, vom prezenta procedura de testare funcțională a CS pentru un set de vectori de test. În cadrul fazei de recunoaștere funcționarea CS este descrisă de ecuația (4.56). În figura 4.44 se prezintă programul de testare funcțională iar în figura 4.45 se prezintă setul de vectori de test utilizat. Fișierul *data.dat* din figura 4.45 definește următorii vectori de test:

$$\begin{aligned}PS_{i,j-1} &= 128 \\D_{i-1,j} &= 005 \\w_{i,j} &= 006\end{aligned}$$

În figurile 4.46, 4.47 și 4.48 se prezintă rezultatele obținute în cadrul fazei de testare funcțională. În fig. 4.46 se prezintă starea CS după inițializare (RESET):

$$\begin{aligned}PS_{i,j-1} &= 0 \\D_{i-1,j} &= 0 \\w_{i,j} &= 0\end{aligned}$$

În figura 4.47 se prezintă starea CS în momentul inițierii etapei de calcul (după încărcarea serială a registrelor):

$$\begin{aligned}PS_{i,j-1} &= 128 \\D_{i-1,j} &= 005 \\w_{i,j} &= 006\end{aligned}$$

În figura 4.48 se prezintă starea CS la finele etapei de calcul:

$$\begin{aligned}PS_{i,j} &= 158 \\D_{i,j} &= 005 \\w_{i,j} &= 006\end{aligned}$$

Rezultatele obținute demonstrează corecta funcționare atât a unității de calcul cât și a unității de comunicație din componența CS. Aplicarea acestei proceduri pe un număr mare de seturi de vectori de test a permis testarea completă a CS. Toate testele efectuate au demonstrat funcționarea corectă a CS atât în faza de calcul cât și în faza de comunicație cu CS vecine.

```
*****
|Script file for the neuro cell test
|10/April/1998
|Mihu Z. Ioan, Univesité Libre Bruxelles
|*****

clpr test.out
|delete_signals
restart
set_mode functional

vector R R[0:23]
vector DOWN DOWN[0:7]
vector WEIGHT WEIGHT[0:7]

radix dec R
radix dec WEIGHT
radix dec DOWN
```

Proiectarea Rețelelor Sistolice dedicate algoritmilor conexioniști

```
watch CLK SHIFT_PS SHIFT_W EAST_OUT

|*****
|CLK
|*****
wfm CLK @0=1\d (10ns=inc by 1)*130

low SHIFT_W
low SHIFT_PS
low FIRST8
sim 20ns

|*****
|LOAD D with 0 PS & Weight
|*****

high SHIFT_X
(assign NORTH_IN < data.dat; sim 20ns) * 8

high SHIFT_PS
high SHIFT_W
(assign WEST_IN < data.dat; assign WEIGHT_IN < data.dat; sim 20ns) * 8
low SHIFT_W
(assign WEST_IN < data.dat; sim 20ns) * 16
low SHIFT_PS

sim 20ns

|*****
|Compute
|*****

high SHIFT_PS
high SHIFT_X
high FIRST8
(assign NORTH_IN < data.dat; sim 20ns) * 8
low FIRST8
low SHIFT_X
sim 320ns
low SHIFT_PS
sim 20ns
```

Figura 4.44 Programul de test pentru CS

```
|*****
|Data file for the neuro cell test
|10/April/1998
|Mihu Z. Ioan, Univesité Libre Bruxelles
|*****

|INIT NORTH_IN d REG
0
0
```

SOLUTII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

```
0
0
0
0
0
0
0
| LOAD PS & LOAD D
0
0
|
0
1
|
0
1
|
0
0
|
0
0
|
0
0
|
0
0
|
1
0
|
0
0
0
0
0
0
0
0
0
0
0
0
0
0
0
0
| NORTH_IN
1
0
1
0
0
0
0
0
```

Figura 4.45 Un set de vectori de test

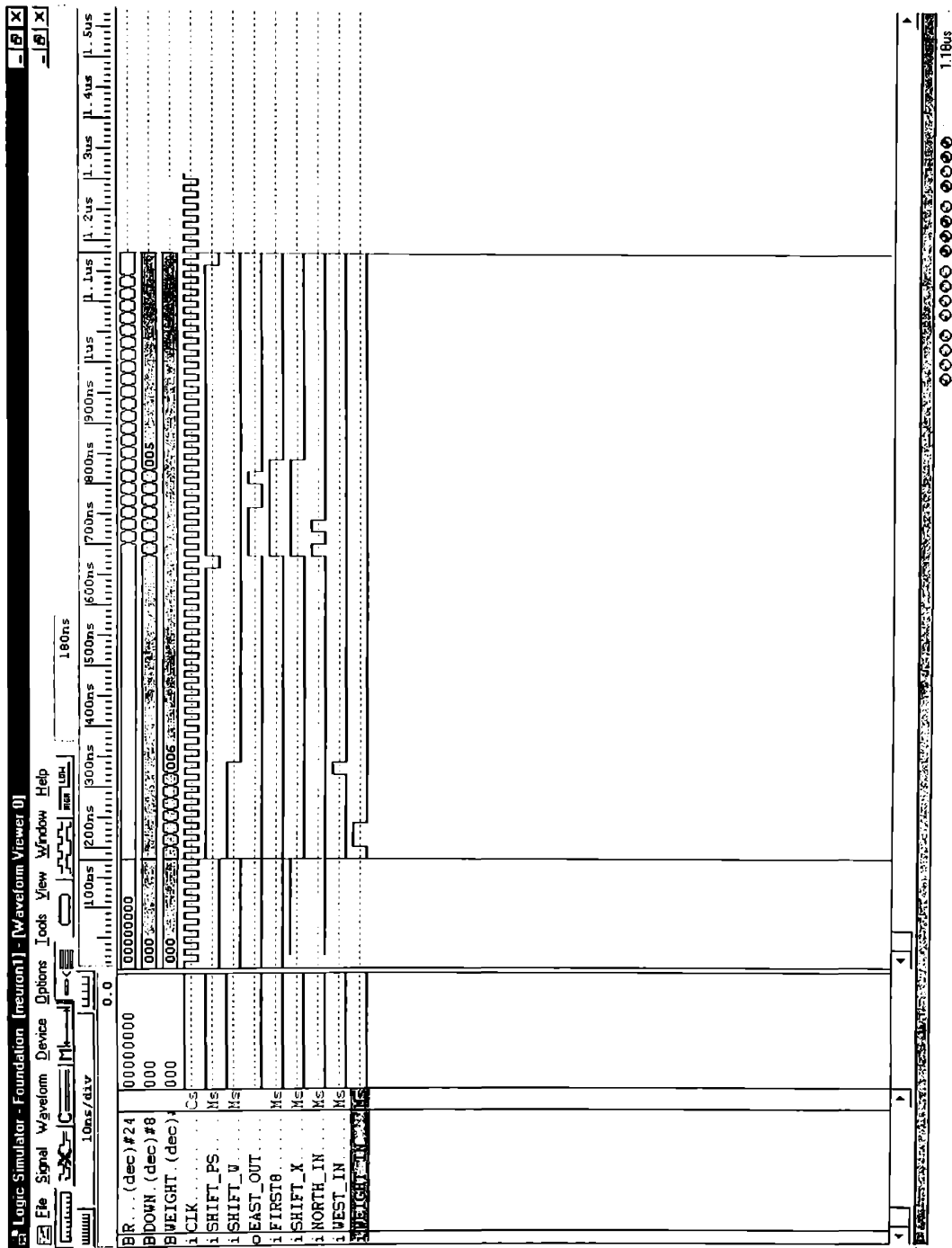


Fig. 4.46. Testarea funcțională a CS. Faza de inițializare: $PS_{i,j-1} = 0$; $D_{i,j} = 0$; $w_{i,j} = 0$

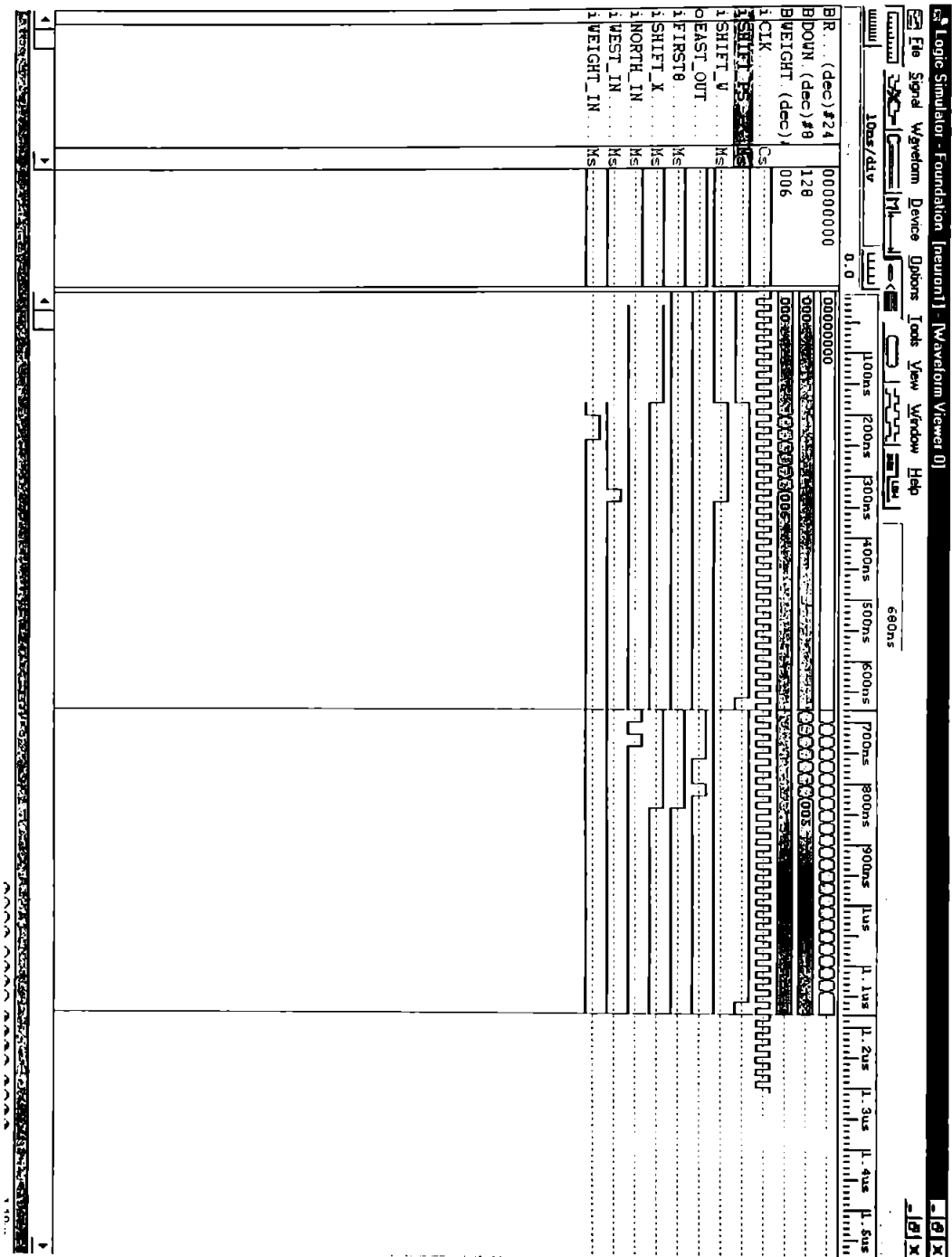


Fig. 4.47. Testarea funcțională a CS. Starea CS în momentul inițierii etapei CALCUL: $PS_{i,j-1} = 128$;
 $D_{i-1,j} = 005$; $w_{i,j} = 006$

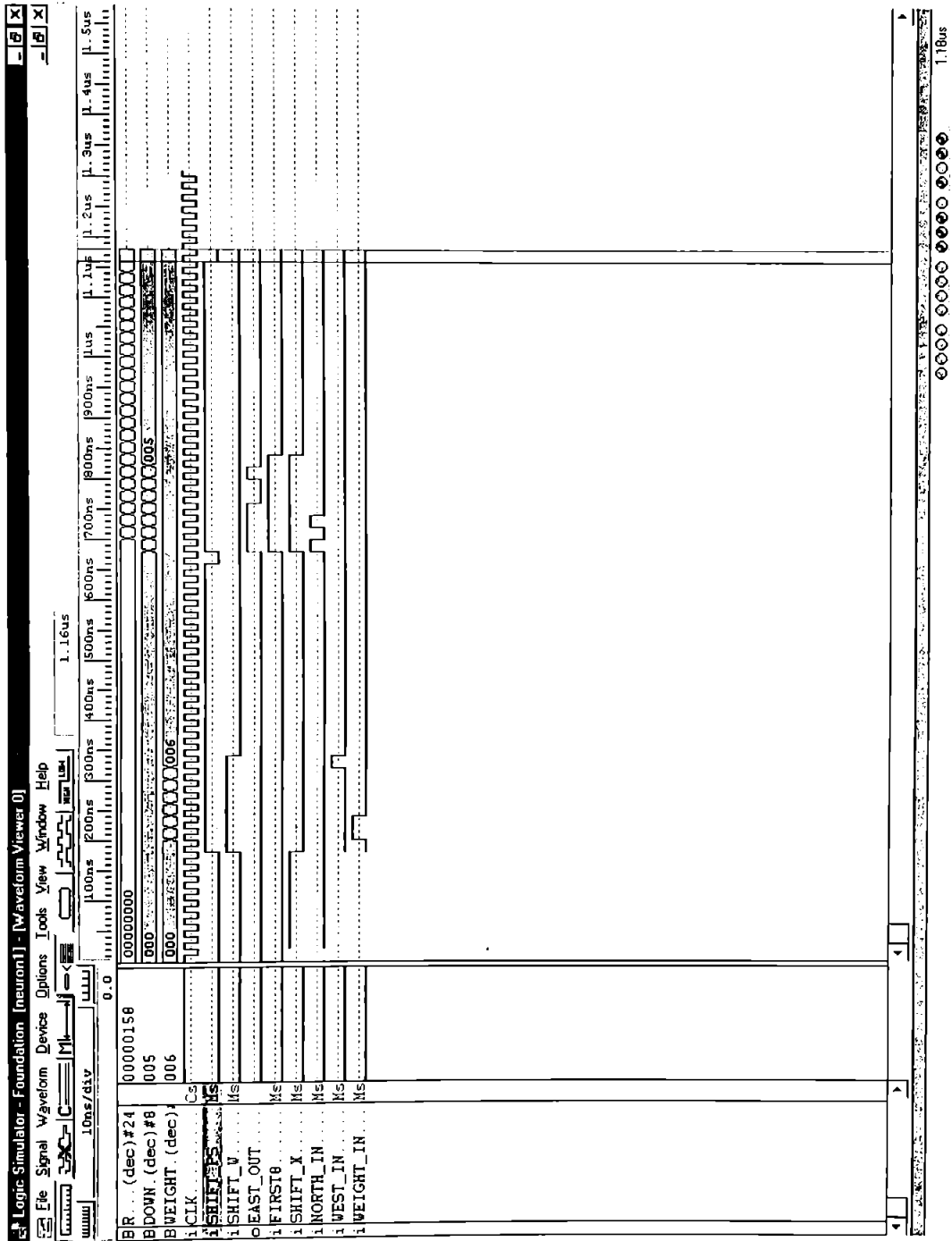


Fig. 4.48. Testarea funcțională a CS. Starea CS la finele etapei CALCUL: $PS_{i,j-1} = 158$; $D_{i-1,j} = 005$; $w_{i,j} = 006$

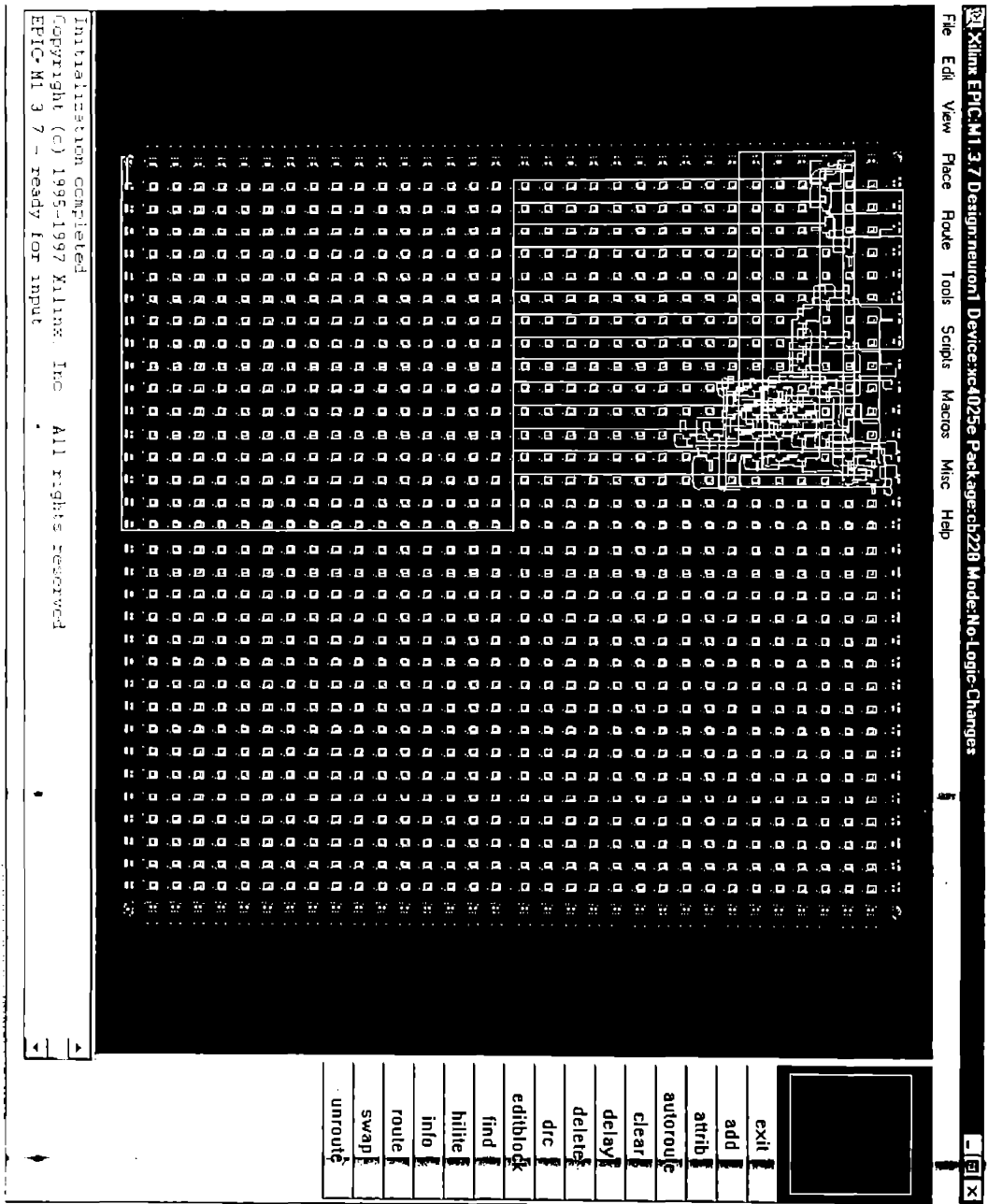


Fig. 4.49. Implementarea CS in circuitul FPGA XC4025E

4.7.2. Implementarea CS

CS proiectată și testată funcțional pe baza procedurilor expuse la 4.7.1. a fost implementată în circuitul FPGA XC 4025 E produs de firma XILINX. Circuitul XC 4025 E conține 1024 CLB organizate într-o matrice de dimensiune 32×32 . Implementarea fizică a CS în acest circuit FPGA este prezentată în figura 4.49. Din cele 1024 CLB existente în circuitul XC 4025 E celula sistolică a ocupat 35 CLB ceea ce înseamnă că un singur circuit XC 4025 E va putea găzdui 29 celule sistolice. Aceasta înseamnă că într-un singur circuit XC 4025 E se poate implementa ușor o RS de 5×5 celule sistolice; extensia RS se poate realiza prin conectarea mai multor circuite XC 4025 E ca în fig. 4.40.

Implementarea în variantă FPGA prezintă un avantaj major: reconfigurabilitate dinamică. Implementarea în FPGA poate constitui soluția practică pentru construcția unei plăci neuronale acceleratoare care urmează a fi integrată într-un sistem gazdă. În aceste condiții RS (placa acceleratoare) nu trebuie să acopere, în orice moment de timp, toți algoritmi conexioniști cunoscuți deoarece o astfel de abordare ar complica peste măsură structura CS din componența RS și ar induce limitări severe în ceea ce privește dimensiunea RS implementabile pe o placă acceleratoare. În funcție de gradele de asemănare dintre diverși algoritmi, se pot forma mai multe clase de algoritmi și pentru fiecare clasă se poate proiecta o RS (CS) adecvată. În funcție de numărul de clase de algoritmi definite se vor obține mai multe arhitecturi de RS, dar toate aceste arhitecturi vor avea o trăsătură comună: simplitatea structurii CS (structura internă a CS diferă de la o RS la alta dar toate CS vor avea structuri interne simple). *Software*-ul de gestiune aferent plăcii acceleratoare va putea reconfigura dinamic RS implementată pe această placă (în funcție de necesități). Pentru o RN de tip SOFM, de exemplu, RS ar putea avea o arhitectură, iar pentru o RN de tip *Hopfield* RS ar putea avea o altă arhitectură. Reconfigurarea RS se va face dinamic în funcție de algoritmul neuronal ce urmează a se executa. Reconfigurabilitatea va conduce la o creștere a eficienței RS (plăcii acceleratoare). De vreme ce structura CS rămâne în permanență simplă (deși se modifică în timp), se vor putea obține (pentru fiecare algoritm neuronal în parte) viteze maxime de procesare în interiorul RS.

Circuitul XC 4025 E nu este un FPGA de densitate mare. În momentul actual firma XILINX produce o diversitate de circuite FPGA de dimensiuni (densități) diverse. De exemplu, circuitul XC 4062 XL conține 48×48 CLB și ar putea găzdui 65 celule sistolice. Aceasta înseamnă că XC 4062 XL ar putea găzdui o RS pătrată de 8×8 celule sistolice. În consecință, o placă acceleratoare cu 36 circuite XC 4062 XL ar putea constitui suportul fizic pentru o RS de 48×48 celule sistolice.

Circuitul XC 40250 conține 20.000 CLB și ar putea găzdui fără probleme o RS de 22×22 celule sistolice. Placa acceleratoare cu 36 circuite XC 40250 ar putea constitui suportul fizic pentru o RS de 132×132 celule sistolice. Toate aceste elemente demonstrează viabilitatea implementării RS cu structuri de tip FPGA.

4.8. Concluzii

Am arătat că RS bidimensională permite implementarea operațiilor vectoriale și matriciale revendicate de algoritmi conexioniști. Implementarea sistolică prezintă mai multe avantaje:

- permite implementarea fazei de recunoaștere pentru toți algoritmi clasici;
- permite implementarea celor mai diverse reguli de învățare. Această capacitate rezultă din faptul că RS realizează o înlănțuire de calcule matriciale - vectoriale fără a specializa mașina suport pe funcții particulare. Structura RS este flexibilă și se bazează pe

SOLUȚII DE IMPLEMENTARE A REȚELELOR NEURONALE PE ARHITECTURI SISTOLICE

comunicații strict locale. Forma pătrată a RS facilitează sincronizarea acestor comunicații locale;

- procesarea în tehnică *pipeline* a vectorilor de intrare este perfect posibilă deoarece toate operațiile se efectuează cu date strict locale. Eficiența RS va crește spectaculos, putând atinge valori de 100 % în fazele de procesare staționară (recurențe);
- extensibilitatea RS se poate realiza printr-un simplu pavaj de circuite VLSI. În cazul comunicațiilor seriale, numărul de pini necesari interconectării circuitelor VLSI este foarte mic.

5. CONSIDERAȚII PRIVIND PERFECȚIONAREA RS DEDICATE ALGORITMILOR CONEXIONIȘTI

5.1. Introducere

Modelul neuronului artificial este reprezentat de un dispozitiv cu N intrări și o ieșire. Ieșirea y_i a neuronului i se calculează:

$$y_i = \sigma(p_i) = \sigma\left(\sum_{j=1}^N w_{i,j} \cdot x_j\right); \quad i = 1, 2, \dots, M \quad (5.1)$$

unde:

w_{ij} reprezintă ponderea sinaptică asociată intrării j a neuronului i
 p_i - potențialul neuronului i
 x_j - componentele vectorului de intrare
 σ - funcția de transfer

Sub formă matricială ecuația (5.1) devine:

$$\mathbf{y} = \sigma(\mathbf{p}) = \sigma(\mathbf{W} \cdot \mathbf{x}) \quad (5.2)$$

Într-o RN, neuronii sunt organizați în straturi; neuronii din componența unui strat vor recepționa același vector pe intrări și vor avea ieșirile conectate la intrările neuronilor din următorul strat. Matricile sinaptice vor fi $\mathbf{W}^{[q]}$, unde q reprezintă indicele sau numărul stratului.

În faza de antrenament (învățare), matricile $\mathbf{W}^{[q]}$ sunt ajustate iterativ pe baza unei reguli de învățare. Una dintre cele mai simple este regula lui *Hebb*:

$$\mathbf{W} := \mathbf{W} + \alpha \cdot (\mathbf{y} \cdot \mathbf{x}^T) \quad (5.3)$$

unde α reprezintă constanta de învățare. Cele mai multe reguli practice de învățare se bazează pe reguli derivate din (5.3).

RN multistrat fac posibilă implementarea oricărei funcții $y = \phi(x)$, x reprezentând vectorul de intrare în primul strat iar $y = y^{[L]}$ vectorul de ieșire din ultimul strat (L). Cea mai uzuală funcție de transfer σ este tangenta hiperbolică sau funcția sigmoidală ((1.3) și (1.4)). RN învață funcția ϕ printr-un proces iterativ de prezentare a perechilor $\{\mathbf{x}_k, \mathbf{d}_k\}$, numite prototipuri sau perechi de antrenament. Vectorul \mathbf{d}_k reprezintă vectorul dorit la ieșire atunci când pe intrarea RN se aplică vectorul \mathbf{x}_k .

Algoritmul *Error Back Propagation* (EBP) realizează o ajustare a ponderilor sinaptice cu scopul minimizării unei funcții de eroare (eroarea medie pătratică); se asigură astfel o coborâre pe suprafața de eroare spre un minim (local sau global) pe direcția și în sens invers gradientului funcției de eroare.

În acest scop se calculează un semnal de eroare asociat fiecărui strat:

$$\delta_i^{[L]} = (d_i - y_i^{[L]}) \cdot \sigma'(p_i^{[L]}) \quad (5.4)$$

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

$$\delta_i^{(q)} = \left(\sum_{k=1}^{M_{q+1}} w_{k,i}^{(q+1)} \cdot \delta_k^{(q+1)} \right) \cdot \sigma'(p_i^{(q)}); \quad q = 1, 2, \dots, L-1 \quad (5.5)$$

$$\text{unde: } \sigma'(v) = \frac{d\sigma(v)}{dv}$$

Ecuatiile (5.4) și (5.5) sunt valide pentru toți neuronii $i = 1, 2, \dots, M_q$ aferenți stratului q . După ce eroarea este propagată dinspre ieșire spre intrarea RN prin toate straturile, ponderile sinaptice vor fi actualizate astfel:

$$\mathbf{W}^{(q)} := \mathbf{W}^{(q)} + \alpha \cdot \delta^{(q)} \cdot (\mathbf{y}^{(q-1)})^T \quad q = 1, 2, \dots, L \quad (5.6)$$

$$\text{unde: } \mathbf{y}^{(0)} = \mathbf{x}.$$

Rețelele recurente complet conectate, cum este cazul modelului *Hopfield*, pot memora și apoi detecta modelele de intrare. După aplicarea unui vector de intrare, ieșirea va fi calculată și reclusată pe intrare până când rețeaua va atinge o stare stabilă. Starea stabilă va reprezenta ieșirea \mathbf{y} asociată vectorului de intrare \mathbf{x} . O regulă de învățare posibilă pentru acest model ar fi regula lui *Hebb* (ecuația (5.3)), forțând $\mathbf{y} = \mathbf{x}$.

Un alt model de RN este rețeaua SOFM (*Self Organizing Feature Map*) a lui *Kohonen*. Există mai multe variante pentru algoritmul SOFM. Neuronii sunt organizați sub forma unei grile regulate pe care se definesc anumite relații topologice (de exemplu, distanțe). Toți neuronii vor recepționa același vector de intrare. Pentru fiecare vector de intrare \mathbf{x} , în spațiul n -dimensional de intrare, se vor calcula distanțele de la \mathbf{x} până la fiecare neuron i (vectorul ponderilor sinaptice

aferent neuronului i). Cele mai uzuale distanțe sunt produsul scalar $\left(\sum_{j=1}^N w_{i,j} \cdot x_j \right)$ și distanța euclidiană:

$$y_i = \Delta_{in}(\mathbf{x}, \mathbf{W}_i) = \sqrt{\sum_{j=1}^N (x_j - w_{i,j})^2} \quad ; \quad i = 1, 2, \dots, M \quad (5.7)$$

unde \mathbf{W}_i reprezintă rândul i al matricii \mathbf{W} (vectorul de ponderi sinaptice aferent neuronului i). Adesea, spațiul de intrare are o dimensiune mult mai mare decât spațiul topologic al rețelei SOFM. Neuronul învingător $I \in \{1, 2, \dots, M\}$ va fi cel care deține vectorul de ponderi sinaptice \mathbf{W}_I cel mai apropiat de vectorul de intrare \mathbf{x} , și va fi identificat prin:

$$\Delta_{in}(\mathbf{x}, \mathbf{W}_I) \leq \Delta_{in}(\mathbf{x}, \mathbf{W}_i) \quad \forall i \in \{1, 2, \dots, M\} \quad (5.8)$$

ponderile sinaptice vor fi apoi ajustate:

$$\mathbf{W}_i := \mathbf{W}_i + \alpha \cdot \lambda(\Delta_{topo}(i, I)) \cdot (\mathbf{x}^T - \mathbf{W}_i); \quad i = 1, 2, \dots, M \quad (5.9)$$

unde $\Delta_{topo}(i, I)$ reprezintă distanța de la neuronul i la neuronul învingător I în spațiul topologic al rețelei SOFM. Funcția de vecinătate λ va permite ajustarea ponderilor sinaptice numai pentru neuronii aflați în vecinătatea învingătorului I . Rezultatul obținut, după prezentarea repetată a vectorilor de intrare, va fi o auto-organizare a rețelei SOFM în populații de neuroni a căror selecție prin ecuația (5.8) este echiprobabilă. Într-o altă versiune a algoritmului SOFM, operațiile de detecție

Considerații privind perfecționarea RS dedicate algoritmilor conexioniști

a neuronului învingător I și de actualizare a ponderilor sinaptice pe baza funcției de vecinătate $\lambda(\Delta_{\text{topo}}(i, I))$, sunt înlocuite cu vectorul de ieșire aferent unei rețele recurente fixe care a fost lăsată să evolueze până la convergență. Această ultimă variantă a algoritmului SOFM este descrisă de ecuația (5.1) în care vectorul x aferent iterației curente va fi reprezentat de vectorul y aferent iterației anterioare.

5.2. Arhitectura RS

În capitolul 4 am proiectat o RS bidimensională dedicată algoritmilor conexioniști; funcționalitatea rețelei a fost demonstrată pentru faza de recunoaștere a rețelei *Hopfield* (RN recurentă care în faza de recunoaștere revendică un anumit număr de iterații până la convergență). Am arătat de asemenea că o mare parte a algoritmilor conexioniști pot fi descompuși în operații matriciale - vectoriale executate cu mare eficiență pe RS bidimensională. În cele ce urmează vom extinde schema de bază a RS proiectate în capitolul 4, în ideea integrării mai multor algoritmi conexioniști (inclusiv cu faza de învățare). Figura 5.1 prezintă arhitectura de bază a RS, în care se evidențiază următorii operanzi:

- matricea ponderilor sinaptice \mathbf{W} , rezidentă în cele $N \times N$ celule sistolice
- secvența de vectori de intrare \mathbf{I}^v (în partea superioară)
- secvența de vectori de intrare \mathbf{I}^h (pe latura de vest)

Rezultatele vor consta în actualizarea ponderilor sinaptice (faza de învățare) sau în vectorii de ieșire \mathbf{O}^v sau \mathbf{O}^h (faza de recunoaștere).

