

UNIVERSITATEA "POLITEHNICA"
TIMIȘOARA

BIBLIOTECA CENTRALĂ

Nr. Inv. 32.107

Dulap 16 Lit. 8

Facultatea "Politehnica" din Timișoara
Departamentul de Calculatoare

Nicolae Szirbik

CONTRIBUȚII LA METODELE DE INTEGRARE SIMBOLIC-
CONEXIONISTE CU APLICAȚII ÎN SISTEMELE DE
RAȚIONAMENT COMPOZIȚIONALE

conducător științific:
Prof.dr.ing. Ioan Jurca

-1998-

Mamei,

BIBLIOTECA CENTRALĂ
UNIVERSITATEA "POLITEHNICA"
TIMIȘOARA

UNIVERSITATEA "POLITEHNICĂ"
TIMIȘOARA
BIBLIOTECA CENTRALĂ
Nr. volum _____
Data: 66 LA C

CUPRINS

Partea I. Motivația acestei cercetări, suportul teoretic existent, soluțiile propuse

Introducere

11. Strategii și metodologii
12. Viabilitatea strategiei hibride
13. Structura tezei

1. Suportul teoretic al unui model neurosimbolic unificat

- 1.1. Modelul neuronului artificial, rețeaua feedforward
- 1.2. Modelul Towell
- 1.3. Concluzii

2. Structurile simbolice și conexioniste utilizate

- 2.1. Grafuri conceptuale
 - 2.1.1. Reguli de construire și utilizare
 - 2.1.2. Structurarea universului discursului
 - 2.1.3. Consistența logică a universului discursului
- 2.2. Rețele conexioniste recurente
 - 2.2.1. Modele conexioniste inspirate din fizica statistică
 - 2.2.2. Rețelele Hopfield
 - 2.2.3. Modelul Hopfield continuu
 - 2.2.4. Algoritm de învățare pentru un sistem neurodinamic
 - 2.2.5. Algoritm de propagare înapoi a erorilor pentru rețele recurente discrete

3. Realizarea legăturii între simbolic și conexionist

- 3.1. Procesul cognitiv privit ca și traiectorie
- 3.2. Translatarea informației simbolice în format conexionist
 - 3.2.1. Paradigma sub-simbolică
 - 3.2.2. Reprezentarea intrării în rețeaua neuronală aleasă
 - 3.2.3. Implementarea conexionistă a reprezentării tensoriale
- 3.3. Metoda propusă pentru translatarea grafurilor conceptuale
- 3.4. Concluzii după prima etapă

Partea II-a. Dezvoltarea modelului propus

4. Modelul hibrid translațional

- 4.1. Modelul cognitiv conexionist în buclă închisă
- 4.2. Principul divergenței
- 4.3. Renunțarea la uniformitate și folosirea unei structuri hibride

- 4.4. Antrenarea unei rețele neurodinamice cu grafuri conceptuale
 - 4.4.1. Dimensiunea rețelei folosite
 - 4.4.2. Timpul necesar antrenării și paralelizarea algoritmului
 - 4.4.3. Testarea generalizării
- 4.5. Fenomenul de bifurcație la antrenare - căi de rezolvare
 - 4.5.1. Analiza fenomenului și corecțiile posibile
 - 4.5.2. Influența formei tiparelor de antrenare
 - 4.5.3. O schemă pentru evitarea bifurcațiilor
 - 4.5.4. Concluzii
- 5. Modelul multimodular omogen
 - 5.1. Considerații dimensionale
 - 5.2. Comunicarea între module
 - 5.3. Alegerea vocabularelor modulelor, experimente
 - 5.3.1. Scopul unui astfel de sistem
 - 5.3.2. Stabilirea vocabularului
 - 5.3.3. Experimentele efectuate
 - 5.4. Concluzii
- 6. Modelul multimodular eterogen
 - 6.1. Modelul compozițional aplicat în inteligența artificială
 - 6.2. Schema pe două nivele
 - 6.2.1. Funcționarea părții simbolice a sistemului
 - 6.2.2. Funcționarea nivelului conexiunist al sistemului eterogen
 - 6.2.3. Posibile moduri de comunicare între cele două nivele
 - 6.3. Antrenarea modulelor nivelului conexiunist
 - 6.3.1. Vocabularul folosit la nivel conexiunist
 - 6.3.2. Traiectoriile de antrenare și alocarea lor la module
 - 6.3.3. Echilibrarea efortului de antrenare al rețelelor neuronale
 - 6.4. Experimente
 - 6.4.1. Faza sistemului minimal
 - 6.4.2. Faza sistemului extins
 - 6.4.3. Faza de urmărire a generalizărilor
 - 6.4.4. O partiționare bazată pe întregul vocabular al sistemului
 - 6.5. Relevanța experimentelor

Concluzii

- C1. Un sumar al rezultatelor
- C2. Contribuții originale
- C3. Posibile direcții de continuare a cercetării

Bibliografie

Notă autobiografică a autorului

Partea I
Motivația acestei cercetări,
suportul teoretic existent,
soluțiile propuse

Introducere

În timpul scurtei sale istorii, domeniul de cercetare numit inteligență artificială a fost arena unei dispute, uneori exagerate, între două arii de cercetare, doi *frères ennemis*, pe care îi putem denota "simbolicism" și "conexionism" [Hilario96]. Inteligența artificială își propune dezvoltarea unor modele care pot fi aplicate la construirea unor sisteme fizice capabile să manifeste însușiri care pot fi considerate inteligente. Pentru că nu există o definiție foarte clară a conceptului de inteligență, singurul procedeu care poate releva "inteligența" unui sistem artificial este testul Turing (1950), care implică o comparație cu o entitate inteligentă umană. Este important de observat că acest test se bazează pe analiza pur comportamentală a sistemului observat și nu pe o analiză formală a mecanismelor interne care guvernează funcționarea lui. Explicația disputei amintite poate fi și o înțelegere diferită a conceptului de inteligență.

Există două școli în domeniul inteligenței artificiale, care au păreri diferite în privința stabilirii "inteligenței" unui sistem. Prima, numită și "clasică", admite doar modele matematice care pot fi demonstrate că produc rezultate corecte în toate circumstanțele. Această școală, derivată din grupul de la Dartmouth (vezi istoria grupului prezentată în [Newell&Simon72]), formată din matematicieni cu o puternică pregătire formală, urmași ai lui Alan Turing și John Von Neumann, își propunea nici mai mult nici mai puțin, decât realizarea unor mașini inteligente (în sensul curent al conceptului de inteligență) care să poată întrece omul în performanțele sale cognitive. Acesta era deja parțial întrecut de mașinile de calcul ale vremii respective, care puteau rezolva anumite probleme, dar doar pe cele care puteau reduce metoda aplicată la un algoritm format dintr-un set de operații de bază. Dificultățile legate de complexitatea sistemelor formale necesare pentru reprezentarea unor probleme mai apropiate de "inteligență" și probleme legate de intensitatea computațională a algoritmilor exhaustivi de căutare folosiți, au făcut ca această abordare să nu aibă rezultate de succes, chiar dacă dezvoltarea sistemelor fizice de calcul (hardware) a fost deosebită în perioada scursă din 1956 până astăzi. Sunt însă de apreciat realizările colaterale ale acestui tip

de cercetare cum ar fi: programarea funcțională, programarea logică, meta-euristicile, logicile non-monotonice, modelele clasice pentru sisteme expert, modelele pentru învățarea automată (machine-learning). A doua școală este bazată pe ideea testului Turing și se numește "comportamentalistă" (behaviorist) [McCulloch88]. Ea stabilește că un sistem este inteligent doar din analiza exterioară a funcționării sistemului. Dacă prin testare (cât de exhaustivă posibil) sistemul dă rezultate corecte, și comportamentul său se apropie în anumite instanțe de cel inteligent uman, atunci sistemul poate fi acceptat, indiferent de mecanismele sale interioare. Se poate observa că lipsa demonstrării unei funcționări invariabil corecte este slăbiciunea acestui tip de abordare. Cele două abordări sunt oarecum extreme, una concentrându-se pe întrebarea CUM se poate simula inteligența, iar a doua pe CE anume trebuie sistemul să poată să facă.

După părerea mea, insuccesul în realizarea unui model matematic al sistemelor inteligente este în primul rând legat de lipsa unei definiții clare, matematice a inteligenței. Există o mulțime de definiții, dar niciuna nu este formală. Definiția mea preferată este următoarea [Clark89]:

Inteligența este tocmai acel *sumum* de însușiri umane al căror conținut nu este exprimat formal. Orice însușire umană, care poate fi exprimată algoritmic este rutinieră. Dacă o însușire este formalizată la un moment dat, ea părăsește grupul însușirilor care formează inteligența și devine rutină.

Această definiție statuează inteligența artificială ca fiind un concept elastic, care se modifică în timp, odată cu găsirea de către om, **prin metode creative, inteligente** a unor explicații pentru propriile sale însușiri. Odată cu aceste explicații, omul poate realiza și simularea acestor activități, care până în momentul respectiv au fost considerate inteligente. Și această definiție este discutabilă, pentru că am văzut cum mașina a învins omul într-un joc despre care era unanim acceptat că necesită "intelență". Conform definiției, șahul nu ar mai reprezenta o activitate inteligentă, lucru care în mod cert nu este adevărat. Dacă vor putea fi construite modele și mașini care să simuleze creativitatea, intuiția, dialogul inter-uman, ar însemna tot conform definiției că inteligența ca și concept ar fi un capitol epuizat, dat fiind faptul că toate manifestările ei sunt exprimabile formal. Dar trebuie să fim atenți, definiția are legătură cu conceptul de inteligență artificială. Până în perioada anilor 60, capacitățile de rezolvare a unor probleme simbolice erau considerate atribute de inteligență, dar odată cu găsirea metodelor matematice, acestea nu au mai fost considerate ca fiind activități "inteligente" **pentru calculator**. Pentru om însă, aceste capacități vor rămâne întotdeauna un atribut al inteligenței, așa cum și calitățile de șahist vor indica calificativul de "om inteligent". De remarcat este că metodele matematice folosite de programatorii mașinilor de jucat șah sunt bazate și pe o serie de metode euristice, mai ales la estimarea valorii poziției, lucru care nu duce la găsirea celei mai avantajoase mutări, sau per ansamblul partidei, la cea mai scurtă victorie. Din acest motiv, este exagerat să afirmăm că jocul de șah este complet formalizat. Pur și simplu, forța brută a mașinii a devenit comparabilă cu forța complet diferită ca natură a minții umane,

despre care trebuie să amintim că ocupă un spațiu de aproximativ 1500 cm cubi și consumă în jur de 20 wați. Chiar dacă mașinile vor avea în viitor dimensiuni și consumuri mai mici, mai rămâne de demonstrat formal, că algoritmul pe care acestea le folosesc, duc la cea mai bună soluție. Interesant este faptul că discuțiile critice care se duc despre conceptul de inteligență artificială eșuează totdeauna în spațiul filozofic, existând multe păreri diferite în ceea ce privește conținutul acestui domeniu.

Revenind la cele două școli, se poate observa că și modelele diferă. În timp ce școala clasică folosea în special metodele logicii matematice, ale sistemelor formale și ale limbajelor formale, școala behavioristă folosea mai mult metode ad-hoc, rezultate din observații făcute asupra comportamentului uman și a sistemului biologic (neuro-fiziologic) implicat în activitățile cognitive. Ajungând la nivelul anilor '60, s-a putut observa o cristalizare a metodelor folosite, școala clasică adoptând procesarea simbolică bazată pe formalismul matematic iar școala behavioristă folosind cu precădere un instrument matematic numit "rețea neuronală artificială" [Hecht90] [Hertz&al91] [Haykin94]. Din păcate, urmașele cele două școli, devenite între timp arii distincte de cercetare denumite "simbolicism" și "conexionism", au avut o evoluție interesantă din punctul de vedere al concurenței științifice. Merită a fi menționată cartea "The Perceptrons" [Minsky&Papert69] care stabilea formal și irefutabil limitările rețelei neuronale feedforward cu un singur strat. Această lucrare a oprit finanțarea cercetărilor din domeniu pentru aproape 20 de ani, deși paradoxal, în același an Bryson și Ho descopereau (întâmplător și în alt domeniu) algoritmul care se numește astăzi "learning by error backpropagation", al cărui redescoperire în '85 (separat de către Parker, Le Cun, Rumelhart&al - plus Werbos în '74) a făcut ca rețelele neuronale să fie din nou o arie de cercetare cu finanțare bogată.

O altă "lovitură" din direcția "simbolicistă" este lucrarea lui Fodor și Pylyshin din 1988 [Fodor&Pylyshyn88] care argumentează că metodele conexioniste nu pot să fie un suport al unei teorii care să modeleze capacitățile cognitive, pentru că nu pot să țină cont de complexitatea combinatorială a structurilor sintactice și semantice ale reprezentărilor mentale. În cel mai bun caz, în opinia celor doi autori, conexionismul este doar o nouă metodă de implementare a acestor reprezentări și a procesării lor, doar o alternativă aproximativă și nesigură la metodele simbolice folosite de școala clasică. Contralovituri există și din direcția grupului de cercetători conexioniști, de remarcat fiind una cât se poate de "extremistă" a lui Pinker și Prince, tot în 1988 [Pinker&Prince88], care susține că teoriile simbolice nu sunt altceva decât aproximări grosolane și foarte incomplete ale fenomenelor care au loc în mintea umană. Există și o a treia părere, observabilă în ambele tabere, că cele două abordări sunt potrivite fiecare pentru un anumit set distinct de probleme și că nu este necesar amestecul lor (direcție numită eliminativism [Smolensky88]).