Evaluarea funcției de transfer σ din ecuația (5.1) se va face cu ajutorul unei tabele (*look - up table*) în exteriorul RS. Similar, scăderea și înmulțirea din ecuația (5.4) se vor executa în exteriorul RS, cu ajutorul unei rețele liniare formată din unități aritmetice auxiliare.

5.2.1. Fluxul feedforward de date

Prima clasă de operații executate în RS va produce ca rezultat un vector de ieșire \mathbf{O}^h sau \mathbf{O}^v ; aceste operații nu vor modifica matricea de ponderi sinaptice \mathbf{W} .

Operațiile executate de celula sistolică $CS_{i,j}$ vor fi:

$$h_{i,j} := \psi(w_{i,j}, v_{i-1,j}, h_{i,j-1}) \quad (5.10)$$

$$v_{i,j} := v_{i-1,j} \quad \text{pentru } i, j = 1, 2, \dots, N \quad (5.11)$$

Condițiile de margine stabilite pe cele 4 laturi ale RS vor fi:

$$h_{i,0} := I_i^h \quad (5.12)$$

$$v_{0,j} := I_j^v \quad (5.13)$$

$$O_i^h := h_{i,N} \quad (5.14)$$

$$O_j^v := v_{N,j} \quad \text{pentru } i, j = 1, 2, \dots, N \quad (5.15)$$

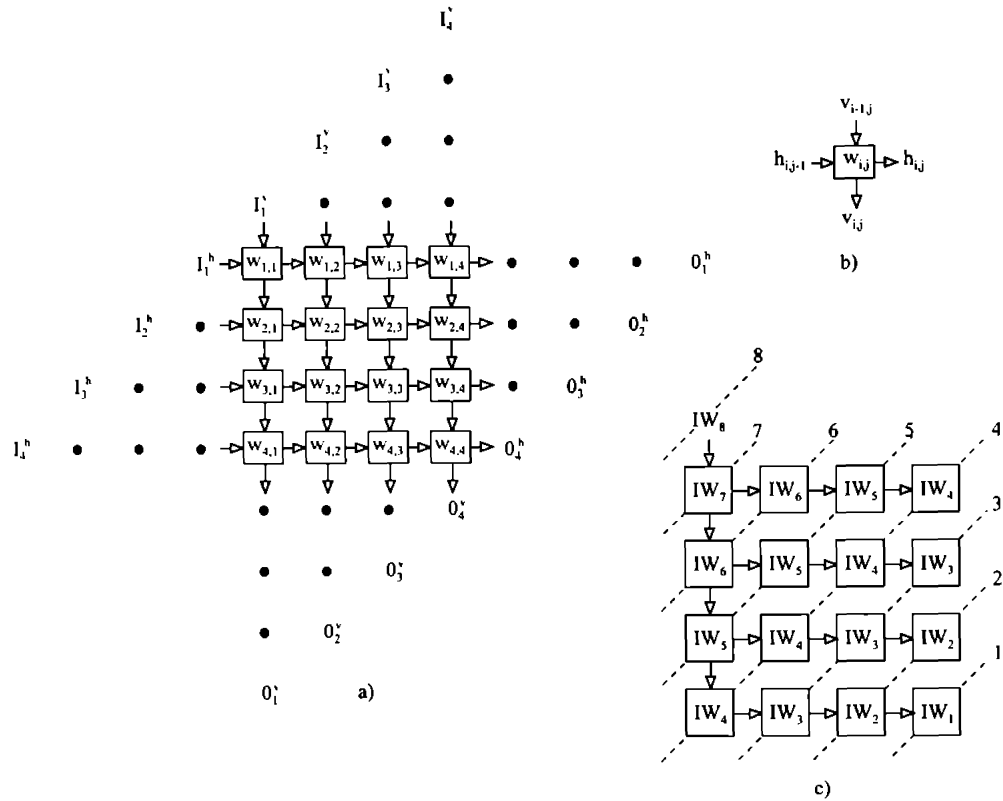


Fig. 5.1. a) Arhitectura RS. b) Operațiile de intrare / ieșire aferente unei CS.
c) Fluxul sistolic de instrucțiuni în RS (IW - instruction word)

În tabelul 5.1 se prezintă setul de funcții ψ necesare implementării principalilor algoritmi neuronali; toate aceste funcții utilizează fluxul sistolic de date descris matematic prin relațiile (5.10) ÷ (5.15) și grafic în figura 5.1. Simbolurile $+\infty$ și $-\infty$ specifică valoarea maximă și respectiv minimă reprezentabilă (limita maximă și respectiv minimă a domeniului de reprezentare); aceste limite depind de formatul de reprezentare (virgulă fixă, virgulă flotantă) și de numărul de biți alocați pentru reprezentare. Operația *mprod* reprezintă o operație clasică de înmulțire a unei matrici cu un vector. Operația *euclidean* este similară unei operații *mprod* în care produsul este înlocuit cu pătratul unei diferențe; această operație va permite calculul distanțelor euclidiene dintre un vector și fiecare rând aferent matricii sinaptice \mathbf{W} . Deoarece vectorul de distanțe rezultat este utilizat numai în ecuația (5.8) pentru identificarea unei valori minime, se poate renunța la calculul rădăcinii pătrate. Includerea termenului adițional I^h în cele două operații, *mprod* și *euclidean*, este absolut necesară pentru cumulara rezultatelor obținute pe matrici \mathbf{W} diferite. Această cumulare de rezultate devine necesară atunci când RS fizică este mai mică decât RN implementată pe RS. În acest caz, matricea \mathbf{W} va fi subdivizată într-un set de matrici care se vor încărca succesiv în RS iar rezultatele obținute se vor cumula pentru a genera un rezultat global.

Considerații privind perfecționarea RS dedicate algoritmilor conexioniști

Denumire operație	$\psi(a,b,c)$	Rezultatul obținut în modul normal de operare	Utilizare
<i>mprod</i>	$a \cdot b + c$	$\mathbf{O}^h = \mathbf{W} \cdot \mathbf{I}^v + \mathbf{I}^h$	ec. (5.1)
<i>euclidean</i>	$(b - a)^2 + c$	$O_i^h = \sum_{j=1}^N (I_j^v - w_{i,j})^2 + I_i^h$	ec. (5.7)
<i>min</i>	c dacă $c \leq b$ $+\infty$ dacă $c > b$	Dacă $I_i^h = \min_i (I_i^h)$ atunci $O_i^h = I_i^h$, toate celelalte elemente din vectorul \mathbf{O}^h fiind setate la $+\infty$	ec. (5.8)
<i>max</i>	c dacă $c \geq b$ $-\infty$ dacă $c < b$	Dacă $I_i^h = \max_i (I_i^h)$ atunci $O_i^h = I_i^h$, toate celelalte elemente din vectorul \mathbf{O}^h fiind setate la $-\infty$	ec. (5.8)

Tabelul 5.1. Operațiile impuse de faza *feedforward* (recunoaștere) în modul de lucru normal

Operațiile *min* și *max* sunt operații de căutare al celui mai mic, respectiv mai mare, element dintr-un vector. Algoritmul utilizat impune ca vectorul investigat să fie introdus în RS simultan atât din direcția vest (\mathbf{I}^h) cât și din direcția nord (\mathbf{I}^v). Aceasta înseamnă că toate componentele introduse orizontal vor întâlni succesiv toate componentele introduse vertical. La nivelul fiecărei CS, pentru operația *min* de exemplu, componenta care intră în CS dinspre vest, va putea tranzita CS orizontal (nemodificată) numai dacă este mai mică decât componenta care intră în CS din direcția nord; în caz contrar componenta care tranzitează orizontal CS va fi saturată la $+\infty$. Cum toate elementele care tranzitează RS orizontal vor întâlni (succesiv) toate elementele care tranzitează RS vertical, pe latura de est a RS vom obține nesaturată numai cea mai mică componentă din vectorul investigat. Operația *min* va fi utilizată pentru calculul distanțelor euclidiene în timp ce operația *max* va fi utilizată la normalizarea ponderilor sinaptice și la calculul produsului scalar (produsul scalar maxim va reveni neuronului al cărui vector de ponderi sinaptice normalizat formează unghiul minim cu vectorul de intrare normalizat). Acest neuron va răspunde preferențial la vectorul de intrare respectiv.

Să notăm că retro-propagarea erorii (faza *backward*) revendică tot înmulțirea unei matrici cu un vector (ec. 5.5), dar matricea sinaptică utilizată în (5.5) este transpusa matricii utilizate în (5.1) (faza *feedforward*). Efectuarea calculelor care implică transpusa matricii \mathbf{W} reprezintă o sarcină simplă care se poate realiza prin simpla interschimbare a rolurilor celor două fluxuri de date din RS (fluxul orizontal și respectiv vertical). Pentru operațiile care implică \mathbf{W}^T , relațiile (5.10) și (5.11) devin:

$$h_{i,j} := h_{i,j-1} \quad (5.16)$$

$$v_{i,j} := \psi(w_{i,j}, h_{i,j-1}, v_{i-1,j}) \quad \text{pentru } i, j = 1, 2, \dots, N \quad (5.17)$$

Rolul intrărilor va fi de asemenea interschimbabil, iar ieșirea va fi:

$$\mathbf{O}^v = \mathbf{W}^T \cdot \mathbf{I}^h + \mathbf{I}^v \quad (5.18)$$

Deși toate celelalte operații din tabelul 5.1 (*euclidean*, *min*, *max*) pot fi executate în modul transpus (arhitectura RS permite), aceste noi operații care implică \mathbf{W}^T nu se vor dovedi necesare.

5.2.2. Fluxul de date impus de ajustarea ponderilor sinaptice

A doua clasă de operații executate pe RS va realiza ajustarea ponderilor sinaptice (\mathbf{W}) fără a produce nici un alt rezultat pe ieșirile \mathbf{O}^h sau \mathbf{O}^v . Operațiile executate de celula sistolică $CS_{i,j}$ vor fi:

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

$$h_{i,j} := h_{i,j-1} \quad (5.19)$$

$$v_{i,j} := v_{i-1,j} \quad (5.20)$$

$$w_{i,j} := \psi(w_{i,j}, v_{i-1,j}, h_{i,j-1}) \quad \text{pentru } i, j = 1, 2, \dots, N \quad (5.21)$$

Condițiile de margine stabilite pe laturile RS sunt cele fixate prin ecuațiile (5.12) și (5.13). Tabelul 5.2 prezintă setul de funcții ψ necesare pentru implementarea fazei de învățare (ajustarea ponderilor sinaptice).

Denumire operație	$\psi(a,b,c)$	Rezultatul obținut în modul normal de operare	Utilizare
<i>hebbian</i>	$a + c \cdot b$	$\mathbf{W} = \mathbf{W} + \mathbf{I}^h \cdot (\mathbf{I}^v)^T$	ec. (5.3) și (5.5)
<i>Kohonen</i>	$a + c \cdot (b - a)$	$\mathbf{W}_i = \mathbf{W}_i + \mathbf{I}_i^h \cdot ((\mathbf{I}^v)^T - \mathbf{W}_i)$	ec. (5.9)

Tabelul 5.2. Operațiile impuse de faza de învățare (*backward*) în modul de lucru normal

Pentru ambele reguli de învățare prezentate în tabelul 5.2, vectorul \mathbf{I}^h este deja înmulțit cu constanta de învățare α (în unități aritmetice auxiliare situate în exteriorul RS). Termenul $\lambda(\Delta_{\text{topo}}(i, j))$ din ecuația (5.9) va fi calculat astfel:

1. Un vector care va conține cele m distanțe (m fiind numărul de neuroni din RN) va fi obținut printr-o operație *min* (distanțe euclidiene) sau o operație *max* (produs scalar). O tabelă externă RS (*look - up table*) va fi utilizată pentru conversia tuturor valorilor $+\infty$ (în cazul *min*) și respectiv $-\infty$ (în cazul *max*) la 0; doar valoarea minimă detectată (în cazul *min*) și respectiv maximă (în cazul *max*) va fi convertită la valoarea 1.
2. Acest vector obținut (cu $m-1$ componente 0 și o componentă 1) va fi apoi multiplicat cu o matrice Λ care exprimă în același timp și distanțele topologice dintre neuroni și funcția de vecinătate utilizată. Elementele acestei matrici simetrice sunt $\Lambda_{i,j} = \lambda(\Delta_{\text{topo}}(i, j))$. Această formulare matricială a distanței topologice are avantajul că permite efectuarea calculului pentru orice tip de vecinătate (paragraful 4.5.3. A).

5.3. Fluxul sistolic al instrucțiunilor

RS este capabilă să efectueze o procesare în tehnică *pipeline*. Dacă numărul de nivele *pipeline* (marcate cu linii punctate în fig. 5.1. c) este T (incluzând nivelele impuse de tabelele *look - up* și de unitățile aritmetice auxiliare externe RS), atunci T vectori de intrare pot fi procesați concurrent ajungându-se astfel la o eficiență maximă a RS. Dacă un grup de T vectori se aplică succesiv pe intrarea RS, primul rezultat va fi disponibil pe ieșirea RS în ciclul care urmează ciclului de injecție în RS a ultimului vector de intrare. RS va genera în continuare următoarele $T-1$ rezultate în timp ce la intrarea RS va trebui setat un nou mod de operare. Procesarea SIMD (*Single instruction flow, multiple data flow*) clasică presupune golirea RS înainte de startarea unui nou mod de operare (a unei noi instrucțiuni). Pentru a îmbunătăți performanțele RS, fiecărui vector de intrare i se poate atașa un cod de instrucțiune (*IW - instruction word*) care va fi propagat sistolic (tehnică *pipeline*) în interiorul RS. În figura 5.1. c) se indică modul de propagare a IW în RS împreună cu nivelele *pipeline* rezultate. Aceasta permite schimbarea progresivă a modului de operare în

Considerații privind perfecționarea RS dedicate algoritmilor conexioniști

interiorul RS, fără a fi necesară golirea RS la trecerea de la un mod de operare (instrucțiune) la alt mod (instrucțiune).

5.4. Structura celulei sistolice

Arhitectura RS și structura CS sunt prezentate figura 5.2.

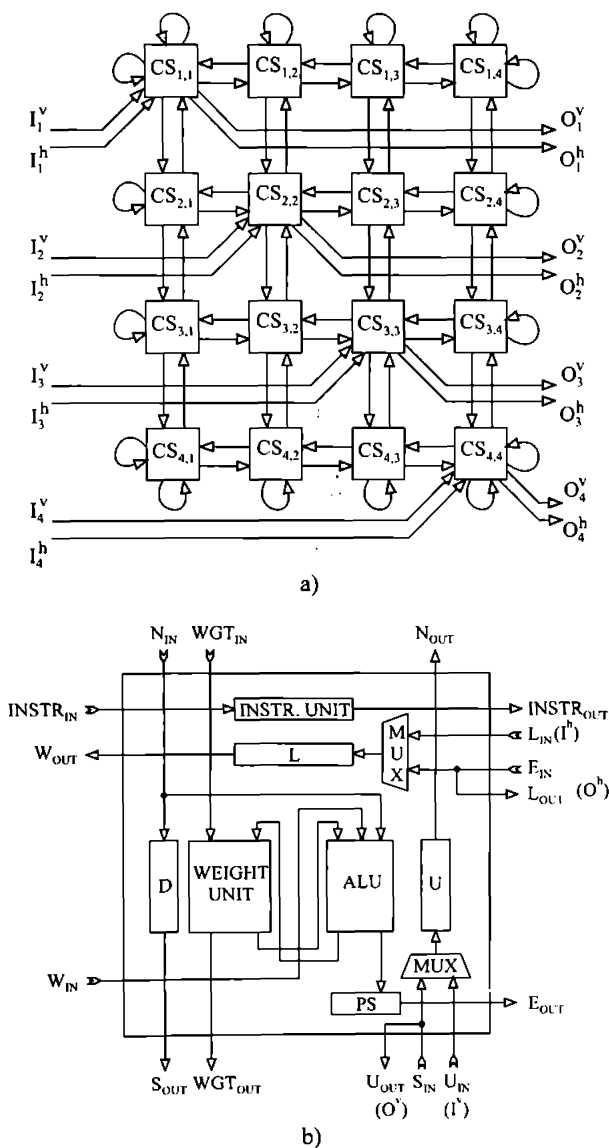


Fig. 5.2. a) Arhitectura RS 4 x 4. b) Structura CS

Spre deosebire de arhitectura de principiu prezentată în figura 5.1 a), operațiile de intrare / ieșire sunt implementate cu ajutorul unor porturi locațe în CS diagonale. Intrările U_{IN} și L_{IN} și ieșirile U_{OUT} și L_{OUT} vor fi prezente în toate celulele sistolice (structură uniformă) dar vor fi utilizate numai în cazul celulelor diagonale; în celulele non - diagonale vor rămâne neconectate. Avantajul major al acestui sistem de intrare / ieșire constă în faptul că nu mai este nevoie de un *hardware* adițional pentru diagonalizarea vectorilor de intrare și de-diagonalizare a vectorilor de ieșire (capitolul 4). Toate componentele unui vector sunt introduse sau extrase simultan.

Mecanismul de transpunere, descris de ecuațiile (5.16) și (5.17), poate fi implementat cu ajutorul a două semnale de control furnizate de INSTRUCTION UNIT. Primul semnal va interschimba intrările N_{IN} și W_{IN} , iar al doilea va interschimba ieșirile S_{OUT} și E_{OUT} (figura 5.3).

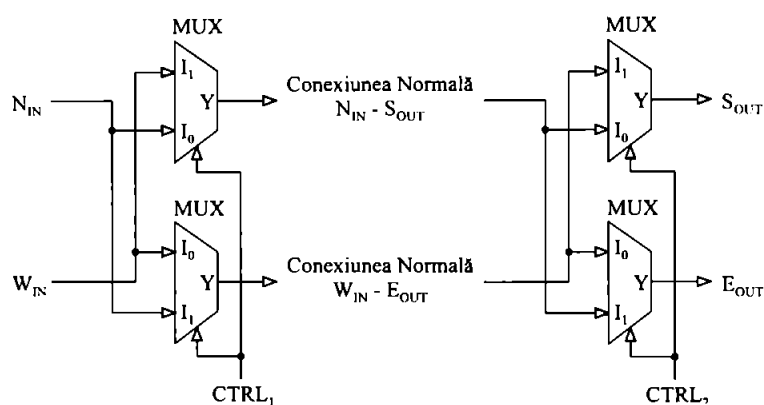


Fig. 5.3. Mecanismul care permite operarea cu transpusa matricii W

Când startează o operație care implică W^T , intrările vor fi interschimbate (se activează $CTRL_1$). Registrele interne CS care conțin valori aferente operației anterioare (care implică W) sunt emise pe căile normale de ieșire din CS. În următorul ciclu vor fi interschimbate și ieșirile (se activează $CTRL_2$). Aceeași succesiune de operații va avea loc când RS revine la modul direct de operare (matricea W).

Comunicațiile din interiorul RS sunt locale (între CS vecine) și vor fi seriale (primul bit transmis va fi bitul cel mai puțin semnificativ). Protocolul serial se impune din două considerente:

- un înmulțitor paralel ar avea o logică foarte complexă și va reduce numărul de CS implementabile într-un circuit (paragraful 4.5.4)
- comunicația paralelă ar revendica un număr mare de pini pentru circuitul VLSI.

5.4.1. Dublarea registrului W

În figura 5.2. b) sunt prezentate două căi de comunicație între celulele sistolice care nu sunt reprezentate și figura 5.2. a):

- legătura dedicată ponderilor sinaptice (nord → sud)
- legătura dedicată instrucțiunilor (vest → est)

Considerații privind perfecționarea RS dedicate algoritmilor conexioniști

Legătura dedicată instrucțiunilor a fost deja discutată în corelație cu figura 5.1. c). Registrele dedicate ponderilor sinaptice (*weight unit*) sunt încărcate prin intermediul unor linii de comunicație dedicate (verticale) și independente de liniile de intrare standard. Cheia succesului pentru un *hardware* accelerator dedicat RN constă în abilitatea acestuia de a emula RN mai mari decât dimensiunea fizică pe care RS o pune la dispoziție. Unitatea *weight unit* conține două registre dedicate ponderilor sinaptice, unul conectat pe linia WGT_{IN} - WGT_{OUT} și celălalt utilizat în calcule (furnizează ponderea sinaptică pentru calculele efectuate în CS respectivă). Multiplexoarele din fig. 5.4 vor realiza interschimbarea celor două registre între două instrucțiuni consecutive.

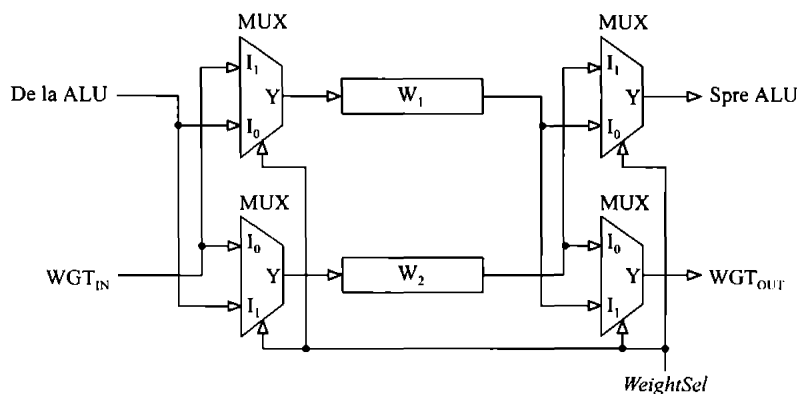


Fig. 5.4. Conectarea celor două registre din unitatea *Weight Unit*

În acest fel, o nouă matrice sinaptică poate fi încărcată în RS (de exemplu în registrele W_2) în timp ce RS procesează matricea anterioară (preluată din registrele W_1). În cazul RN mari, care depășesc capacitatea RS, cele 2 matrici sinaptice vor proveni din aceeași matrice care a trebuit să fie partiționată. În figura 5.5 se prezintă timpii de interschimbare pentru situația în care fiecare nouă epocă (care procesează T vectori de intrare) necesită încărcarea unei noi matrici sinaptice.

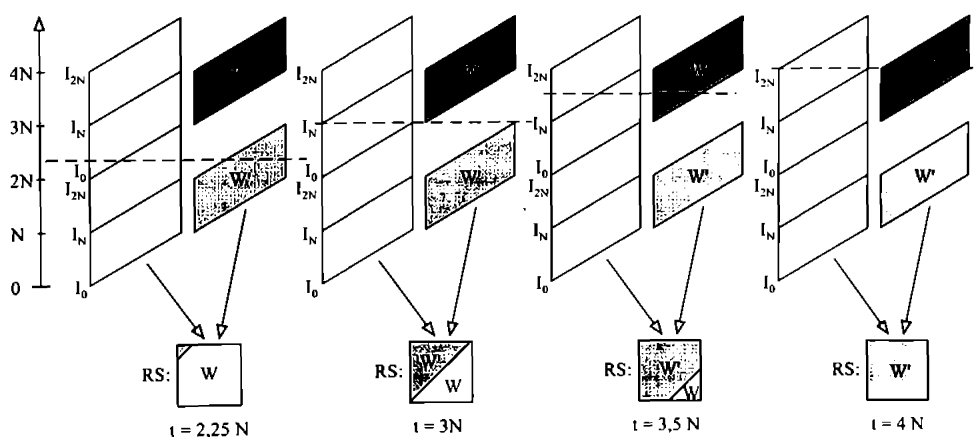


Fig. 5.5. Mecanismul de interschimbare a matricilor sinaptice

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

Pătratele din subsolul figurii reprezintă RS ($N \times N$) iar zonele delimitate în interior indică matricea sinaptică utilizată în calcule de către celulele sistolice din zona respectivă. Fiecare epocă va conține $2N$ vectori de intrare. $2N$ reprezintă și latența *pipeline*-ului la nivelul instrucțiunii propagate în RS (fig. 5.1. c)). Grupurile de $2N$ vectori de intrare, fiecare compus din N elemente, sunt reprezentate prin intermediul unor paralelograme în figura 5.5. Vectorii de intrare sunt figurați ca și cum ar fi injectați în RS din direcția nord; în realitate aceștia sunt injectați în celulele sistolice de pe diagonala RS. Toate componentele unui vector intră simultan în RS (capitolul 4). Până în momentul $t = 2N$ RS utilizează în calcul matricea sinaptică W (injectată anterior). De la $t = 0$ până la $t = 2N$ se încarcă în RS (în registrele W libere) matricea sinaptică W' . În momentul $t = 2,25 \cdot 3/4$ din celulele sistolice aferente primei linii din RS utilizează încă în calcul matricea W ; matricea W' va intra în calcul progresiv începând cu celula sistolică situată în colțul de nord - vest al RS. În momentul $t = 3N$ (la jumătatea celei de a doua epoci reprezentate în figură), jumătate din celulele sistolice utilizează în calcul matricea W' . Linia oblică împarte RS în două părți: partea nehașurată reprezintă celulele sistolice care utilizează în calcul matricea W și cea hașurată reprezintă celulele care utilizează în calcul matricea W' . Linia oblică va intersecta în permanență setul de celule în care se realizează interschimbarea registrelor W (semnalul *Weight Sel* din figura 5.4). La momentul $t = 4N$ toate celulele sistolice utilizează în calcul matricea W' . Încărcarea matricii W'' demarează în momentul $t = 3N$ iar în calcul aceasta va fi utilizată începând din momentul $t = 4N$. Acest proces revendică o epocă cu lungimea $2N$ vectori de intrare (egală cu latența *pipeline*-ului).

Matricile sinaptice diagonalizate sunt introduse din direcția nord în RS. Trei registre W implementate în fiecare celulă sistolică ar fi permis încărcarea normală a matricilor sinaptice (nediagonalizate); un registru suplimentar crește însă complexitatea CS și reduce numărul de celule integrabile într-un singur circuit VLSI.

5.4.2. Precizia calculelor

Stabilirea preciziei care se impune în calcul este o problemă delicată. O primă variantă, motivată de simplitatea *hardware*-ului, constă în reprezentarea tuturor variabilelor în virgulă fixă (complement de 2). Cât despre numărul de biți necesari în reprezentare, analizele teoretice conduc de regulă spre valori minime care însă nu garantează convergența pentru orice aplicație. În absența unor studii analitice clarificatoare, proiectanții recurg adesea la simulări experimentale. În [ASA 90] și [HOL 91] sunt prezentate simulări pentru determinarea preciziei revendicate de anumite aplicații tipice ale algoritmului EBP. Astfel de simulări au fost realizate și pentru rețeaua SOFM a lui *Kohonen*. Concluziile acestor investigații pot fi sintetizate astfel:

- Ponderile sinaptice pentru operațiile *mprod* și *euclidean* (faza *feedforward*) oferă suficientă precizie dacă sunt reprezentate pe 16 biți. Simulările au arătat că pentru aplicațiile rețelei SOFM sunt necesari mai mult de 9 - 10 biți în reprezentare. Este deci oportună alegerea primului număr multiplu de 8 superior acestei valori. În concluzie $N_W = 16$ biți. În mod similar, pentru vectorii de intrare se poate conveni $N_D = 16$ biți. $N_W = N_D$ deoarece, în cazul rețelei SOFM, ponderile sinaptice și componentele vectorilor de intrare aparțin aceluiași spațiu.
- Notăm cu n numărul maxim de intrări în RN. Parametrul n este independent de dimensiunea RS. În aplicațiile reale în general n nu depășește valoarea $256 = 2^8$. Numărul de biți necesari reprezentării sumelor parțiale PS (în cadrul operației *mprod*) va fi:

$$N_{PS} = N_W + N_D + \lceil \log_2(n) \rceil \quad (5.22)$$

Considerații privind perfecționarea RS dedicate algoritmilor conexioniști

Obținem valoarea $N_{PS} = 16 + 16 + 8 = 40$ biți. Această valoare permite implementarea RN cu până la 256 intrări.

- Practic, datorită mecanismului de transpunere a matricii W și datorită operațiilor *min* și *max*, este de dorit ca toate registrele de date din structura CS să conțină 40 biți. Pentru a evita depășirile de capacitate în timpul calculelor este de dorit ca inițializarea RS să se facă cu $N_W = N_D = 16$ biți.

5.4.3. Unitatea ALU

ALU execută operațiile listate în tabelele 5.1 și 5.2 și va conține 3 blocuri: un înmulțitor, un sumator și un scăzător. Stabilirea structurii circuitului dedicat înmulțirilor reprezintă o chestiune delicată. Operația *mprod* nu pune probleme deosebite. Din acest punct de vedere se poate implementa un înmulțitor serie - paralel deoarece coeficientul sinaptic este rezident în celula sistolică și poate fi preluat în paralel, în timp ce celălalt factor este preluat de pe linia serială de comunicație nord - sud (sub formă serială). Din păcate această schemă nu poate fi utilizată și în operațiile de ajustare a ponderilor sinaptice (*hebbian* și *Kohonen*) și nici pentru operația *euclidean*, deoarece ambii factori sunt recepționați sub formă serială de pe liniile de comunicație intercelulare. O primă soluție ar consta în implementarea următoarei secvențe de operații:

- recepția serială a unui factor și memorarea sa într-un *buffer* în interiorul CS
- recepția serială a celui de-al doilea factor și execuția concomitentă a operației de înmulțire cu ajutorul înmulțitorului serie - paralel.

Această soluție permite utilizarea unui înmulțitor cu structură simplă (serie - paralel) dar implică o pierdere de performanță (un ciclu de transfer suplimentar pentru transferul și memorarea primului operand).

O a doua soluție ar consta în utilizarea unui înmulțitor serie - serie. Înmulțitoarele serie - serie tipice prezintă de regulă un număr de cicluri de latență între momentul aplicării primilor biți pe intrare și momentul obținerii primului bit pe ieșire [DAD 89]. Această latență conduce de asemenea la degradarea performanțelor RS. Problematika înmulțitoarelor serie - serie performante (zero cicluri latență) a fost tratată și în [IEN 94]. Fiecare CS revendică un înmulțitor serie - serie pe 17 biți; în cele ce urmează vom propune un astfel de înmulțitor care prezintă 2 avantaje majore:

- zero cicluri latență;
- extensia automată a semnului produsului atâta timp cât pe intrare cei doi factori își mențin extensia de semn.

Să notăm că ALU recepționează operanzii direct de pe liniile de comunicație externe (fig. 5.2) și încarcă rezultatele serial în registre (PS, W_1 , W_2) plasate pe ieșirile ALU. Acest fapt simplifică implementarea operațiilor *min* și respectiv *max*, deoarece permite, pe ultimul ciclu de tact (la încheierea operației de comparare), saturarea la $+\infty$ sau $-\infty$ a valorii propagate în registrul PS.

A. Notații și simboluri utilizate

Vom considera că cei 2 factori aplicați pe intrarea înmulțitorului, X și Y , au N biți și sunt reprezentați în cod complementar (virgulă fixă). Produsul P obținut la ieșire va avea K biți, unde $K \geq 2N$. În figura 5.6 sunt reprezentate simbolurile logice utilizate. Simbolurile a) până la e) reprezintă elemente de întârziere (bistabile D care întârzie semnalul de la intrare cu un ciclu de tact). Pentru simplificare nu s-a mai reprezentat și intrarea de *clock*.