Între aceste poziții, oarecum radicale, au apărut la sfârșitul anilor '80 [Reeke&Edelman88] [Touretzky&Hinton88] [Hendler89] [Goodman&al89], lucrări care se plasau la interfața dintre procesarea simbolică și conexionistă. Originea lor este legată de încercările de a realiza anumite aspecte ale procesării simbolice cu metode

specific conexioniste și de a realiza legătura între rețele neuronale cu rol de senzor sau rol motor cu un nivel de decizie superior, realizat prin metode simbolice (în special în domeniul conducerii roboților). Această nouă arie de cercetare a fost denumită *integrare neuro-simbolică* (neuro-symbolic integration) [Sun91]. Se pot identifica și aici două direcții de lucru: una *unificată* și una *hibridă*. Prima încearcă să implementeze toate operațiile dintr-un sistem inteligent prin rețele neuronale, pe când a doua încearcă să combine cele două paradigme, punând de exemplu în același sistem, rețele neuronale și grafuri de decizie. Merită a fi remarcat că Marvin Minsky, realizatorul lui SNERT (1954), primul "calculator neuronal" și adversarul ideilor conexioniste în perioada '65-'75, este unul din cei mai importanți susținători ai cercetării în domeniul sistemelor hibride [Minsky86] [Minsky91].

Această teză arată calea străbătută pentru identificarea unei metodologii de realizare a unor sisteme de raționament automat hibride, simbolic-conexioniste.

II. Strategii și metodologii

Integrarea simbolic-conexionistă este motivată de natura bivalentă a capacităților cognitive ale minții umane. S-a observat că procesele mentale pot fi clasificate în diverse moduri, toate aceste procese luate în considerare având un rol important în obținerea unui comportament inteligent. Un mod de clasificare este de a grupa procesele în două tipuri: *sinetice* și *analitice*. Clasificarea nu se poate face în mod rigid, majoritatea proceselor având și o componentă analitică și una sintetică, într-o măsură mai mică sau mai mare. Procesele pot fi dispuse pe o axă limitată la capete, care să conțină la o extremitate procesele pur sinetice, iar la cealaltă extremitate pe cele pur analitice.

Procesele sintetice sunt de regulă asociate cu operații de nivel scăzut, care nu pot fi descompuse în operații mai simple, cu caracteristică de non-conștiență (omul nu le execută în mod voit). Tipice sunt operațiile de percepție, cum sunt vederea și recunoașterea vorbirii (a cuvintelor). Ele sunt folosite de procese superioare, cum sunt identificarea unor obiecte din scena văzută, sau extragerea unui conținut de idei din fluxul verbal ascultat. Ceea ce este cel mai caracteristic, cunoștințele implicate în acest tip de operații sunt non-propoziționale și este aproape imposibil să le exprimăm formal, pentru că ele sunt mai degrabă ilustrative.

În contrast, procesele analitice sunt asociate cu operații conștiente, care pot fi descompuse în operații mai simple. Exemple: rezolvarea problemelor, luarea unor decizii în mod conștient, generarea unui text, dialogul (la un nivel propozițional). Cunoștințele implicate folosite au o exprimare propozițională, legată de un limbaj, care poate fi limbajul natural, sau o serie de limbaje artificiale, mai mult sau mai puțin formale. [Newell90] oferă o clasificare graduală (pe scara complexității analitice) a acestor procese, folosind ca factor de discriminare, timpul necesar pentru efectuarea acestor operații. Ceea ce rezultă foarte clar este că majoritatea proceselor nu sunt pur

sintetice sau pur analitice. Un exemplu tipic este calculul mental: 4 înmulțit cu 3 este pur sintetic (pentru cine știe pe de rost tabla lui Pitagora), pe când 34 înmulțit cu 13 este analitic în mare măsură (doar operațiile obținute prin secvențierea și descompunerea necesară sunt sintetice). Dar cum sunt 12 înmulțit cu 7 sau 5 înmulțit cu 15? Posibil ca undeva între cele două extreme.

Un caz interesant este învățarea cântatului la un instrument. Un începător la pian trebuie să fie conștient de fiecare mișcare pe care o face, cunoștințele sale sunt explicite, pe când un pianist mai experimentat este mai puțin conștient, concentrându-se mai mult asupra întregului și fluenței, pe când pentru profesioniști, cântatul în sine este un automatism, ei fiind mai mult conștienți de nuanțele și tipul de interpretare pe care o doresc. Se poate observa cum experiența contribuie la deplasarea aceleiași operații executate pe axa imaginată, de la analitic la sintetic. Sunt însă și operații care implică mai multe procese mentale în același timp, procese care pot fi plasate pe porțiuni foarte diferite ale axei. Operația mentală de a recunoaște un membru al familiei după voce este pur sintetică și fără nevoia de conștientizare. Operația mentală de a decide dacă iau umbrela colegului de birou implică o analiză conștientă a implicațiilor și consecințelor posibile. Aceste două operații sunt la extremele axei, dar să ne imaginăm ce fel de procese sunt implicate în timpul în care cineva conduce o mașină în trafic intens (vezi figura).

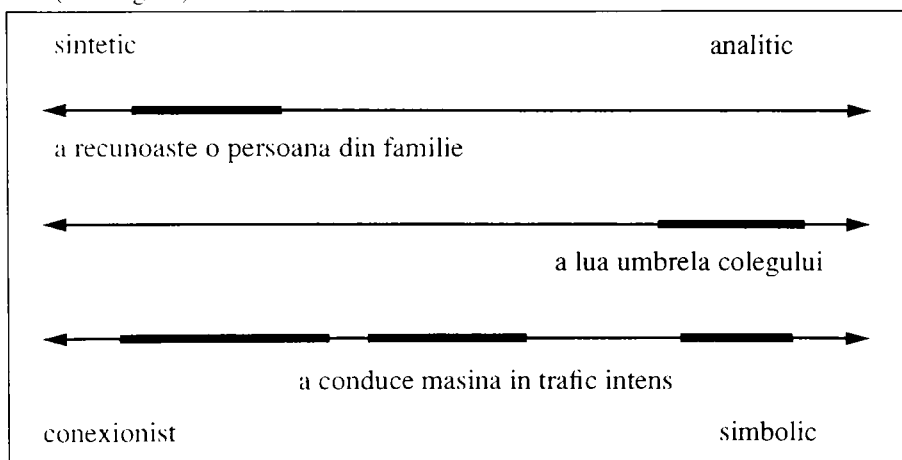


Figura I-1.

Este cât se poate de limpede că simularea proceselor mentale sintetice se poate face cu succes mai ales cu metode conexioniste, iar cele analitice cu metode simbolice. Se poate imagina o axă paralelă cu cea a proceselor mentale, dedicată metodelor de realizarea artificială folosite. La extremitatea sintetică să va oglindi extremitatea algoritmilor pur numerici ai rețelelor neuronale, iar la extremitatea analitică, va apare o analogie cu procesarea simbolică. Aici trebuie făcută observația crucială: **modelele oferite de domeniul inteligenței artificiale sunt la cele două extremități ale axei**

metodelor folosite. Nu există la ora actuală nici un fel de modele bine aprofundate, de succes, care să se afle într-un punct al axei isomorf cu punctul unde procesele mentale sunt în aceeași măsură sintetice și analitice (deși astfel de procese există, mai mult, ele reprezentând majoritatea operațiilor "inteligente" efectuate de om [Gutknecht&al91]).

Un exemplu al unui astfel de proces mental este stabilirea unui diagnostic medical de către un specialist. Identificarea unui model matematic corespunzător este necesară pentru transpunerea acestei operații într-un sistem expert. Se pot folosi mai multe metode, de exemplu inferența Bayesiană [Todd91] care are dezavantajul că cunoștințele propoziționale nu pot fi reprezentate, sau logica fuzzy, care are dezavantajul că sistemul nu poate să fie antrenat prin consultarea unei populații statistice de cazuri medicale (vezi sistemele fuzzy CaDiag și FAULT [Leung&Lam89]). Folosirea modelelor computaționale conexioniste se poate face pe două căi. Fie se găndește tot sistemul ca fiind conexionist, în acest caz reprezentarea și procesarea cunoștințelor propoziționale facându-se prin metode conexioniste, fie se implementează o parte a sistemului prin metode simbolice, o altă parte prin metode conexioniste, accentul căzând în acest caz pe comunicarea dintre cele două părți. Putem vorbi de o strategie unificată sau de una hibridă.

Strategia unificată se bazează pe ideea că în astfel de sisteme inteligente, nu este nevoie de structuri și reprezentări simbolice precum și de procesare simbolică explicită. Există o explicație pentru această presupunere, cei care o propun sunt de regulă neurobiologi, foarte îndepărtați de matematica simbolică și ei caută de regulă modele care au o identificare cât mai apropiată de sistemele neuronale naturale (mai ales că ei studiază structuri primitive relativ la operațiile analitice - orientarea șobolanilor de exemplu). Există bineînțeles și argumentul că procesarea cea mai analitică posibil a simbolurilor la nivelul minții umane se realizează tot prin rețele neuronale. O analiză mai detaliată a acestei probleme este prezentată în subcapitolul 4.3 al tezei.

Un caz special al strategiei unificate este **procesarea conexionistă a simbolurilor** (CSP - Connexionist Symbolic Processing). Aceasta nu trebuie confundată cu **procesarea neuronală a simbolurilor** (NSP - Neuronal Symbolic Processing) [Edelman92], a cărui obiectiv este identificarea unor modele de neuroni cât mai apropiați de neuronii naturali și care prin funcționarea lor colectivă pot emula operații de tip analitic, cu caracter simbolic. CSP este o cercetare care încearcă realizarea unor operații de regulă simbolice (raționamentul automat, planificarea, analiza sintactică) prin metode conexioniste unitare (un singur tip de rețea neuronală, pe cât posibil mai simplă). Exemple de modele reușite sunt: reprezentările și clasificarea conceptelor [Feldman&Ballard82] [Shastri88], procesarea unor liste prin primitive gen LISP cu sistemul BOLTZCONS [Touretzky&Hinton88], extragerea unor proprietăți prin reprezentarea ierarhică a cunoștințelor [Sun91]. S-au obținut unele rezultate și în domeniul inferenței logice [Hölldobler&Kurfess91] [Towell91], analiza limbajului natural [Bookman87] [Dyer91] și implementarea conexionistă a sistemelor expert [Gallant88] [Dragomirescu&al93b].

12. Viabilitatea strategiei hibride

Rostul implementatorului de sisteme inteligente (sisteme expert, sistem de conducere a roboților, sisteme de dialog mașină-om) este de a realiza aceste sisteme "acum și aici". Pentru a fi capabil de acest lucru, este necesar ca el să aibă la dispoziție modele și metodologii rezultate din cercetarea în domeniul inteligenței artificiale. Toate sistemele înșiruite mai sus necesită modele cu un grad de complexitate destul de ridicat, care face imposibilă realizarea lor prin utilizarea unei singure paradigme computaționale. Pentru a cuprinde diversele operații care trebuie realizate, este necesară utilizarea simultană a mai multor modele, atât din spectrul simbolic, cât și din cel conexionist. Complexitatea, dar mai ales diversitatea proceselor cognitive implicate chiar în activități aparent simple pentru om, a condus la ideea că un singur model unificat nu este suficient pentru atingerea unor scopuri imediate, lucru accentuat și în [Minsky91]: "*To solve really hard problems, we'll have to use several different models and representations*". Există firește și scopuri cu termen lung. Cel mai important este firește emularea prin mașină a principalelor capacităților cognitive umane. Probabil că nu va conta în acel moment dacă acest lucru a fost obținut prin sisteme hibride. La fel de importantă este și realizarea unor modele plauzibile din punct de vedere biologic, dar mai importantă este plauzibilitatea psihologică, legată mai mult de comportamentul sistemului și mai puțin de detaliile de implementare. Contraargumentul invocat de direcția de cercetare unificată este că succesul sistemelor de calcul existente este bazat pe arhitectura Von Neumann, care are un caracter puternic unificat. Dar pentru a realiza **acum** sisteme viabile, este necesar să se folosească din arsenalul existent, metodele cele mai avantajoase din punct de vedere computațional. Folosirea unei strategii unificate duce de regulă la necesitatea unor resurse computaționale (memorie, timp de calcul) exagerate. S-a observat că folosind o strategie hibridă, pentru aceeași problemă, necesarul de resurse este apreciabil mai mic [Miikulainen94].

Cercetarea în direcția integrării neuro-simbolice (perioada 1986-1996) a încercat să ofere **modele** pentru realizarea unor sisteme economic viabile. Pentru că nici unul dintre acestea nu are la baza o teorie demonstrată a corectitudinii funcționării, doar un indice de performanță procentual, este puțin exagerată terminologia de model (deși ea este folosită în majoritatea lucrărilor de specialitate). Eu consider însă că termenul de *metodologie* este mai potrivit. În momentul actual este imposibil de oferit un model corect funcțional, care să aibă în componența lui rețele neuronale, pentru că nu s-a demonstrat că aceste rețele au o funcționare corectă (indiferent de natura informației de intrare). Bineînțeles, în această discuție, contează și ce înțelegem prin *corectitudine*. Dacă o rețea recunoaște toate tiparele cu care a fost antrenată și în plus recunoaște un set de probă alcătuit dintr-o mulțime mai mare de tipare decât baza de antrenare (fără a conține însă tipare foarte diferite de cele învățate), mulți cercetători afirmă că rețeaua funcționează corect, mai ales dacă poate rejecta tipare foarte diferite de cele antrenate. Este doar un punct de vedere, pentru că un analist mai riguros ar clasa ca funcționare incorectă, orice rejectare sau interpretare greșită a oricărui tipar, indiferent de gradul său de similitudine cu cele învățate. Din acest motiv, atunci când în teză este folosit

cuvântul model, cititorul riguros îl poate interpreta ca **metodologie**.

Metodologia propusă în teză este rezultatul unor idei personale, dar în același timp ea reflectă și observațiile făcute asupra unor sisteme propuse de alții. Există o serie de metodologii propuse pentru sisteme hibride și am încercat dealungul tezei să compar ideile și soluțiile mele cu acestea. Pentru a-mi plasa metodologia în contextul altor propuneri, este necesară o sistematizare a acestora. Modelele hibride pot fi grupate în două mari familii: modele hibride translaționale și modele hibride funcționale [Hilario86].

- **modelele translaționale** sunt o clasă intermediară între modelele strategiei unificate și modelele funcționale. Ca și modelele unificate (numite și modele uniforme) cele translaționale se bazează pe rețele neuronale pe post de procesoare de informație, dar datele de start și rezultatele sunt structuri simbolice. De cele mai multe ori acestea sunt reguli propoziționale [Fu&Fu90] [Towell91] (dacă este privit ca un generator de reguli - vezi discuția despre acest model din capitolul 1) [Kunicky&al92]. În alte cazuri [Crucianu&Memmi92] [Sun92] [Giles93] [Miikulainen96] [Wermter&Weber97], limbajul natural (sau mai bine zis o componentă simplificată a lui) este cel pe baza căruia se construiesc structurile simbolice. Acestea nu sunt procesate ca atare, prin metode simbolice, ci sunt traduse fie în rețele neuronale, care își modifică parametrii sau chiar structura prin adăugarea unor noi componente prin antrenare, fie în vectori de activare neuronală, care sunt intrări pentru o rețea [Medsker94]. Rezultatul se obține prin traducerea rețelei finale sau a ieșirii ei într-o structură simbolică. Merită remarcat faptul că sistemele construite pe baza acestui model, de exemplu CONSYDERR [Sun94], NP [Giles95] și SCREEN [Wermter&Weber97] au performanțe foarte reduse, mai ales datorită limitărilor dimensionale ale rețelelor folosite. Unul din sistemele care au rezultate bune (în faza experimentală) este DISCERN [Miikulainen96], dar în acest caz, cu toate că toate procesările se fac prin rețele neuronale, există mai multe elemente de procesare, cu topologie și funcționalitate diferită. Primele mele propuneri de model [Szirbik&al95] [Szirbik95b] [Szirbik&Babii96] se pot înscrie în această clasă a modelelor translaționale.

- **modelele hibride funcționale** încorporează componente *complete* simbolice și conexiuniste. Acestea sunt caracterizate pe de o parte de reprezentările (structurile de date) folosite și pe de altă parte de schemele de procesare. Pe lângă diferite tipuri de rețele neuronale putem avea: interpretoare de reguli, analizoare sintactice, automate de raționament și demonstratoare de teoreme. Aceste sisteme pot fi diferențiate privind dealungul mai multor dimensiuni imaginare, cum sunt:

- clasa de aplicații pentru care se potrivește modelul,
- gradul de generalitate,
- stabilirea componentei dominante (conexiunistă sau simbolică), numită și *mod de integrare* [Minsky91],
- gradul de integrare între componentele simbolice și cele conexiuniste.

Cele mai importante aspecte pentru mine au fost ultimele două. O taxonomie bazată pe gradul de integrare oferă două variante: prin cuplare puternică (tight coupling) și cuplare slabă (loose coupling). În modelele cu cuplare slabă, interacțiunea între cele două tipuri de componente este strict localizată în spațiu și timp. Informația se transferă prin intermediul unui intermediar (agent, supervisor sau un "blackboard") iar interacțiunea este inițiată întotdeauna explicit, fie pe baza unui protocol, fie datorită unui agent de control exterior. În sistemele cu cuplare puternică, datele nu mai sunt transferate de la componentele simbolice la cele conexioniste și invers, ci sunt partajate prin structuri comune de date. Astfel, o modificare făcută de una din componente are repercursiuni imediate asupra celeilalte componente, fără să existe necesitatea inițierii unei interacțiuni explicite. Unul din sistemele propuse este SEMMARK [Hendler89], unde unități ale unei rețele neuronale aveau legături directe cu acțiunile semantice asociate unei analize sintactice hibride. O alta propunere este "filtrul" neurosimbolic, care asigură legătura reprezentărilor conexioniste distribuite cu structurile simbolice din SYNTHESIS [Giacometti92]. Astfel de sisteme puternic cuplate (integrate) nu au avut rezultate deosebite, explicația fiind incompatibilitatea naturală a reprezentărilor simbolice și conexioniste. Se poate însă observa o revenire a lor [Fahlman97] [Farhat97] [Shavlik97], deși nu se poate vorbi încă de rezultate utilizabile în sistemele viabile economic (în plus, toate propunerile fiind legate mai ales de utilizarea rețelelor nesupervizate și mai puțin revenirea la sistemele puternic cuplate).

Anii 92-96 au fost caracterizați mai ales prin dezvoltarea modelelor cu cuplare slabă. Rezultate încurajatoare au fost obținute mai ales în domeniul medical ca FIBRIL [Jabri&al92], REASONER [Piramuthu&Shaw94] sau NST-EXPERT [Alonso&al95], dar și în alte domenii, de menționat fiind controlul roboților (pe larg, subiectul este tratat în [Leonard&Durrant94]). Unul din sistemele folosite efectiv, chiar pe scară largă este WASTEWATER [Krovvidy&Wee92]. Modelele mai recente, cum sunt DUAL [Kokinov96], NESSY3L [Orsier&Labbi96] sau GENUES (Generic Neuro-Expert System) [Khosla&Dillon96], propun noi arhitecturi, făcând o analiză prealabilă, puternic experimentală și mai puțin teoretică a ideilor propuse. Cu toate că au existat și încercări notabile de formalizare a integrării neurosimbolice [Smolensky86,88,90,92], complexitatea covârșitoare a problemei a făcut ca domeniul să fie unul eminentamente experimental. Există critici serioase relativ la acest tip de abordare [Touretzky90], [Smolensky92], la un moment dat folosindu-se chiar epitetul de "alchimie" [Crevier93], dar există și argumentări pro care pleacă de la imposibilitatea unei alternative [Honavar94], [Alexandre96].

Modul de integrare sau *schema ierarhică* se referă la modul în care componentele simbolice și conexioniste sunt configurate relativ unele la altele și relativ la întregul sistem. Patru scheme pot fi identificate:

- procesare secvențială,
- sub-procesare,
- meta-procesare,
- co-procesare.

În cazul procesării secvențiale (numită și chainprocessing), una din cele două

componente are rolul principal (fie cea conexionistă fie cea simbolică) și cealaltă se ocupă de operațiile de pre/post procesare. Schema este folosită mai ales de sistemele expert medicale care au și intrare analogică [Hayes&al92], [Alonso95] (ambele sisteme referite având componenta principală simbolică - un caz invers fiind WASTEWATER [Krovvidy&Wee92]). În cazul sub-procesării, una din componente este interioară celeilalte, care acționează ca și componentă principală. De obicei componenta simbolică este cea principală, ca în INNATE/QUALMS [Becraft&al91], unde un generator neuronal de situații candidate este apelat în mod repetat de procesorul simbolic. În domeniul roboticii, schema cu meta-procesare este cea mai folosită. Un nivel simbolic de meta-procesare care monitorizează mai multe componente neuronale este realizat în RSAA (Robotic Skill Acquisition Architecture) [Handelman&al89] și COOHYS [Gutknecht&al91].

Cea mai interesantă și mai plauzibilă psihologic și biologic este co-procesarea. Componenta simbolică și cea neuronală sunt parteneri egali, ambele pot avea interacțiuni directe cu exteriorul, ambele pot transmite și recepționa informație de la cealaltă componentă. Există două posibilități: fie cele două funcționează fără supervizare (FIBRIL [Jabri&al91]), fie există un al treilea nivel de metaprocăsare care monitorizează cooperarea celor două componente ca în NESSY3L [Orsier&Labbi96]. O variantă pe care am încercat-o eu a fost schema cu doar două nivele [Szirbik96] [Szirbik97], unde nivelul simbolic are și rolul implicit de control.

La nivelul componentei conexioniste, este important de precizat ce fel de metode se folosesc pentru procesare. Putem avea rețele nesupervizate sau supervizate. Din punct de vedere al dinamicii rețelelor, putem avea rețele incrementale, care se antrenează în timpul utilizării sistemului, sau rețele statice, care sunt antrenate doar în timpul construirii lui. Se pot folosi rețele feedforward sau recurente, cu valori continue sau discrete, sau combinații ale acestora. În cazul sistemelor expert, se folosesc rețele feedforward supervizate (sau nesupervizate în cazul INNATE/QUALMS [Becraft&al91]) cu valori de intrare reale și ieșiri binare. Pentru aplicațiile în robotică, se folosesc mai mult rețele recurente dinamice.

Ca ultim aspect taxonomic, trebuie menționate aria de aplicabilitate și generalitate a modelului. Deși primele sisteme expert hibride dezvoltate de mine au fost în domeniul medical, mai exact în domeniul bolilor infecțioase [Dragomirescu&al93a] [Dragomirescu&al93b] [Dragomirescu&al94], căutarea unui model nou, m-a făcut să aleg un domeniu cu un grad mai mare de generalitate, și anume raționamentul automat, fără a aplica imediat metodologia la o problemă specifică. Am încercat diverse modele de raționament automat cu rețele feedforward [Szirbik&al93] [Szirbik93], aplicate apoi la sistemele expert medicale.

Modelul propus de mine în finalul acestei teze, poate fi clasificat ca fiind un **model hibrid funcțional, cu cuplare slabă între componentele simbolice și cele conexioniste, cu co-procesare și meta-procesare deopotrivă, care folosește rețele recurente supervizate** (iar într-o variantă îmbunătățită folosește și rețele multistrat).

Această combinație este oarecum asemănătoare cu cea propusă de Fahlman [Fahlman97], care are însă o cuplare mai puternică și rețele recurente nesupervizate (rezultatele fiind însă neconcludente până în prezent).

13. Structura tezei

Majoritatea rezultatelor prezentate în teză au fost publicate în volumele unor conferințe sau în periodice de specialitate. Există și o serie de lucrări publicate care sunt doar referite dar nu și prezentate pe larg în teză (perioada 1991-94). Ele se referă la cercetarea din domeniul sistemelor expert pentru studierea procesului infecțios. Este prezentată însă metodologia folosită în aceste sisteme expert [Szirbik&al93] [Szirbik93] și anume o variantă specifică a modelului Towell [Towell91]. Deasemenea, rezultatele din 1994 și 1995 sunt prezentate la un grad mult mai mare de detaliu decât în lucrările publicate pe această temă [Szirbik&al95] [Szirbik95a] [Szirbik95b]. Metodologia finală propusă a fost prezentată în [Szirbik96] și [Szirbik97], dar în teză se face în plus o prezentare exhaustivă a rezultatelor experimentale și a interpretării lor.

Ca și confirmare a viabilității metodologiei propuse și a aplicabilității ei, invitația de a discuta o serie de rezultate și idei la US Air Force Wright Laboratory (Dayton, Ohio) și la Neuro-Lab al Carnegie-Mellon University (Pittsburgh, Pennsylvania) au condus la două seminarii cu colective de cercetători din aceste instituții (aprilie 1997). Rezultatele prezentate în Statele Unite și interpretările specialiștilor de acolo sunt prezentate în teză.

- Capitolul 1: Prezintă un model de raționament implementat cu ajutorul rețelelor neuronale feedforward, propus de Towell în 1991. Se analizează performanțele, scăderile și limitările acestui tip de procesare. Se introduc o serie de termeni care vor fi pe larg utilizați în teză.

- Capitolul 2: Prezintă cele două instrumente de lucru folosite, sistemul neurodinamic discret (instrument pur conexiionist) și grafurile conceptuale (instrument de reprezentare pur simbolică a cunoștințelor). Se insistă în special pe unele particularități care au avut o mare importanță relativ la experimentele făcute cu modelul propus.

- Capitolul 3: Prezintă modul în care am implementat transformarea structurilor simbolice în vectori de activare conexiionști. Sunt discutate și o serie de considerații dimensionale.

- Capitolul 4: Prezintă prima arhitectură implementată și testată, modelul monomodular hibrid translațional. Este prezentată și o posibilă aplicație pentru planificare robustă. Sunt pe larg prezentate rezultatele experimentale care au arătat puterea de generalizare a unui astfel de sistem. Problema bifurcațiilor la antrenare este explicată în detaliu, precum și o soluție de evitare a ei. **Un fenomen descoperit**

experimental și demonstrat formal, numit de mine principiul divergenței este deasemenea prezentat.