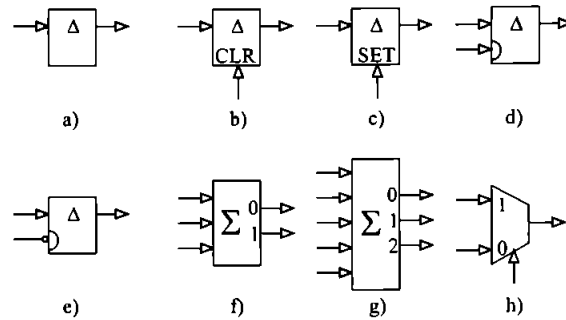


Fig. 5.6. Simbolurile aferente elementelor utilizate. a) Element de întârziere (bistabil D).
 b) Element de întârziere cu intrare de ștergere sincronă activă *high*.
 c) Element de întârziere cu intrare de set sincronă activă *high*.
 d) Element de întârziere cu intrare de validare (*enable*) activă *high*.
 e) Element de întârziere cu intrare de validare (*enable*) activă *low*.
 f) (3,2) counter (full-adder). g) (5,3) counter. h) multiplexor 2:1

Simbolurile d) și e) au un semnal de validare activ *high* și respectiv *low*. Pentru a încărcă valoarea de pe intrare în element, semnalul de validare trebuie să fie activ în momentul tranziției semnalului de *clock*.

Termenul (*n, k*) counter desemnează un dispozitiv cu *n* intrări și *k* ieșiri; codul binar furnizat pe cele *k* ieșiri reprezintă numărul de intrări active. O definiție alternativă a acestui dispozitiv ar fi aceea de sumator care însumează *n* cuvinte de un bit aplicate la intrare și generează suma reprezentată pe *k* biți la ieșire. Numărul de ieșiri va fi: $k = \lceil \log_2(n + 1) \rceil$

B. Algoritmi fundamentali

În figura 5.7 este redat algoritmul de înmulțire fundamental rezultat din definiția codului complementar (complementul față de 2).

				x_3	x_2	x_1	x_0
				y_3	y_2	y_1	y_0
+					x_2y_0	x_1y_0	x_0y_0
-				x_3y_0			
+				x_2y_1	x_1y_1	x_0y_1	
-			x_3y_1				
+			x_2y_2	x_1y_2	x_0y_2		
-		x_3y_2					
-		x_2y_3	x_1y_3	x_0y_3			
+	x_3y_3						
	p_7	p_6	p_5	p_4	p_3	p_2	p_1
					p_3	p_2	p_1
						p_2	p_1
							p_1
							p_0

Fig. 5.7. Algoritmul de înmulțire dedicat numerelor reprezentate în cod complementar

Considerații privind perfecționarea RS dedicate algoritmilor conexioniști

Algoritmul calculează produsul final pe baza unor produse parțiale. Primii $N - 1$ biți ai produsului parțial se adună la suma de produse anterior obținută, iar bitul N (cel mai semnificativ) va fi scăzut din suma de produse anterior obținută. Excepție de la această regulă face ultimul produs parțial unde operațiile de adunare și respectiv scădere sunt inversate. Acest algoritm revendică un control complex (alternarea operațiilor de adunare și respectiv scădere) și din acest motiv nu reprezintă o bună soluție pentru implementarea în *hardware*.

O posibilă simplificare a algoritmului constă în extensia de semn la nivelul produselor parțiale până la lungimea K (lungimea produsului final), și execuția "normală" a operațiilor de adunare și respectiv scădere a produselor parțiale, păstrând de fiecare dată doar cei mai puțini semnificativi K biți. Algoritmul rezultat este reprezentat în figura 5.8 unde s-a considerat $K = 9$. Acest algoritm poate fi privit și ca pornind de la un operand X cu semn extins la lungimea de K biți.

				x_3	x_3	x_3	x_3	x_3	x_3	x_2	x_1	x_0
									y_3	y_2	y_1	y_0
+				x_3y_0	x_3y_0	x_3y_0	x_3y_0	x_3y_0	x_3y_0	x_2y_0	x_1y_0	x_0y_0
+			x_3y_1	x_3y_1	x_3y_1	x_3y_1	x_3y_1	x_3y_1	x_2y_1	x_1y_1	x_0y_1	
+		x_3y_2	x_3y_2	x_3y_2	x_3y_2	x_3y_2	x_3y_2	x_2y_2	x_1y_2	x_0y_2		
-	x_3y_3	x_3y_3	x_3y_3	x_3y_3	x_3y_3	x_3y_3	x_2y_3	x_1y_3	x_0y_3			
				p_7	p_7	p_6	p_5	p_4	p_3	p_2	p_1	p_0

Fig. 5.8. Algoritmul de înmulțire dedicat numerelor reprezentate în cod complementar cu X reprezentat cu extensie de semn

Totuși algoritmul din figura 5.8 prezintă dezavantajul că ultimul produs parțial trebuie scăzut în timp ce toate celelalte produse parțiale trebuie adunate. Acest dezavantaj poate fi eliminat dacă se realizează o extensie de semn prealabilă și la operandul Y . Se va intra în operație cu ambii operanzi extinși la K biți și după obținerea produsului final se vor reține doar cei mai puțini semnificativi K biți ai acestuia. Acest algoritm, prezentat în figura 5.9, este unul dintre cei mai utilizați algoritmi. Atât algoritmul din figura 5.8 cât și cel din figura 5.9 generează rezultatul cu extensie de semn.

					x_3	x_3	x_3	x_3	x_3	x_3	x_2	x_1	x_0
					y_3	y_3	y_3	y_3	y_3	y_3	y_2	y_1	y_0
+					x_3y_0	x_3y_0	x_3y_0	x_3y_0	x_3y_0	x_3y_0	x_2y_0	x_1y_0	x_0y_0
+				x_3y_1	x_3y_1	x_3y_1	x_3y_1	x_3y_1	x_3y_1	x_2y_1	x_1y_1	x_0y_1	
+			x_3y_2	x_3y_2	x_3y_2	x_3y_2	x_3y_2	x_3y_2	x_2y_2	x_1y_2	x_0y_2		
+		x_3y_3	x_3y_3	x_3y_3	x_3y_3	x_3y_3	x_3y_3	x_2y_3	x_1y_3	x_0y_3			
+		x_3y_3	x_3y_3	x_3y_3	x_3y_3	x_3y_3	x_3y_3	x_2y_3	x_1y_3	x_0y_3			
+		x_3y_3	x_3y_3	x_3y_3	x_3y_3	x_3y_3	x_3y_3	x_2y_3	x_1y_3	x_0y_3			
+		x_3y_3	x_3y_3	x_3y_3	x_3y_3	x_3y_3	x_3y_3	x_2y_3	x_1y_3	x_0y_3			
+		x_3y_3	x_3y_3	x_3y_3	x_3y_3	x_3y_3	x_3y_3	x_2y_3	x_1y_3	x_0y_3			
+	x_3y_3	x_3y_3	x_3y_3	x_3y_3	x_3y_3	x_3y_3	x_2y_3	x_1y_3	x_0y_3				
					p_7	p_7	p_6	p_5	p_4	p_3	p_2	p_1	p_0

Fig. 5.9. Algoritmul de înmulțire a 2 numere reprezentate în cod complementar cu extensie de semn

C. Înmulțitorul serie - serie

Utilizând algoritmul din figura 5.9, numărul de termeni care trebuie adunați pentru a genera un bit p_i al produsului (numărul de termeni de pe o coloană) crește liniar cu ponderea i . Pentru rangurile p_i de pondere mai mare sau egală cu $2^{2^{N-1}}$ cei mai mulți termeni care trebuie adunați se datorează extensiei de semn și au forma $x_{N-1} \cdot y_{N-1}$. Se poate arăta că regularitatea acestor termeni face posibilă ignorarea lor fără a afecta corectitudinea rezultatului.

Procedând într-o manieră similară celei expuse în [STR 82], produsul parțial obținut după aplicarea bitului i poate fi exprimat sub o formă recursivă:

$$P_i = X_i \cdot Y_i = P_{i-1} + x_i \cdot Y_{i-1} \cdot 2^i + y_i \cdot X_{i-1} \cdot 2^i + x_i \cdot y_i \cdot 2^{2i} \quad (5.23)$$

unde X_i și Y_i reprezintă valorile celor 2 operanzi considerând doar biții de la 0 la i (deci, $X_i = X \bmod 2^{i+1}$). Valorile inițiale sunt: $P_{-1} = X_{-1} = Y_{-1} = 0$. Dacă extensia semnului este menținută până în momentul declanșării următoarei înmulțiri, numai primii $N-1$ biți aferenți celor 2 operanzi vor trebui memorati (biții de semn vor rămâne în continuare stabili pe intrări până la declanșarea următoarei operații de înmulțire). Această idee poate sugera un algoritm în care, în ciclurile $i \geq N$, operanzii parțiali se consideră estinși numai până la bitul $N-2$ și astfel termenii simetrici nu vor mai intra în sumă:

$$\bar{P}_i = \bar{P}_{i-1} + x_{N-1} \cdot Y_{N-2} \cdot 2^i + y_{N-1} \cdot X_{N-2} \cdot 2^i \quad \text{pentru } i \geq N \quad (5.24)$$

$$\bar{P}_i = P_i \quad \text{pentru } i < N \quad (5.25)$$

Eroarea indusă prin utilizarea relației (5.24) în locul relației (5.23), în ciclul $j \geq N$ va fi:

$$\begin{aligned} (P_j - P_{j-1}) - (\bar{P}_j - \bar{P}_{j-1}) &= x_j \cdot Y_{j-1} \cdot 2^j - x_{N-1} \cdot Y_{N-2} \cdot 2^j + y_j \cdot X_{j-1} \cdot 2^j - \\ &- y_{N-1} \cdot X_{N-2} \cdot 2^j + x_j \cdot y_j \cdot 2^{2j} = x_{N-1} \cdot y_{N-1} \cdot \left(2^{2j} + 2 \cdot 2^j \cdot \sum_{k=N-1}^{j-1} 2^k \right) \end{aligned} \quad (5.26)$$

Ultima expresie se obține considerând că $x_j = x_{N-1}$, pentru $j \geq N$, care implică de asemenea $Y_{j-1} = Y_{N-2} + y_{N-1} \cdot \sum_{k=N-1}^{j-1} 2^k$. Relații similare pot fi deduse pentru y_j și X_{j-1} . Prin urmare, după ciclul i , eroarea totală va fi:

$$\begin{aligned} E_i &= x_{N-1} \cdot y_{N-1} \cdot \sum_{j=N}^i \left(2^{2j} + 2 \cdot 2^j \cdot \sum_{k=N-1}^{j-1} 2^k \right) = \\ &= x_{N-1} \cdot y_{N-1} \cdot \sum_{j=N}^i \left(2^{2j} + \sum_{k=N}^j 2^{j+k} \right) = \\ &= x_{N-1} \cdot y_{N-1} \cdot \sum_{j=N}^i \left(2^{2j} + \sum_{k=j}^i 2^{j+k} \right) = \\ &= x_{N-1} \cdot y_{N-1} \cdot 2^{i+1} \cdot \sum_{j=N}^i 2^j \end{aligned} \quad (5.27)$$

Considerații privind perfecționarea RS dedicate algoritmilor conexiști

Din relația (5.27) rezultă că: $E_i \bmod 2^{i+1} = 0$, pentru orice i . Aceasta indică faptul că partea ignorată din suma parțială nu va afecta ieșirea serială a înmulțitorului (produsul) deoarece, până la ciclul $i = K - 1$ (la sfârșitul operației de înmulțire), vor fi generați doar cei mai puțini semnificativi K biți ai produsului.

Mai mult, dacă se înmulțesc numere fără semn (extinse cu zero), E_i este identică cu zero. Prin urmare expresia \bar{P}_i poate fi utilizată la proiectarea unui înmulțitor care va opera corect atât pentru operanzi reprezentați în cod complementar (cu extensie de semn) cât și pentru operanzi fără semn (extinși cu zero).

Tehnica clasică *add & shift* poate fi utilizată la implementarea algoritmului dacă se introduce:

$$Q_i = \frac{\bar{P}_i}{2^{i+1}} = \frac{1}{2} (Q_{i-1} + x_i \cdot Y_{i-1} + y_i \cdot X_{i-1} + 2^i \cdot x_i \cdot y_i) \quad \text{pentru } i < N \quad (5.28)$$

$$Q_i = \frac{\bar{P}_i}{2^{i+1}} = \frac{1}{2} (Q_{i-1} + x_{N-1} \cdot Y_{N-2} + y_{N-1} \cdot X_{N-2}) \quad \text{pentru } i \geq N \quad (5.29)$$

Ecuția (5.28) arată că fiecare produs parțial poate fi generat cu o întârziere de tip combinațional (de propagare printr-o logică combinațională) utilizând valorile memorate ale celor 2 operanzi dar numai până la nivelul bitului precedent (X_{i-1} și Y_{i-1}) și biții de intrare curenți (x_i și y_i). Termenul simetric $x_i \cdot y_i$ trebuie adunat *on-line* la bitul de pe rangul i . Termenul $1/2$ echivalează cu o deplasare aritmetică la dreapta și partea fracțională (bitul) care rezultă din această deplasare reprezintă noul bit aferent produsului final. Înainte de deplasare, valoarea care se adună produsului parțial obținut în ciclul anterior va fi o valoare întreagă. Deci, dacă restul împărțirii (deplasării) este extras din înmulțitor la fiecare iterație, vom putea construi un sistem echivalent care va memora numai partea întreagă a rezultatului:

$$\bar{Q}_i = \left\lfloor \frac{1}{2} \cdot (Q_{i-1} + x_i \cdot Y_{i-1} + y_i \cdot X_{i-1} + 2^i \cdot x_i \cdot y_i) \right\rfloor \quad \text{pentru } i < N \quad (5.30)$$

O expresie similară poate fi derivată din (5.29) și pentru $i \geq N$.

Algoritmul poate fi înțeles mai ușor dacă se prezintă prin comparație cu algoritmul de bază din figura 5.9. În figura 5.10 se prezintă un exemplu de înmulțire a 2 operanzi pe 4 biți, care generează un produs pe 9 biți.

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

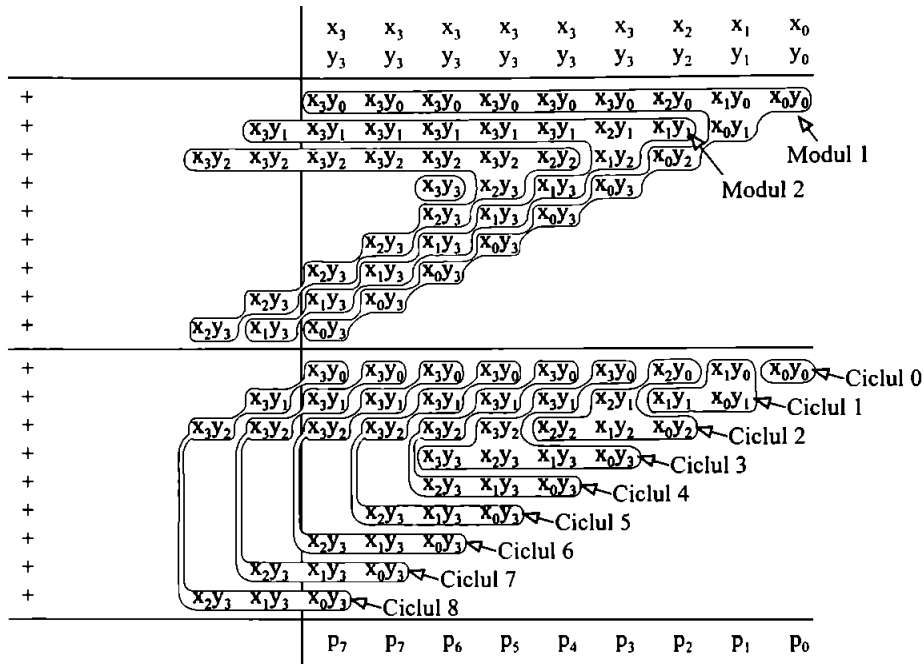


Fig. 5.10. Algoritm de înmulțire serie - serie pentru numere reprezentate în cod complementar

Termenii prezenți în figura 5.9 și absenți în figura 5.10 corespund termenilor generați de ultima sumă din relația (5.27).

Un *bit - slice* aferent înmulțitorului bazat pe algoritmul (5.30) este redat în figura 5.11.

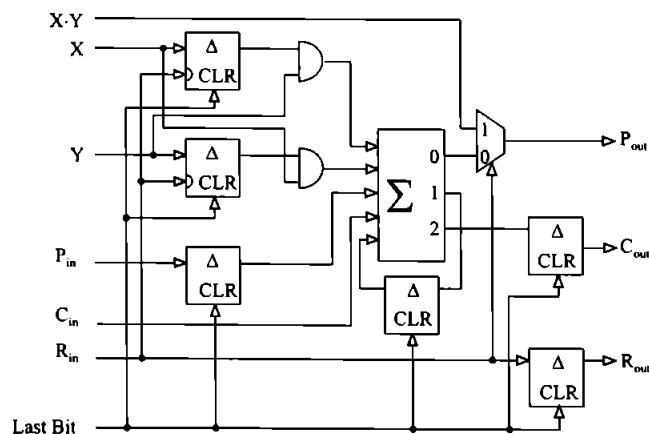


Fig. 5.11. Un *bit - slice* aferent înmulțitorului serie - serie

În acord cu structura sa, acest înmulțitor revendică operanzi de aceeași lungime. Dacă operanzii sunt de lungimi diferite, atunci operandul mai scurt va suporta o extensie de semn până va egala

Considerații privind perfecționarea RS dedicate algoritmilor conexioniști

lungimea celuilalt operand. $N-1$ astfel de *bit - slice* -uri vor fi necesare pentru înmulțirea a 2 numere reprezentate pe N biți. Cele $N-1$ *bit - slice* -uri trebuie interconectate ca în figura 5.12.

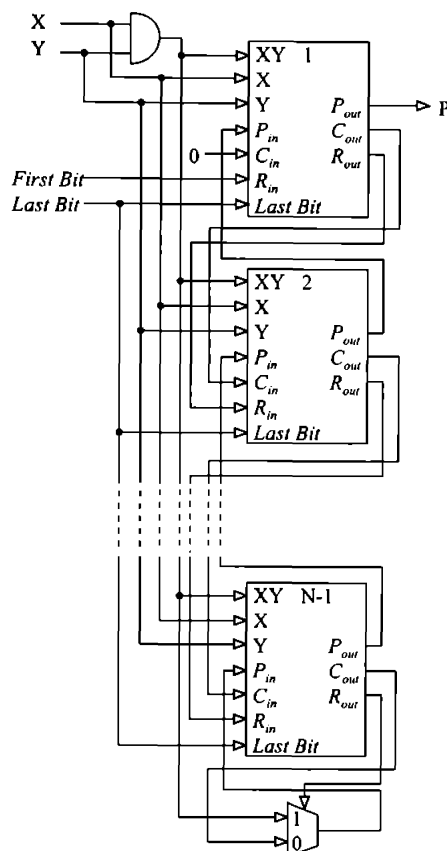


Fig. 5.12. Înmulțitor serie - serie pentru numere reprezentate în cod complementar

Prin multiplicarea porții AND din figura 5.12 și includerea ei în interiorul fiecărui modul și dacă se înlocuiește multiplexorul cu un al N -lea modul care are intrarea P_{in} conectată la zero, se va obține o modularitate perfectă (cu costul unui *hardware* redundant). Ambii operanzi vor intra serial în circuit, începând cu cel mai puțin semnificativ bit; preluarea operanzilor se va realiza în N cicluri de tact. Așa cum am menționat deja, dacă se operează cu numere fără semn, acestea trebuie extinse cu zerouri pe durata ultimelor $K - N$ cicluri (cu $K \geq 2N$ pentru a asigura un rezultat corect); dacă se operează cu numere reprezentate în complement față de 2 acestea vor suporta extensia semnului. Comanda *FirstBit* va fi setată *high* pe durata primului ciclu după care va fi menținută la zero. Similar, comanda *LastBit* va fi activată pe durata ultimului ciclu din cadrul operației de înmulțire sau, mai precis, pe durata ciclului ce precede o nouă înmulțire. Ieșirea (produsul) va fi disponibilă sub formă serială începând cu primul ciclu. Termenii încadrați în seturile din jumătatea superioară a figurii 5.10 corespund termenilor procesați (însumați) de modulul specificat; termenii încadrați în seturile din jumătatea inferioară corespund celor care vor fi adăugați sumei parțiale în același ciclu de tact.

SOLUȚII DE IMPLEMENTARE A REȚELELOR NEURONALE PE ARHITECTURI SISTOLICE

În general doi noi termeni sunt adunați la suma parțială în fiecare *bit - slice*, pe fiecare ciclu de tact. Aceasta sugerează utilizarea unui element $(5, 3)$ counter care va genera un *far carry* (pondere $\times 4$) spre *bit - slice*-ul următor. Apare deci un element *delay* pe ieșirea C_{out} . Un alt element *delay* este prevăzut pentru memorarea semnalului *carry* local. Multiplexorul de pe ieșirea *bit - slice* - ului $(i + 1)$ este utilizat pentru adunarea termenului simetric $x_i \cdot y_i$ la bitul corespunzător din cadrul sumei parțiale. Elementele *delay* cu intrări de validare (*enable*) memorează progresiv operanzii (X_i, Y_i) și sunt utilizate pentru a genera cei 2 termeni noi care se vor aduna sumei parțiale $(x_i \cdot Y_{i-1}$ și $y_i \cdot X_{i-1})$. În sfârșit, elementele *delay* conectate pe traseul $P_{in} - P_{out}$ compun registrul care memorează și deplasează suma parțială \bar{Q}_i .

Registrul de deplasare obținut prin înlănțuirea semnalelor $R_{in} - R_{out}$ aferente diverselor *slice*-uri activează succesiv modulele înmulțitorului; aceasta permite adunarea termenului simetric $x_i \cdot y_i$ și validează memorarea biților aplicați pe intrarea *slice*-urilor. Înmulțitorul serie - serie prezentat are câteva avantaje majore (foarte importante pentru RS proiectată în această lucrare):

- grad înalt de modularitate (extrem de important pentru implementările VLSI);
- latență redusă între biții aplicați pe intrare și bitul corespondent obținut la ieșire. Latența este cauzată doar de un circuit combinațional care conține pe calea critică numai un $(5, 3)$ counter și câteva porți logice;
- permite extensia produsului obținut pe ieșire la orice lungime; acesta rămâne corect (extensie semn) atâta timp cât se menține extensia de semn la operanzii aplicați la intrare.

5.5. Concluzii

Simplitatea celulei sistolice, comunicațiile strict locale (limitate la nivelul vecinilor cei mai apropiați), numărul extrem de redus de conexiuni și regularitatea fac ca rețeaua sistolică prezentată să reprezinte o candidată cu mari șanse la integrarea în structuri VLSI. Plăcile acceleratoare integrate într-un calculator de tip universal (*workstation*) constituie o soluție alternativă în munca de cercetare în domeniul conexiunilor. Comparativ cu un supercomputer de mare capacitate, calculatorul de tip *workstation* dotat cu un *hardware* accelerator reprezintă un bun compromis între cost și performanțe. RS prezentată reprezintă propunerea noastră pentru acest *hardware* accelerator.

6. O PROPUNERE ALTERNATIVĂ: REȚEAUA SISTOLICĂ LINIARĂ

RS bidimensională se bazează pe schema de principiu prezentată în figura 6.1.

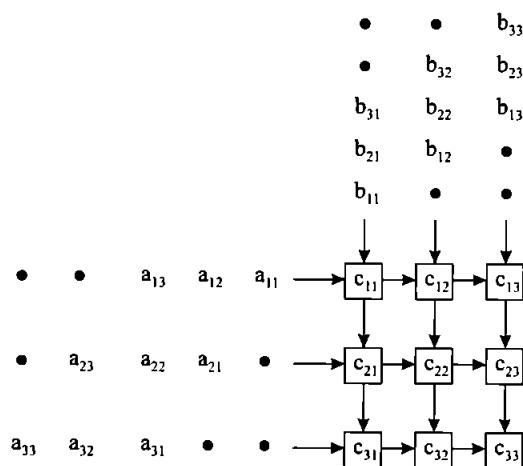


Fig. 6.1. Implementarea produsului matricial pe RS bidimensională

Pe fiecare impuls de tact, celulele sistolice din rețeaua rectangulară vor efectua operațiile prezentate în figura 6.2.

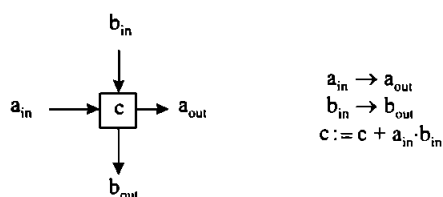


Fig. 6.2. Funcționarea celulei sistolice

Elementele matricii produs c_{ij} se vor acumula ca mărimi de stare în celulele rețelei. În cazul în care matricile A și B vor fi $N \times N$ va fi necesară o RS cu N^2 celule iar matricea rezultat $C(N \times N)$ se va obține în $3N - 2$ perioade tact. În cazul general ($A(N \times M)$, $B(M \times K)$) va fi necesară o RS dreptunghiulară de dimensiuni $N \times K$, iar matricea produs $C(N \times K)$ se va obține în $N + M + K - 2$ perioade tact.

Când se aplică în domenii care revendică procesarea unor cantități mari de date în timp real, cum ar fi de exemplu procesarea imaginilor, rețelele neuronale revendică implementări *hardware* speciale. În capitolele precedente am investigat implementarea RN pe RS bidimensionale rectangulare. În acest capitol vom investiga implementarea RN pe rețele sistolice liniare (RSL). Vom investiga posibilitatea implementării RN *feedforward* multistrat pe RSL, ca soluție alternativă la RS bidimensională.

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

După cum s-a arătat în capitolul introductiv, implementările sistolice se pot realiza în două variante:

- cu grad înalt și fin de paralelism (paralelism la nivelul sinapselor). În acest caz celula sistolică este alocată unei sinapse. RS bidimensională studiată în capitolele precedente face parte din această categorie;
- cu grad grosier de paralelism (paralelism la nivelul neuronilor). În acest caz celula sistolică va fi alocată unui neuron; va avea o structură mai complexă și va avea o memorie locală în care va stoca toți coeficienții sinaptici aferenți neuronului respectiv. Pentru a evidenția această diferență de complexitate, vom denumi celula sistolică aferentă RSL element de procesare (EP).

Au fost propuse diverse variante de RSL dedicate algoritmilor conexioniști; s-au efectuat simulări și s-au dezvoltat modele analitice pentru identificarea principiilor care trebuie să stea la baza implementărilor algoritmilor conexioniști pe RSL [NAY 95]. Pentru analiza ce urmează vom utiliza drept suport RSL HANNIBAL care implementează fazele *feedforward* (recunoaștere) și respectiv *backward* (învățare) aferente RN *feedforward* multistrat bazate pe algoritmul EBP. Proiectul HANNIBAL a fost dezvoltat de British Telecom în colaborare cu Loughborough University of Technology [MYE 91].

În figura 6.3 se prezintă o RSL care execută un produs matricial.

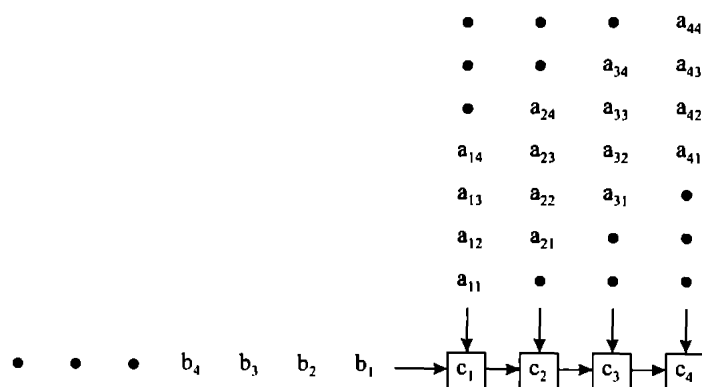


Fig. 6.3. Implementarea produsului matricial pe RSL

Fiecare CS va calcula un element al matricii (vectorului) produs pe baza operațiilor prezentate în figura 6.4.

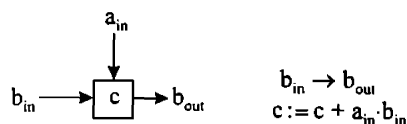


Fig. 6.4. Funcționarea celulei sistolice aferentă RSL

În cazul general în care matricea A este de dimensiuni $N \times M$, va fi necesară o RSL cu N celule iar matricea produs C de dimensiuni $N \times 1$ se va obține în $N + M + 1$ perioade de tact.

6.1. Arhitectura RSL dedicată RN feedforward multistrat

6.1.1. Corespondența CS - neuron

HANNIBAL este un circuit VLSI care poate fi cascadat în scopul construirii unei RSL dedicate algoritmului EBP. Semnalele de control și de tact sunt generate centralizat de către un *controller* RSL, iar celulele sistolice adiacente pot transfera date prin intermediul unui bus de 16 biți. Multe studii care vizează estimarea performanțelor RSL dedicate algoritmilor conexioniști iau ca referință arhitectura HANNIBAL pe care o simulează pe calculatoare universale.

Modul de implementare a unei RN *feedforward* multistrat pe o astfel de RSL este redat în figura 6.5.

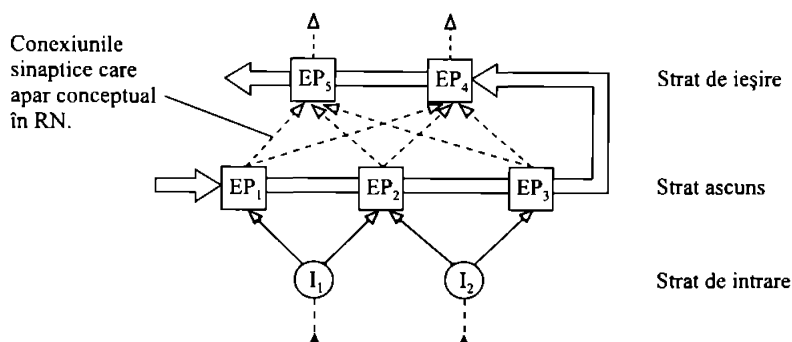


Fig. 6.5. O RN *feedforward* implementată pe RSL. RSL formează o cascadă *pipeline* în care datele de intrare se aplică direct primului element procesor din stratul ascuns

Figura 6.5 prezintă principalul mod de mapare în RSL a unei RN *feedforward* cu 3 straturi (2 - 3 - 2); 2 intrări, 3 neuroni în stratul ascuns și 2 în stratul de ieșire. Conexiunile sinaptice conceptuale indică o RN complet conectată; o rețea incomplet conectată se poate obține prin simpla setare a unor *flag-uri* asociate conexiunilor sinaptice inexistente. Datele de intrare sunt aplicate direct la începutul stratului ascuns. Neuronii din stratul ascuns și stratul de ieșire sunt mapați în elementele de procesare (EP) conținute în aria sistolică HANNIBAL. Un circuit HANNIBAL conține 4 EP fiecare dotat cu o memorie RAM locală care permite memorarea a 256 ponderi sinaptice reprezentate pe 16 biți. Fiecare EP conține o singură unitate de înmulțire - acumulare (MAC) pentru execuția calculelor matriciale - vectoriale, comune celor mai mulți algoritmi neuronali. Unitatea MAC lucrează cu operanzi în virgulă fixă pe 32 biți. O operație MAC (înmulțirea a doi factori și acumularea produsului la rezultatul obținut în operația MAC anterioară) se execută într-o singură perioadă de tact. Funcția de activare neliniară din cadrul fiecărui EP se obține prin aproximarea funcției sigmoideale. Busul care transferă date între două EP adiacente are 16 biți. El poate fi configurat în structură *feedforward* și respectiv în structură *feedback* (fig. 6.6).

Structura *feedforward* divizează cascada *pipeline* de elemente EP în două busuri de date de 8 biți. Busul superior (cascada *pipeline* superioară) este utilizat (pentru fiecare grup de procesoare EP aferente unui anumit strat neuronal) pentru transmiterea datelor de intrare, în timp ce busul inferior extrage rezultatele obținute de acest strat neuronal și le transmite primului EP din următorul strat. Granița dintre două straturi neuronale (între EP₃ și EP₄ în exemplul nostru) este configurată la inițializarea rețelei. Avantajul acestei structuri *pipeline* liniare constă în faptul că datele de ieșire aferente fiecărui strat neuronal pot fi transmise următorului strat de îndată ce devin disponibile. Un

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

nou vector de intrare poate fi transmis unui start neuronal fără a aștepta ca datele de ieșire să fie emise din cascada *pipeline* aferentă acestui strat.