- Capitolul 5: Prezintă modul în care au fost dezvoltate două sisteme de test, bazate pe un model distribuit, numit modelul (translațional hibrid) multimodular omogen, precum și limitările sale.

- Capitolul 6: Modelul final, denumit modelul multimodular hibrid eterogen (sau modelul pe două nivele) este prezentat ca și metodologie de dezvoltare pentru sisteme cognitive. Sunt prezentate experimentele care demonstrează puterea de generalizare a componentei conexiunilor și capacitatea acestora pentru anticiparea unor situații neantrenate sau implementate în sistem.

- Un ultim capitol, este dedicat concluziilor, contribuțiilor originale ale autorului tezei și posibilelor linii de continuare ale acestei cercetări.

1. Suportul teoretic al unui model neurosimbolic unificat

Modelul neuronilor artificiali, propus de McCulloch și Pitts în 1943, îmbunătățit de Hebb și rafinat de Rosenblatt în anii '50 (o prezentare detaliată a istoriei domeniului se poate găsi în [Hecht90]), este pe de o parte foarte depărtat de modelul neuronilor naturali (o excelentă prezentare "inginerescă" a lor în [DeFelipe97]) și este foarte diferit ca și natură matematică de operațiile folosite în logică și sistemele formale. Făcând această constatare, este normal să ne întrebăm dacă este posibilă realizarea unor operații simbolice cu ajutorul rețelelor neuronale formate din neuroni artificiali. Că se pot realiza operații simbolice cu neuroni naturali, este mai mult decât evident, din moment ce omul poate să facă operații simbolice la nivel mental. Întrebarea este dacă rețele neuronale artificiale, formate dintr-un număr mult mai mic de neuroni decât cei disponibili în creierul uman pot procesa simboluri și mai ales care sunt limitările impuse de modelul artificial. Există și modele matematice pentru neuroni artificiali mai apropiate de schemele neurofiziologice ale neuronilor naturali (există foarte multe tipuri de celule nervoase, mai mult sau mai puțin specializate), dar acestea sunt încă în fază experimentală și au dezavantaje: complexitate teoretică mare și un necesar de resurse de implementare ridicat. Din acest motiv, trebuie să utilizăm vechiul model, foarte simplu și ușor de implementat. Din punct de vedere teoretic, este destul de bine știut ce se poate face și ce nu se poate face cu astfel de rețele [Hertz&al91], [Haykin94].

În continuare, se prezintă succint un model foarte simplu de procesare simbolică cu rețele neuronale (unificat), și se introduc câteva noțiuni de bază, care sunt folosite pe larg în teză.

1.1. Modelul neuronului artificial, rețeaua feedforward

Neuronul artificial este un procesor foarte simplu de informație, care primește la intrare un set de valori, grupate de regulă într-un vector, și generează o singură ieșire. Valorile pot fi întregi, reale sau grupuri de valori de tip matrice de diferite

tipuri. De regulă, întrările sunt numere reale. Întotdeauna intrările sunt însumate, obținându-se *nivelul intern de activare* (notat în teză cu h). Aspectul esențial care definește un neuron și îl deosebește de ceilalți cu care se află în rețea este *vectorul de ponderi* w (de aceeași lungime ca și vectorul de intrare pentru neuronul respectiv). Din acest motiv, în loc de o simplă sumă a intrărilor se folosește ca nivel intern de activare, produsul scalar dintre cei doi vectori. Asupra acestei mărimi se aplică o funcție de transfer (notată în teză cu g), iar valoarea acestei funcții reprezintă ieșirea neuronului matematic, numită și *nivel de activare* al neuronului. Dacă funcția g este funcția identitate sau o simpla multiplicare, atunci neuronul se numește *liniar*. Dacă notăm cu x vectorul de intrare și cu o ieșirea, ecuația care descrie un asemenea neuron liniar este următoarea (γ fiind factorul de multiplicare):

$$o = \gamma \cdot \sum_{i=1}^n w_i \cdot x_i \quad (1-1)$$

Acest tip de neuron nu are prea multe aplicații [Smolensky86]. Pentru neuroni *neliniari* se folosesc alte funcții de transfer, de tip treaptă, continue sau discontinue. Funcția pe care am folosit-o cu predilecție în teză este tangenta hiperbolică, care are domeniul de definiție în \mathbf{R} , și ia valori în mulțimea reală $(-1,1)$. Ar fi mai simplu să folosim funcții discontinue, de tip treaptă:

$$o = \begin{cases} 1 & \text{daca } h \geq T \\ 0 & \text{daca } h < T \end{cases} \quad (h = \sum w_i x_i) \quad (1-2)$$

dar există considerente matematice care impun folosirea funcțiilor continue (derivabilitatea). Funcția treaptă (care are valori în $\{0,1\}$) este însă mai apropiată de ideea că un neuron are o funcționare asemănătoare cu a unui bistabil, neuronul având doar două stări posibile: activ și non-activ. Dar chiar folosind ca funcție de transfer tangenta hiperbolică, se poate suplimentar transforma ieșirea dintr-o valoare reală într-o valoare binară:

$$o' = \begin{cases} 1 & \text{daca } o \in [0.5, 1) \\ 0 & \text{daca } o \in (-1, -0.5] \end{cases} \quad (1-3)$$

și se poate observa că apare posibilitatea unei a treia stări (incerte, atunci când $o \in (-0.5, 0.5)$). Există însă tehnici pentru a evita apariția acestei stări nedorite.

O astfel de stare binară se poate asocia valorilor de TRUE și FALSE din logică. Unui neuron i se poate asocia un simbol (cuvânt). Dacă avem mai mulți astfel de neuroni cu simboluri asociate într-un sistem neuronal (rețea), ieșirea lor poate fi interpretată ca valoarea de adevăr a simbolurilor asociate. Sunt însă și situații când se dorește valoarea exactă a funcției *tanh*, când folosim de exemplu un sistem logic multivalent, sau factori de certitudine [Shortliffe76]. De exemplu, în sistemele expert dezvoltate de mine și de colegii mei în colaborare cu medici [Dragomirescu&al94], am folosit o variantă de logică cuadrivalorică, ale cărei valori de adevăr se obțineau printr-

o transformare simplă din valorile de ieșire ale neuronilor (care aveau funcția de transfer *tanh*). Modelul care va fi prezentat în continuare (numit și model Towell [Towell91]) lucrează doar cu două valori de adevăr.

Foarte rar se întâmplă să existe aplicații pentru un singur neuron. Când sunt folosiți mai mulți, sunt legați între ei și formează rețele. Modul în care sunt legați definește *topologia* rețelei. Există multe topologii posibile [Hecht90], dar două sunt cele mai folosite: feedforward și complet recurentă (aceasta din urmă va fi comentată în detaliu în capitolul următor). Modelul unificat Towell folosește o rețea cu topologie feedforward. Cea mai simplă rețea feedforward este rețeaua cu un singur strat. Aceasta este formată dintr-un număr de *m* neuroni care au în comun faptul că toți au același număr de intrări *n*. Funcționarea acestei rețele este foarte simplă: ca intrare se furnizează un vector (de obicei real) de lungime *n*, vector care este intrare pentru fiecare neuron. Ieșirea va fi un vector de lungime *m*, binar în cazul în care funcțiile de transfer sunt toate de tipul descris de formula (1-2), sau real, când cel puțin una din funcții este continuă. Un vector real se poate transforma în unul binar, aplicând formula (1-3), cu condiția ca valorile sale să fie în $(-1,0.5] \cup (0.5,1)$.

Rețeaua descrisă mai sus poate fi folosită ca și element constructiv pentru rețelele feedforward multistrat. Dacă avem o rețea cu *n1* intrări și *n2* ieșiri, putem cupla la ieșirea acesteia o nouă rețea de același tip, singura restricție fiind ca aceasta să aibă *n2* intrări. Se poate lega o altă rețea în continuare și așa mai departe, numărul de rețele legate în lanț definind numărul de straturi al rețelei finale. Teorema Vapnik-Chervonenkis [Hertz&al91] arată că pentru marea majoritate a aplicațiilor este suficient să folosim două straturi cu unități continue, dar pentru rețelele Towell, numărul de straturi este dat de forma cunoștințelor simbolice folosite pentru construirea sistemului. Pentru o rețea multistrat, care are toți neuronii dotați cu aceeași funcție de transfer *g*, funcția de transfer globală este:

$$o = \Gamma(W_s \times \Gamma(W_{s-1} \times \dots \Gamma(W_1 \times x) \dots)) \quad (1-4)$$

unde *s* este numărul de straturi, iar Γ este o funcție vectorială care aplică funcția de transfer a neuronilor folosiți în rețea (*g*) peste un vector, element cu element. W_i (*i*=1,*s*) sunt matricile care definesc rețeaua, ale căror linii sunt formate din vectorii de ponderi transpuși ai neuronilor de pe fiecare strat. Se poate observa că rețeaua se poate implementa foarte simplu, atât software cât și hardware [Beiu96].

O astfel de rețea poate primi la intrare un vector de o lungime fixată și poate genera la ieșire un alt vector de altă lungime, deasemenea fixă. Timpul de calcul necesar este foarte mic, răspunsul unei astfel de structuri computaționale fiind unul dintre cele mai scurte, comparativ cu alte implementări care fac aceeași transformare a vectorului de intrare în vector de ieșire. Acest tip de transformare poate fi descris prin două metode cât se poate de diferite. Prima, riguroasă, descrie maparea între cele două spații vectoriale (de intrare și ieșire) printr-o formulă matematică. A doua, experimentală, descrie maparea printr-o mulțime de perechi de vectori (cât mai mare),

fără a oferi nici o exprimare formală. Foarte multe din problemele din tehnică sunt exprimate prin a doua variantă, fiind extrem de dificil a se construi o funcție care să exprime legătura dintre intrarea și ieșirea unui sistem.

Fiind dată o mulțime de perechi de vectori $\{(x_p, o_p), p=1, P\}$, există metode prin care se poate realiza maparea între spațiul de intrare unde sunt fixați vectorii x_p și spațiul de ieșire, unde sunt fixați vectorii o_p , folosind o rețea neuronală feedforward multistrat. Acestea realizează o aproximare locală a funcției de mapare (care nu poate fi identificată formal) prin așa ziiși *algoritmi de antrenare* care calculează prin metode bazate pe minimizarea funcțiilor pe bază de gradient, valorile matricilor de ponderi ale rețelei. Pentru o rețea de tipul descris, cel mai folosit algoritm este clasicul *backpropagation* (o foarte bună descriere, în [Haykin94]). O versiune adaptată folosită pentru topologiile complet recurente este discutată pe larg în capitolul 2.

1.2. Modelul Towell

Cea mai răspândită formă de reprezentare a cunoștințelor la nivel simbolic sunt *regulile*. Acestea pot fi descrise la nivel propozițional sau la diferite nivele predicative [Luger89]. Forma generală a unei reguli propoziționale conjunctive este:

$$C :- P_1 \text{ and } \dots \text{ and } P_n.$$

unde C și P_i sunt propoziții atomice. Semantica unei astfel de piese de cunoaștere se poate exprima prin propoziția: " C este o propoziție adevărată, dacă și numai dacă toate propozițiile P_i sunt adevărate". Există și reguli disjunctive, unde este suficient ca o singură propoziție P_i să fie adevărată. O a treia categorie este cea a regulilor combinate, dar este simplu de arătat [Towell91] că un set de reguli combinate poate fi redus la un set de reguli pur disjunctive și un set de reguli pur conjunctive. O astfel de colecție de reguli poate fi privită ca o bază de cunoștințe foarte simplă. Există două metode de a utiliza static această bază de cunoștințe, folosind două mecanisme de inferență diferite, *forward-chaining* sau *backward-chaining*. Modelul Towell precizează calea de implementare a unei astfel de baze de cunoștințe care este interogată prin mecanismul de *forward-chaining*, folosind o rețea neuronală feedforward cu mai multe straturi.

Dacă avem o regulă ca cea de mai sus (conjunctivă), asociem valorile de adevăr TRUE și FALSE cu valorile numerice 1 și -1 și putem folosi ca automat de evaluare a consecinței C a regulii un neuron cu n intrări asociate premiselor P_i . Acest neuron va avea de acum încolo un simbol asociat (C). Dacă folosim următoarea funcție de transfer:

$$o = \begin{cases} 1 & \text{daca } h \geq 1 \\ -1 & \text{daca } h < 1 \end{cases} ; h = \sum_{i=1}^n w_i \cdot x_i ; w_i = \frac{1}{n} \quad (i=1, n) \quad (1-5)$$

unde:

$$x_i = \begin{cases} 1 & \text{daca } P_i = \text{TRUE} \\ -1 & \text{daca } P_i = \text{FALSE} \end{cases} \quad (1-6)$$

se observă că făcând transformarea inversă față de (1-6) pentru ieșirea o se poate obține o valoare de adevăr pentru simbolul C , dacă sunt cunoscute valorile de adevăr pentru simbolurile P_i , iar acest calcul este corect din punct de vedere logic [Shavlik&Towell91]. Pentru regulile disjunctive (1-5) se schimbă în:

$$o = \begin{cases} 1 & \text{daca } h \geq 1 \\ -1 & \text{daca } h < 1 \end{cases}; \quad h = \sum_{i=1}^n w_i \cdot x_i; \quad w_i = 1 \quad (i=1, n) \quad (1-7)$$

unde diferă doar metoda de setare a ponderilor. Neuronii astfel obținuți pot fi transformați în neuroni cu funcție de transfer continuă prin înlocuirea funcțiilor de mai sus prin tangenta hiperbolică. Această modificare presupune mutarea treptei în origine:

$$\sum_{i=1}^n w_i x_i \geq 1 \Rightarrow \sum_{i=1}^n w_i x_i - 1 \geq 0 \Rightarrow \sum_{i=0}^n w_i x_i \geq 0 \quad (w_0 = -1; x_0 = 1)$$

fapt care introduce o intrare suplimentară, setată mereu pe TRUE, și o pondere suplimentară, care are din start valoarea -1, indiferent dacă regula este disjunctivă sau conjunctivă. Pentru a avea o treaptă cât mai abruptă, asemănătoare cu funcția discretă, funcția de transfer $o = \text{tanh}(h)$ se înlocuiește cu $o = \text{tanh}(\beta h)$. Valoarea β se alege în funcție de numărul de intrări inițiale în neuronul respectiv, o valoare propusă prin experiment [Towell91] fiind $\beta = n/1.65$. Pentru a avea o funcționare cât mai apropiată de funcția treaptă, ponderea de indice zero (cea suplimentară) se setează pe valoarea -0.85, ca să fim siguri că în cazul în care h este mai mare decât 1, la ieșire vom avea o valoare apropiată de 1.0, transformată prin (1-3) în 1 întreg, iar în continuare, de inversa ecuației (1-6) în TRUE.

Fiecare regulă din baza de cunoștințe inițială va constitui în final un neuron, care va avea un număr de intrări egal cu numărul de premise din partea dreaptă a regulii plus o intrare suplimentară setată mereu pe TRUE (pe aceasta o vom ignora în continuarea prezentării, ea considerându-se implicită). Pentru ca avem două structuri neuronale care instanțiază reguli pur disjunctive sau pur conjunctive, este necesară o metodă de transformare a regulilor inițiale *combine* în reguli care sunt doar de aceste două tipuri. Ilustrată printr-un exemplu, metoda este următoarea:
-dacă avem regula

A :- (B and C and D) or (E and F and G).

regula va fi rescrisă în trei reguli care au aceeași semantică:

A :- A' or A''.

A' :- B and C and D.

A'' :- E and F and G.

rezultând trei neuroni, care vor avea ponderile egale cu $1/3$ pentru regulile cu partea stângă A' și A", iar pentru regula cu partea stângă A cele două ponderi vor fi egale cu 1. Având o bază de cunoștințe astfel transcrisă și care nu are reguli izolate (nu este formată din simboluri care nu aparțin și altor reguli), transformând toate regulile în neuroni, aceștia pot fi legați între ei. Ieșirea unui neuron care are asociat un simbol X (partea stângă a regulii este X) va fi o intrare a altor neuroni care sunt asociați unor reguli care îl conțin pe X în partea dreaptă. Exemplul din figura 1-1 (a., b., c. și d.) arată cum se leagă neuronii între ei, pornind de la un mic set de reguli. În [Towell91] este dată și soluția pentru legarea în rețea a unor simboluri precedate de operatorul logic unar "NOT", prin folosirea unor ponderi cu valoare negativă și prin schimbarea semnului valorilor de ieșire din neuroni, dar nu voi intra aici în astfel de amănunte.

Unitățile (neuroni) din figură care sunt desenați ca triunghiuri, au doar rol de a transforma valorile logice de adevăr ale simbolurilor asociate în valori numerice și invers. Se poate observa că rezultă o rețea feedforward, care are ca intrări valorile asociate acelor simboluri care apar doar în partea dreaptă a unei reguli din baza inițială de cunoștințe (posibile premise de raționament). Ieșirile acestei rețele sunt asociate unor simboluri care apar numai în parte stângă în setul de reguli (posibile concluzii finale de raționament). Dacă știm despre un subset de posibile premise că sunt adevărate și pornim rețeaua cu aceste valori, vom obține la ieșire un subset de simboluri care sunt concluzii finale și care au valoarea logică de adevăr TRUE. În [Towell91] se demonstrează izomorfismul funcțional dintre un sistem de raționament bazat pe *modus-ponens* (forward-chaining) și această rețea. Modelul Towell mai este cunoscut în literatură și sub denumirea de KBANN (Knowledge Based Artificial Neural Networks). Cele două modele de raționament (unul pur simbolic și unul conexionist), care pot fi implementate în mod foarte diferit, au fiecare o serie de avantaje și de dezavantaje.

Modelul pur simbolic poate fi ușor implementat printr-un mecanism de pattern-matching (de exemplu în LISP - chiar pe o mașină dedicată). Regulile pot fi exprimate ușor ca liste, manipularea lor și interfața sistemului cu utilizatorul fiind facile. Marele avantaj al unui astfel de sistem este însă *plasticitatea* [Shavlik&Towell91]. Este foarte simplu să adăugăm sau să ștergem reguli într-un astfel de sistem. Dezavantajul major al modelului simbolic este o caracteristică numită "crispness" (termen propus în [Pearl88] - și pe care nu am reușit să-l traduc). Aceasta se manifestă prin următoarele două probleme:

- un astfel de sistem nu poate fi construit decât plecând de la cunoștințe care pot fi exprimate prin reguli care sunt în orice context adevărate.
- Dacă lipsește un singur element de intrare pentru un anumit raționament, acesta este invalidat (ceea ce poate fi considerat și o lipsă de robustețe).

Aceste probleme apar mai ales atunci când:

- cunoștințele care sunt folosite la construirea sistemului au un anumit grad de incertitudine, sau nu pot fi reprezentate prin reguli.

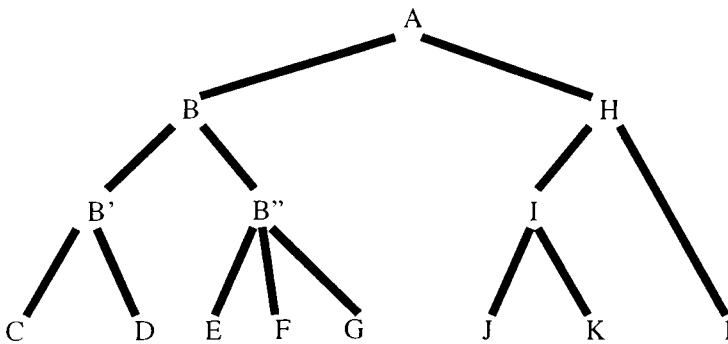
figura 1-1.

A :- B and H.
B :- (C and D) or (E and F and G).
H :- G and I.
G :- J and K.

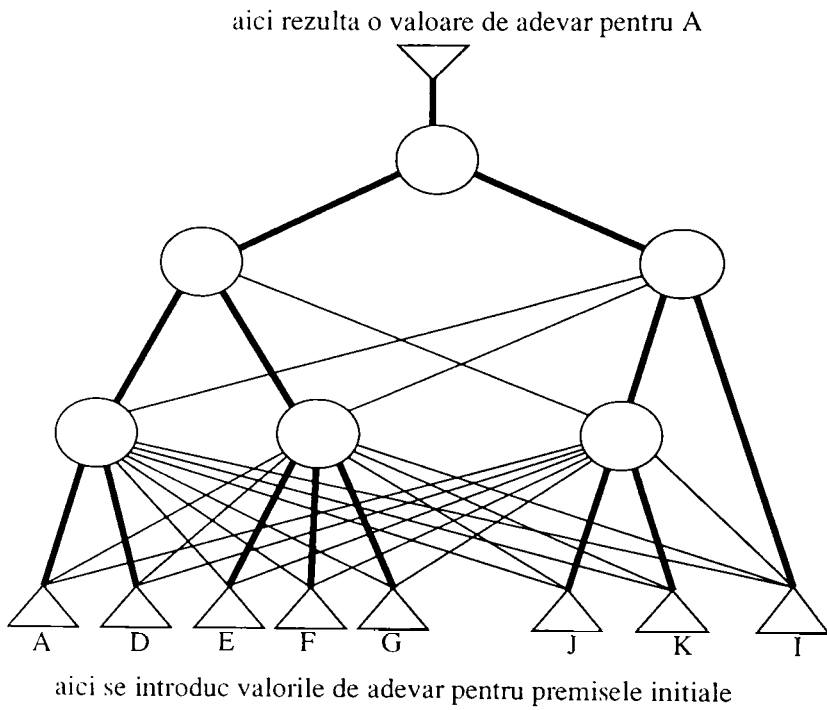
a. baza de cunostinte initiala

A :- B and H.
B :- B' and B''.
B' :- C and D.
B'' :- E and F and G.
H :- G and I.
G :- J and K.

b. baza de cunostinte rescrisa
(are doar reguli pur conjunctive si pur disjunctive)



c. grafurile de decizie construite pe baza regulilor



d. reseaua finala

Legenda:

- ponderi setate initial prin formulele 1-5 si 1-7
- ponderi practic nule (valori mici aleatoare)

figura 1-1 (continuare)

- datele de intrare au un grad de incertitudine sau sunt incomplete.

Pe lângă abordarea KBANN, au mai existat și alte metodologii care au încercat surmontarea acestor probleme, ca MYCIN [Shortliffe76], ODYSSEUS [Wilkins88], RTLS [Ginsberg90] și KR-FOCL [Pazzani&Brunk91]. Nici una până la acea dată (1991) nu a încercat însă o abordare neuronală.

Pentru KBANN, principalul dezavantaj este lipsa de plasticitate. Inserarea unor noi reguli în rețeaua neuronală este dificilă, modificându-se de fiecare dată matricile de ponderi ale rețelei rezultate. Dar există o serie de avantaje față de modelul simbolic. Dacă lucrăm cu funcții de transfer continue, valorile care pot apărea la intrarea și ieșirea sistemului nu vor fi valori "crisp" de adevăr (TRUE și FALSE) ci valori reale în intervalul $[-1,1]$ care pot reprezenta un anumit grad de incertitudine în premise și rezultate. Cel mai mare avantaj este însă posibilitatea inserării în sistem a unor cunoștințe empirice, exprimate prin perechi de seturi intrare/ieșire, care pot avea și o exprimare incertă (cunoștințe imposibil de exprimat prin reguli). Putem antrena rețeaua (obținută prin algoritmul prezentat în figura 1-1) cu aceste perechi, folosind backpropagation, modificând ponderile astfel încât sistemul să poată reproduce și raționamente empirice. Schema de antrenare [Towell91] presupune și introducerea de noi noduri și ponderi. Trebuie păstrat însă un echilibru între cunoștințele inițiale, exprimate prin regulile de start și cunoștințele empirice, pentru că o antrenare cu prea multe cunoștințe empirice poate distruge structura inițială a rețelei.

Există două posibilități de utiliza un astfel de sistem conexionist (care este unificat, pentru că odată implementat, nu mai are nevoie de regulile inițiale sau de alt suport simbolic). Prima posibilitate este de a folosi sistemul ca și sistem expert, construit pe baza unor cunoștințe simbolice (chiar incerte, vezi [Dragomirescu&a94]) și a unor cunoștințe cu aspect statistic. A doua posibilitate este de a crea noi reguli, extrăgând din rețeaua finală, neuroni care au fost atașați fără a li se asocia inițial un simbol, dar care au dobândit pe parcursul procesului de antrenament cu cunoștințe empirice, ponderi semnificative [Towell91]. Indiferent de stilul de utilizare, KBANN este o metodologie care arată că se pot integra cu succes cunoștințe simbolice și empirice (statistice), folosind o structură computațională conexionistă relativ simplă.

1.3. Concluzii

Modelul Towell a avut succes nu numai ca "methodology demonstrator" ci și în aplicații efective, de exemplu în conducerea proceselor chimice [Scott&a92] și în cardiologie (sistemul SOAR -[Watrous&a92]). Pentru ca modelul inițial era relativ simplist, s-au relevat o serie de probleme și dezavantaje, în special imposibilitatea de a modifica rețeaua după ce cunoștințele empirice au fost inserate prin antrenare neuronală. Altă problemă este folosirea rețelei ca estimator Bayesian, pornind de la un set de cunoștințe statistice (de exemplu pacienți pentru o aplicație medicală). S-a demonstrat că o astfel de rețea se comportă anormal ca estimator Bayesian în cazul în

care rețeaua are o structură predefinită și are mai mult de două straturi (un estimator Bayesian perfect este o rețea cu numai două straturi de calcul, primul strat (cel "ascuns") având un număr cât mai mare de unități - infinit pentru cazul ideal - vezi [MacKay92]). Aceste probleme au fost parțial rezolvate de îmbunătățiri ale modelului, cele mai importante fiind utilizarea unor metode statistice pentru calculul ponderilor [Koppel&al94] și utilizarea unor algoritmi genetici pentru modificarea dinamică a topologiei rețelei prin adăugarea unor noi unități în timpul antrenării, ca în modelul REGENT [Opitz&Shavlik97].