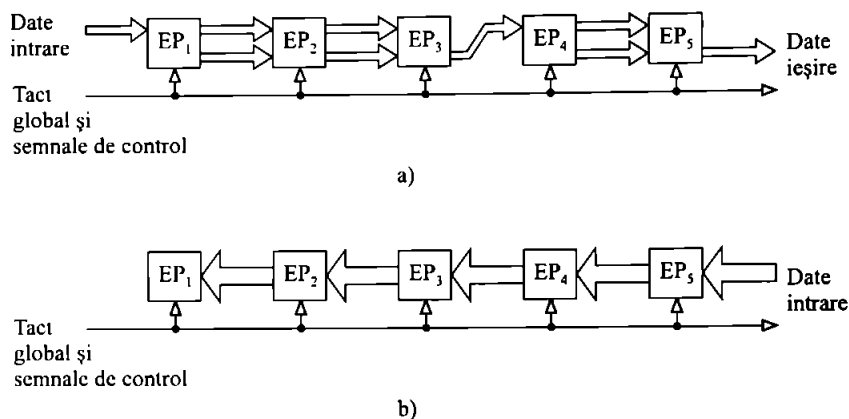


Fig. 6.6. Interconectarea *pipeline* a procesoarelor EP

- Structura *feedforward* dispune de 2 busuri de 8 biți pentru transmiterea simultană a datelor de intrare și respectiv ieșire
- Structura *feedback* utilizează un singur bus de 16 biți pentru propagarea înapoi a erorilor în timpul procesului de învățare

Structura *feedback* prezentată în fig. 6.6. b) este utilizată numai în timpul antrenării RN (învățării) și reconfigurează RSL sub forma unui bus *pipeline* de 16 biți inversând direcția datelor. Această structură este utilizată pentru propagarea erorii înapoi prin lanțul de EP aferente fiecărui strat neuronal spre EP-urile corespunzătoare din stratul neuronal anterior. Există studii care arată că reprezentarea erorilor pe 16 biți asigură acuratețea necesară procesului de ajustare a ponderilor sinaptice și este suficientă pentru a asigura convergența algoritmului de învățare [VIN 92].

6.1.2. Implementarea algoritmului EBP pe RSL

Implementarea algoritmului EBP pe RSL HANNIBAL se realizează cu doar câteva instrucțiuni sau moduri. *Controller*-ul fiecărui EP generează informația de inițializare cum ar fi poziția în cadrul stratului neuronal, numărul stratului neuronal etc., ceea ce permite execuția aproape autonomă a algoritmului. *Controller*-ul centralizat al RSL trebuie să controleze operațiile mai puțin frecvente și să sincronizeze aplicarea datelor de intrare și memorarea rezultatelor.

A. Faza *feedforward*

Faza *feedforward* utilizează *pipeline*-ul de date din fig. 6.6. a) pentru transferul datelor de intrare și a valorilor de activare aferente EP-urilor din RSL. Valoarea de activare, A_j^l , aferentă unui EP (neuron) din stratul ascuns sau de ieșire este dată de ecuația fazei *feedforward* caracteristică RN multistrat:

$$A_j^l = \sigma \left(\theta_k^l + \sum_{k=0}^{K-1} (w_{jk}^l \cdot A_k^{l-1}) \right) \quad (6.1)$$

unde:

- l - indexul stratului neuronal (l = 0 strat intrare, l = 1 strat ascuns etc.)
- j - indexul EP din stratul curent
- k - indexul EP din stratul anterior stratului curent
- K - numărul de EP-uri (neuroni) conținute în stratul anterior celui curent
- w - ponderea sinaptică pentru EP-ul indicat de indecși
- θ - valoarea de prag aferentă EP-ului indicat de indecși
- $\sigma(\cdot)$ - funcția sigmoidală

Structura *feedforward pipeline* duală (fig. 6.6. a)) permite transferul spre stratul de ieșire a valorilor de activare generate de către stratul ascuns în paralel cu aplicarea unui nou vector de intrare în stratul ascuns. În consecință, mai mulți vectori de intrare parțial procesați pot fi prezenți simultan în RSL, deși fiecare strat poate procesa doar elementele unui singur vector la un moment dat.

Controller-ul ariei RSL selectează modul ACTIVARE pentru a iniția secvența de stări RSL caracteristică fazei *feedforward*. Diagrama spațiu - timp pentru rețeaua 2 - 3 - 2 din fig. 6.5 care operează în modul ACTIVARE este redată în fig. 6.7.

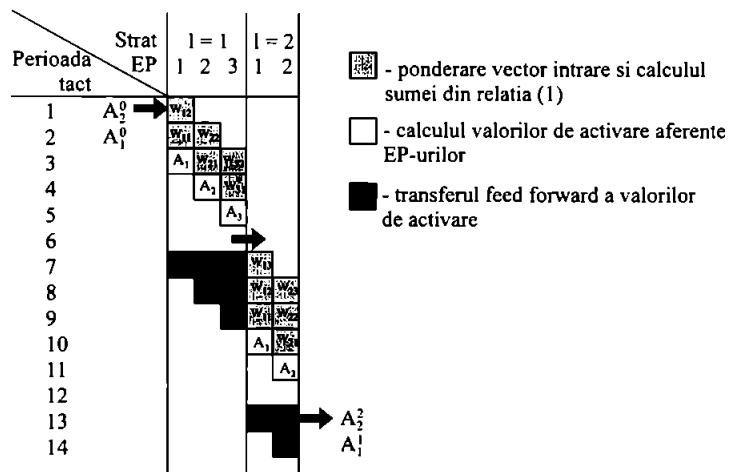


Fig. 6.7. Diagrama spațiu - timp corespunzătoare modului ACTIVARE aferentă RN 2 - 3 - 2. $A_1 - A_3$ sunt calculate în acord cu relația (6.1), înainte de a fi transferate stratului de ieșire (perioada de tact 7) pe busul *pipeline* inferior din fig. 6.6. a). Același proces se repetă și pentru stratul de ieșire și rezultatul (vectorul de ieșire) se generează în perioada de tact nr. 13.

Fig. 6.7 indică procesele care au loc în fiecare EP din RSL, în fiecare ciclu de tact, în cazul aplicării unui vector la intrarea RN (vectorul $A^0 = A_1^0, A_2^0$). Spațiile libere indică cicluri WAIT. EP₁ recepționează primul element al vectorului de intrare (A_2^0) în primul ciclu, execută prima operație MAC și transferă data de la intrare prin cascada *pipeline* spre EP₂. Aceste operații se repetă pentru toate EP-urile din stratul neuronal respectiv și pentru toate elementele vectorului de intrare. Așa cum indică fig. 6.6. a), doar busul de date superior este utilizat pentru transferul datelor de intrare prin cascada *pipeline* ce formează stratul neuronal respectiv. În momentul în care EP₁ a recepționat

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

toate elementele vectorului de intrare, va calcula valoarea de activare A_j^1 și apoi va intra în WAIT. Când toate EP-urile stratului respectiv au definitivat aceste operații, într-un ciclu suplimentar, toate valorile de activare sunt transferate busului inferior aferent cascadei *pipe*. Același proces repetă în stratul de ieșire care va genera vectorul de ieșire din RN, A^2 , asociat vectorului de intrare, A^0 .

Fiecare EP din RSL va repeta operațiile aferente fazei *feedforward* procesând noi vectori de intrare atâta timp cât *controller*-ul centralizat menține modul ACTIVARE activ. Frecvența cu care noii vectori de intrare pot fi aplicați la intrarea RSL este determinată de structura RN. Aceasta va determina în consecință caracteristicile de productivitate, latență și eficiență aferente RSL. Deci, structura RN este foarte importantă pentru performanțele RSL în faza *feedforward*.

B. Învățarea

Procesul de învățare sau antrenare a RN constă în ajustarea ponderilor sinaptice în ideea reducerii erori globale calculată la ieșirea RN. În cazul algoritmului EBP aceasta se realizează prin regula de învățare delta generalizată. Diferența dintre valoarea dorită la ieșire și respectiv valoarea reală obținută la ieșire generează eroarea aferentă stratului de ieșire care este propagată înapoi prin rețea spre straturile ascunse (stratul ascuns). Ponderile sinaptice sunt ajustate pentru a reduce gradientul erorii spre un minim global. Acest proces trebuie aplicat pentru fiecare pereche de vectori intrare - ieșire pe care RN trebuie să-o învețe. Există 3 stadii în procesul învățării:

1. propagarea înainte (*forward*) a valorilor de activare
2. propagarea înapoi (*backpropagation*) a erorilor
3. ajustarea ponderilor sinaptice și a valorilor de prag.

Primul stadiu reprezintă tocmai faza *feedforward* care se obține în cadrul RSL prin aplicarea modului ACTIVARE. Va conține aceeași secvență de operații prezentată deja în fig. 6.7, până la nivelul ciclului de tact nr. 11 moment în care acest prim stadiu se încheie. Valorile vectorului de ieșire nu sunt transferate la ieșirile RSL în timpul învățării (așa cum se întâmplă în cazul fazei *feedforward*), deoarece aceste valori sunt utilizate doar intern la calculul erorilor.

Următorul stadiu al procesului de învățare poate să înceapă imediat după ciclul nr. 11 din fig. 6.7. Acest stadiu implică calculul erorii dintre vectorul de ieșire așteptat și cel real. EP-urile stratului de ieșire calculează erorile locale δ_j^2 pe baza ecuației:

$$\delta_j^1 = A_j^1(1 - A_j^1)(D_j - A_j^1) \quad (6.2)$$

unde D_j reprezintă valoarea dorită la ieșirea neuronului j din stratul de ieșire, iar A_j^1 reprezintă valoarea reală obținută pe ieșirea aceluiași neuron. Ecuația (6.2) este valabilă în cazul funcției de activare unipolară (1.3).

EP-urile din stratul ascuns generează erorile locale pe baza ecuației:

$$\delta_j^i = A_j^i(1 - A_j^i) \sum_{l=0}^{I-1} (\delta_j^{l+1} \cdot w_{jl}^{l+1}) \quad (6.3)$$

unde: i - indexul EP-urilor din stratul următor celui curent
 I - numărul de EP-uri din stratul următor celui curent

Relația (6.3) reprezintă în fapt propagarea erorii înapoi într-o rețea *feedforward* de forma celei din fig. 6.8:

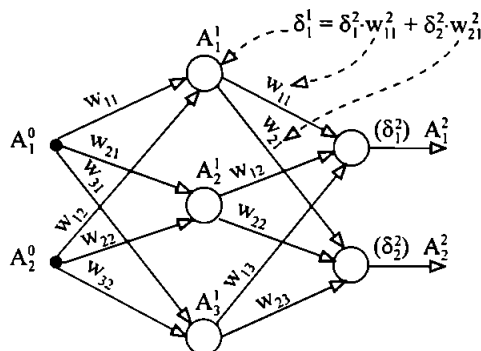


Fig. 6.8. Propagarea erorii înapoi în RN *feedforward* 2 - 3 - 2

Valorile erorilor de la nivelul stratului de ieșire (valori pe 16 biți) sunt propagate înapoi prin cascada RSL utilizând structura *pipeline* prezentată în fig. 6.6. b). Ecuațiile (6.2) și (6.3) sunt implementate în arhitectura HANNIBAL în modul CALCUL. Diagrama spațiu - timp aferentă acestui mod pentru o RN *feedforward* 2 - 3 - 2 este prezentată în fig. 6.9:

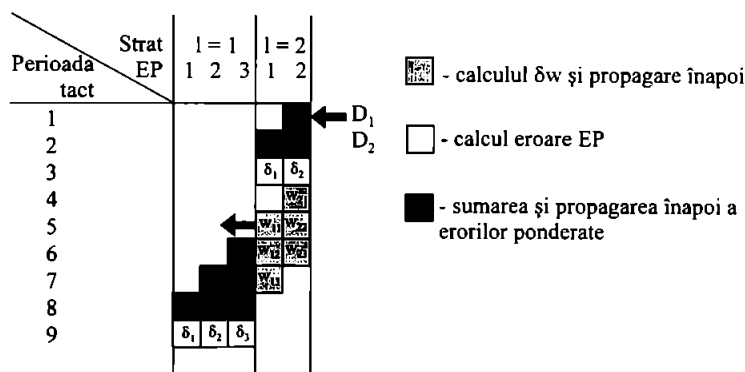


Fig. 6.9. Diagrama spațiu - timp aferentă modului CALCUL pentru RN 2 - 3 - 2. Vectorul de ieșire așteptat D (8 biți) este încărcat în RSL în sens invers (fig. 6.6. b)) și se calculează erorile aferente stratului de ieșire în ciclul 3 (relația (6.2)). Termenii obținuți prin ponderarea erorilor de ieșire (16 biți) sunt transferați înapoi spre stratul ascuns unde se calculează erorile aferente acestui strat (relația (6.3))

Vectorul dorit la ieșire, $D = (D_1, D_2)$ conține 2 componente reprezentate pe 8 biți. Acestea sunt încărcate succesiv în EP-urile corespunzătoare din stratul de ieșire. Vectorii de eroare, δ_1^2 și δ_2^2 , sunt apoi calculați pe baza relației (6.2). Derivata întâi a funcției de activare sigmoiale unipolare, $A_j^1(1 - A_j^1)$, este precalculată inițial de către fiecare EP în momentul în care se determină A_j^1 . Pentru a genera eroarea aferentă procesorului EP₁ din stratul ascuns, suma erorilor ponderate, $\delta_1^2 \cdot w_{11}^2 + \delta_2^2 \cdot w_{21}^2$, trebuie propagată înapoi prin EP-urile 2 și 3 aferente acestui strat

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

(fig. 6.8 și fig. 6.6. b)). După această propagare ecuația (6.3) se poate aplica pentru calculul erorii δ_j^l . Același proces se repetă întocmai pentru erorile EP-urilor 2 și 3 din stratul ascuns.

Exemplul dat prezintă o RN cu un singur strat ascuns. Pentru RN cu mai multe straturi ascunse ciclurile pentru generarea erorilor pe baza ecuației (6.3) se vor repeta pentru fiecare strat ascuns.

Stadiul final aferent procesului de învățare implică ajustarea ponderilor sinaptice și a valorilor de prag în acord cu eroarea memorată acum în fiecare EP. Actualizarea ponderilor sinaptice pentru fiecare EP se face pe baza relației:

$$w_{jk}^l \leftarrow w_{jk}^l + \alpha \cdot \delta_j^l \cdot A_k^{l-1} \quad (6.4)$$

Valorile de prag vor fi actualizate pe baza relației:

$$\theta_j^l \leftarrow \theta_j^l + \alpha \cdot \delta_j^l \quad (6.5)$$

unde α reprezintă constanta de învățare.

În arhitectura HANNIBAL, ajustarea ponderilor sinaptice și a valorilor de prag se realizează în modul ACTUALIZARE, utilizând structura *pipeline* din fig. 6.6. a). Fig. 6.10 reprezintă diagrama spațiu - timp aferentă acestui mod pentru o RN 2 - 3 - 2.

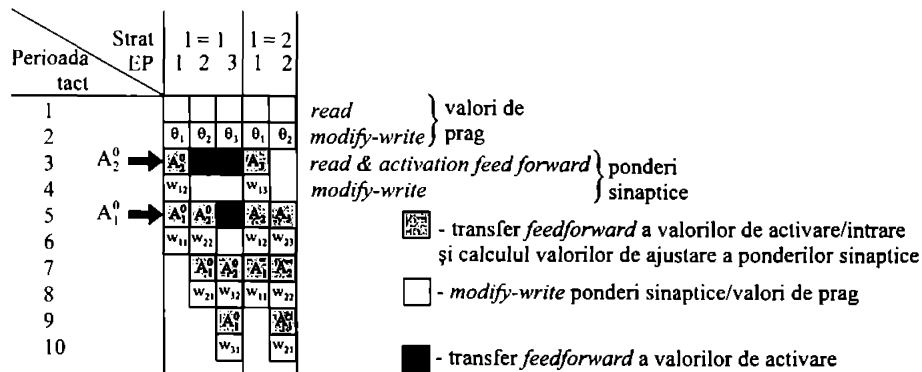


Fig. 6.10. Diagrama spațiu - timp aferentă modului ACTUALIZARE pentru RN 2 - 3 - 2. Fiecare strat execută în paralel două cicluri *read - modify - write* asupra ponderilor sinaptice. Datele de intrare pentru fiecare strat (vectorul de intrare în RN sau valorile de activare aferente stratului anterior) sunt transferate de-a lungul structurii *pipeline* din două în două cicluri. Ecuațiile (6.4) și (6.5) descriu procesul de actualizare a ponderilor sinaptice și respectiv a valorilor de prag

Procesul de actualizare este în fapt o operație *read - modify - write* cu durata a două cicluri care se aplică mai întâi valorilor de prag. Relația (6.4) arată că ajustarea ponderilor sinaptice respectă următoarea regulă: pentru EP_j, valoarea cu care se modifică ponderea sinaptică aferentă sinapsei ce conectează EP_j la EP_k din stratul anterior, este o funcție care depinde de valoarea de activare aferentă procesorului EP_k și de eroarea din EP_j. În consecință, vectorul de intrare trebuie aplicat din

nou pe intrarea RSL și valorile de activare aferente stratului ascuns transferate stratului de ieșire. Structura *pipeline* a RSL din fig. 6.6 a) permite execuția în paralel a acestor transferuri. Frecvența cu care se aplică datele pe intrarea EP din stratul ascuns și respectiv de ieșire trebuie înjumătățită datorită duratei extinse (2 cicluri) a operațiilor *read - modify - write*. În primul ciclu se citește valoarea ponderii sinaptice din memorie și se calculează valoarea cu care trebuie ajustată. În ciclul 2 ponderea sinaptică este ajustată și rescrisă în RAM.

Învățarea este, prin urmare, un proces secvențial care permite numai unei perechi de antrenament (vectori intrare - ieșire) să fie procesată la un moment dat. Consecința acestui fapt va fi că latența și productivitatea procesului de învățare vor fi mai puțin afectate de structura RN decât faza *feedforward*. Totuși, stratul care va necesita cele mai multe cicluri pentru actualizarea ponderilor sinaptice aferente, va determina durata globală a procesului de învățare. Prin urmare, optimizând structura RN se poate obține o creștere a performanțelor procesului de învățare bazat pe algoritmul EBP.

6.2. Analiza performanțelor

În cele ce urmează ne vom concentra asupra performanțelor pe care arhitectura HANNIBAL le obține în procesul emulării RN *feedforward* multistrat cu algoritmul de învățare EBP. Pentru caracterizarea performanțelor *hardware*-ului vom utiliza 3 metrici:

1. productivitatea (modele procesate per 1000 cicluri de tact)
2. latența (numărul de cicluri necesare pentru procesarea unui model)
3. eficiența ariei RSL

Prin model înțelegem un vector de intrare în cazul fazei *feedforward* și respectiv o pereche de antrenament (o pereche de vectori intrare - ieșire) în cazul procesului de învățare. Eficiența ariei reprezintă o măsură a gradului de utilizare a procesoarelor EP din cascada *pipeline* pe durata de timp necesară procesării unui model. Pentru a cuantifica eficiența, procesele care apar pe durata unui ciclu de tact în fiecare EP au fost evaluate și asignate uneia dintre următoarele 4 categorii:

- a) calcul (*calc*) - procesări în ALU, citire / scriere memorie, transfer între registre
- b) comunicație (*com*) - transferurile de date între diversele EP care vor afecta ieșirea rețelei
- c) calcul și comunicație (*calc & com*) - procesele a) și b) executate simultan. Un procentaj ridicat de astfel de operații caracterizează o arie sistolică eficientă.
- d) wait (*wait*) - un procesor (EP) în stare de așteptare sau o operație care nu are efect asupra ieșirilor rețelei.

Evaluările de performanță au fost făcute luând în considerare un singur parametru: dimensiunea stratului ascuns în cazul RN cu un singur strat ascuns (fig. 6.5). Dimensiunea stratului ascuns este considerată mult mai flexibilă comparativ cu cea a stratului de intrare și respectiv ieșire. Dimensiunea acestor straturi este adesea fixată și restricționată de formatul datelor (vectorilor de intrare - ieșire) utilizate în aplicația RN. A fost de asemenea arătat [KUN 88] că schimbarea dimensiunii stratului ascuns nu va conduce în mod obligatoriu la deteriorarea performanțelor algoritmice ale aplicației. În consecință, un număr de EP redundante pot fi adăugate în stratul ascuns dacă acest lucru se dovedește benefic pentru performanțele RSL. O dimensiune tipică de RN utilizată în recunoașterea imaginilor a fost utilizată pentru evaluarea performanțelor. Numărul de EP din stratul ascuns (*h*) a fost variat între 50 și 1000, iar dimensiunile stratului de intrare și respectiv ieșire au fost fixate (*I* = 512, *O* = 10).

O propunere alternativă: Rețeaua Sistolice Liniară

tipul, frecvența și comutarea operațiilor depind nu numai de h ci și de mărimea sa în raport cu I și O . Pentru RN cu $I > O$, apar trei cazuri în funcție de valoarea parametrului h :

- $0 < h \leq I - O$ (fig. 6.11. a)). Stratul ascuns procesează vectorii de intrare cu eficiență maximă. Prin urmare productivitatea fiecărui strat și productivitatea globală a RSL este determinată de acest strat. Stratul de ieșire poate executa toate operațiile MAC ce-i revin și poate descărca vectorul de ieșire aferent modelului de intrare $p-1$ înainte ca valorile de activare aferente modelului p (următor) să fie generate de stratul ascuns și evident transmise. În consecință, un număr de cicluri *wait* vor apărea în stratul de ieșire.
- $I - O < h \leq I$ (fig. 6.11. b)). Dimensiunea stratului ascuns determină suprapunerea descărcării vectorului de ieșire aferent modelului $p-1$ cu calculul vectorului de ieșire aferent modelului p . Pentru că cele două operații utilizează părți diferite din *pipeline*-ul RSL, eficiența procesoarelor EP din stratul ascuns nu suferă și productivitatea globală a RSL rămâne nemodificată.
- $h > I$ (fig. 6.11. c)). Va exista cel puțin un EP în stratul ascuns care va revendica mai multe cicluri de tact pentru transferul spre stratul de ieșire a valorilor de activare aferente modelului $p-1$ decât ar fi obligat să revendice pentru execuția operațiilor MAC aferente modelului p . Prin urmare, stratul de ieșire poate recepționa în mod continuu valorile de activare transmise de stratul ascuns în timp ce simultan descarcă vectorul de ieșire curent. Deci, stratul de ieșire operează cu eficiență maximă, în timp ce stratul ascuns va trebui să introducă periodic un număr de stări *wait*.

Cele două straturi ale RN funcționează în cadrul RSL ca două module *pipeline* cascade. În consecință, latența se obține simplu cumulând întârzierile de procesare cauzate de cele două straturi.

$$\begin{aligned} \text{Latența} &= \text{Întârziere strat ascuns} + \text{Întârziere strat ieșire} = \\ &= (I + h) + (h + O) = I + 2h + O \end{aligned} \quad (6.6)$$

Productivitatea este determinată de dimensiunea individuală a straturilor neuronale și nu de dimensiunea RN în ansamblu. Productivitatea este reprezentată de numărul de modele procesate în unitatea de timp (perioadă de tact, 1000 perioade de tact etc.). Figurile 6.11. a) și b) indică:

$$\text{Productivitatea} = \frac{1}{I} \quad \text{pentru } h \leq I \quad (6.7. a)$$

iar fig. 6.11. c) indică:

$$\text{Productivitatea} = \frac{1}{h} \quad \text{pentru } h > I \quad (6.7. b)$$

Relațiile (6.7) indică faptul că productivitatea este determinată de cel mai mare strat neuronal din componența RN, fiind invers proporțională cu dimensiunea acestuia. Dacă dimensiunea acestui strat este fixă (impusă de către aplicație) atunci un grad de flexibilitate va exista pentru celelalte straturi fără a afecta productivitatea. În analiza noastră de exemplu ($I = 512 = \text{fix}$), productivitatea maximă se va obține atâta timp cât este îndeplinită condiția $h \leq I$, indiferent de dimensiunea O . Pentru $h > I$, productivitatea descrește invers proporțional cu h .

Fig. 6.12 prezintă rezultatele obținute în baza acestor simulări pentru productivitatea RSL și latență.

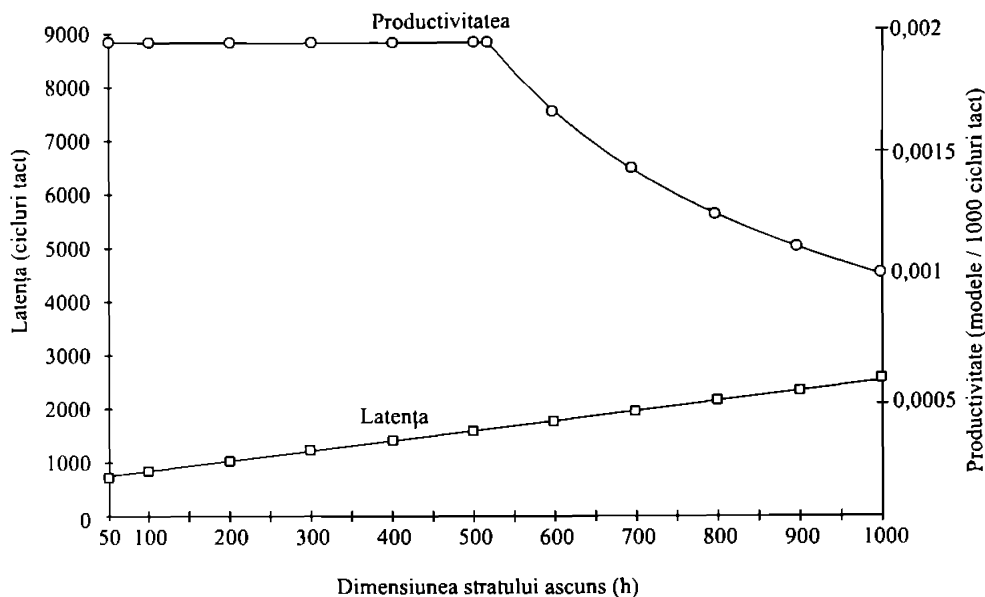


Fig. 6.12. Latența și productivitatea aferente fazei *feedforward* pentru o RN 512 - h - 10. Productivitatea este invers proporțională cu dimensiunea celui mai mare strat din structura RN. Latența este dată de suma numărului de intrări și ieșiri aferente tuturor straturilor din structura RN (straturile de intrare și respectiv ieșire participă la sumă cu numărul de EP ce le conțin iar straturile ascunse cu dublul numărului de EP ce apar în componența lor)

În regiunea $h < I$ productivitatea este constantă. Latența este independentă de productivitate pentru orice valoare a lui h .

Pentru a evalua eficiența RSL trebuie examinate operațiile executate de fiecare EP pe durata unui ciclu *feedforward*. Un astfel de ciclu începe cu aplicarea primei componente a unui vector de intrare și se încheie după obținerea ultimei componente din vectorul de ieșire asociat. Pe baza ariilor delimitate în fig. 6.11 corelate cu diagrama spațiu - timp din fig. 6.7 se pot deduce următoarele relații cantitative pentru cele 4 categorii de procese executate de către procesoarele EP pe durata unui ciclu *feedforward*:

$$calc = h + O \quad \text{pentru } \forall h \quad (6.8)$$

Cea mai mare pondere în cadrul unui ciclu *feedforward* o au procesele *calc & com*. Cantitativ aceste procese prezintă o dependență liniară în raport cu h și relația poate fi dedusă analitic pe baza ariilor delimitate în fig. 6.11:

$$calc \ \& \ com = h \cdot (I + O) \quad \text{pentru } \forall h \quad (6.9)$$

Conform figurii 6.11 ultimele două procese prezintă evaluări cantitative dependente de h . Stratul mic de ieșire (O) utilizat în această analiză a limitat influența regiunii $I - O < h \leq I$ în raport cu răspunsul global al RSL.

Ecuțiile aferente fiecărei regiuni sunt:

$$com = \frac{O^2}{2} \quad \text{pentru } 0 < h \leq I - O \quad (6.10)$$

$$com = \frac{(I-h)^2}{2} \quad \text{pentru } I - O < h \leq I \quad (6.11)$$

$$com = \frac{(h-I)^2}{2} \quad \text{pentru } h > I \quad (6.12)$$

$$wait = O(I-h) - \frac{O^2}{2} \quad \text{pentru } 0 < h \leq I - O \quad (6.13)$$

$$wait = O(I-h) - \frac{(I-h)^2}{2} = \frac{I-h}{2}(2O - I + h) \quad \text{pentru } I - O < h \leq I \quad (6.14)$$

$$wait = \frac{h^2}{2} - \frac{I^2}{2} \quad \text{pentru } h > I \quad (6.15)$$

Fig. 6.13 prezintă rezultatele obținute în baza acestei simulări.

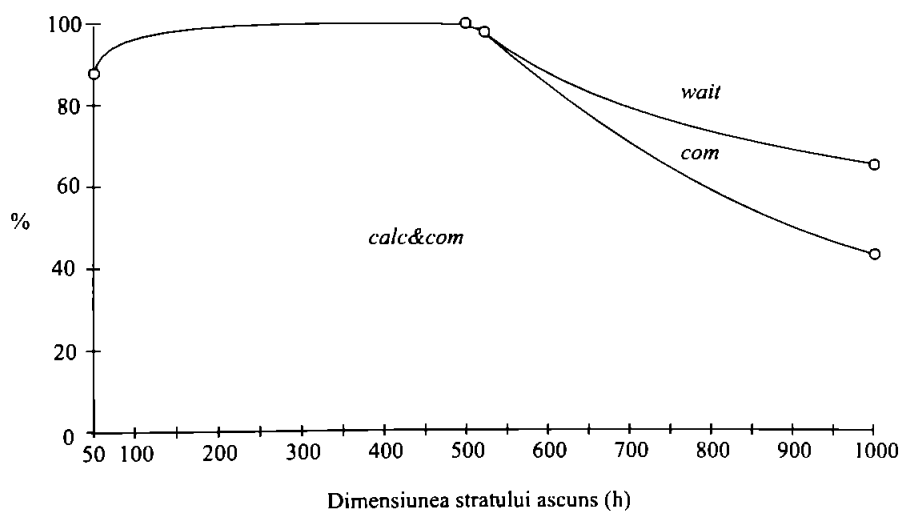


Fig. 6.13. Structura proceselor din RSL pentru RN 512 - h - 10 pe durata fazei *feedforward*. Eficiența RSL se apropie de 100 % când h = I = 512. Peste această valoare performanțele se degradează. Un alt punct de discontinuitate apare la valoarea h = I - O = 502.

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

Zona *com* este neglijabilă și constantă pentru $h \leq I - O$ și tinde la zero pe măsură ce h se apropie de valoarea I . Un strat de ieșire de dimensiune mai mare ar provoca creșterea influenței operațiilor *com* în aceste zone. Pentru $h > I$ o ecuație pătratică descrie creșterea ponderii proceselor *com*. Procesele *wait* sunt dependente de h în toate regiunile. Regresia liniară din zona $h \leq I - O$, (6.13), și dependența pătratică în raport cu h în zona $I - O < h \leq I$, (6.14), sugerează că RSL lucrează cu eficiență maximă când $h = I$. Aceasta se poate observa și din fig. 6.13. Pentru $h > I$, procesele *wait* cresc cu h^2 și atât procesele *wait* cât și procesele *com* sunt independente de O . În consecință, dimensiunea stratului de ieșire este flexibilă în această regiune și ar putea constitui un parametru ajustabil pentru creșterea eficienței RSL.

6.2.2. Învățarea

Diagrama spațiu - timp din fig. 6.14 reprezintă o formă simplificată a celor prezentate deja în fig. 6.7, 6.9 și 6.10 pentru modurile ACTIVARE, CALCUL și respectiv ACTUALIZARE.

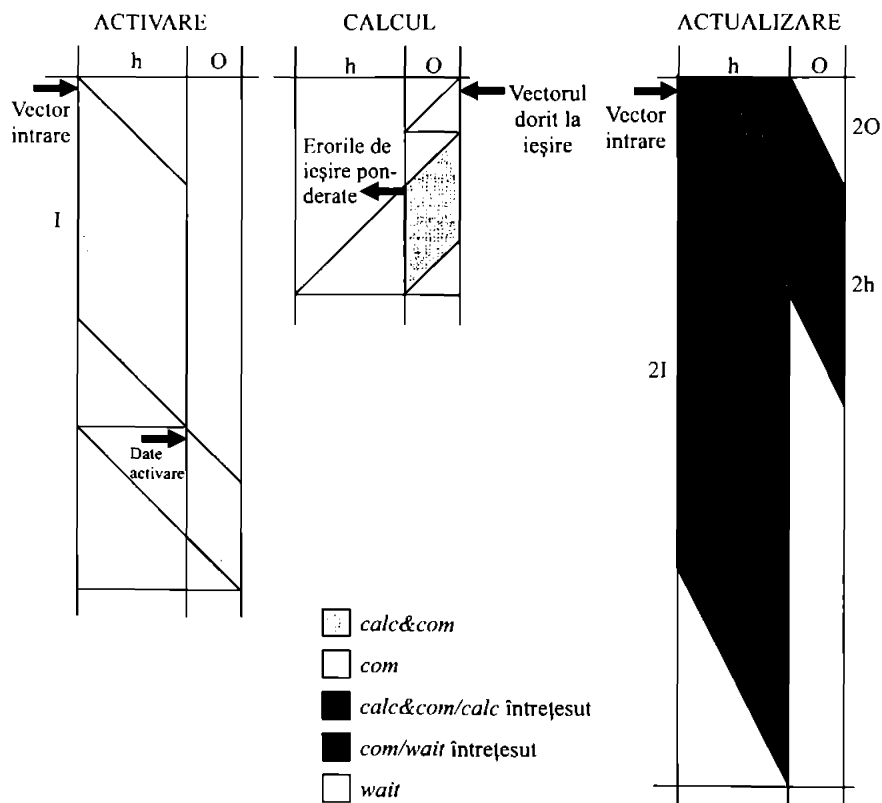


Fig. 6.14. Diagrama spațiu - timp simplificată aferentă celor 3 stadii (moduri) ce compun învățarea RN. Modul ACTIVARE nu mai conține zona *com* aferentă stratului de ieșire ca în cazul fazei *feedforward* (la învățare vectorul de ieșire nu mai trebuie comunicat în exterior; el va fi utilizat în modul calcul la evaluarea erorilor). Modul ACTUALIZARE conține operații întreșesute datorită ciclurilor *read - modify - write* care se execută pe două perioade de tact succesive.

O propunere alternativă: Rețeaua Sistolică Liniară

Se poate observa din fig. 6.14 că există un procentaj semnificativ de cicluri *wait* și respectiv *com* în cadrul fazei de învățare. Acestea vor afecta serios eficiența RSL și latența.

$$\text{Latența} = (I + 2h + O) + (2O + h) + (2 \max_{l=0}^2 (L_l + L_{l+1}) + 2h + 2O) \quad (6.16)$$

unde L_l reprezintă dimensiunea stratului l .

Relația (6.16) indică faptul că dimensiunile relative ale straturilor RN joacă un rol pentru latență în cazul învățării așa cum s-a întâmplat și în faza *feedforward*. Din (6.16) se poate deduce expresia latenței pentru RN cu un singur strat ascuns:

$$\text{Latența} = 3I + 7h + 5O \quad \text{pentru } O \leq I \quad (6.17)$$

$$\text{Latența} = I + 7h + 7O \quad \text{pentru } O > I \quad (6.18)$$

În consecință, dacă două RN conțin același număr de neuroni, cea care îndeplinește condiția $O \leq I$ va prezenta cea mai mică latență. Pentru rețeaua analizată de noi ($512 - h - 10$) latența este dată de relația (6.17). Pentru că perechile de antrenament sunt procesate secvențial putem afirma că în acest caz este valabilă relația care exprimă legătura dintre productivitate și latență în cazul proceselor secvențiale:

$$\text{Productivitate} = \frac{1}{\text{Latența}} \quad (6.19)$$

Fig. 6.15 prezintă rezultatele obținute pentru latență și productivitate în urma acestor simulări ale procesului de învățare.

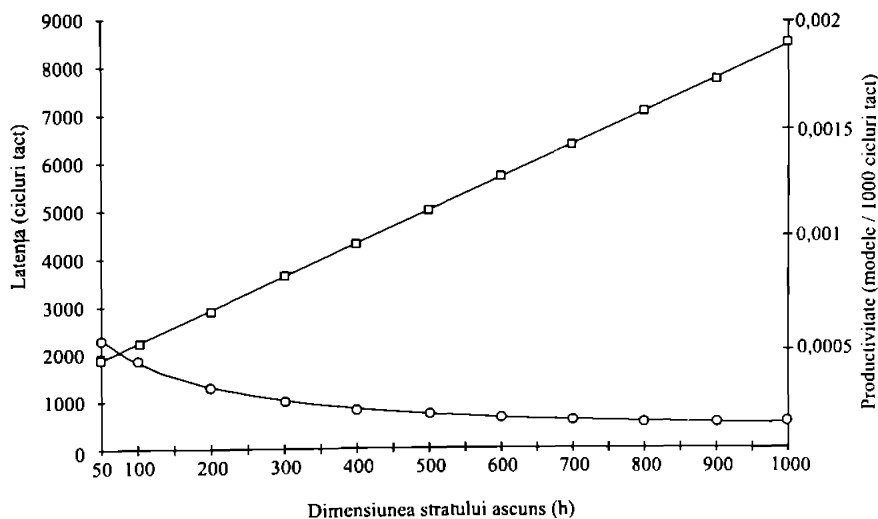


Fig. 6.15. Latența și productivitatea aferente fazei de învățare pentru o RN $512 - h - 10$. Modelele sunt procesate secvențial și ciclul de învățare aferent modelului curent trebuie încheiat înainte de startarea ciclului aferent modelului următor. În această situație productivitatea este invers proporțională cu latența.

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

Pentru a evalua eficiența RSL trebuie examinate operațiile executate de fiecare EP pe durata unui ciclu de învățare. Pe baza ariilor delimitate în fig. 6.14 și a diagramelor spațiu - timp din figurile 6.7, 6.9 și 6.10 se deduc următoarele relații:

$$\begin{aligned} calc &= \text{ACTIVARE}_{calc} + \text{CALCUL}_{calc} + \text{ACTUALIZARE}_{calc} = \\ &= (h + O) + (h + O) + (Ih + hO + 2h + 2O) = \\ &= h(4 + I) + O(4 + h) \end{aligned} \quad (6.20)$$

Procentul de procese *calc* care apar în modurile ACTIVARE și CALCUL este neglijabil și nici nu s-au mai evidențiat aceste procese în fig. 6.14, dar au fost incluse în relația (6.20) pentru o evaluare cât mai riguroasă.

$$calc \& \ com = (Ih + hO) + (hO) + (Ih + hO) = h(2I + 3O) \quad (6.21)$$

Procentul de procese *com* este mult mai mare în cazul învățării decât în cazul fazei *feedforward*. Aceasta se datorează propagării înainte a valorilor de activare urmată de propagarea înapoi a erorilor. Propagarea înapoi are un efect mult mai amplu.

$$com = \left(\frac{h^2}{2}\right) + \left(\frac{O^2}{2} + \frac{h^2}{2}\right) + \left(\frac{h^2}{2}\right) = \frac{3h^2}{2} + \frac{O^2}{2} \quad (6.22)$$

După cum indică ecuația (6.22) procentul de procese *com* este independent de I.

Ca și în cazul latenței, procentul de cicluri *wait* este determinat de dimensiunile relative ale straturilor adiacente. Stratul care va avea cele mai multe ponderi sinaptice de ajustat va determina cât de mult vor rămâne în *wait* toate celelalte straturi. De asemenea, stări *wait* adiționale sunt introduse în modul ACTUALIZARE datorită faptului că valorile de activare sunt transferate spre EP următor prin cascada *pipeline* RSL doar din două în două cicluri de tact. Aceste stări *wait* pot fi observate în ciclurile de tact 4 și 6 din fig. 6.10.

$$\begin{aligned} wait &= \left(\frac{3h^2}{2} + O^2 + O(I + 2h)\right) + \left(\frac{3O^2}{2} + 2Oh + \frac{h^2}{2}\right) + \left(\frac{3h^2}{2} + 2IO\right) = \\ &= \frac{7h^2}{2} + \frac{5O^2}{2} + O(3I + 4h) \quad \text{pentru } O \leq I \end{aligned} \quad (6.23)$$

$$\begin{aligned} wait &= \left(\frac{3h^2}{2} + O^2 + O(I + 2h)\right) + \left(\frac{3O^2}{2} + 2Oh + \frac{h^2}{2}\right) + \left(\frac{3h^2}{2} + 2O^2 + 2h(O - I)\right) = \\ &= \frac{7h^2}{2} + \frac{9O^2}{2} + 6Oh + I(O - 2h) \quad \text{pentru } O > I \end{aligned} \quad (6.24)$$

Pentru RN 512 - h - 10 pe care o analizăm este valabilă relația (6.23) care cuantifică cel mai mic procent de stări *wait*. Fig. 6.16 ilustrează eficiența RSL pe durata învățării. Răspunsul RSL este foarte diferit pe durata învățării în raport cu faza *feedforward* și nu numai prin faptul că cantitatea de procese *calc* devine mult mai semnificativă în faza de învățare. Deși toate cele patru categorii de procese cresc cantitativ cu h, fig. 6.16 arată că termenii de gradul 2 care apar în expresiile

proceselor *com* și *wait* au un efect devastator asupra eficienței RSL. În faza de învățare nu va fi niciodată posibil să se atingă nivelul înalt de eficiență obținut în faza *feedforward*.

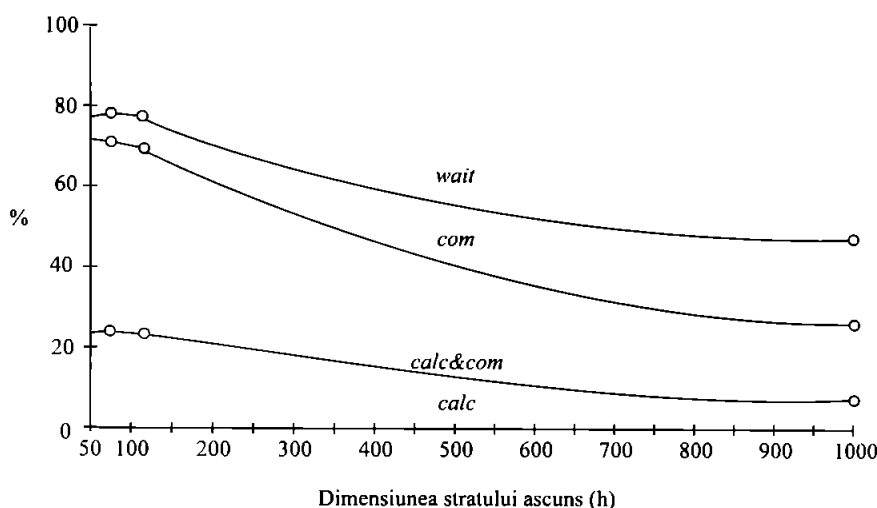


Fig. 6.16. Structura proceselor din RSL pentru RN 512 - h - 10 pe durata învățării. Eficiența RSL este mult mai mică la învățare decât în cadrul fazei *feedforward*; eficiența scade pe măsură de h crește, indiferent de structura rețelei

6.3. Considerații privind optimizarea implementării

Analiza pe care am efectuat-o poate fi utilizată pentru a determina opțiunile posibile atunci când se pune problema optimizării structurii RN mapate în RSL. Performanțele algoritmice ale aplicației trebuie însă menținute permanent în atenție, deoarece atunci când rețeaua este restructurată cu scopul de a optimiza performanțele *hardware*-ului, se poate ajunge la variante de RN incapabile să învețe pentru respectiva aplicație.

Adesea, cel mai important parametru utilizat în evaluarea performanțelor RN este productivitatea aplicației neuronale. În consecință, strategiile care optimizează înainte de toate acest parametru sunt larg aplicate. Așa cum am arătat prin analiza pe care am efectuat-o, eficiența RSL, latența și productivitatea nu sunt întotdeauna dependente între ele și ca atare o creștere a eficienței RSL nu este automat urmată de o creștere similară a productivității sau de o descreștere a latenței. Aceasta este în special adevărat în faza *feedforward*. Prin urmare, optimizarea nu trebuie să însemne întotdeauna simpla minimizare a dimensiunii globale a RN.

6.3.1. Faza *feedforward*

Analiza fazei *feedforward* a demonstrat că productivitatea este determinată de dimensiunea celui mai mare strat din structura RN. Dacă dimensiunea acestui strat nu se mărește pe parcursul restructurării, atunci productivitatea poate fi menținută constantă, în timp ce celelalte straturi sunt modificate pentru a optimiza latența și eficiența RSL. Pe baza analizei făcute, aceasta înseamnă a se opera în zona $h \leq I$. Obținem astfel:

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

Regula 1: Abilitatea de a optimiza productivitatea unei RN implementată într-o RSL crește dacă stratul cel mai puțin flexibil din punct de vedere dimensional (constrâns de către aplicație la o anumite dimensiune) este și cel mai mare.

Această regulă poate fi aplicată cel mai adesea stratului de intrare a cărui dimensiune este foarte puțin flexibilă (impusă de către aplicație).

Latența minimă se atinge cu cea mai mică RSL posibilă. Dacă latența este foarte importantă pentru aplicație, atunci reducând dimensiunile straturilor ascunse vom obține performanțe maxime în ceea ce privește performanțele *hardware*-ului.

Regula 2: Reducerea latenței se obține prin reducerea dimensiunii globale a RSL, și în particular prin reducerea numărului de EP din stratul ascuns (straturile ascunse). Structura optimă care generează latența minimă nu produce în mod necesar și productivitate maximă. Dacă productivitatea și / sau eficiența sunt considerate mai importante, atunci latența trebuie reevaluată după optimizarea rețelei din punct de vedere al productivității și / sau eficienței.

Eficiența RSL trebuie îmbunătățită fără degradarea productivității și fără a cauza o creștere semnificativă a latenței. Setând $h \leq I$ pentru a obține productivitatea optimă, structura RSL devine parțial determinată. Fig. 6.13 și ecuațiile care definesc eficiența indică faptul că un strat mic de ieșire asigură menținerea unei eficiențe mari în regiunea $h \leq I$. Dacă se operează în afara acestei regiuni, sau cu un strat de ieșire mare, se vor degrada productivitatea, latența și eficiența.

Regula 3: Menținerea unei productivități maxime și creșterea eficienței se obțin prin descreșterea monotonă a dimensiunilor straturilor RN, $I > h > O$. Dacă $h = I$ și $O < I$, se poate obține o eficiență de 100 % a RSL fără a degrada productivitatea.

În cazul $h = I$; $O < I$ trebuie acceptată totuși o concesie: creșterea latenței. Este un caz tipic de acceptare a creșterii latenței în favoarea creșterii eficienței RSL.

Ultima regulă se va referi la necesitatea unui strat de ieșire mic. Dacă cel mai mare strat este l, atunci eficiența maximă se va obține în procesoarele EP din stratul $l + 1$. Fig. 6.11 demonstrează această afirmație pentru $l = 0$ și $l = 1$. Cazul $l = 2$ (stratul de ieșire) nu a fost tratat în analiza pe care am efectuat-o datorită valorii fixe alocate parametrilor O și I . Ecuația care exprimă eficiența RSL este validă pentru $I > O$, dar efectul vizibil în fig. 6.11. c) unde stările *wait* se cantonează în stratul ascuns, ar fi vizibil și în stratul de ieșire dacă condiția $I > O$ ar fi inversată. În concluzie, rețelele în care stratul de ieșire este cel mai extins dimensional trebuie evitate.

Regula 4: Stratul de ieșire nu trebuie să fie cel mai mare din cadrul rețelei. Dacă această condiție nu poate fi satisfăcută, eficiența RSL nu poate atinge niciodată limita 100 %, indiferent de valorile dimensiunilor I și h .

6.3.2. Învățarea

Considerațiile privind optimizarea implementării pentru învățare revendică alte principii comparativ cu faza *feedforward*. În acest caz productivitatea este invers proporțională cu latența. Drept consecință, cel mai important criteriu de optimizare vizează reducerea latenței. Există de asemenea o mai strânsă relație între eficiența RSL și productivitate (latență) comparativ cu faza *feedforward*. Această dependență este reliefată și de rezultatele simulării prezentate în fig. 6.15 și 6.16, unde se poate observa că atât eficiența cât și productivitatea scad pe măsură ce latența crește. Totuși, simpla reducere a dimensiunii rețelei nu este suficientă, deoarece, pentru orice dimensiune

O propunere alternativă: Rețeaua Sistolică Liniară

globală particulară a rețelei va exista o structură optimă (o repartizare optimă a procesoarelor EP pe straturi).

Cea mai importantă condiție pentru învățare este asigurarea că $O \leq I$. Aceasta va garanta că latența se va evalua pe baza relației (6.17) și ca atare va fi minimizată.

Regula 5: Pentru orice rețea și independent de valoarea lui h , pentru asigurarea latenței minime (productivitate maximă) trebuie îndeplinită condiția $O \leq I$.

O altă observație importantă care rezultă din relațiile (6.17) și (6.18) este influența puternică pe care stratul ascuns o are asupra latenței. Această dependență se manifestă independent de relația dintre dimensiunile I și O , și este cauzată de strânsa interdependență dintre eficiență și latență. Eficiența scade odată cu creșterea valorii termenilor în h^2 din procesele *com* și *wait* (relațiile (6.22), (6.23) și (6.24)). Dimensiunea stratului de ieșire are de asemenea o puternică influență asupra volumului de operații *com* și *wait*. Prin urmare, asigurând condiția ca stratul de intrare să fie cel mai mare (dimensiunea I este impusă de regulă de aplicație) va rezulta o creștere a eficienței și o reducere a latenței.

Regula 6: Este indicat ca O și h să fie mai mici decât I . Dacă stratul de intrare este și cel mai mare, eficiența RSL va crește și latența va scădea datorită scăderii timpului alocat operațiilor de comunicație (între procesoarele EP) și de *wait*.

Aceste reguli pot fi utilizate pentru a adopta structura RN atunci când aceasta este implementată în RSL pentru a ne asigura că posibilitățile RSL sunt exploatate pentru performanțe maxime.

6.4. Optimizarea performanțelor aplicațiilor

O aplicație neuronală tipică dezvoltată și demonstrată de către mulți cercetători este recunoașterea caracterelor. RN utilizată preia imaginea unui caracter de regulă printr-o matrice de 8×8 pixeli. O RN *feedforward* tipică utilizată în clasificarea celor 26 caractere din setul Roman ar avea structura 64 - 20 - 26. Această structură și altele similare au fost considerate tipice pentru problema recunoașterii caracterelor utilizând RN *feedforward* cu algoritmul EBP sau alți algoritmi [VIN 92], [SAM 91], [REY 91]. Pentru implementarea pe RSL această structură de RN necesită o optimizare pentru obținerea unor performanțe deosebite deoarece câteva dintre regulile enunțate la 6.3. sunt violate de către structura inițială. Stratul de intrare este și cel mai extins (regula 6 îndeplinită), dar $O > h$ (regula 3 neîndeplinită). Satisfacerea regulii 3 reprezintă cea mai stringentă necesitate pentru orice RN implementată în RSL. Una dintre principalele probleme care afectează această structură este că atât productivitatea cât și eficiența vor fi reduse datorită lungimii deosebite a vectorului de intrare și a neîndeplinirii condiției $h = I$ (regula 3). De asemenea, deși stratul de ieșire nu este cel mai mare (regula 4 îndeplinită), el este totuși destul de mare pentru a cauza o creștere a procentului de operații *com* în cadrul fazei *feedforward* (regulile 3 și 4) și a celei de învățare (regula 6). Ambele dimensiuni, I și O , trebuie reduse pentru a îmbunătăți performanțele acestei RN mapată în RSL.

O metodă de reducere a dimensiunii vectorului de intrare ar fi compresia. Aceasta impune introducerea unui *hardware* (sau *software*) adițional ceea ce ar putea afecta performanțele aplicației. O altă soluție ar consta în utilizarea unei matrici de pixeli de rezoluție mai redusă; o matrice de 6×8 pixeli ar revendica un vector de intrare cu doar 48 componente. O metodă simplă de reducere a dimensiunii vectorului de ieșire constă în utilizarea unui vector de ieșire codificat. Cu numai cinci EP în stratul de ieșire și un decodificator extern care nu va cauza o creștere semnificativă a latenței vor putea fi decodificate cele 26 caractere. Simulările funcționale realizate pentru rețeaua originală

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

(64 - 20 - 26) și pentru rețeaua modificată (48 - 20 - 5) au arătat că nu sunt necesare mai multe procesoare EP în stratul ascuns al rețelei modificate pentru a produce caracteristici de eroare similare; rețeaua modificată realizează suplimentar o convergență mai rapidă.

O creștere semnificativă a performanțelor s-a constatat odată cu trecerea de la rețeaua originală la cea modificată. Productivitatea RSL a crescut cu 32 % în cadrul fazei *feedforward* și cu 42 % în cadrul fazei de învățare. Latența a scăzut cu 27 % și respectiv 29 % iar eficiența a crescut semnificativ (procentul de operații *calc & com* a crescut de la 58,4 % la 83,1 % în cadrul fazei *feedforward* și de la 23,9 % la 33,3 % în cazul fazei de învățare).

6.5. Concluzii

Regulile care au fost definite la 6.3. sunt relativ generale dar nu sunt întotdeauna aplicabile pentru toate RN *feedforward* multi-strat mapate în RSL. Există câțiva factori legați de condițiile de investigare fixate care pot influența eficiența acestor reguli.

În primul rând modelul analitic elaborat s-a limitat la RN cu un singur strat ascuns cu $I > O$. Această restricție nu afectează, totuși, atât de mult generalitatea analizei elaborate. Foarte importante sunt dimensiunile straturilor neuronale și pozițiile acestora în cadrul RN. Prin urmare, dacă apar mai multe straturi ascunse acestea pot fi luate în considerare și, de exemplu, regula 3 poate fi extinsă în cazul RN cu n straturi ascunse la:

$$I > h_1 > h_2 > \dots > h_n > O.$$

Cazul $I < O$ nu a fost acoperit și prin simulări, deși analiza teoretică a contribuit la obținerea regulilor elaborate. Concluzia generală care rezultă din aceste reguli privitor la stratul de ieșire este aceea că acest strat ar fi preferabil să fie cel mai mic din cadrul RN.

Un alt factor care trebuie luat în considerare este utilizarea arhitecturii HANNIBAL ca suport de investigare. În modul ACTIVARE se execută operații matriciale - vectoriale care sunt comune tuturor algoritmilor neuronali de tip *feedforward* implementați pe RSL. În consecință, regulile obținute pentru faza *feedforward* au un caracter de generalitate și sunt valabile pentru cele mai multe RN *feedforward* implementate pe RSL.

Faza de învățare aferentă algoritmului EBP are o mult mai mare specificitate pentru arhitectura HANNIBAL. Calculul erorilor se poate realiza mult mai eficient în arhitecturile de tip RSL decât în arhitecturile multiprocesor unibus. Degradarea performanțelor este cauzată totuși de inversarea sensului de transfer al informațiilor (propagarea erorilor înapoi). Ar fi necesar totuși un studiu mai larg care să investigheze mai multe tipuri de arhitecturi RSL dedicate RN pentru a asigura o certă generalitate regulilor stabilite pentru faza de învățare.

7. ALGORITMI NEURONALI PERFORMANȚI ADAPTAȚI IMPLEMENTĂRII PE REȚELE SISTOLICE

În capitolele 3, 4 și 5 am proiectat o RS bidimensională dedicată algoritmilor conexioniști. RS obținută este foarte performantă în execuția operațiilor specifice algoritmilor conexioniști și reprezintă o candidată extrem de bine cotată pentru implementările *hardware* (plăci acceleratoare, circuite VLSI pentru aplicații particulare etc.).

Pentru maximizarea performanțelor globale (la nivelul aplicației), trebuie acționat pe două direcții:

1. Implementarea RN dedicată aplicației respective pe o structură suport cât mai performantă (maximizarea performanțelor mașinii suport).
2. Maximizarea performanțelor algoritmului neuronal utilizat în aplicația respectivă.

Prima direcție a constituit obiectul capitolelor 3, 4 și 5 din această lucrare. A doua direcție vizează performanțele algoritmice ale aplicației și va constitui obiectul acestui capitol. Vom prezenta în cele ce urmează un exemplu de optimizare algoritmică; algoritmul vizat va fi algoritmul SOFM elaborat de *Kohonen*.

Algoritmul *Self - Organizing Feature Map* (SOFM) a fost pentru prima dată introdus de *Kohonen* [KOH 97] și de atunci a fost unul dintre cei mai populari algoritmi utilizați atât în studiile teoretice cât și în aplicații cum ar fi procesarea vorbirii. RN cu auto-organizare sunt motivate biologic prin abilitatea creierului de a mapa lumea externă pe cortex, unde stimuli externi apropiați sau asemănători sunt codați pe arii corticale adiacente. Algoritmul SOFM elaborat de *Kohonen* este un algoritm simplu care realizează o astfel de mapare. O secvență aleatoare de vectori de intrare este aplicată rețelei, iar ponderile sinaptice sunt astfel ajustate încât să reproducă cât mai fidel posibil distribuția de probabilitate de la intrare. Ponderile sinaptice se auto-organizează în sensul că neuroni vecini răspund la vectori de intrare asemănători și tind către valori asimptotice care cuantizează spațiul de intrare cât mai fidel posibil.

Sistemele *hardware* dedicate implementării RN încearcă să exploateze paralelismul implicit al algoritmilor neuronali. Implementările digitale cele mai indicate sunt arhitecturile sistolice (RS) formate dintr-un număr mare de elemente de procesare (EP) care operează concurrent. Maximizarea utilizării acestor resurse distribuite sugerează anumite modificări a algoritmilor neuronali clasici. Astfel de modificări (care vizează creșterea performanțelor RN de tip SOFM implementate în RS) vizează și algoritmul clasic elaborat de *Kohonen*. În cele ce urmează vom prezenta algoritmul *Kohonen* clasic precum și două versiuni alternative care devin mai eficiente în cazul implementării SOFM pe RS.

7.1. Algoritmul SOFM clasic și algoritmi SOFM modificați

Rețelele SOFM sunt formate fie din M neuroni dispuși într-o structură liniară (unidimensională) fie, mai general, din $M \times M$ neuroni dispuși într-o structură bidimensională de tip grilă. Vom discuta rețeaua sub forma bidimensională. Răspunsul obținut pe ieșirea unui neuron $i = (i_1, i_2)$, cu $1 \leq i_1, i_2 \leq M$, la un vector n -dimensional aplicat la intrare, $\mathbf{x} = (x_1, x_2, \dots, x_n)$, este dat de norma euclidiană $\|\mathbf{x} - \mathbf{w}_i\|_2$, unde $\mathbf{w}_i = (w_{i1}, w_{i2}, \dots, w_{in}) = (w_{(i_1, i_2)1}, \dots, w_{(i_1, i_2)n})$ este vectorul ponderilor sinaptice aferent neuronului $i = (i_1, i_2)$. Vom nota cu \mathbf{W} setul complet de ponderi sinaptice (vectorul ponderilor sinaptice):

$$\mathbf{W} = (w_{(1,1)}, w_{(2,1)}, \dots, w_{(M,1)}, w_{(1,2)}, \dots, w_{(M,M)})$$

7.1.1. Algoritmul SOFM clasic

Fie $\mathbf{x}(t)$ vectorul aplicat pe intrarea RN la momentele $t = 0, 1, 2, \dots$. Răspunsul $o_i(t)$ obținut pe ieșirea neuronului i va fi o măsură a distanței dintre vectorul său de ponderi sinaptice și vectorul $\mathbf{x}(t)$ aplicat la intrare.

$$o_i(t) = \|\mathbf{x}(t) - \mathbf{w}_i(t)\|_2 \quad (7.1)$$

Regula de selecție a neuronului învingător (*winner*) selectează neuronul $z(t) = (z_1(t), z_2(t))$ care prezintă răspuns minim la ieșire, ca învingător:

$$z(t) = z(\mathbf{x}(t), \mathbf{w}(t)) = \arg \min_i o_i(t) \quad (7.2)$$

Ponderile sinaptice sunt apoi ajustate în acord cu următoarea ecuație:

$$\mathbf{w}_i(t + 1) = \mathbf{w}_i(t) + \alpha(t) \cdot V(z(t) - i, t) \quad (7.3)$$

unde $\alpha(t)$ reprezintă factorul de adaptare sau constanta de învățare ($0 < \alpha(t) < 1$) și $V(z - i, t)$ este o funcție de vecinătate care ia valoarea 1 pentru $i = z$. O formă frecvent aleasă pentru funcția $V(z - i, t)$ este forma rectangulară:

$$V(z - i, t) = \begin{cases} 1 & \text{pentru } \|z - i\|_p \leq R(t) \\ 0 & \text{pentru } \|z - i\|_p > R(t) \end{cases} \quad (7.4)$$

unde $R(t)$ reprezintă mărimea vecinătății (în general dependentă de timp). Adesea norma $\|z - i\|_p$ este $\|z - i\|_\infty = \max\{|z_1 - i_1|, |z_2 - i_2|\}$.

Când un vector de intrare \mathbf{x} este aplicat rețelei, regula de ajustare (7.3) "apropie" vectorii ponderilor sinaptice aferenți neuronului învingător și celor din vecinătatea (7.4) de vectorul de intrare \mathbf{x} . Vectorii ponderilor sinaptice vor tinde, prin urmare, să reproducă cât mai fidel posibil distribuția de probabilitate de la intrare. Rețeaua SOFM realizează o corespondență între topologia spațiului vectorilor de intrare și topologia funcțiilor de ieșire din matricea SOFM, în sensul că neuroni vecini vor răspunde la stimuli (vectori de intrare) apropiați.

Ecuația (7.2) definește selecția *on-line* a neuronului câștigător, deoarece neuronul câștigător la momentul t , $z(t)$, este determinat comparând vectorul de intrare $\mathbf{x}(t)$ cu vectorul ponderilor sinaptice $\mathbf{W}(t)$. Relația (7.3) realizează actualizarea *on-line* a ponderilor sinaptice în sensul că, după aplicarea vectorului $\mathbf{x}(t)$ pe intrarea rețelei, vectorul ponderilor sinaptice este ajustat în acord cu (7.3). Ponderile sinaptice la momentul $t + 1$, $\mathbf{W}(t + 1)$, depind numai de intrarea $\mathbf{x}(t)$ și de ultima actualizare a vectorului $\mathbf{W}(t)$. Algoritmul definit de ecuațiile (7.2) și (7.3) este algoritmul SOFM clasic. Acesta actualizează ponderile sinaptice după fiecare aplicare a unui vector pe intrarea rețelei.

7.1.2. Algoritmul SOFM Batch

În algoritmul *batch*, ponderile sinaptice nu sunt actualizate după aplicarea fiecărui vector de intrare ci numai după aplicarea unui grup de T vectori. Acest nou parametru, T , reprezintă lungimea epocii. Timpul va fi exprimat în acest caz prin relația $t = kT + t'$ unde $k = 0, 1, 2, \dots$ reprezintă numărul epocii iar t' reprezintă numărul iterației din cadrul unei anumite epoci ($0 \leq t' \leq T - 1$).