Există un alt domeniu de cercetare, numit ML (machine learning), care include și o parte din teoriile din rețele neuronale, mai exact, abilitatea lor de a stoca cunoștințele. Metodologia KBANN și derivatele ei au două laturi din punct de vedere al acumulării de cunoștințe. Prima este dată de inserarea regulilor în sistem prin definirea unităților și este numită într-un context mai general, *stocare constructivistă* [Kibler&Langley88] (care se poate face și prin alte metodologii, nu neapărat cu rețele neuronale). Ceea ce este caracteristic pentru acest tip de învățare este că cunoștințele de plecare (*Prior Knowledge*) sunt întotdeauna reprezentate simbolic. Restul cunoștințelor folosite pentru construirea sistemului pot fi reprezentate simbolic (dar nu sub formă de reguli), sau numeric (statistic) și singura condiție impusă asupra lor este translabilitatea în tipare de antrenare neuronală. Aceste cunoștințe empirice sunt de regulă exemple de funcționare ale sistemului dorit, fără a se cunoaște fenomenologia care duce la comportamentul acestui sistem. Învățarea din exemple se mai numește și învățare *inductivă* (inductive learning) [Kibler&Langley88]. Ea poate fi realizată nu numai prin metode conexioniste, existând mai multe căi posibile, parte din ele fiind deja puternic teoretizate [Dietterich&al95].

Remarcabil în modelul prezentat este faptul că stocarea constructivistă și învățarea inductivă sunt folosite împreună. Acest lucru este benefic, pentru că într-o problemă tehnică complexă, suportul de cunoștințe este în general format dintr-un set restrâns de cunoștințe teoretice și un set mai bogat de cunoștințe empirice. Pentru rezolvarea problemei se folosesc în primul rând cunoștințele teoretice existente, apoi încercându-se extragerea de noi cunoștințe teoretice din cele empirice (acesta fiind și scopul *in sine* al cercetării în general). Folosirea **în sistem** a cunoștințelor empirice este un pas important, care pe de o parte simplifică munca celui care implementează sistemul și pe de altă parte, îmbunătățește robustețea sistemului. De aceea am considerat că modelul Towell, cu toate slăbiciunile sale (descrise pe larg în [Koppel&al94]) este un pas de cotitură în integrarea simbolic-conexionistă.

Am folosit modelul Towell la primele mele aplicații practice ale sistemelor expert cu componentă neuronală [Dragomirescu&al94], iar rezultatele au fost bune. Încercând însă să aplic acest model la probleme ceva mai complicate (tot din domeniul medical și anume analiza procesului infecțios), am constatat o scădere drastică a performanțelor, în termeni de rezultate obținute și resurse computaționale necesare. Observația mea a fost atunci a fost că de vină sunt:

- rețelele folosite, care sunt mult prea simple,
- paradigma că structura rețelei este cea care guvernează raționamentul,
- structurile simbolice de tip regula "if-then", care sunt mult prea sărace pentru a captura semantica unor cunoștințe care aveau un grad mult mai puternic de asociere și ierarhizare.

Pornind de la aceste observații, am încercat să gășesc soluții pentru aceste trei "slăbiciuni":

- pentru că "puterea" unei rețele este dată de gradul ei de conectare [Hecht90], am hotărât să folosesc rețelele complet recurente, pentru că ele au cel mai mare grad de conectare,
- structura unei rețele fiind de regulă fixă, iar raționamentul fiind un concept dinamic, am stabilit că **starea** rețelei este cea care guvernează raționamentul,
- am trecut la un nivel superior în materie de structuri simbolice folosite, mai exact spus, la *grafuri conceptuale* (vezi capitolul următor).

Important de notat este că noul model pe care mă pregăteam să-l dezvolt era unul eminamente unificat (având o componentă ușor hibridă - și anume o parte translațională). Problemele întâmpinate, descrise pe larg în capitolele următoare, au făcut ca acest model să se transforme într-unul hibrid, conform taxonomiei propuse în Introducere.

2. Structurile simbolice și conexiunile utilizate

În cazul unui sistem bazat pe metodologia KBANN, utilizatorul trebuie să fie un bun cunoscător a două instrumente folosite des în inteligența artificială: sistemul expert cu mecanism de inferență tip forward chaining, construit din reguli propoziționale și rețeaua neuronală feedforward, care are de obicei un număr mic de unități, dar poate avea un număr mare de straturi. Pentru că în metodologia propusă de mine se folosesc două instrumente mai puțin cunoscute, care au un grad relativ mare de complexitate, am prezentat în acest capitol, pe larg, o serie de aspecte teoretice și experimentale legate de acestea. În plus, am exprimat și o serie de considerații personale, derivate nu atât din bibliografia folosită, ci mai mult din experiența dobândită.

2.1. Grafuri conceptuale

Regulile "if-then" sunt o transcriere a cunoștințelor exprimate în limbaj natural. Necesitatea formei lor rigide este datorată mecanismului de inferență folosit pentru raționamente. Din păcate, la translația unor cunoștințe din limbajul natural în reguli, se pierde foarte mult din conținutul semantic al propozițiilor inițiale. Din punct de vedere al reprezentării cât mai corecte a semanticii propozițiilor, este recunoscut faptul [Allen87] [Roberts91] [Lendaris92] că cea mai bună alegere sunt grafurile conceptuale, o propunere de structură simbolică făcută de John Sowa în 1984 [Sowa84], derivată din grafurile de dependență conceptuală ale lui Roger Schank [Schank81] [Luger&Stubblefield89]. Pe lângă informația semantică purtată de aceste structuri, se păstrează, mai ales datorită metodelor de translație, și o parte din informația *sintactică* asociată propoziției. Nu vom insista asupra metodelor de translație, ele fiind tehnici de mult cunoscute [Riesbeck86] [Sowa95].

Poate ar fi fost util un studiu de caz pentru alegerea structurii de reprezentare

intermediare, dar există un argument foarte puternic pentru alegerea grafurilor conceptuale. Sowa [Sowa95] demonstrează că toate structurile de reprezentare a cunoștințelor general acceptate și utilizate pot fi transformate prin metode deterministe în grafuri conceptuale și invers. Dar se poate observa că întotdeauna transformarea **graf conceptual** → **altă structură** → **graf conceptual**, implica o pierdere de conținut (mai ales semantic), pierdere care se înregistrează în primul pas.

2.1.1. Reguli de construire și utilizare

Ce este un graf conceptual?

Definiție

*Un graf conceptual este un graf finit, orientat, conectat și bipartit. Arcele nu sunt etichetate și nodurile sunt fie **concepte**, fie **relații conceptuale**.*

Restricție

Nodurile concept pot avea arce doar spre noduri relație, iar nodurile relație pot avea arce doar spre noduri concept.

Conceptele reprezintă obiecte concrete sau abstracte din universul discursului.

Definiție

Se numește universul discursului, mulțimea formată din ierarhia de noțiuni, idei, cuvinte folosită pentru susținerea unei teorii, unde numele teoriei se află în vârful acestei ierarhii (ex. discursul fizicii, discursul filozofic etc).

Pentru a avea un punct de plecare solid, ar trebui formalizăm matematic această noțiune de univers al discursului. Formalizare care ar implica o serie de definiții topologice, pentru că relevant pentru limbajul folosit sunt relațiile care există între cuvinte și modul de ordonare al acestora. Dat fiind faptul că lucrăm cu cuvinte ale limbajului natural, acestea având în majoritatea cazurilor sensuri legate de contextul în care sunt folosite, este foarte dificil să realizăm o astfel de formalizare. O tentativă de a construi un astfel de sistem de relații este "dicționarul lui Quillian" [Quillian69], care surprinde foarte multe relații semantice între cuvinte, modul de organizare fiind foarte simplu, dar care nu este utilizabil într-un sistem automat, fiind o structură eminentemente declarativă și deloc procedurală. Scrierea unor algoritmi care să utilizeze acest dicționar, pentru rezolvarea unor probleme practice, s-a dovedit a fi deosebit de dificilă, modelul fiind abandonat. Interesant este de remarcat faptul că, spre deosebire de toate abordările anterioare, care insistau pe structura sintactică a informației și rolurile sintactice ale cuvintelor, acest model a abandonat complet problemele legate de sintaxă, concentrându-se numai asupra relațiilor care există între cuvinte și definiții.

Motivul principal pentru care acest model nu a putut fi folosit niciodată, este că suferă de sindromul "solving the world". Autorii lui s-au gândit la o soluție universal valabilă pentru toate aplicațiile posibile. Datorită gradului prea mare de generalitate, modelul s-a dovedit inutilizabil în aplicațiile practice. Din acest motiv, este indicat ca în orice aplicație care implică folosirea limbajului natural, modelul matematic folosit să fie pragmatic, adică să slujească idea aplicației și nu să dea soluții generale pentru tratarea limbajului natural.

Din acest motiv, formalizarea universului discursului pe care am ales-o este gândită în așa fel încât aceasta să fie cât mai simplu de folosit într-o aplicație în care structura de bază pe care o folosim sunt grafurile conceptuale.

Decizia principală în proiectarea acestui model pleacă de la următoarea observație: *Nodurile concept pot conține două tipuri de concepte: generale (clase) și individuale (instanțe).*

2.1.2. Structurarea universului discursului

Sistemele expert care folosesc ca și formă de reprezentare a cunoștințelor, grafurile conceptuale, au un mecanism de inferență bazat pe regula de formare canonică "join" (compunere) [Sowa84] [Sowa95]. Spre deosebire de mecanismele de inferență bazate pe logică, inferența prin compoziție nu oferă certitudinea păstrării riguroase a valorilor de adevăr.

Operația de de compunere (join) se aplică asupra a două grafuri conceptuale. Ea devine posibilă doar în cazul în care cele două grafuri conceptuale luate în considerare au cel puțin un nod concept comun. Dacă vom lua în considerare grafurile din figura 2-1a, rezultă prin compunere graful conceptual din figura 2-1b. Conceptul care este comun celor două grafuri, devine calea de legătură dintre ele, graful rezultat fiind o piesă de cunoaștere cu un conținut semantic mai bogat decât cele două grafuri inițiale, și care din punct de vedere pragmatic este utilizabil în procesul eventual de generare a unor propoziții (acestea vor avea în mod sigur sens, însă nu putem să facem prea multe presupuneri despre valorile lor de adevăr).

Dacă grafurile conceptuale inițiale au mai multe concepte comune, toate acestea vor fi reunite, în graful rezultat neputând apărea același concept de două ori. Se poate observa din exemplul următor (figura 2-2), că operația de combinare duce uneori la apariția unor noduri relație redundante, fiind necesară o operație specială, numită simplificare. Relația r3 apare de două ori și acest lucru nu este necesar.

Cazurile în care două grafuri conceptuale se pot compune nu sunt chiar atât de frecvente pe cum ar părea la prima vedere. Acest lucru este și mai evident dacă mulțimea completă a grafurilor conceptuale folosite într-un sistem este relativ redusă, iar grafurile conțin concepte comune care nu apar des. Din acest motiv este necesară o operație care să faciliteze operația de combinare. Pentru a introduce această

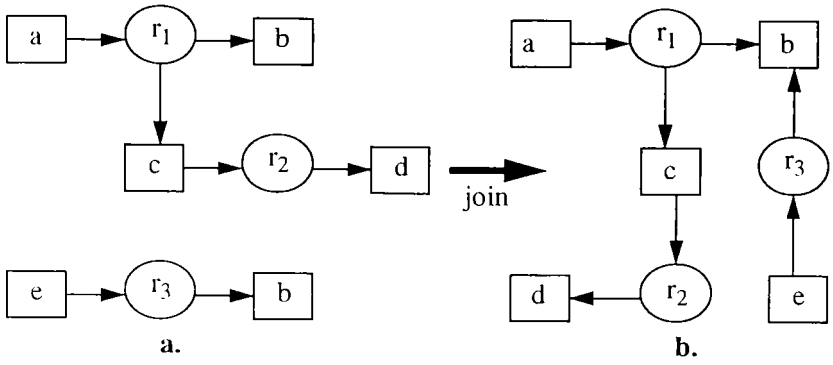


figura 2-1.

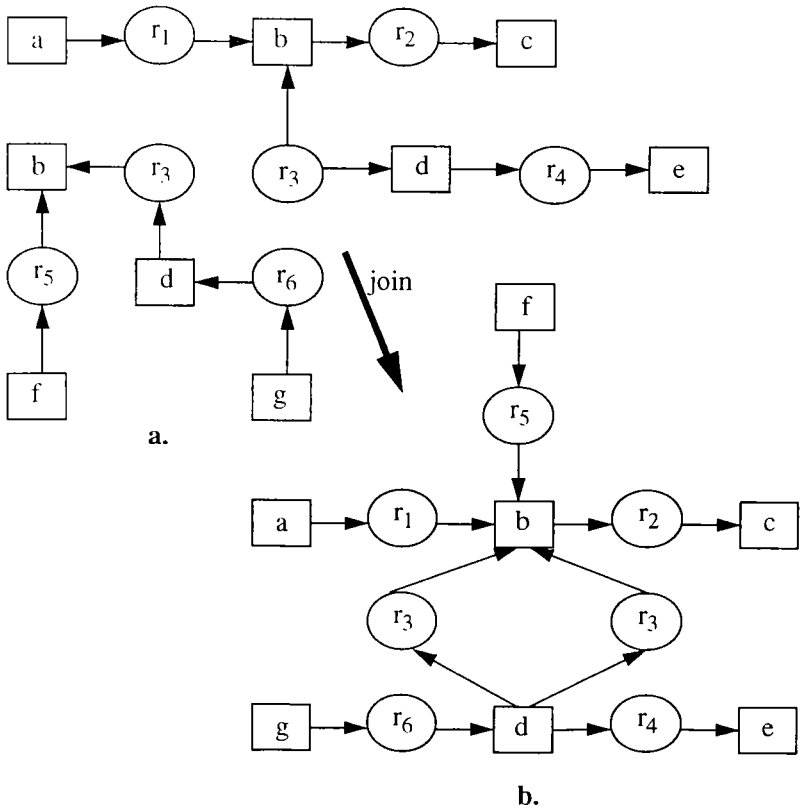


figura 2-2.