Algoritmi neuronali performanți adaptați implementării pe rețele sistolice

Algoritmul *batch* se definește prin selecția unui învingător de grup, deoarece neuronul învingător $z(kT + t')$ va fi selectat prin compararea vectorilor de intrare $\mathbf{x}(kT + t')$ cu vectorul ponderilor sinaptice $\mathbf{W}(kT)$ existent la debutul epocii k .

$$z(kT + t') = z(\mathbf{x}(kT + t'), \mathbf{W}(kT)) = \arg \min_i \|\mathbf{x}(kT + t') - \mathbf{w}_i(kT)\|_2 \quad (7.5)$$

Actualizarea ponderilor sinaptice este o actualizare de grup deoarece ajustarea acestora se calculează pornind de la vectorul $\mathbf{W}(kT)$ existent la începutul epocii k .

$$\mathbf{w}_i(kT + t' + 1) = \mathbf{w}_i(kT + t') + \alpha(kT) \cdot V(z(kT + t') - i, kT) \cdot (\mathbf{x}(kT + t') - \mathbf{w}_i(kT)) \quad (7.6)$$

Dacă facem $T = 1$ se obține algoritmul SOFM clasic.

7.1.3. Algoritmul Batch - M

Algoritmul *batch - M* reprezintă o combinație intermediară între algoritmul SOFM clasic și algoritmul *batch*. Algoritmul *batch - M* va defini un neuron învingător de grup, așa cum o face și algoritmul *batch* prin relația (7.5), dar va realiza o actualizare *on-line* a ponderilor sinaptice:

$$\mathbf{w}_i(kT + t' + 1) = \mathbf{w}_i(kT + t') + \alpha(kT) \cdot V(z(kT + t') - i, kT) \cdot (\mathbf{x}(kT + t') - \mathbf{w}_i(kT + t')) \quad (7.7)$$

Dacă facem $T = 1$ se obține din nou algoritmul SOFM clasic deoarece lungimea epocii se reduce la o unitate.

7.1.4. Implementarea algoritmilor clasic, batch și batch - M

În funcție de strategia aplicată pentru detecția neuronului învingător și respectiv pentru ajustarea ponderilor sinaptice, caracteristicile celor 3 algoritmi sunt sintetizate în tabelul 7.1.

Algoritm	Selecție învingător	Actualizare ponderi sinaptice
SOFM clasic	<i>on-line</i>	<i>on-line</i>
SOFM <i>batch</i>	<i>batch</i>	<i>batch</i>
SOFM <i>batch - M</i>	<i>batch</i>	<i>on-line</i>

Tabelul 7.1. Caracteristicile fundamentale ale algoritmilor SOFM clasic, *batch* și *batch - M*

Vom justifica în continuare atenția acordată versiunilor *batch* și *batch - M* și vom arăta că acestea prezintă avantaje importante în raport cu versiunea clasică atunci când se pune problema implementării algoritmului SOFM în *hardware*.

Bucula principală a algoritmului SOFM clasic conține următorii pași:

1. Extrage aleator un vector din setul de antrenament
2. Calculează distanțele dintre vectorul de intrare pe de o parte și toți neuronii ce compun rețeaua SOFM pe de altă parte (ecuația (7.1))
3. Caută neuronul învingător (ecuația (7.2))
4. Actualizează ponderile sinaptice aferente neuronilor din vecinătatea învingătorului (ecuația (7.3))

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

Este evident că pasul 2 este cel mai mare consumator de timp și deci primul candidat la paralelizare. Simplist, neuronii pot fi partiționați în grupe și distanțele aferente diverselor grupe de neuroni pot fi calculate concurrent pe elemente de procesare diferite. Apoi se determină neuronul câștigător și, în final, unitățile de procesare responsabile cu neuronii din vecinătatea învingătorului vor actualiza ponderile sinaptice tot concurrent. Acest prim pas de paralelizare este posibil fără a aduce modificări algoritmului clasic. Creșterea de performanță se va limita în momentul în care se va atinge pragul de un unic neuron repartizat pe o unitate de procesare. Arhitecturile obișnuite, care dispun de un număr relativ mic de procesoare, nu pot depăși această limită de paralelizare.

Anumite structuri *hardware* speciale urmăresc creșterea productivității prin utilizarea unui număr foarte mare de elemente de procesare (EP) și deci printr-o paralelizare și mai accentuată a proceselor. Tipic, această creștere de productivitate se poate obține aplicând 2 strategii:

1. Procesarea vectorilor de intrare în tehnică *pipeline*. Într-o arhitectură cu număr mare de EP, mai multe unități pot partaja calculul unei singure distanțe. Această partajare se poate implementa în tehnică *pipeline*; fiecare EP calculează o distanță parțială (pe o anumită coordonată) și o va aduna la suma parțială transmisă din unitate în unitate. Pentru ca o mașină *pipeline* să fie utilizată eficient, procesoarele care și-au încheiat activitatea pe un set de date (pe un vector de intrare) trebuie să-și starteze activitatea pe următorul set; toate procesoarele EP vor procesa concurrent pe seturi de date diferite. Principala caracteristică a structurii *pipeline* este aceea că are un ciclu (timpul minim între 2 vectori aplicați succesiv pe intrări) mai mic decât latența (timpul minim dintre aplicarea vectorului de intrare și obținerea vectorului rezultat corespunzător). Pentru a menține structura *pipeline* plină, mai mulți vectori de intrare vor trebui să intre în pasul 2 (calculul distanțelor) înainte ca primul rezultat să devină disponibil pe ieșire și deci pasul 4 (actualizarea ponderilor sinaptice) să poată fi inițiat. Acest tip de procesare *batch* este utilizat în sistemul SYNAPSE - 1 la calculul produsului scalar în cadrul rețelelor de percepțiuni multistrat [RAM 92]. Caracterul strict *on-line* al algoritmului SOFM clasic împiedică proiectanții să aplice această tehnică în cazul rețelelor SOFM, irosind astfel o parte din potențialul de paralelism al mașinii disponibile. Să notăm că RS proiectată în capitolele 3 și 4 și perfecționată în capitolul 5 asigură procesarea în tehnică *pipeline* a vectorilor de intrare. Fiecărui vector de intrare i se asociază un cod de instrucțiune care permite schimbarea progresivă a modului de operare în interiorul RS (paragraful 5.3.3). Procesarea în tehnică *pipeline* a vectorilor de intrare înseamnă de fapt implementarea pe RS a unui algoritm SOFM de tip *batch / batch - M*.
2. Procesarea vectorilor de intrare în paralel. O altă posibilitate de a exploata eficient numărul mare de procesoare EP constă în multiplicarea RN pe *hardware*-ul existent; toate copiile RN vor partaja aceleași ponderi sinaptice și deci vor procesa simultan vectori de intrare diferiți. Această variantă implică de asemenea o procesare de grup (*batch*), deoarece aceeași matrice este utilizată pentru a calcula efectele produse de toți vectorii de intrare procesați concurrent iar actualizările ponderilor sinaptice vor fi efectuate la sfârșit, cu contribuția tuturor copiilor RN implementate în matricea EP. Această tehnică este utilizată de exemplu tot în cazul MLP pentru exploatarea arhitecturii de tip dublu inel a sistemului SYNAPSE - 1. Desigur, spre deosebire de costurile introduse de structura *pipeline* din cazul 1., această metodă revendică un pas adițional pentru actualizarea ponderilor sinaptice partajate (cu contribuția tuturor copiilor RN multiplicat în aria EP). În anumite cazuri ar putea fi ineficientă partajarea RN pe rețeaua EP și ar putea fi mai benefică partajarea datelor având în fiecare procesor EP o copie completă a RN [MAN 90].

După implementarea RN în una din cele 2 variante descrise, se poate afirma că paralelismul RN a fost complet exploatat; paralelizarea în continuare ar putea viza numai nivelul intim al operațiilor aritmetico - logice (funcții de activare sau unități aritmetice *pipeline*-izate).

Implementarea SOFM pe sisteme de mare productivitate nu reprezintă singura motivație pentru algoritmi *batch* și *batch - M*. În majoritatea sistemelor există anumite costuri fixe (în ceea ce privește timpul de procesare) asociate cu anumite operații care sunt independente de numărul de vectori procesați. Cel mai la îndemână exemplu ar fi costul încărcării matricii de ponderi sinaptice în sistem. Această operație poate fi foarte costisitoare în anumite sisteme atunci când se implementează RN a căror matrice sinaptică depășește capacitatea memoriei locale la nivelul fiecărui EP și când matricea este descompusă în submatrici. În loc să se reîncarce fiecare submatrice pentru fiecare vector de intrare, acest cost poate fi distribuit pe un grup de vectori de intrare. Costuri similare apar atunci când se transferă anumite tabele în sistem (de exemplu tabela care implementează funcția de activare) sau la golirea / încărcarea structurii *pipeline* atunci când se trece de la pasul curent la pasul următor în cadrul algoritmului.

Atât algoritmul *batch* cât și algoritmul *batch - M* acceptă procesarea de grup (grupuri de vectori de intrare) și deci vor fi foarte indicați pentru implementare pe arhitecturile care vizează un grad înalt de paralelism. Versiunea *batch*, în comparație cu versiunea *batch - M*, diferă doar prin faptul că permite *pipeline*-izarea relației (7.6). Dacă această *pipeline*-izare nu este exploatată, rezultatul $w(kT + t' + 1)$ va deveni disponibil după procesarea fiecărui vector de intrare și ajustarea ponderilor sinaptice poate fi realizată ca în (7.7). Diferențele între versiunile *batch* și *batch - M* sunt minore în comparație cu distanța de la algoritmul clasic la cei doi algoritmi modificați. Totuși, diferența poate deveni mai evidentă dacă se implementează strategia 2.

Studiile teoretice realizate privind auto - organizarea ponderilor sinaptice și convergența algoritmilor *batch* și *batch - M*, au arătat că aceste 2 versiuni respectă proprietățile fundamentale ale algoritmului clasic [IEN 97]. Primele abordări au avut ca obiect algoritmul clasic. Studiile privind auto - organizarea ponderilor sinaptice (ordonarea ponderilor) consideră algoritmul clasic descris de ecuațiile (7.2) și (7.3) un lanț Markov cu număr infinit de stări (configurațiile ponderilor sinaptice $\mathbf{W} = (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_M)$). Nu este dificil de arătat că cele 2 configurații strict monotone ale ponderilor sinaptice, F^+ (configurația strict crescătoare) și F^- (configurația strict descrescătoare), reprezintă 2 clase de absorbție ale acestui lanț Markov (o clasă de absorbție a unui lanț Markov este reprezentată de un set de stări a căror probabilitate de a rămâne în aceeași clasă este 1). Pentru a demonstra că orice configurație de ponderi inițial dezordonată va deveni ordonată într-un timp finit cu o probabilitate egală cu 1, se utilizează o proprietate a lanțului Markov, care stipulează că, dacă există cel puțin un set de secvențe de vectori de intrare $x(t)$ care prezintă o probabilitate diferită de zero de a ordona orice configurație de ponderi sinaptice, atunci o secvență aleatoare de vectori de intrare va ordona orice configurație de ponderi într-un timp finit cu probabilitatea 1. Cu alte cuvinte, dacă $P_{\mathbf{w}(0)}$ reprezintă legea de probabilitate a lanțului Markov $\mathbf{W}(t)$, care pornește din starea inițială $\mathbf{W}(0)$ de coordonate distincte ($\mathbf{w}_1(0) \neq \mathbf{w}_2(0) \neq \dots \neq \mathbf{w}_M(0)$), atunci $P_{\mathbf{w}(0)}(t \in F^+ \cup F^- < +\infty) = 1$. O astfel de demonstrație a fost conturată de *Kohonen* [KOH 88], a fost stabilită în detaliu de *Cottrell* și *Fort* [COT 87] pentru o distribuție de probabilitate uniformă pe intrare ($P(x) = 1$ pe intervalul (0,1)), și de către *Bouton* și *Pages* [BOU 93] pentru o distribuție $P(x)$ neuniformă dar continuă pe intervalul (0,1). Alți autori [ERW 92], [FOR 93] au extins aceste rezultate la alte clase de funcții de vecinătate. Algoritmul *batch - M* respectă în totalitate proprietățile fundamentale de auto - organizare și de convergență ale algoritmului clasic. În condițiile în care factorul de adaptare $\alpha(t)$ are o valoare fixă pozitivă care satisface relația $0 < \alpha(kT) < \frac{1}{T}$ pentru $\forall k \geq 0$, atunci și algoritmul

batch respectă aceste proprietăți de auto - organizare și convergență. Totuși, studiile teoretice lasă și anumite probleme deschise și acestea se referă la viteza de convergență și la robustețea algoritmului în raport cu valorile alese pentru parametrii de învățare. În acest scop, au fost realizate simulări comparative pe probleme reale pentru a demonstra și completa rezultatele teoretice.

7.2. Simulări experimentale

7.2.1. Inițializare

Pentru a separa efectul procesării de grup (*batch*) de efectul celorlalte particularități algoritmice, acestea din urmă au fost menținute la minim. Toate valorile (vectori de intrare și de ieșire, ponderi sinaptice) au fost procesate în virgulă flotantă. Funcția de vecinătate utilizată a fost cea specificată prin relația (7.4). Distanțele în rețeaua SOFM au fost exprimate:

$$\|z - i\|_p = \|z - i\|_1 = |z_1 - i_1| + |z_2 - i_2|$$

Toate simulările au fost realizate cu $T = 50$ (lungimea epocii). Această valoare poate fi considerată, pe de o parte, suficient de mare pentru a scoate în evidență potențialele deficiențe ale algoritmilor *batch* și *batch - M* și, pe de altă parte realistă pentru cazul unor implementări *hardware*. RS de dimensiune $N \times N$ perfecționată în capitolul 5 revindică o epocă $T = 2N$ (egală cu latența *pipeline*-ului). $T = 50$ înseamnă o RS de dimensiune 25×25 ; este o dimensiune rezonabilă pentru tehnologiile actuale dar care poate fi extinsă fără mari eforturi în viitor.

A. Setul de antrenament

Datele de test utilizate provin dintr-o aplicație de recunoaștere a vorbirii (clasificator de cod de carte prin recunoașterea vorbirii). Rețeaua SOFM utilizată a fost o rețea bidimensională 10×10 iar vectorii de intrare utilizați au fost 12 - dimensional. Setul de antrenament a fost compus din 10000 de prototipuri (vectori de intrare distincți) pentru fiecare experiment. Fișierul complet de date a fost compus din 112274 vectori. Fiecare rezultat raportat reprezintă media a 10 experimente, prin urmare utilizând 100000 de vectori. Fișierul de test, aplicat în cadrul fiecărui pas al algoritmului de învățare (epocă), a conținut 1000 de vectori din cei 12274 neutilizați în procesul de învățare. Vectorii de intrare au fost selectați aleator (un vector o singură dată) și aceleași secvențe de antrenament și de test au fost aplicate în cazul fiecăruia dintre cei 3 algoritmi (aceleași condiții de antrenament). Fiecare dintre cele 10 experimente au startat cu ponderi sinaptice inițiale aleatoare dar, ca și în cazul vectorilor de intrare, ponderile inițiale au fost aceleași în cazul tuturor celor 3 algoritmi.

B. Criteriul de convergență

Criteriul de eroare utilizat pentru a evalua convergența algoritmilor a fost eroarea cuantificată prin relația:

$$e(t) = \frac{1}{E \cdot P_{\text{test}}} \cdot \sum_{e=0}^{E-1} \sum_{i=0}^{P_{\text{test}}-1} \left(\|x_i^e - w_{z(i)}(t)\|_2 \right)^2 \quad (7.8)$$

unde:

- z - indică indexul neuronului învingător
- i - parcurge toate prototipurile (vectorii) din setul de test P_{test} .
- e - realizează o mediere a erorilor înregistrate în cele E experimente.

Algoritmi neuronali performanți adaptați implementării pe rețele sistolice

Eroarea a fost calculată la finele fiecărei epoci ($t = kT$). Acest criteriu de eroare nu este singurul criteriu de evaluare a performanțelor rețelelor SOFM în aplicațiile de recunoaștere a vorbirii; alte criterii cum ar fi rata de recunoaștere pot fi luate în considerare. Totuși, criteriul de eroare specificat aici va fi cel mai indicat pentru evaluarea diferențelor dintre algoritmul SOFM clasic și cele 2 versiuni modificate. La compararea rezultatelor experimentale nu a fost utilizată valoarea finală a erorii ci valoarea medie a acesteia pe ultimele 10 epoci (valoarea e). Această strategie elimină hazardul pe care l-ar putea introduce anumite experimente particulare care înregistrează o foarte slabă convergență dar care, printr-un joc al hazardului, înregistrează o foarte mică eroare finală.

C. Parametrii de învățare

Rata de învățare $\alpha(t)$ din ecuațiile (7.3), (7.6) și (7.7) a fost variată pe parcursul simulărilor conform cu relația:

$$\alpha(t) = \frac{\alpha_0}{1 + K_\alpha \cdot t}$$

În toate experimentele mărimea vecinătății $R(t)$ din relația (7.4) a fost stabilită prin relația:

$$R(t) = 1 + \left\lfloor \frac{R_0}{1 + K_R \cdot t} \right\rfloor$$

În algoritmi *batch* și *batch - M* ambii parametri ($\alpha(t)$ și $R(t)$) s-au recalculat la începutul fiecărei epoci ($t = kT$) și apoi au fost menținuți constanți în interiorul epocii, așa cum indică relațiile (7.6) și (7.7).

D. Limitarea ponderilor sinaptice

În secțiunea (7.1.4) s-a arătat că, în cazul algoritmului *batch*, pentru a asigura ordonarea ponderilor sinaptice este necesară îndeplinirea condiției:

$$0 < \alpha(kT) < \frac{1}{T} \text{ pentru } \forall k \geq 0 \quad (7.9)$$

În caz contrar, nu numai ordonarea ponderilor nu poate fi garantată ci chiar și menținerea acestor ponderi în domeniul valorilor de intrare. Cu anumite secvențe particulare de vectori de intrare, procesul de învățare poate conduce la divergența anumitor ponderi, cu rezultate catastrofice în ceea ce privește convergența algoritmului *batch*.

Pentru a putea starta algoritmul cu parametrul $\alpha(t)$ plecând de la valori inițiale mari și fără a exista riscul divergenței unor ponderi sinaptice, ponderile pot fi limitate la anumite valori limită minime și respectiv maxime. Deci, la sfârșitul fiecărei epoci ponderile care depășesc limitele prescrise vor fi setate pe aceste limite. Limitarea ponderilor a fost aplicată în cazul algoritmului *batch* în toate experimentele în care valoarea inițială α_0 a fost mai mare decât $1/T$ ($\alpha_0 > 0,02$). Această tehnică nu a fost aplicată în cazul algoritmilor clasic și *batch - M* deoarece constrângerea cuprinsă în relația (7.9) nu se aplică.

7.2.2. Rezultate obținute

Tabelul 7.2 sintetizează valorile parametrilor utilizați în toate experimentele.

Parametru	Valoare
Dimensiunea rețelei SOFM	10 x 10
Număr intrări	12
Domeniul de intrare	[-2,01; 1,89]
Limitarea ponderilor	-2,5; +2,5
Lungimea epocii (T)	50
Număr epoci	200
Număr experimente considerate la mediere	10
Vectori de test	1000

Tabelul 7.2 Configurația parametrilor comună tuturor experimentelor

Mărimea vecinătății și evoluția acesteia a fost identică în cadrul celor 3 algoritmi; valoarea inițială a fost setată la $R_0 = 5$ în toate simulările. Această valoare a fost menținută constantă deoarece ea este determinată de regulă de mărimea matricii SOFM (menținută constantă în toate experimentele) și nu depinde de diferențele existente între diversele versiuni ale algoritmului de învățare. Efectul parametrului K_R a fost investigat utilizând valorile 0,2 ; 0,02 ; 0,002 ; 0,0002.

Pentru toți algoritmi rata inițială de învățare α_0 a fost setată pe 5 valori distribuite relativ uniform pe intervalul (0, 1). Aceste valori au fost: 0,1 ; 0,3 ; 0,5 ; 0,7 și 0,9. Rata de descreștere K_α a fost setată pe 6 valori: 0,25 ; 0,025 ; 0,0025 ; 0,00025 ; 0,000025 și 0,0000025.

Toate combinațiile acestor valori au fost aplicate celor 3 algoritmi; adițional pentru versiunea *batch* au fost aplicate câteva valori speciale. În acest caz α_0 trebuie să rămână în intervalul (0, 1/T), în acord cu condiția (7.9). Valorile utilizate au fost 0,002 ; 0,006 ; 0,010 ; 0,014 și 0,018. Valori foarte mici au fost utilizate și pentru coeficientul K_α pentru a facilita o bună convergență a algoritmului *batch*.

A. Evoluția erorii

Fig. 7.1 indică evoluția erorii pe durata primelor 60 epoci pentru cele mai bune 10 rulări ale fiecăruia dintre cei 3 algoritmi. Axa orizontală pentru versiunile *batch* și *batch - M* indică numărul epocii; pentru a obține aceeași scară pe abscisă la versiunea normală eroarea a fost marcată după fiecare set de T vectori aplicați la intrare. Erorile finale, în cazul celor mai bune rulări pentru cele 3 versiuni clasic, *batch* și *batch - M*, au fost aproape identice: 0,176227 ; 0,178781 și respectiv 0,174834.

B. Viteza de convergență

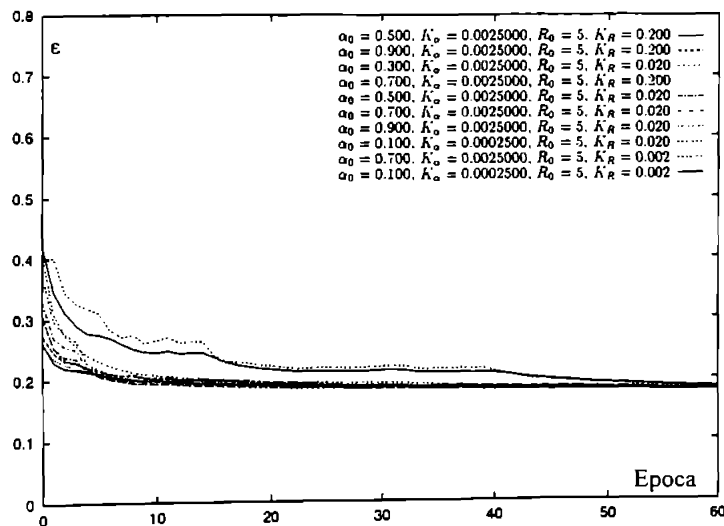
Scopul esențial al noilor versiuni introduse (*batch* și *batch - M*) este acela de a crește viteza procesului de învățare. Acest deziderat s-ar putea atinge printr-o implementare *hardware* cu un grad mai înalt de paralelism, grad pe care cele 2 versiuni *batch* îl pot permite. Pentru a menține acest avantaj este important ca cele 2 versiuni *batch* să nu fie semnificativ mai lent convergente decât versiunea clasică.

Pentru o comparație cât mai exactă, o rulare bună și cu aceleași valori pentru K_R și K_α trebuie selectată pentru fiecare versiune a algoritmului SOFM. Acest lucru este foarte important deoarece cei 2 parametri (primul în special) afectează profund viteza de organizare a rețelei SOFM.

Algoritmi neuronali performanți adaptați implementării pe rețele sistolice

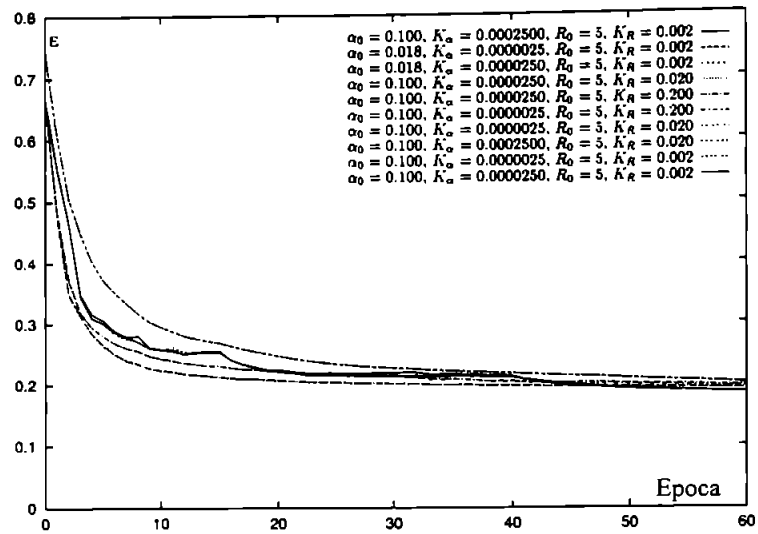
Amplerea acestui efect poate fi observată în figura 7.1. a) unde 2 rulări ale algoritmului original converg sensibil mai lent datorită valorilor mici ale coeficientului K_R . În afara acestui efect, putem observa că diferențierea între cele mai bune rulări este mică, ceea ce justifică compararea celor mai bune rulări (câte una pentru fiecare versiune).

Au fost aleși coeficienții celei mai bune rulări a algoritmului clasic ($K_\alpha = 0,0025$ și $K_R = 0,2$); utilizând aceste 2 valori au fost alese și cele mai bune rulări aferente versiunilor *batch*. Prin contrast cu valorile comune alese pentru K_α și K_R , pentru coeficientul α_0 a fost aleasă cea valoare (diferită de la caz la caz) care corespunde comportamentului optimal al fiecărei versiuni. Rata de descreștere, K_α , aleasă pentru algoritmul clasic s-a dovedit a fi optimă și pentru algoritmul *batch - M*. În cazul algoritmului *batch*, luând în considerare cel mai bun experiment nondivergent, convergența este rapidă la început dar valoarea excesivă a ratei de descreștere asociată cu valoarea mică a constantei inițiale α_0 , încetinește ulterior convergența (viteza de scădere a erorii) comparativ cu celelalte 2 versiuni. Acest lucru este vizibil în figura 7.2 unde sunt ilustrate primele 60 epoci din cadrul algoritmului de învățare pentru fiecare dintre cele 3 versiuni (cu $K_\alpha = 0,0025$ și $K_R = 0,2$); tot în fig. 7.2 se asociază și cea mai bună rulare pentru algoritmul *batch* ($K_\alpha = 0,000025$ și $K_R = 0,002$). Această cea mai bună rulare care este relativ mai lent convergentă în prima parte a intervalului, este mai apropiată ca performanțe de celelalte rulări și de aceea a fost utilizată în comparațiile care urmează. Algoritmul clasic este cel mai rapid convergent: după cea de-a șaptea epocă atinge deja pragul de $(1 + 20\%) \cdot e_f$. Am notat cu e_f - eroarea finală. Algoritmul *batch - M* este cu foarte puțin mai lent convergent decât cel clasic: după a 20-a epocă rezultatele acestuia sunt aproape identice cu cele ale algoritmului clasic. Versiunea *batch*, deși avantajată de acea alegere liberă a parametrilor de învățare, este net mai lent convergentă pe întreg domeniul și necesită 42 epoci pentru a atinge pragul de $(1 + 20\%) \cdot e_f$. Rezultatele sunt sintetizate în tabelul 7.3 care evidențiază la ce număr de epocă eroarea devine mai mică sau egală cu valorile de prag considerate. Valorile de prag sunt considerate ca procente peste cea mai mică eroare absolută obținută (versiunea *batch - M*: 0,174834).

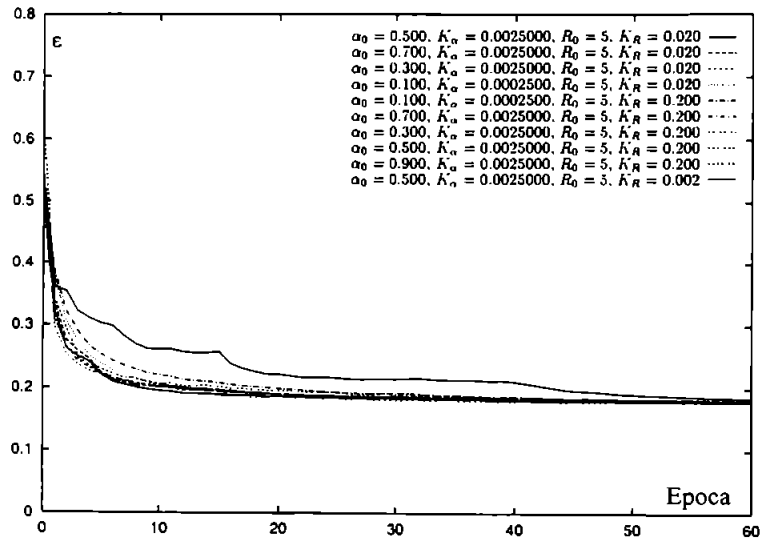


a) algoritmul clasic

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE



b) algoritmul *batch*



c) algoritmul *batch - M*

Fig. 7.1. Cele mai bune 10 rulări

Algoritmi neuronali performanți adaptați implementării pe rețele sistolice

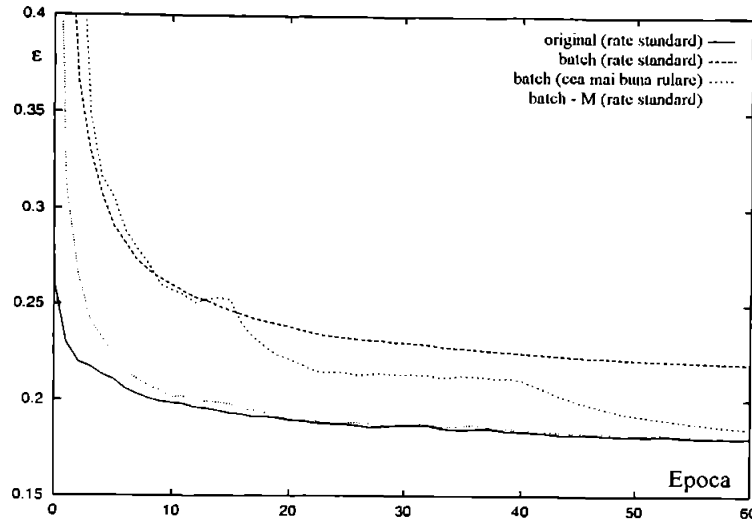


Fig. 7.2. Viteza de convergență comparativă a celor mai bune rulări cu rate standard ($K_{\alpha} = 0,0025$; $K_R = 0,2$)

Prag	Eroare	clasic	<i>batch</i>	<i>batch - M</i>
50 %	0,262	1	10	4
20 %	0,210	7	42	9
10 %	0,192	18	52	21
5 %	0,184	40	68	43
1 %	0,177	136	–	142

Tabelul 7.3. Viteza de convergență a celor 3 algoritmi

C. Robustețea

Rezultatele prezentate în paragraful anterior arată că, pe o problemă reală, viteza de convergență a algoritmului *batch - M* și, cu o oarecare aproximație, și a algoritmului *batch* sunt comparabile cu cea a algoritmului SOFM clasic. Datorită imposibilității practice de a rula algoritmi cu foarte multe combinații ale parametrilor de învățare, este important să se evalueze cât de critică este problema alegerii acestor parametri.

O primă posibilitate de evaluare constă în raportarea numărului de experimente (pentru fiecare versiune în parte) care conduc la o eroare finală cu peste 20 % mai mare decât cea corespunzătoare celei mai bune rulări, la numărul total de experimente efectuate. În cazul algoritmului original, 58 din cele 120 experimente au condus la o eroare cu 0,210 mai mare decât eroarea finală minimă (au prezentat o distorsiune finală mai mare de 0,210). Similar, în cazul algoritmului *batch - M* acest lucru s-a întâmplat cu 54 experimente. Așa cum era de așteptat situația este sensibil mai proastă în cazul algoritmului *batch* cu parametri normali, deoarece multe experimente au fost divergente: 89 experimente n-au atins nivelul de distorsiune finală de 0,210. Mai surprinzător, rezultatele sunt chiar mai slabe și în cazul experimentelor efectuate cu parametri adaptați (micșorați) caz în care numărul de experimente care nu ating nivelul final de distorsiune

SOLUȚII DE IMPLEMENTARE A REȚELELOR NEURONALE PE ARHITECTURI SISTOLICE

impus este 95. Deoarece rulările cu parametri reduși sunt mult mai puțin catastrofice decât cele cu parametri normali (toate experimentele converg în cazul rulărilor cu parametri reduși), este clar că viteza de convergență a algoritmului *batch* este "stânjenită" de necesitatea satisfacerii condiției (7.9).