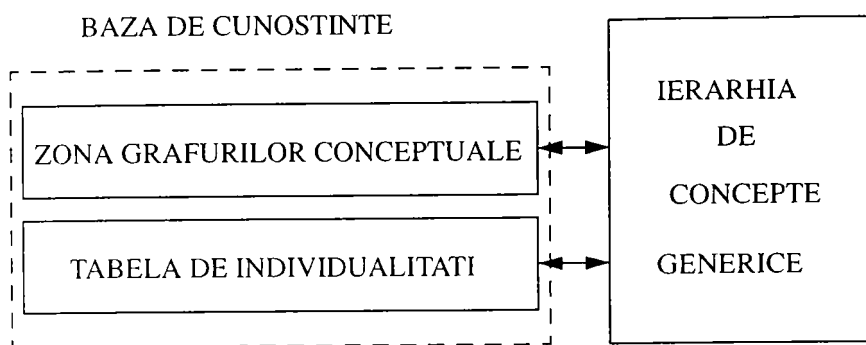


figura 2-3.

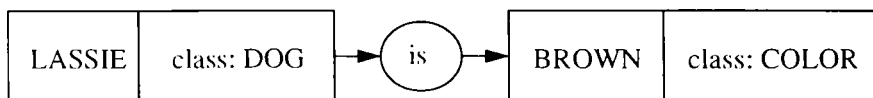


figura 2-4.

nouă operație, este necesară o rafinare suplimentară a nodurilor concept.

Conținutul nodurilor concept este un simbol (șir de caractere) care denotă conceptul respectiv. Observând mulțimea conceptelor care aparțin universului discursului putem aprecia că între acestea există relații de moștenire de tip clasă/superclasă. Forțând această constatare, putem construi în toate cazurile practice, o ierarhie completă a conceptelor. Forțarea se datorește introducerii necesare a unor concepte suplimentare, (care nu aparțin inițial universului discursului), pentru a putea crea toate lanțurile de moștenire clasă/superclasă complete.

Ceea ce urmărim este crearea unei structuri a ierarhiei de moșteniri de forma unui arbore, a cărui rădăcină să fie un supertip universal T, care nu trebuie să aiba o semantică definită. Important este faptul că fiecare concept aparține unei anumite clase și acest lucru este evident, pentru că în acest arbore, cu excepția supertipului universal T, toate nodurile au unul sau mai multe noduri părinte (considerând posibil și cazul moștenirii multiple). Operația nou introdusă, posibilă doar în contextul unui arbore de clase, care facilitează operația de combinare, se numește generalizare/specializare și acționează asupra valorilor conținute în nodurile concept.

Dacă avem un nod concept care conține o clasă (un concept din ierarhia de moșteniri), nod numit și concept generic, putem instanția această clasă, înlocuind simbolul care denotă clasa cu una din instanțele sale (specializare). Dacă avem un nod concept care conține o instanță a unei clase (nod numit și concept individual), putem înlocui simbolul care îl denotă cu simbolul asociat clasei sale.

Pentru a ușura implementarea, este util să separăm conceptele individuale de cele generice. Se conturează astfel două structuri de date diferite: o tabelă de individualități și o structură de concepte generice, care se poate observa că este chiar ierarhia inițială de moșteniri, din care am eliminat toate conceptele individuale. Deci, pentru a putea efectua operații cu grafuri conceptuale, vom folosi schema de principiu a organizării datelor ilustrată în figura 2-3. Aceste trei zone de grupare a informației care reprezintă universul discursului și cunoștințele înglobate în grafurile conceptuale sunt legate între ele în felul următor:

- conținutul fiecărui nod concept aparținând grafurilor conceptuale este format din doi pointeri: unul indicând către clasa din care face parte conceptul și celălalt indicând către o poziție din tabela de individualități. Dacă nodul concept este generic și denotă chiar clasa, atunci pointerul către tabela de individualități va avea o valoare specială, pe care o vom nota cu *. Pentru o implementare cât mai rapidă este util ca toate clasele care aparțin ierarhiei de moșteniri să aiba un pointer și un număr întreg asociat către tabela de individualități, care să indice prima individualitate a acestei clase și numărul de individualități de care această clasă dispune în universul curent al discursului. Evident că în acest caz, organizarea tabelii de individualități trebuie făcută grupat pe clase.

O ultimă îmbunătățire, care la prima vedere, pare redundantă, este atașarea la

fiecare poziție a tabelii de individualități a unui pointer către un nod din ierarhia de moșteniri care să indice clasa din care face parte individualitatea din poziția respectivă. Toate acestea sunt ilustrate în exemplul următor:

- Considerăm grafurile conceptuale din figura 2-4. În universul discursului există clasele DOG și COLOR (concepte generice). În tabelii de individualități există trei nume de câini: LASSIE, SCOOPY, DROOPY și trei culori posibile: BLACK, BROWN, WHITE. O imagine parțială a structurii de date arată ca în figura 2-5. Dacă avem două grafuri conceptuale, ca în figura 2-6, putem aplica operația de specializare asupra conceptului generic DOG, iar în momentul în care ajungem la instanța SCOOPY, putem combina grafurile și vom obține o informație mai completă (ca în figura 2-7a). Dacă generalizăm conceptul SCOOPY vom obține grafurile din figura 2-7b care conține mai puțină informație.

Putem însă ajunge și la rezultate greșite din punct de vedere logic. Dacă vom considera cele două grafuri conceptuale din figura 2-7c, prin specializare și combinare vom obține grafurile conceptuale din figura 2-7d, care reprezintă o piesă de cunoaștere falsă (pentru că avem deja o altă piesă de cunoaștere care afirmă că LASSIE este BROWN).

2.1.3. Consistența logică a universului discursului

Dacă traducem, tripleții de tipul CONCEPT-relație-CONCEPT care constituie grafurile conceptuale rezultate din ultima combinare, în axiome predicative și trebuie comparăm sub forma unui test de consistență logică (nu trebuie să existe două axiome care se contradic) cu toți ceilalți tripleți extrași din grafurile conceptuale existente înainte de operația de combinare efectuată, atunci operația "join" efectuată este corectă logic. Dacă apare o inconsistență logică, operația de combinare respectivă va fi abandonată.

Relațiile care unesc conceptele pot fi și ele standardizate, în primul rând poate fi făcută o clasificare în verbe, locuțiuni și însușiri. Verbele pot fi la rândul lor clasificate, cea mai importantă clasificare fiind cea legată de timp (prezent, trecut, viitor, imperfect, gerunziu etc). O analiză mai serioasă a relațiilor este realment necesară, dar nu constituie pe moment subiectul acestei lucrări. Este o zonă de interes, care trebuie în mod evident investigată, dar am considerat că este prioritar să ne concentrăm asupra problematicii conceptelor.

Totuși o structurare primitivă a relațiilor este necesară, astfel că în final vom avea și un arbore al relațiilor (conceptul tot în genul unei ierarhii de moșteniri, acest lucru fiind necesar pentru un anumit tip de procesare tensorială, folosită mai târziu în această lucrare). În final vom avea universul discursului formalizat în felul următor:

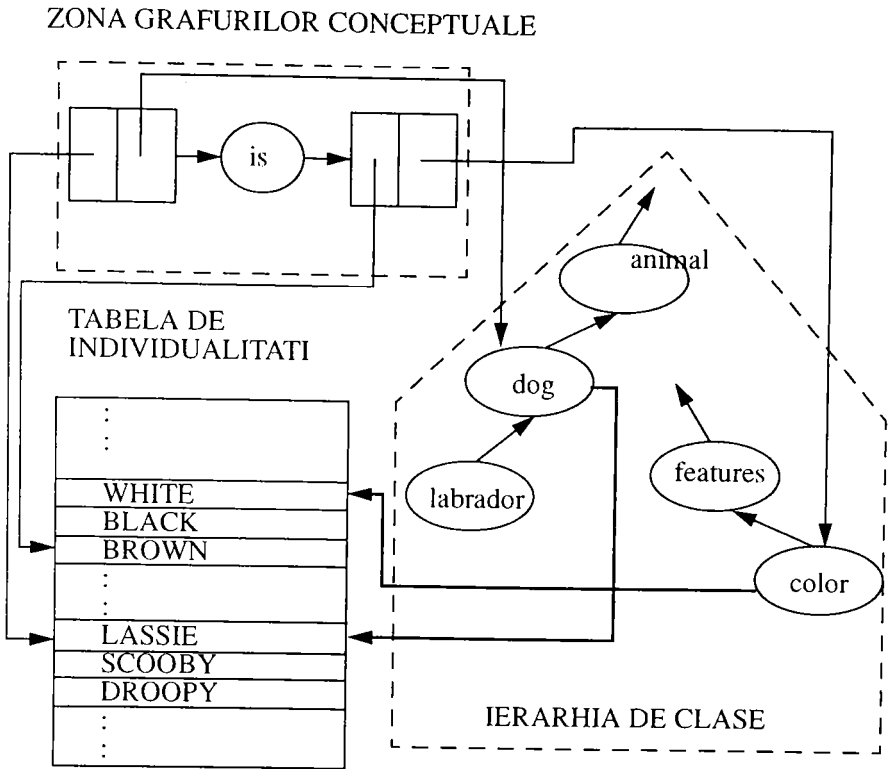


figura 2-5.

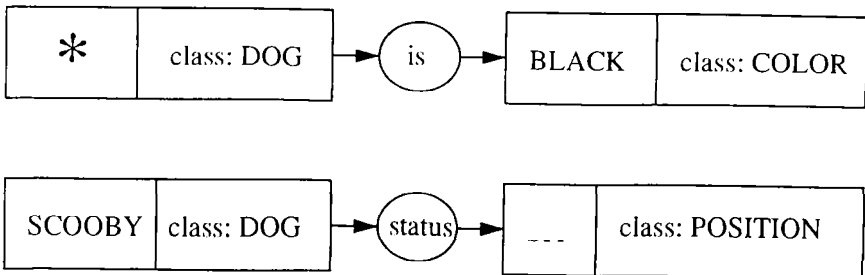
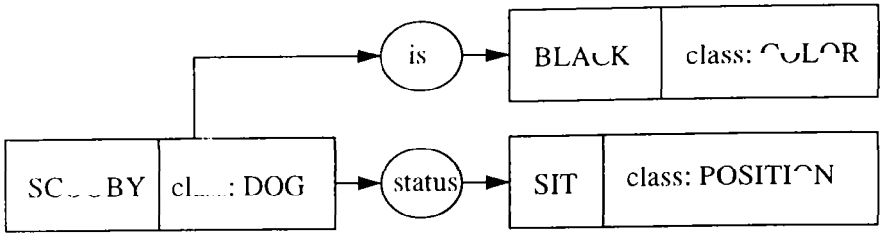
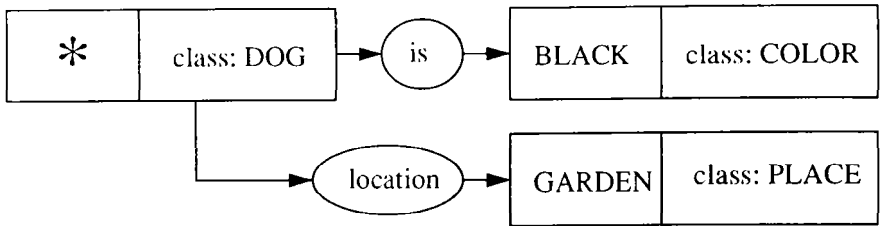


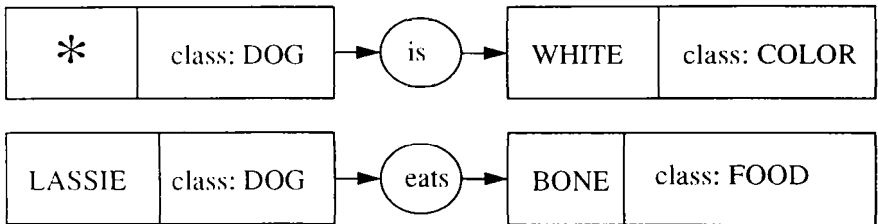
figura 2-6.



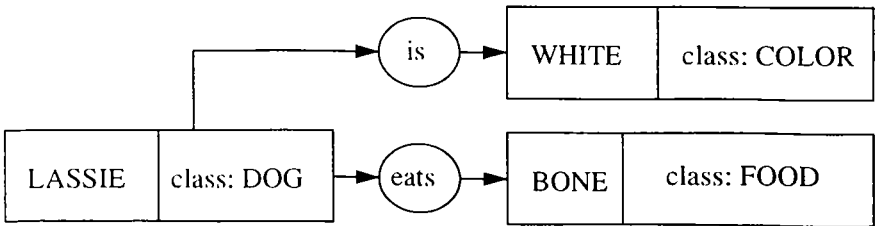
a.



b.



c.



d.

figura 2-7.

$$CG = \{ cg \mid cg \subset Tr, cg \subset P \} \quad (2-1)$$

$$Tr = \{ a(b, c) \mid a \in ; b, c \in CS \} \quad (2-2)$$

$$P = \{ (pC, pI) \mid pC \in IC, pI \in TI \} \quad (2-3)$$

unde:

- R este mulțimea relațiilor, Tr mulțimea tripleților, P mulțimea pointerilor.
- CS este mulțimea tuturor conceptelor.
- IC este ierarhia de clase, o mulțime parțial ordonată (sortată topologic):

$$IC = \{ cl \mid cl \prec T, cl \succ 1, \exists scl \prec cl, \exists scl \succ cl, scl, scl \in IC \} \quad (2-4)$$

- TI este mulțimea individualităților (tabelă ordonată și grupată pe clase).
- În plus, formalizarea mulțimii relațiilor este dată de:

$$R = \{ r \mid r \prec T, \exists sr \succ r, sr \in R \} \quad (2-5)$$

ceea ce definește o structură de arbore.