În final, trebuie menționat că parametrii corespunzători la șapte dintre cele mai bune zece experimente aferente algoritmului clasic se regăsesc și în cele mai bune zece experimente aferente algoritmului *batch - M*: aceasta demonstrează că orice metodă euristică de optimizare a parametrilor aferenți algoritmului clasic poate fi aplicată fără nici o modificare și în cazul algoritmului *batch - M*. Rezultatele sunt sintetizate în tabelul 7.4.

Algoritm	Experimente cu Distorsiune $e \geq 1,2 \cdot e_{\text{optim}}$	Distorsiunea e în experimentul cel mai slab
clasic	58 48 %	0,53
<i>batch</i> (parametri normali)	89 74 %	21,03
<i>batch</i> (parametri adaptați)	95 79 %	1,22
<i>batch - M</i>	54 45 %	0,43

Tabelul 7.4. Robuștețea algoritmilor clasic, *batch* și *batch - M*

D. Concluzii

Simulările experimentale prezentate în paragraful anterior confirmă rezultatele studiilor teoretice. Într-adevăr rezultatele experimentale confirmă că, pe o problemă reală, algoritmul *batch - M* dă rezultate sensibil identice (accidental mai bune) decât algoritmul original. Aceste rezultate au fost evaluate cu referire la:

1. calitatea rezultatului final
2. viteza de convergență
3. sensibilitatea redusă la parametrii de învățare (robustețe).

Numai în ceea ce privește viteza de convergență rezultatele experimentale reliefează o convergență ceva mai lentă a algoritmului *batch - M*; să menționăm totuși că diferența este mică și va putea fi supracompensată și mult depășită de gradul mult mai înalt de paralelism pe care acest algoritm îl permite în implementările *hardware* (procesarea în tehnică *pipeline* a vectorilor de intrare, *pipeline*-izare asigurată și de RS proiectată în capitolele 3, 4 și 5). Situația este mai puțin pozitivă în cazul algoritmului *batch*. În acest caz, restricția de utilizare a parametrilor adaptați, impusă de condiția (7.9) poate afecta (uneori sever) viteza de convergență.

8. CONCLUZII

8.1. Obiectivele lucrării

Conducerea adaptivă în timp real a proceselor rapide și complexe reclamă sisteme de conducere capabile să reacționeze prompt la evenimentele care au loc în proces. *Neurocontroller*-ele implementate în *hardware* propuse în această lucrare satisfac pe deplin aceste cerințe.

Studiile prezentate în lucrare au avut drept scop concepția unor arhitecturi specializate pentru implementarea modelelor conexioniste. Problematika vizată a fost descompusă în mai multe părți:

- În domeniul rețelelor neuronale există o multitudine de modele care diferă prin structură și prin funcțiile pe care le realizează. Pentru o problemă dată, alegerea rețelei optime nu se poate baza încă pe niște criterii bine definite; alegerea se bazează în mare măsură pe un grad de pricepere câștigat prin experiență. În această lucrare s-au studiat principalele modele conexioniste actuale evidențiindu-se elementele lor constitutive: funcția de activare neuronală, structura, regulile de învățare etc. Studiul a permis evidențierea unui ansamblu de operații matriciale - vectoriale comune tuturor modelelor studiate. În faza de implementare *hardware* a algoritmilor conexioniști aceste operații comune (set de primitive) ne-au permis o detașare de constrângerile arhitecturale impuse de diversele modele. Astfel, în locul unei simple integrări de modele, studiul s-a concentrat asupra integrării unor lanțuri de primitive astfel încât să se obțină o implementare de înaltă performanță care să emuleze toate modelele conexioniste studiate.
- Pentru a înțelege și a inventaria constrângerile pe care le creează implementarea *hardware* (VLSI) a rețelelor neuronale, în capitolul 2 am trecut în revistă cele mai reprezentative realizări ale etapei actuale (digitale și analogice); aceasta a permis inventarierea principiilor esențiale care trebuie să stea la baza concepției unei arhitecturi conexioniste. Optând pentru integrarea în variantă digitală a matricii sinaptice, am estimat că implementarea sistolică ne va permite să obținem performanțe deosebite; aceste așteptări au fost confirmate de RS proiectată, implementată, perfecționată și evaluată în capitolele 3, 4 și 5.
- Rețeaua sistolică bidimensională propusă realizează într-o manieră originală funcțiile esențiale aferente algoritmilor conexioniști, atât pentru faza de recunoaștere cât și pentru faza de învățare. Această arhitectură permite implementarea operațiilor conexioniste de bază (setul de primitive): produsul matrice - vector, adunarea și scăderea matricilor, produsul unui vector cu transpusul său, detecția identității dintre 2 vectori etc. Ea permite de asemenea ajustarea ponderilor sinaptice care reprezintă funcția de bază a tuturor sistemelor evolutive. Arhitectura RS este modulară, se bazează numai pe conexiuni locale și permite extensia dimensională printr-un simplu pavaj rectangular de circuite VLSI.
- Un prim aport algoritmic se referă la dezvoltarea mai multor algoritmi care se bazează pe setul de primitive stabilit, atât pentru faza de recunoaștere cât și pentru cea de învățare. Formularea matematică a operațiilor impuse de algoritmi vizati a fost transformată în ecuații de dependență între diferitele puncte de calcul (celule) din structura RS. Atât transferurile de date din interiorul RS cât și cascada diferitelor

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

operații impuse de algoritmi vizati au fost optimizate, exploatând toate avantajele pe care le oferă rețeaua sistolică bidimensională: sincronizarea naturală a transferurilor de date, reducerea numărului de legături fizice dintre celulele din structura RS (comunicații strict locale), *pipeline*-izarea operațiilor, a transferurilor de date și a instrucțiunilor în cadrul RS. Calculele efectuate și comunicațiile de date sunt strict locale, celula sistolică este programabilă (diagonală sau non-diagonală) ceea ce va permite o extensibilitate foarte comodă a RS.

- Avantajele arhitecturii sistolice prezentate pot fi rezumate astfel:

A. Modularitate

Modularitatea este prezentă la mai multe nivele:

- ♦ celula sistolică este construită utilizând 3 tipuri de blocuri funcționale: registru, sumator, multiplexor. Înmulțitorul paralel - serie are la bază un sumator și un registru. Înmulțitorul serie - serie introduce 2 elemente noi: *(3, 2) counter (full - adder)* și *(5, 3) counter*.
- ♦ circuitul VLSI este constituit din celule identice care pot fi programate diagonale sau non-diagonale.
- ♦ rețeaua sistolică este realizată printr-un pavaj regulat de circuite identice, cu semnalele de la pini interconectate direct (fără interfețe specializate).

Este posibilă construcția unor rețele de dimensiuni mari prin simpla adăugare de circuite. Limita practică de extensibilitate este fixată doar de lungimea registrului PS care memorează potențialul neuronului.

B. Simplitate

Celula sistolică este simplă, constituită doar din câțiva operatori și câteva registre. Proiectarea, testarea și evaluarea ei se poate realiza rapid și cu un grad înalt de siguranță. Validarea celulei se poate face examinând și evaluând fiecare bloc independent din structura acesteia. Proiectarea și realizarea celulei sistolice înseamnă de fapt realizarea întregii rețele.

C. Reconfigurabilitate

Implementarea în variantă FPGA prezintă avantajul reconfigurabilității. Rețeaua sistolică nu trebuie să acopere în orice moment de timp, toți algoritmi conexioniști cunoscuți deoarece o astfel de abordare ar complica inutil structura celulei sistolice. Pe baza gradului de asemănare se pot constitui clase de algoritmi conexioniști și pentru fiecare clasă se poate proiecta celula adecvată. *Software*-ul de gestiune aferent plăcii acceleratoare va putea reconfigura dinamic RS în funcție de algoritmul neuronal ce urmează a se executa (aplicația lansată în execuție pe sistemul gazdă). De vreme ce structura CS rămâne în permanență simplă (deși se reconfigurează dinamic în timp), se vor putea obține (pentru fiecare algoritm neuronal în parte) viteze maxime de procesare în interiorul RS. Reconfigurarea RS se reduce la o simplă programare a circuitelor FPGA (modificarea conținutului tabelor LUT).

D. Grad înalt de paralelism

RS propusă reprezintă o implementare care integrează sinapse și nu neuroni. În contextul unei implementări VLSI, paralelismul masiv constituie un factor esențial. Paralelismul trebuie să fie cât mai fin posibil pentru a reproduce

structura distribuită aferentă RN; orice multiplexare a calculelor va conduce la reducerea calităților intrinseci aferente rețelelor conexiunite: reprezentare distribuită și toleranță la defect. RS propusă satisface aceste deziderate deoarece fiecare sinapsă este integrată într-o celulă sistolică proprie (cel mai fin paralelism posibil). Trebuie să precizăm că decizia de a integra sinapse și nu neuroni atrage și un dezavantaj: dimensiunea relativ redusă a RN implementabile într-un VLSI. Aici progresul tehnologic își va spune cuvântul. În momentul actual un FPGA XC 40250, de exemplu, poate găzdui o RN recurentă complet interconectată compusă din 22 neuroni. În condițiile dezvoltărilor arhitecturale propuse la 5.4 o matrice sinaptică de dimensiuni mari ar putea fi partajată în submatrici de dimensiuni egale cu dimensiunea RS. În timp ce RS procesează submatricea curentă, următoarea submatrice poate fi încărcată în RS (fără consum suplimentar de timp). În condițiile în care interfața RS - sistem gazdă asigură transferuri de date în ritmul solicitat de RS, dimensiunea RN emulate pe RS propusă nu mai constituie o problemă.

E. Performanțe

RS propusă este perfect adaptată algoritmilor conexioniști. Performanțele (CPS, UPS) evaluate la 4.6 sunt net superioare în raport cu alte realizări cunoscute. Aceste performanțe estimate ne îndreptățesc să considerăm că RS analizată conduce la implementări VLSI de mare eficiență utilizabile în aplicații reale complexe (timp real).

- Al doilea aport algoritmic se referă la elaborarea unor versiuni alternative ale algoritmilor neuronali clasici. RS reprezintă un *hardware* distribuit care implică un grad masiv și fin de paralelism în procesare. Maximizarea utilizării acestor resurse *hardware* distribuite sugerează anumite modificări ale algoritmilor neuronali clasici. O astfel de problematică este abordată în capitolul 7 al lucrării unde se prezintă algoritmul SOFM clasic precum și două versiuni alternative (*batch* și *batch - M*) care devin mult mai eficiente în cazul implementării rețelei SOFM pe rețeaua sistolică propusă. Spre deosebire de algoritmul clasic, algoritmi *batch* și *batch - M* permit un grad mult mai înalt de paralelism la nivelul implementării *hardware*. Acești algoritmi permit procesarea în tehnică *pipeline* a vectorilor de intrare, *pipeline*-izare asigurată și chiar revendicată de către mașina suport (RS).

8.2. Contribuții

În continuare vom prezenta, pe capitole, contribuțiile originale aduse în cadrul prezentei lucrări:

Capitolul 1

- ♦ o sinteză a principalelor modele conexiunite cu evidențierea elementelor lor constitutive: funcția de transfer, structură, reguli de învățare etc.
- ♦ evidențierea unui set de primitive (operații matricial - vectoriale comune tuturor modelelor studiate). Acest set de primitive a permis o detașare de constrângerile arhitecturale impuse de diversele modele. Implementarea *hardware* vizată nu va fi o simplă integrare de modele; vor fi integrate lanțuri de primitive care vor conduce la o implementare de înaltă performanță capabilă să emuleze toate modelele conexiunite studiate.

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

Capitolul 2

- ◆ un studiu comparativ al structurilor candidate la integrare (integrarea neuronilor și respectiv integrarea sinapselor) cu evidențierea avantajelor / dezavantajelor fiecărei structuri
- ◆ o sinteză a celor mai importante realizări în domeniul neuroprocesoarelor digitale și analogice
- ◆ stabilirea principiilor de proiectare a neuroprocesoarelor sistolice.

Capitolul 3

- ◆ o analiză a metodelor și descrierilor sistolice de tip SFG și DG care permit localizarea spațială a algoritmilor conexioniști
- ◆ elaborarea unei arhitecturi sistolice generale din care vor deriva 3 arhitecturi exploatabile:
 - inel sistolic cu componentele vectorului fixe
 - inel sistolic cu sume parțiale fixe
 - arhitectura cu coeficienți sinaptici ficși

Capitolul 4

- ◆ definirea arhitecturii sistolice optime și proiectarea acestei arhitecturi
- ◆ definirea algoritmului sistolic și demonstrarea validității acestuia (atât pentru faza de recunoaștere cât și pentru faza de învățare) pe baza Sistemului de Ecuații Recurente Uniforme
- ◆ proiectarea celulei sistolice și a rețelei sistolice caracterizată de următoarele avantaje:
 - permite implementarea fazei de recunoaștere pentru toți algoritmi neuronali clasici și a celor mai diverse reguli de învățare
 - permite procesarea în tehnică *pipeline* a vectorilor de intrare ceea ce conduce la o creștere spectaculoasă a eficienței RS; aceasta poate atinge valori de 100 % în fazele de procesare staționară (recurențe)
 - permite extensibilitatea RS printr-un simplu pavaj rectangular de circuite VLSI
- ◆ proiectarea și implementarea RS în structuri FPGA; testarea funcțională și evaluarea performanțelor rețelei proiectate.

Capitolul 5

- ◆ elaborarea unei strategii de perfecționare a rețelei sistolice dedicate algoritmilor conexioniști prin *pipeline*-izarea fluxurilor de date și de instrucțiuni în interiorul RS, prin perfecționarea structurii celulei sistolice și prin implementarea unui mecanism performant de interschimbare a matricilor sinaptice
- ◆ proiectarea unui înmulțitor serie - serie de înaltă performanță dedicat celulei sistolice perfecționate. Avantajele majore ale înmulțitorului proiectat sunt:
 - grad înalt de modularitate (extrem de important pentru implementările VLSI)
 - latență extrem de redusă între biții aplicați la intrare și bitul corespondent obținut la ieșire
 - permite extensia produsului obținut la ieșire la orice lungime (extensie de semn)

Capitolul 6

- ◆ analiza unei propuneri alternative: Rețeaua sistolică liniară (RSL) cu integrare de neuroni și nu de sinapse
- ◆ definirea unei strategii și stabilirea unui set de reguli privind optimizarea și creșterea eficienței RSL.

Capitolul 7

- ◆ elaborarea unor versiuni alternative pentru algoritmi neuronali clasici, versiuni care permit un grad mult mai înalt de paralelism la nivelul implementării *hardware* (procesarea în tehnică *pipeline* a vectorilor de intrare)
- ◆ analiza performanțelor și a vitezei de convergență aferente algoritmilor propuși prin metode teoretice și prin simulări experimentale.

8.3. Dezvoltări viitoare

Pornind de la studiul prezentat în această lucrare, pot fi angajate mai multe dezvoltări. O primă axă abordată deja parțial în lucrare ar viza realizarea unui neurosistem bazat pe un computer gazdă (de uz general) și o placă acceleratoare bazată pe RS propusă. Prin puterea de calcul pusă la dispoziție, un astfel de sistem ar constitui un suport esențial pentru cercetările desfășurate în domeniul algoritmilor conexioniști.

Realizarea neurosistemului presupune abordarea următoarei problematice:

- Proiectarea și construcția plăcii acceleratoare. Este o sarcină complexă care implică rezolvarea următoarelor probleme:
 - ◆ definitivarea structurii celulei sistolice
 - ◆ alegerea tehnologiei VLSI care se va utiliza pentru implementarea RS
 - ◆ definirea dimensiunii RS și matricii de circuite VLSI aferentă RS
 - ◆ proiectarea *hardware*-ului care implementează funcțiile sigmoideale. Funcția sigmoidală nu a făcut obiectul cercetărilor prezentată în această lucrare. Ea a fost lăsată în sarcina sistemului gazdă dar ar putea fi de asemenea implementată într-un circuit specializat. Circuitul VLSI dedicat s-ar putea baza pe tabele *look - up* sau pe implementări *hardware* speciale [ZHA 96]
 - ◆ proiectarea memoriei locale aferentă plăcii acceleratoare. Memoria locală are sarcina de a reține și a pune la dispoziția RS datele de lucru (matricea sinaptică, vectorii de intrare etc.). Interfața dintre memoria locală și RS reclamă o atenție deosebită; transferurile de date între cele două subsisteme trebuie să asigure un ritm normal de procesare în interiorul RS (fără stări *wait*)
 - ◆ proiectarea interfeței dintre placa acceleratoare și sistemul gazdă. Lărgimea canalului de comunicație dintre sistemul gazdă și placa acceleratoare trebuie să fie suficient de mare pentru a asigura procesarea în ritm normal în interiorul RS (fără stări *wait* cauzate de indisponibilitatea anumitor date).
- Proiectarea *software*-ului de gestiune aferent plăcii acceleratoare. Performanțele globale depind esențial de eficiența acestui *software*.

O a 2-a axă de dezvoltare vizează perfecționarea algoritmilor conexioniști. Această direcție a fost abordată în capitolul 7 dar a vizat doar algoritmi SOFM. Cercetarea trebuie continuată și extinsă și asupra celorlalte modele conexioniste.

SOLUȚII DE IMPLEMENTARE A REȚELOR NEURONALE PE ARHITECTURI SISTOLICE

O a 3-a axă de dezvoltare ar viza concepția unor circuite *Wafer Scale Integration* reconfigurabile care ar putea găzdui rețele complexe compuse dintr-un număr foarte mare de celule sistolice. Aceste implementări ar urma să fie utilizate în aplicațiile reale cele mai complexe (timp real).

Prescurtări utilizate

- ART - *Adaptive Resonance Theory* (arhitectură specializată de RN)
BAM - *Bidirectional Associative Memory* (memorie asociativă bidirecțională)
CLB - *Configurable Logic Block* (bloc logic configurabil)
CNAPS - *Connected Network of Adaptive Processors* (rețea interconectată de procesoare adaptive)
CP - *Control Processor* (procesor de control)
CPIF - *Control Processor Interface* (interfața procesorului de control)
CPS - *Connections Per Second* (conexiuni pe secundă)
CS - Celulă Sistolică
CSC - *CNAPS Sequencer Chip* (circuitul de secvențiere aferent ariei CNAPS)
DAP - *Distributed Array Processor* (arie de procesoare distribuită)
DG - *Dependence Graph* (graf de dependență)
DRAM - *Dynamic Random Access Memory* (memorie RAM dinamică)
DSP - *Digital Signal Processor* (procesor de semnal digital)
EBP - *Error Back Propagation* (propagarea erorii înapoi); algoritmul de învățare utilizat în cazul RN *feedforward* multistrat
EP - Element de Procesare
FPGA - *Field Programmable Gate Array* (circuite VLSI programabile)
GCPS - *Giga-Connections Per Second* (giga conexiuni pe secundă)
ICANN - *Proceedings of International Conference on Artificial Neural Networks*
ICNN - *Proceedings of the IEEE International Conference on Neural Networks*
IW - *Instruction Word* (instrucțiune)
LUT - *Look - Up Table* (tabelă look - up)
LVQ - *Learning Vector Quantization* (arhitectură specializată de RN)
MAC - *Multiplication & ACumulation unit* (unitate de înmulțire și acumulare); realizează înmulțirea a 2 factori urmată de acumularea produsului la rezultatul obținut în operația MAC anterioară.
MCPS - *Mega-Connections Per Second* (mega conexiuni pe secundă)
MIMD - *Multiple Instruction - Multiple Data streams* (fluxuri multiple de instrucțiuni și date)
MLP - *MultiLayer Perceptrons* (rețele de perceptroni multistrat)
MMU - *Memory Management Unit* (management de memorie)
PS - *Partial Sum* (sumă parțială)
RBF - *Radial Basis Function* (rețele neuronale cu funcții radiale de bază)
RCE - *Restricted Coulomb Energy* (arhitectură specializată de RN)
RISC - *Reduced Instruction Set Computer* (calculator cu set redus de instrucțiuni)
RN - Rețea Neuronală
RS - Rețea Sistolică
RSL - Rețea Sistolică Liniară
SERU - Sistem de Ecuații Recurente Uniforme
SIMD - *Single Instruction flow, Multiple Data flow* (flux unic de instrucțiuni și multiple fluxuri de date)
SFG - *Signal Flow Graph* (graf de mișcare a datelor)
SOFM - *Self Organizing Feature Map* (rețeaua cu auto-organizare a lui Kohonen)
SRAM - *Static Random Access Memory* (memorie RAM statică)
TDNN - *Time Delay Neural Networks* (arhitectură specializată de RN)
UPS - *Updates Per Second* (actualizări pe secundă)
VLSI - *Very Large Scale Integration* (circuite integrate pe scară largă)
WCNN - *World Congress on Neural Networks*

BIBLIOGRAFIE

- [ACI 93] Antonio d'Acierno, Roberto Vaccaro. "A Parallel Implementation of the Back - Propagation of Errors Learning Algorithm on a SIMD Parallel Computer", ICANN 1993, pag. 1074 - 1077.
- [ALI 91] C. Alippi, G. Storti - Gajani. "Simple Approximation of Sigmoid Functions: Realistic Design of Digital VLSI Neural Networks", Proc. IEEE Int. Symp. Circuits and Systems, 1991, pag. 1505 - 1508.
- [ALM 89] L. Almeida. "Back Propagation in perceptrons with feedback", Neural Computers, Springer Verlag, R. Eckmiller, C.v.d. Malsburg Eds., 1989.
- [AMI 85] D. Amit, H. Gutfreund, H. Sompolinski. "Storing infinite number of patterns in a spin - glass model of neural networks", Physical Review Letters, vol. 55, nr. 14, 1985, pag. 1530 - 1533.
- [AND 77] J. A. Anderson, J. W. Silverstein, S. A. Ritz, R. S. Jones. "Distinctive features, categorical perception on probability learning: some applications of a neural model", Psychological Review, Nr. 84, 1977, pag. 413 - 451.
- [ASA 90] K. Asanovic, N. Morgan. "Experimental determination of precision requirements for back - propagation training of artificial neural networks", 1st International Conference on Microelectronics for Neural Networks, Dortmund, 1990, pag. 9 - 15.
- [BAN 94] S. H. Bang, B. J. Sheu. "A Quick Search of optimal solutions for Cellular Neural Networks", WCNN, vol. II, 1994, pag. 549 - 554.
- [BAR 97] D. Baratta, G. Marco Bo, D. D. Caviglia, M. Valle, G. Canepa, R. Parenti and C. Penno, "A Hardware Implementation of Hierarchical Neural Networks for Real - Time Quality Control Systems in Industrial Applications", ICANN 1997, pag. 1229 - 1234.
- [BEN 90] Y. Bennani, F. Fogelman - Soulie, P. Gallinari. "Text - dependent speaker identification using LVQ", Proc of the Int. Neural Networks Conf., vol. 2, Paris, July 1990, pag. 1087 - 1090.
- [BER 96] A. Bermak, D. Martinez. "A Variable - Precision Systolic Architecture for ANN Computation", IEEE Proc. of MicroNeuro, 1996, pag. 347 - 354.
- [BLA 92] F. Blayo, M. Verleysen. "Setting initial conditions for the RCE model", Proc. of the First IFIP Working Group - 10.6 Workshop, Grenoble, 1992.
- [BLA 95] F. Blayo, A. G. Dugué, N. Maria. "Implementing Radial Basis Functions Neural Networks on the Systolic MANTRA Machine", Proceedings of International Workshop on Artificial Neural Networks, Malaga - Torremolinos, Spain, June 1995, pag. 781 - 788.
- [BOA 93] K. A. Boahen, A. G. Andeon. "Design of a bidirectional associative memory chip", Associative Neural Memories: Theory and Implementation, M. Hasson, Ed. New York: Oxford Univ. Press, 1993.
- [BOS 92] B. Boser, E. Sackinger, J. Bromley, Y Le Cun, L. D. Jackel. "Hardware requirements for neural network pattern classifiers", IEEE Micro, Feb. 1992, pag. 32 - 40.
- [BOU 93] C. Bouton, G. Pages. "Self - organization and a.s. convergence of the one - dimensional Kohonen algorithm with nonuniformly distributed stimuli", Stochastic Processes Applicat., vol. 47, 1993, pag. 249 - 274.
- [CHA 94] J. C. F. Chang. "An Efficient Digital VLSI Neural Processor Design for Image Processing", Signal and Image Proc. Institute Report, University of Southern California, 1994.
- [CHO 93a] W. Choi, H. Burlison, D. Phatak. "Fixed -point round off error analysis of large feed forward neural networks", Proc. IJCNN '93, 1993, pag. 1947 - 1950.
- [CHO 93b] A. N. Choudhary, J. H. Patel, N. Ahuja. "NERTA: A hierarchical and Partionable Architecture for Computer Vision Systems", IEEE Trans. on Parallel and Distributed Systems, vol. 4, nr. 10, Oct. 1993, pag. 1092 - 1104.

- [CHU 91] Chung Yoon. "A systolic array exploiting the inherent parallelism of artificial neural networks", International Conference on Parallel Processing, 1991.
- [CIM 96] L. Ciminiera, P. Montuschi. "Carry - Save Multiplication Schemes Without Final Addition", IEEE Trans. on Computers, vol. 45, nr. 9, Sept. 1996, pag. 1050 - 1055.
- [COM 90] P. Comon, C. Jutten, J. Héroult. "Separation of sources: statistical interpretation", Annexe 2 du rapport DRI - E3, ESPRIT - BRA NERVES Project, June 1990.
- [COM 94] P. Comon, J. L. Voz, M. Verleysen, "Estimation of performance bounds in supervised classification", In M. Verleysen, editor, ESANN: European Symposium on Artificial Neural Networks, Bruxelles, April 1994, pag. 37 - 42.
- [COR 96] T. Cornu, P. Jenne, D. Niebur, P. Thiran, M. A. Viredaz. "Design, implementation, and test of a multimodel systolic neural - network accelerator", Scientific Programming, vol. 5, nr.1, Spring 1996, pag. 47 - 61.
- [COS 89] M. Cosnard, Y. Robert. "Une introduction a l'algorithmique systolique", Actes de l'Université d'été CIRILLE, Lyon, Juillet, 1989.
- [COT 87] M. Cottrel, J. C. Fort. "Etude d'un algorithme d'auto - organisation", Annales de l'Institut Henri Poincare, vol. 23, nr. 1, 1987, pag. 1 - 20.
- [COT 88] M. Cottrel. "Stability and attractivity in Associative Memory Networks", Biol. Cyb. 58, 1988, pag. 129 - 139.
- [COT 89] M. Cottrel, J. C. Fort. "Aspects théoriques de l'algorithme d'auto - organisation de Kohonen", Annales du groupe CARNAC nr. 2, Lausanne, 1989.
- [DAD 89] L. Dadda. "On serial - input multipliers for two's complement numbers", IEEE Trans. On Computers, vol. 38, nr. 9, Sept. 1989, pag. 1341 - 1345.
- [DAD 96] L. Dadda, V. Piuri. "Pipelined Adders", IEEE Trans. on Computers, vol. 45, nr. 3, March 1996, pag. 348 - 356.
- [DIE 87] S. Diederich, M. Oppel. "Learning of Correlated patterns in spin - glass network by local learning rules", Physical Review Letters, vol. 58, nr. 9, 1987, pag. 949 - 952.
- [DOW 93] E. M. Dowling, Z. Fu, R. S. Drafz. "HARP: An Open Architecture for Parallel Matrix and Signal Processing", IEEE Trans. on Parallel and Distributed Systems, vol. 4, nr. 10, Oct. 1993, pag. 1081 - 1091.
- [DUD 73] R. O. Duda, P. E. Hart. "Pattern classification and scene analysis", Wiley and Sons, 1973.
- [EDE 82] G. M. Edelman, G. N. Reeke. "Selective networks capable of representative transformation, limited generalizations and associative memory", Proc. Natl. Acad. Sci., vol. 79, March 1982, pag. 2091 - 2095.
- [ERW 92] E. Erwin, K. Obermayer, K. Schulten. "Self - organizing maps: Ordering, convergence properties, and energy functions", Biol. Cybern., vol. 67, 1992, pag. 47 - 55.
- [FAN 92] W. C. Fang et al. "A VLSI neural processor for image data compression using self - organizing networks", IEEE Trans. on Neural Networks, vol. 3, May 1992, pag. 506 - 518.
- [FAT 96] Sherif Kassem Fathy, Mostafa Mohamed Syiam. "A Parallel design and Implementation for Backpropagation Neural Network Using MIMD Architecture", ICNN 1996, pag. 1361 - 1366.
- [FEN 96] S. T. J. Fenn, M. Benaissa, D. Taylor. "GF (2^m) Multiplication and Division Over the Dual Basis", IEEE Trans. on Computers, vol. 45, nr. 3, March 1996, pag. 319 - 327.
- [FIS 97] T. Fischer, W. Eppler, H. Gemmeke, G. Kock and T. Becker. "The SAND Neurochip and Its Embedding in the MIMD System", ICANN 1997, pag. 1234 - 1240.
- [FOR 87] J. A. B. Fortes, B. W. Wah. "Systolic arrays: from concepts to implementation", IEEE Computer, July 1987.

- [FOR 93] J. C. Fort, G. Pages. "Sur la convergence p.s. sure de l'algorithme de Kohonen Généralisé", Comptes Rendus de l'Académie des Sciences de Paris, vol. 317, 1993, pag. 389 - 394.
- [FUK 75] K. Fukushima. "Cognitron: a self organizing multilayered neural network", Biol. Cyb. Nr. 20, 1975, pag. 121 - 136.
- [GRA 88] D. H. Graf, W. R. Lalonde. "A neural controller for collision free movement of general robot manipulator", Int. Conf. on Neural Networks, San Diego, July 24 - 27, vol. 1, 1988, pag. 77 - 84.
- [GRI 91] M. Griffin, G. Tahara, K. Knorpp, R. Pinkham, B. Riley, "An 11-million transistor neural networks execution engine", Tech. Digest. IEEE Inter. Solid - State Circuits Conference, San Francisco, CA, Feb. 1991, pag. 180 - 181.
- [GRO 86] S. Grossberg. "The adaptive Brain", Elsevier, North Holland, Amsterdam, 1986.
- [GRO 87] S. Grossberg. "Competitive learning: from Interactive Activation to Adaptive Resonance Theory", Cognitive Science, vol. 11, 1987, pag. 23 - 63.
- [GUI 89] M. Guivarch. "Réseaux: techniques de base et apprentissage", 3-eme université d'été CIRILLE, Lyon, 3 - 7 Juillet, 1989.
- [GUY 88] I. Guyon. "Réseaux de neurones pour la reconnaissance de formes: architectures et apprentissage", These de Doctorat de l'Université Paris VI, Décembre, 1988.
- [GUY 90] I. Guyon. "Neural networks and applications", Internal Report AT & T Bell Labs. September 1990.
- [HAM 87] S. E. Hampson, D. J. Volper. "Disjunctive models of boolean category learning", Biol. Cyb. Nr 56, 1987, pag. 121 - 137.
- [HAM 91] D. Hammerstrom. "A highly parallel digital architecture for neural network emulation", VLSI for Artificial Intelligence and Neural Networks, J. G. Delgado, Frias, W. Moore, Eds., New York: Plenum, 1991, pag. 357 - 366.
- [HEB 49] D. Hebb. "The organisation of behavior", John Wiley and sons, New York, 1949.
- [HEC 87] R. Hecht - Nielsen. "Counterpropagation networks", Applied Optics, nr. 26, Dec. 1987, pag. 4976 - 4984.
- [HOL 91] J. L. Holt, T. E. Baker. "Back propagation simulation using limited precision calculation", International Joint Conference on Neural Networks, Seattle, Wash., 1991.
- [HOL 93] J. Holt, J. N. Hwang. "Finite Precision Error Analysis of Neural Network Hardware Implementations", IEEE Trans. on Computers, vol. 42, 1993, pag. 281 - 290.
- [HOL 94] P. W. Hollis, J. J. Paulos. "A Neural Network Learning Algorithm Tailored for VLSI Implementation", IEEE Trans. on Neural Networks, vol. 5, nr. 5, Sept 1994, pag. 748 - 791.
- [HOP 82] J. J. Hopfield. "Neural Networks and physical systems with emergent collective computational abilities", Proc. of the National Academy of Science, USA, vol. 79, 1982 pag. 2554 - 2558.
- [HUB 77] D. Hubbel, T. Wiesel. "Functional architecture of the macaque monkey visual cortex", Ferrier lecture Proc. Royal Society, London. B. 198, 1977, pag. 1 - 59.
- [HUG 92] F. M. Hugen, Z. Houkes. "Systolic arrays for real - time recursive linear least squares parameter estimation", Proceedings of the Workshop on Circuits, Systems and Signal Processing, Houthalen, 1992, pag. 283 - 289.
- [IEN 94] P. Ienne, M. A. Viredaz. "Bit - serial multipliers and squarers", IEEE Trans. on Computers, vol. 43, nr. 12, Dec. 1994, pag. 1445 - 1450.
- [IEN 95] P. Ienne, M. A. Viredaz. "GENEZ IV: A bit - serial processing element for multimodel neural - network accelerator", Journal VLSI Signal Processing, vol. 9, nr. 3, Apr. 1995, pag. 257 - 273.

- [IEN 97] P. Ienne, P. Thiran, N. Vassilas. "Modified Self - Organizing Feature Map Algorithms for Efficient Digital Hardware Implementation", IEEE Trans. on Neural Networks, vol. 8, nr. 2, March 1997, pag. 315 - 330.
- [INT 93] Intel Neural Network Group. "Design and Implementation of a recognition accelerator", Proc. of the Canadian Conference on Very Large Scale Integration, 1993.
- [JAH 97] A. Jahnke, T. Schönauer, U. Roth, K. Mohraz and H. Klar. "Simulation of Spiking Neural Networks on Different Hardware Platforms", ICANN 1997, pag. 1187 - 1192.
- [JEA 94] J. S. N. Jean, J. Wang. "Weight Smoothing to Improve Network Generalization", IEEE Trans. on Neural Networks, vol. 5, nr. 5, Sept. 1994, pag. 752 - 763.
- [JOH 92] A. Johannet, L. Personnaz, G. Dreyfus, J. D. Gascuel, M. Weinfeld. "Specification and implementation of a digital Hopfield - type associative memory with on - chip training", IEEE Transaction on Neural Networks, vol. 3 (4), July 1992, pag. 529 - 539.
- [JOH 93] K. T. Johnson, A. R. Hurson. "General - purpose systolic arrays", IEEE Comp., November 1993, pag. 20 - 31.
- [JOH 94] S. R. Johns, K. Sammut. "Learning in linear systolic neural network engines: Analysis and implementation", IEEE Trans on Neural Networks, vol. 5, nr. 4, 1994, pag. 584 - 593.
- [JUT 87] C. Jutten. "Calcul neuro - mimétique et traitement du signal. Analyse en composantes indépendantes", These de Doctorat d'etat, Grenoble, Juin 1987.
- [JUT 88] C. Jutten, J. Héroult. "Une solution neuromimétique au probleme de séparation de sources", Tratement du Signal, vol. 5, nr. 6, 1988, pag. 389 - 403.
- [JUT 90] C. Jutten, A. Guérin, J. Héroult. "Simulation machine and integrated implementation of neural networks", Lecture Notes in Computer Science: Neural Networks, Almeida and Wellekens Eds. Springer Verlag, 1990.
- [KAN 93a] V. Kantabutra. "Accelerated Two - level carry - Skip Adders - A Type of Very Fast Adders", IEEE Trans. on Computers, vol. 42, nr. 11, Nov. 1993, pag. 1389 - 1393.
- [KAN 93b] V. Kantabutra. "A Recursive Carry - Lookahead / Carry - Select Hybrid Adder", IEEE Trans. on Computers, vol. 42, nr. 12, Dec. 1993, pag. 1495 - 1499.
- [KIN 94] F. S. King, P. Saratchandran, N. Sundararajan. "A Theoretical Study of Training Set Parallelism for Back propagation Networks on a Traspouter Array", Proceed. of WCNN, 1994, vol. II, pag. 519 - 524.
- [KNO 88] S. C. Knowles, R. F. Wood, J. G. McWriter, J. V. McCanny. "Bit level systolic arrays for the IIR filtering", Int. Conf. on Systolic Arrays, San Diego, May 25 - 27, 1988, pag. 653 - 663.
- [KOH 86] T. Kohonen. "Self organisation and associative memory", Springer, Berlin, 1986, 2nd Edition.
- [KOH 88a] T. Kohonen. "The "neural" phonetic typewriter", Computer, March 1988, pag. 11 - 22.
- [KOH 88b] T. Kohonen. "Self - Organization and Associative Memory", 2nd edition, Berlin, Springer - Verlag, 1988.
- [KOH 97] T. Kohonen. "Self - Organizing Maps", Springer - Verlag, Berlin Heidelberg, 1997.
- [KOL 97] Pasi Kolinummi, Timo Härmäläinen and Jukka Saarinen. "Mapping of Radial Basis Function Networks to Partial Tree Shape Parallel Neurocomputer", ICANN 1997, pag. 1259 - 1264.
- [KON 93] A. König, X. Geng, M. Glesner. "Hardware Implementation of Kohonen's Feature Map by Scalar and SIMD - Array Processors", ICANN 1993, pag. 1046 - 1049.
- [KOR 90] I. Koren, O. Zinaty. "Evaluating Elementary Functions in a Numerical Co-Processor Based on Rational Approximations", IEEE Trans. on Computers, vol. 39, nr. 8, Aug. 1990, pag. 1030 - 1037.

- [KOR 93] I. Koren. "Computer Arithmetic Algorithms", Englewood Cliffs, New Jersey, Prentice Hall, 1993.
- [KOS 87] B. Kosko, "Bidirectional associative memories", BYTE Magazine, Sep. 1987.
- [KUN 78] H. T. Kung, C. E. Leiserson. "Algorithm for VLSI processor arrays", Symposium on Sparse Matrix Computations and Their Applications, Knoxville, Nov. 2 - 3, 1978.
- [KUN 88a] S. Y. Kung. "VLSI array processors", T. Kailath Ed., Prentice Hall, Englewood Cliffs, New Jersey, 1988.
- [KUN 88b] S. Y. Kung, J. N. Hwang. "An algebraic projection for optimal hidden units size and learning rates in backpropagation learning", Proc. IEEE International Conference on Neural Networks, San Diego, vol. 1, 1988, pag. 363 - 370.
- [LAW 93] J. C. Lawson, N. Maria, H. J. Smart. "A neuro - computer using sparse matrices", In Euro Micro PDP Workshop, Canari (Spain), January 1993.
- [LeC 85] Y. Le Cun. "Une procédure d'apprentissage pour réseau a seuil asymetrique", Actes de COGNITIVA 85, Paris, Juin 1985, pag. 599 - 604.
- [LEH 93] C. Lehmann, M. Viredaz, F. Blayo. "A generic systolic array building block for neural networks with on - chip learning", IEEE Trans. on Neural Networks, vol. 4, nr. 3, May 1993, pag. 400 - 407.
- [LEN 97] P. Lenders, S. Rajopadhye. "Multirate VLSI Arrays and Their Synthesis", IEEE Trans. on Computers, vol. 46, nr. 5, May 1997, pag. 515 - 529.
- [LI 85] G. H. Li, B. W. Wah. "The design of optimal systolic arrays", IEEE Transaction on Computers, vol. 34, nr. 1, 1985, pag. 66 - 67.
- [LIP 87] R. Lippmann. "An introduction to computing with neural nets", IEEE ASSP Magazine, April 1987, pag. 4 - 22.
- [MAL 94] Q. M. Malluhi, M. A. Bayounmi, T. R. N. Rao. "A parallel Algorithm for Neural Computing", WCNN, vol. II, 1994, pag. 563 - 569.
- [MAN 90] R. Mann, S. Haykin. "A parallel implementation of Kohonen feature maps on the Warp systolic computer", Proc. Int. Joint Conf. on Neural Networks, vol. II, Washington, Jan. 1990, pag. 84 - 87.
- [MAR 90] P. W. Markstein. "Computation of Elementary Functions on the IBM RISC System/6000 Processor", IBM J. Research and Development, vol. 34, nr. 1, Jan. 1990, pag. 111 - 119.
- [MAR 94] N. Maria, A. G. Dugué, F. Blayo. "1D and 2D Systolic Implementations for Radial Basis Function Networks", IEEE Proc. of MicroNeuro, 1994, pag. 34 - 45.
- [MAS 92] R. Mason et al. "An hierarchical VLSI neural - network architecture", IEEE J. Solid - State Circuits, vol. 27, Jan. 1992, pag. 106 - 108.
- [MAS 97] E. I. El - Masry, H. K. Yang, M. A. Yakout. "Implementations of Artificial Neural Networks Using Current - Mode Pulse Width Modulation Technique", IEEE Trans. on Neural Networks, vol. 8, nr. 3, May 1997, pag. 532 - 548.
- [MAU 92] N. Mauduit. "Lneuro 1.0: a piece of hardware LEGO for building neural network systems", IEEE Trans. on Neural Networks, March 1992.
- [McC 43] W. S. McCulloch and W. Pitts. "A logical calculus of the ideas immanent in nervous system", Bulletin of Mathematical Biophysics, nr. 5, 1943, pag. 115 - 133.
- [McG 87] R. J. Mac Gregor. "Neural and Brain modeling", Neuroscience Series Academic Press Inc. London, 1987.
- [McL 97] S. McLoone, G. W. Irwin. "Fast Parallel Off - Line Training of Multilayer Perceptrons", IEEE Transactions on Neural Networks, vol. 8, nr. 3, May 1997, pag. 646 - 652.
- [MEL 92] M. S. Melton, T. Phan, D. S. Reeves, D. E. Van den Bout, "The TinMANN VLSI chip", IEEE Trans. Neural Networks, vol. 3, nr. 3, May 1992, pag. 375 - 384.
- [MIH 93] I. Z. Mihiu. "Some aspects concerning the markovian analysis of overlap systems performance", Acta Universitatis Ciniensis Technical series, vol. X(1), Sibiu, 1993. pag. 52 - 59.

- [MIH 95] I. Z. Mihu, V. Tulai. "O implementare hard în tehnică pipeline a rețelelor neuronale feed forward", Acta Universitatis Ciniensis, Seria Tehnică, vol. XXII, Sibiu 1995, pag. 72 - 77.
- [MIH 96a] I. Z. Mihu. "Tendențe în ingineria sistemelor asistată de calculator", Referat doctorat, Universitatea Tehnică Timișoara, 1996.
- [MIH 96b] I. Z. Mihu. "Sisteme bazate pe ingineria cunoașterii în conducerea proceselor", Referat doctorat, Universitatea Tehnică Timișoara, 1996.
- [MIH 97] I. Z. Mihu. "Arhitectura sistemelor de calcul. Paralelismul în sistemele microprocesor", JEP 08012/94, Univ. Tehnică Cluj - Napoca, 1997.
- [MIN 88] M. Minsky, S. Pappert. "Perceptrons", MIT Press, Expanded Edition, 1988.
- [MOR 92] N. Morgan, J. Beck, P. Kohn, J. Bilmes, E. Allman, J. Beer, "The Ring Array Processor: A multiprocessing peripheral for connexionist applications", Journal of Parallel and Distributed Comput., vol. 14, nr. 3, March 1992, pag. 248 - 259.
- [MOR 95] J. M. Moreno, J. Madrenas, S. S. Anselmo, F. Castillo, J. Cabestany. "Digital Hardware Implementation of ROI Incremental Algorithms", Proceedings of International Workshop on Artificial Neural Networks, Malaga - Torremolinos, Spain, June 1995, pag. 761 - 770.
- [MUE 93] D. Mueller, D. Hammerstrom. "A neural network systems component", Proc. IEEE Inter. Conf. Neural Networks, vol. 3, San Francisco, CA, Mar. 1993, pag. 1258 - 1264.
- [MUL 92] U. A. Müller, B. Bäuml, P. Kohler, A. Gunzinger, W. Guggenbühl. "Achieving supercomputer performance for neural - net simulation with an array of digital signal processors", IEEE Micro, Oct. 1992, pag. 55 - 65.
- [MUR 92] P. Murtagh, A. C. Tsoi. "Implementation Issues of Sigmoid Function and its Derivative for VLSI Digital Neural Networks", IEE Proc. - E, vol. 139, nr.3, May 1992.
- [MYE 91] D. J. Myers, J. M. Vincent, D. A. Orrey, "HANNIBAL: A VLSI building block for neural networks with on - chip back propagation learning", Proc. of 2nd Int. Conf. Microelectronics for Neural Networks, Munich, 1991, pag. 171 - 181.
- [MYO 96] Jean Frédéric Myoupo, David Semé. "A single - layer Systolic Architecture for Back propagation Learning", ICNN 1996, pag. 1329 - 1333.
- [NAS 92] N. M. Nasrabadi, C. Y. Choo, "Hopfield Network for Stereo Vision Correspondence", IEEE Trans. on Neural Networks, January 1992.
- [NAY 95] D. Naylor, S. Jones, D. Myers. "Backpropagation in Linear Arrays - A performance Analysis and Optimization" IEEE Trans. on Neural Networks, vol. 6, nr. 3, May 1995, pag. 583 - 594.
- [OJA 82] E. Oja. "A simplified neuron model as a principal component analyzer", J. Math. Biol. Nr. 15, 1982, pag. 267 - 273.
- [ORR 91] D. A. Orrey, D. J. Myers, J. M. Vincent, "A high performance digital processor for implementing large artificial neural networks", Proc. IEEE Custom Integrated Circuits Conference, San Diego, May 1991, pag. 16.3.1. - 16.3.4.
- [OKL 96] V. G. Oklobdzija, D. Villeger, S. S. Liu. "A Method for Speed Optimized Partial Product Reduction and Generation of Fast Parallel Multipliers Using an Algorithmic Approach", IEEE Trans. on Computers, vol. 45, nr. 3, March 1996, pag. 294 - 306.
- [PAR 96] R. Parisi, E. D. Di Claudio, A. Rapagnetta, G. Orlandi. "Recursive Least Squares Approach to Learning in Recurrent Neural Networks", ICNN 1996, pag. 1350 - 1354.
- [PAT 96] Minesh I. Patel, N. Ranganathan. "PANTHER: A Parallel Neuro Systolic Architecture for Real - Time Processing", ICNN 1996, pag. 1006 - 1011.
- [PER 86] L. Personnaz, I. Guyon, G. Dreyfus. "Collective computational proprieties of neural networks", Physical Review A, vol. 34, nr.5, Nov. 1986, pag. 4217 - 4228.
- [PER 87] L. Personnaz, I. Guyon, G. Dreyfus. "High order neural networks: information storage without errors", Europhysics letters, nr. 4(8), Oct. 1987, pag. 863 - 867.

- [PER 97] A. Pérez. - Uribe and E. Sanchez, "FPGA Implementation of a Network of Neurolike Adaptive Elements", ICANN 1997, pag. 1246 - 1252.
- [PIU 94] V. Piuri, S. Bettola. "High - performance digital Neural Networks: the use of redundant binary representation for concurrent error detection", WCNN, vol. II, 1994, pag. 531 - 535.
- [POR 91] W. Porter. "Neuronic Arrays: Design and Performance", Proceedings of the Second Workshop on Neural Networks: Academic / Industrial / NASA / Defense, WNN - AIND 91, pag. 119 - 128.
- [PRZ 91] K. W. Przytula, W. M. Lin, V. K. P. Kumar, "Partitioned implementation of neural networks on mesh connected array processors", VLSI Signal Processing IV, IEEE Pres.: New York, Aug. 1991, pag. 106 - 115.
- [PUL 96] Ville Pulkki, Taneli Harju. "An implementation of the Self - Organizing Map on the CNAPS Neurocomputer", ICNN 1996, pag. 1345 - 1349.
- [QUI 84] P. Quinton. "Automatic synthesis of systolic arrays from Uniform Recurrent Equations", Proc. of 11th Annual Symposium on Computer Architecture, Ann Arbor, Michigan, 1984, pag. 208 - 214.
- [QUI 89] P. Quinton, Y. Robert. "Algorithmes et architectures systoliques", Masson, Paris, 1989.
- [RAM 91] U. Ramacher. "Guidelines to VLSI Design of Neural Networks", Silicon Archs. For Neural Nets, Eds. M. Sami, J. Calzadilla - Daguerre, Elsevier Science Publisher B. V. (North - Holland), 1991.
- [RAM 91a] U. Ramacher, J. Beichter, N. Brüls. "Architecture of a general - purpose neural signal processor", Proc. IEEE / INNS Inter. Joint Conf. Neural Networks, vol. 1, Seattle, WA, July 1991, pag. 443 - 446.
- [RAM 91b] U. Ramacher, J. Beichter, W. Raab, J. Anlauf, N. Brüls, U. Hachmann, M. Wesseling. "Design of a 1st generation neurocomputer", in VLSI Design of Neural Networks (U. Ramacher, U. Rückert, eds.), Kluwer Academic Publisher, Boston, 1991, pag. 271 - 310.
- [RAM 91c] U. Ramacher, B. Schürmann. "Unified description of neural algorithms for time - independent pattern recognition", in VLSI Design of Neural Networks (U. Ramacher, U. Rückert, eds.), Kluwer Academic Publisher, Boston, 1991, pag. 255 - 270.
- [RAM 92] U. Ramacher. "SYNAPSE - A neurocomputer that synthesizes neural algorithms on a parallel systolic engine". Journal of Parallel and Distributed Comput., vol. 14, nr. 3, Mar. 1992, pag. 306 - 318.
- [RAM 93] U. Ramacher, W. Raab, J. Anlauf, J. Beichter, U. Hachmann, N. Brüls, M. Weßeling, E. Sicheneder. "Multiprocessor and Memory Architecture of the Neurocomputer SYNAPSE - 1", ICANN 1993, pag. 1034 - 1038.
- [REZ 97] S. M. Rezaul Hasan, N. K. Siong. "A Parallel Processing VLSI BAM Engine", IEEE Trans. on Neural Networks, vol. 8, nr. 2, Mar. 1997, pag. 424 - 436.
- [REY 91] L. M. Reyneri, E. Philippi, "An analysis on the performance of silicon implementations of back propagation algorithms for artificial neural networks", IEEE Trans. Comput., vol. 40, nr. 12, Dec. 1991, pag. 1380 - 1389.
- [ROB 92] M. E. Robinson, H. Yoneda, E. Sanchez - Sinencio, "A Modular CMOS Design of a Hamming Network", IEEE Trans. on Neural Network, vol. 3, No. 3, May 1992.
- [ROS 58] F. Rosenblatt, "The Perceptron; a probabilistic model for information storage and organization in the brain", Psychological Review, Nr. 65, 1958, pag. 386 - 408.
- [ROS 97] M. Rossmann, A. Bühlmeier, G. Manteuffel and K. Goser. "Short - and Long - Term Dynamics in a Stochastic Pulse Stream Neuron Implemented in FPGA", ICANN 1997, pag. 1241 - 1246.

- [ROT 97] Ulrich Roth, Axel Jahnke and Heinrich Klar. "On - line Hebbian Learning for Spiking Neurons: Architecture of the Wecht - Unit of NESPINN", ICANN 1997, pag. 1217 - 1222.
- [RUE 90] P. F. Ruedi, E. Sorouchyari. "Commande d'un bras de robot redondant par reseaux de neurones", Centre Suisse d'Electronique et de Microtechnique, Rapport technique Nr. 354, Mai 1990.
- [RUE 93] S. Rueping, U. Rueckert, K. Goser. "Hardware Design for Self - Organizing Feature Maps with Binary Input Vectors", Proceedings of the IWANN '93, Sigtes, Spain, 1993, pag. 488 - 493.
- [RUE 94] S. Rueping, U. Rueckert, K. Goser. "A Chip for Selforganizing Feature Maps", Proc. of the 4th Int. Conf. of microelectronics for Neural Networks and Fuzzy Systems, IEEE Computer Society Press, 1994, pag. 26 - 33.
- [RUE 95] S. Rueping, U. Rueckert, K. Goser. "A Chip for Selforganizing Feature Maps", IEEE MICRO, vol. 15, Nr. 3, June 1995, pag. 57 - 59.
- [RUI 95] V. Ruiz de Angulo, C. Torras. "On-Line Learning with Minimal Degradation in Feedforward Networks", IEEE Trans. on Neural Networks, vol. 6, Nr. 3, May 1995, pag. 657 - 668.
- [RUM 86] D. E. Rumelhart, J. McClelland. "Parallel distributed processing, vol. 1: Foundations, vol. 2: Psychological and Biological models", Institute for Cognitive Science, Bradford Book, The MIT press, J. A. Feldman, P. J. Hayes, D. E. Rumelhart Eds., 1986.
- [RUM 90] D. E. Rumelhart. "Brain Style Computation: Learning and Generalization", Introduction to Neural and Electronic Networks, New York, Academic Press, 1990.
- [RUP 96] S. Rüping, U. Rückert. "A Scalable Processor Array for Self - Organizing Feature Maps", IEEE Proc. of MicroNeuro, 1996, pag. 285 - 291.
- [SAC 96] E. Sackinger, H. P. Graf. "A board System for High - Speed Image Analysis and Neural Networks", IEEE Trans. on Neural Networks, vol. 7, nr. 1, Jan. 1996, pag. 214 - 221.
- [SAE 94] R. Sacks, K. Priddy, K. Schneider, S. Stowell. "On the Design of an MIMD Neural Network Processor", WCNN, vol. II, 1994, pag. 590 - 595.
- [SAM 91] K. Sammut, S. R. Jones. "Implementing the adaptive resource theory in linear array", Proc. of 2nd Int. Conference Microelectronics for Neural Networks, Munich, 1991, pag. 69 - 76.
- [SAN 92] E. Sanchez - Sinencio, C. Lau (Eds.). "Artificial Neural Networks: Paradigms, Applications and Hardware Implementations", IEEE Press: Piscataway, 1992.
- [SAT 92] S. Satyanarayan et al. "A reconfigurable VLSI neural network", IEEE J. Solid - State Circuits, vol. 27, Jan. 1992, pag. 67 - 81.
- [SCH 92] I. Scherson, D. Kramer, B. Alleyne, "Bit - parallel arithmetic in massively - parallel associative processor", IEEE Trans. on Computers, vol. 41 (10), 1992, pag. 1201 - 1210.
- [SER 93] M. J. Serrano, B. Parhani. "Optimal Architectures and Algorithms for Mesh - Connected Parallel Computers with Separable Row / Column Buses", IEEE Trans. on Parallel and Distributed Systems, vol. 4, nr. 10, Oct. 1993, pag. 1073 - 1080.
- [SHE 95] Bing J. Shen, J. Choi. "Neural Information Processing and VLSI", Kluwer Academic Publishers, Boston, 1995.
- [SIG 91] A. Siggelkow, J. Nijhuis, at. al. "Influence of Hardware Characteristics on the Performance of a Neural System", Artificial Neural Networks, vol. 1, North - Holland, Eds. Kohonen, K. Mäkisaro, O. Simula, J. Kangas, Amsterdam, 1991.
- [SPE 93] H. Speckmann, P. Thole, W. Rosenstiel. "COKOS: A COprocessors for Kohonen's Selforganizing Map, ICANN, 1993, pag. 1040 - 1044.
- [SPE 94] H. Speckmann, P. Thole, M. Bogdan, W. Rosenstiel. "Coproprocessors for special neural networks KOKOS and KOBOLD", WCNN, vol. II, 1994, pag. 612 - 617.

- [STE 97] M. Stevenson. "The Effects of Limited - Precision Weights on the Threshold Adaline", IEEE Trans. on Neural Networks, vol. 8, nr. 3, May 1997, pag. 549 - 552.
- [STR 82] N. R. Strader, V. T. Ryhne. "A canonical bit - sequential multiplier", IEEE Trans. on Computers, vol. C - 31, nr. 8, Aug. 1982, pag. 791 - 795.
- [STR 95] A. Strey, N. Avellana, R. Holgado, J. A. Fernandez, R. Capillas, E. Valderrama. "A Massively Parallel Neurocomputer with a Reconfigurable Arithmetical Unit", Proceedings of International Workshop on Artificial Neural Networks, Malaga - Torremolinos, Spain, June 1995, pag. 800 - 806.
- [SUN 93] M. H. Sunwoo, J. K. Aggarwal. "A Sliding Memory Plane Array Processor", IEEE Trans. on Parallel and Distributed Systems, vol. 4, nr. 6, Jun. 1993, pag. 601 - 612.
- [TAM 97] S. Tamura, M. Tateishi. "Capabilities of a Four - Layered Feedforward Neural Network: Four Layers Versus Three", IEEE Trans. on Neural Networks, vol. 8, nr. 2, Mar. 1997, pag. 251 - 255.
- [THI 94] P. Thiran, M. Hasler. "Self - organization of a one - dimensional Kohonen network with quantized weights and inputs", Neural Networks, vol. 7, nr. 9, 1994, pag. 1427 - 1439.
- [THI 95] P. Thissen, M. Verleysen, J. D. Legat. "An associative processor architecture for pattern recognition", Proc. of the ProRISC / IEEE Benelux Workshop on Circuits, Systems and Signal Processing, Martie 1995, pag. 309 - 315.
- [TRY 89] V. Tryba, S. Metzen, C. Goser. "Designing basic integrated circuits by self - organizing feature maps", Proc. of Neuro Nimes 89, Int. Workshop, Nimes, Nov. 13 - 16, 1989, pag. 225 - 234.
- [TSA 95] J. C. Tsay, P. Y. Chang. "Some New Designs of 2-D Array for Matrix Multiplication and Transitive Closure", IEEE Transactions on Parallel and Distributed Systems, vol. 6, nr. 4, April 1995, pag. 351 - 362.
- [UCH 92] K. Uchimura, O. Saito, Y. Amemiya. "An 8G connections - per - second 54mW digital neural network chip with low - power chain - reaction architecture", Tech. Digest IEEE Inter. Solid - State Circuits Conference, San Francisco, Feb. 1992, pag. 134 - 135.
- [VAS 96a] Nikolaos Vassilas, Patrick Thiran, Paolo lenne. "On Modifications of Kohonen's Feature map Algorithm for an Efficient Parallel Implementation", ICNN 1996, pag. 932 - 937.
- [VAS 96b] S. Vassiliadis, S. Cotozana, K. Bertels. "2 - 1 Addition and Related Arithmetic Operations with Threshold Logic", IEEE Trans. on Computers, vol. 45, nr. 9, Sept. 1996, pag. 1062 - 1067.
- [VIN 92] J. M. Vincent. "Finite wordlength, integer arithmetic multi - layer perceptron modeling for hardware realization", Neural Networks for Vision, Speech and Natural Language, R. Linggard, D. J. Myers and C. Nightingale - editors, London, Chapman & Hall, 1992, pag. 293 - 311.
- [VIR 92] M. A. Viredaz, C. Lehmann, F. Blayo, P. lenne. "MANTRA: A multi - model neural - network computer", In Proc. of the 3rd International Workshop on VLSI for Neural Networks and Artificial Intelligence, Oxford, UK, September 1992.
- [VIR 93] M. A. Viredaz. "MANTRA I: An SIMD processor array for neural computation", In Proceedings of Euro - ARCH '93, Munich, Germany, October 1993.
- [WAI 89] A. Waibel, T. Hanazawa, G. Hinton, K. Shikano, K. Lang. "Phoneme recognition using time - delay neural networks", IEEE Trans on ASSp, nr. 37, mar. 1989, pag. 328 - 339.
- [WAT 93] T. Watanabe, K. Kimura, M. Aoki, T. Sakata, K. Ito. "A single 1,5 V digital chip for 10⁶ synapse neural network", IEEE Trans. Neural Networks, vol. 4, nr. 3, May 1993, pag. 387 - 393.
- [WAW 93] J. Wawrzynek, K. Asanovic, N. Morgan. "The design of a neuromicroprocessor", IEEE Trans. Neural Networks, vol. 4, nr. 3, May 1993, pag. 394 - 399.

- [WEI 89] M. Weinfield. "A fully digital integrated CMOS Hopfield network including the learning algorithm", VLSI for Artificial Intelligence, Delgado - Frias & Moore Eds., Kluwer Academic Publishers, 1989.
- [WER 88] P. J. Werbos. "Generalization of Back - Propagation with application to Recurrent Market Model", Neural Networks, vol. 1, 1988, pag. 339 - 356.
- [WID 60] B. Widrow, M. E. Hoff. "Adaptive switching circuits", 1960 IRE WESCON Convention Record, New York: IRE, 1960, pag. 96 - 104.
- [WID 85] B. Widrow, S. D. Stearns. "Adaptive Signal Processing", Prentice Hall Signal Processing Series, A. V. Oppenheim Ed., Englewood Cliff, New Jersey, 1985.
- [WRZ 96] A. Wrzyszczyk, D. Milford, E. L. Dagless. "A New Approach to Fixed - Coefficient Inner Product Computation Over Finite Rings", IEEE Trans. on Computers, vol. 45, nr. 12, Dec. 1996, pag. 1345 - 1355.
- [YAS 90] M. Yasunaga, N. Masuda, M. Yagyū, M. Asai, M. Yamada, A. Masaki. "Design, fabrication and evaluation of a 5-inch wafer scale neural network VLSI composed of 576 digital neurons", Proc. IEEE / INNS Inter. Joint Conf. Neural Networks, vol. 2, June 1990, pag. 527 - 535.
- [ZHA 94] D. Zhang, M. I. Elmasry. "High Performance Compressor Building Blocks for Digital N. N. Implementation", WCNN, vol. II, 1994, pag. 607 - 611.
- [ZHA 96] M. Zhang, S. Vassiliadis, J. G. Delgado - Frias. "Sigmoid Generators for Neural Computing Using Piecewise Approximations", IEEE Trans. on Computers, vol. 45, nr. 9, Sep. 1996, pag. 1045 - 1049.
- [ZHA 97] C. X. Zhang, D. A. Mlynski. "Mapping and Hierarchical Self - Organizing Neural Networks for VLSI Placement", IEEE Trans. on Neural Networks, vol. 8, nr. 2, Mar. 1997, pag. 299 - 314.
- [ZUR 92] J. M. Zurada. "Introduction to Artificial Neural Systems", West Publishing Company, St. Paul, 1992.

CURRICULUM VITAE

Nume : MIHU
Prenume : IOAN
Data nașterii : 14 Decembrie 1954
Locul nașterii : sat DEAL, jud. ALBA, ROMÂNIA
Starea civilă : căsătorit
Copii : 2
Adresa : str. LUDOȘ, Nr. 3, Bl. 51, Ap. 45, 2400, SIBIU,
ROMÂNIA, Tel. 481006, E-mail: mihuz@cs.sibiu.ro

Studii:

1961 – 1969, Școala Generală, sat DEAL, jud. ALBA
1969 – 1974, Liceul Industrial Energetic, Sibiu
1976 – 1981, Facultatea de Electrotehnică, specializarea Automatizări și Calculatoare,
Institutul Politehnic "Traian Vuia" Timișoara.
1992 Admis la doctorat în domeniul "Sisteme Automate"
1998 Absolvirea doctoraturii în "Sisteme Automate" (V. T. Timișoara)

Experiența profesională:

1981 – 1990, inginer, Întreprinderea pentru Întreținerea și Repararea Utilajelor de Calcul
București, Filiala Sibiu.
1990 – prezent, șef lucrări, Universitatea "Lucian Blaga" din Sibiu, Facultatea de Inginerie,
Catedra de Calculatoare și Automatică (Discipline: Calculatoare Numerice,
Structura Sistemelor de Calcul, Arhitecturi Paralele).

Activități manageriale:

1992 – 1993, secretar științific al Consiliului Facultății de Inginerie
1992 – prezent, membru al Consiliului Facultății de Inginerie
1996 – prezent, coordonator local Program TEMPUS S_JEP 11005/96

Activitate științifică și didactică:

- 14 lucrări științifice publicate
- 2 cărți publicate în domeniul arhitecturii sistemelor de calcul

Limbi de circulație cunoscute:

- engleza;
- franceza.

Domenii de interes:

- arhitecturi / procesări paralele
- rețele neuronale