Clasa este la rândul ei un triplet:

$$cl(pI, n), pI \in TI, x(pI) \in IC, x(pI+i)_{i=0, n} \in IC \quad (2-6)$$

iar operația de combinare se reduce la :

$$cg1 \circ cg2 = cg3 \quad (2-7)$$

cu proprietatea că:

$$cg3T = \{ a(b, c) \mid \exists (b \wedge c) \in (cg1 \vee cg2) \} \quad (2-8)$$

iar:

$$cg3P = \{ (pC, pI) \mid (pC, pI) \in (cg1P \vee cg2P) \} \quad (2-9)$$

- cg3T este mulțimea tripleților rezultați din graful combinat,
 - cg3P este mulțimea pointerilor,
- iar graful combinat este complet determinat prin mulțimea neomogenă:

$$cg3 = cg3T \cup cg3P \quad (2-10)$$

Mulțimea rezultată din definiția (2-8) poate fi vidă, iar în acest caz se consideră

operația de combinare ca fiind nereușită.

Dacă mulțimea cg_3T nu este vidă, vom aplica asupra tuturor tripleților ei, principiul rezoluției [Robinson65] pentru verificarea consistenței. Acest principiu se poate aplica pentru grafurile conceptuale în felul următor:

1. Vom considera fiecare mulțime cg_iT ca pe o clauză multiplă de tip conjuncție:

$$cg_iT \Delta \bigwedge_{j=1, m_i} a_j(b_j, c_j) \quad i=1, n \quad (2-11)$$

unde n este numărul curent de grafuri conceptuale și m_i numărul de tripleți ai fiecărui graf.

2. Trecem clauzele de forma (2-11) în formă disjunctivă folosind ecuațiile lui De Morgan.

3. Adăugăm pe rând fiecare triplet din mulțimea cg_3T în formă negată la setul de ecuații disjunctive inițial.

4. Dacă aplicând refutația prin rezoluție (vezi ecuația 2-12), obținem în final clauza vidă, înseamnă că nu avem nici o contradicție și operația de combinare verificată este consistentă logic.

$$\begin{aligned} & a_1 \vee \dots \vee a_{i-1} \vee a_i \vee a_{i+1} \vee \dots \vee a_f \\ & b_1 \vee \dots \vee b_{j-1} \vee \neg a_i \vee b_j \vee \dots \vee b_f \quad \rightarrow \\ & a_1 \vee \dots \vee a_{i-1} \vee a_{i+1} \vee \dots \vee a_f \vee b_1 \vee \dots \vee b_{j-1} \vee b_j \vee \dots \vee b_f \end{aligned} \quad (2-12)$$

2.2. Rețele conexiunile recurente

Am avut mai multe alternative pentru alegerea unei topologii potrivite pentru componenta conexiunistă a sistemului. Din start, am eliminat o clasă importantă de rețele neuronale, rețelele cu propagare înainte (feedforward). Am considerat că pentru procesarea neurosimbolică, o idee simplă este de a alege modele cât mai apropiate de modelul natural. În acest caz se pot face două observații.

Prima observație este de natură cantitativă. Creierul uman, singurul capabil în lumea vie cunoscută de a realiza activități de complexitate procesării simbolice, are în medie un număr de $2 \cdot 10^{11}$ neuroni. Ordinul de mărime al celor mai mari rețele neuronale artificiale simulate pe cele mai puternice calculatoare actuale este de 10^5 unități de calcul (neuroni matematici). S-ar putea ca rezultatele pe care le vom obține să fie alterate de această discrepanță dimensională. În acest punct putem face o serie de presupuneri, care sunt confirmate cu un mare factor de certitudine de către

cercetătorii din domeniul neuro-psihologic [McClelland94].

a. Creierul uman folosește doar o mică parte a capacității sale pentru activitățile cognitive care ne interesează.

b. Pentru o problemă foarte simplă, "procesorul" cognitiv uman nu folosește întreaga sa capacitate de procesare simbolică, având o structura modulară, din componente slab cuplate, fiecare componentă fiind specializată în rezolvarea unui alt tip de problemă. Organizarea acestor componente (așa cum sugerează și teoria asociaționistă, vezi [Quillian&Collins69]) este o structură ierarhică impusă de ierarhia conceptelor memorate de mintea umana.

La ora actuală acestea sunt doar presupuneri, dar deși există posibilitatea ca ele să fie infirmate de neuropsihologie și neurofiziologie, nici un experiment sau observație clinică nu le infirmă [Röbel94].

Cel mai puternic argument fizic este separația pe zone corticale a funcțiilor organismului. Nu vom intra în amănunte, dar un exemplu banal care se poate da, este specializarea emisferelor cerebrale pe activități legate de jumătatea opusă de corp. Un alt argument în favoarea acestei specializări a fost dat de comportamentul persoanelor cu "partial brain-damage", care suferă de amnezii locale [Elman90]. Aceste informații experimentale ne indică destul de clar localizarea unor informații cu semantică apropiată, în zone apropiate în creierul uman. Ca și punct de plecare, putem considera ca plauzibilă ideea că un număr relativ mic de neuroni antrenaj corespunzător sunt capabili să execute o procesare de nivel cantitativ mai scăzut, dar calitativ ridicată.

A doua observație este de natură calitativă. Modelul matematic al neuronului izolat este foarte depărtat de modelul natural. Dar pentru o populație suficient de mare de neuroni artificiali comportamentul unei rețele este destul de apropiat de comportamentul unei rețele naturale. De exemplu, la un nivel simplist, retina a putut fi simulată folosind rețele artificiale. Esențială este însă organizarea unei astfel de rețele. Spre deosebire de creierul uman, noi putem crea virtual orice organizare posibilă. Aici rezidă puterea rețelelor neuronale artificiale. Creierul uman în dezvoltarea sa fizică depune un efort imens și foarte puțin eficient de a apropia zone corticale diferite prin procesul de "împachetare" a scoarței cerebrale și generare a circumvoluțiunilor. Zone corticale îndepărtate nu se vor putea uni niciodată, decât indirect. Lungimea axonilor fiind finită (dendritele fiind mobile), un neuron are în medie doar 10^4 vecini de care este legat direct.

2.2.1. Modele conexiunile inspirate din fizica statistică

Putem crea rețele conexiunile în care toți neuronii să fie legați între ei. Inițial toate ponderile vor avea valoarea 0, iar dacă doi neuroni vecini sunt în situația să se asocieze (pentru că zonele din care fac parte trebuie să se "aprope") ponderea legăturii

dintre ei se va modifica corespunzător. Pentru a implementa această idee, avem nevoie de un suport teoretic confirmat. Lucrarea de față nu își propune să caute noi modele matematice pentru arhitecturi conexiuniste dedicate pentru procesarea simbolică neuronală, ci să investigheze modelele cunoscute pentru a detecta care sunt cele mai potrivite din punct de vedere al performanțelor. Este interesant de observat că modelele care vor fi prezentate în continuare au fost inspirate în special din fizică și mai puțin din biologie.

Considerăm un material cu proprietăți magnetice, solid, cu o structură cristalină rectangulară. Atomii care compun structura cristalină sunt echidistanți și au toți aceleași proprietăți. Proprietatea principală este spinul magnetic. Acesta este o mărime vectorială, care indică o direcție în spațiu și are lungime unitară. Pentru un astfel de material, atomii au un număr limitat de direcții posibile pentru spinul magnetic. Un astfel de model a fost propus de Ising [Hertz&al91]. El a propus ca spinul magnetic să poată lua doar direcții posibile, notate generic cu $\langle I \rangle$ și $\langle -I \rangle$.

Câmpul magnetic generat de un atom va fi notat cu b_i și va avea direcția spinului său magnetic. Întreaga structură poate fi supusă unui câmp magnetic extern B_{ext} . Pe lângă acest câmp, fiecare atom al rețelei este supus și acțiunii combinate ale câmpurilor atomilor vecini. Factorul de influență dintre un atom i și atomul j se notează cu p_{ij} și depinde de poziția relativă a celor doi atomi unul față de celălalt (distanță și unghi). Un atom este supus unei influențe interne datorată tuturor atomilor, din rețea, influență care este totdeauna prezentă și o influență externă care poate să lipsească. Ising a formalizat aceste fenomene prin formula:

$$\bar{B}'_i = \left(\sum_{j=1}^N \bar{p}_{ij} \cdot \bar{B}_j \right) + \bar{B}_{ext} \quad (2-13)$$

Vectorul b'_i este rezultanta influenței interne și externe. Se observă că o componentă a rezultantei interne este dată de valoarea "anterioară" a vectorului b_i . Acest lucru ne obligă să dăm și o dimensiune temporală acestei formule. Putem face acest lucru discret sau continuu. Deoarece în contextul discuției noastre această problemă reprezintă doar o problemă de implementare vom amâna discuția despre factorul timp. Important este faptul că din interacțiunea acestor vectori se poate modifica spinul. Se poate demonstra că la un moment dat:

$$s_i = f \left(\sum_{j=1}^N |\bar{p}_{ij}| \cdot s_j \right) \quad (2-14)$$

Unde f este o funcție bipolară care trebuie în mod obligatoriu să aibă o zonă de saturație. Concluzia este că spinul unui atom depinde numai de spinii atomilor vecini din structură [Hertz&al91].

2.2.2. Rețelele Hopfield

Făcând o analogie cu modelul Ising, Hopfield (1982) a descris o rețea neuronală folosind aceleași formule. Dacă vom considera N neuroni, fiecare corespunzând cu un atom din modelul magnetic, putem folosi același principiu de transmitere a interacțiunilor. Într-o astfel de rețea, fiecare neuron își trimite semnalul de ieșire prin "fan-out", către toți ceilalți neuroni din rețea, inclusiv către el însuși și primește ca intrare suma ponderată a ieșirilor tuturor neuronilor. O astfel de rețea care nu respectă topologia clasică a perceptronului (rețeaua feedforward) și permite și legături înapoi, se mai numește și rețea recurentă. Fiind și o rețea complet conectată, rețeaua Hopfield se mai numește și rețea recurentă complet conectată. Dacă vom nota intrarea ponderată într-un neuron i cu h_i și ieșirea cu S_i , vom avea pentru un neuron i :

$$h_i = \sum_{j=1}^N w_{ij} \cdot S_j \quad (2-15)$$

Hopfield a ales în analogie cu spinul magnetic care putea fi $\langle 1 \rangle$ sau $\langle -1 \rangle$, funcția de ieșire a neuronului astfel încât ea să poată lua doar două valori: 1 și -1. O funcție potrivită pentru a genera o astfel de ieșire este funcția *signum* care dă valoarea 1 pentru argumente pozitive și -1 pentru argumente negative.

$$S_i = \text{sgn}(h_i) = \begin{cases} 1 & \text{dacă } h_i \geq 0 \\ -1 & \text{dacă } h_i < 0 \end{cases} \quad (2-16)$$

Se poate observa o analogie cu formula spinului magnetic (2-14). Dacă rescriem formula vom obține:

$$S_i = \text{sgn}\left(\sum_{j=1}^N w_{ij} \cdot S_j\right) \quad (2-17)$$

O astfel de rețea reprezintă cel mai simplu model posibil pentru procesarea simbolică [Hertz&al90]. De ce? După cum am menționat anterior, tocmai o astfel de arhitectură, unde fiecare neuron poate fi legat cu oricare alt neuron ne interesează. Se poate observa însă că numărul de legături între neuroni devine prohibitiv (N^2). Observația este importantă, pentru că ponderile sunt numere reale și ne interesează și precizia lor. În acest caz, dacă simulăm funcționarea rețelei cu ajutorul unui program scris într-un limbaj de programare de nivel înalt, vom reprezenta aceste ponderi ca și variabile în virgulă flotantă. Dacă am încerca să memorăm ponderile unei rețele cu un miliard de neuroni (10^9) care reprezintă cam a două suta parte din numărul de neuroni din creierul uman, am avea nevoie de $(10^9)^2 = 10^{18}$ variabile în virgulă flotantă, adică 8 miliarde de gigaoceteți de memorie (considerând că o variabilă în virgulă flotantă se reprezintă pe 8 oceteți).

Dar dacă ne vom baza pe considerațiile făcute anterior, că pentru o activitate cognitivă simplă putem presupune ca suficienți doar 10^4 neuroni adică 10^8 ponderi,

rezultă un necesar de 800 de megaocteți, care este o valoare rezonabilă dacă vom memora extern ponderile, având în memoria RAM doar o parte din acestea. Se observă că în această situație toate ponderile pot fi grupate într-o matrice de dimensiune $N \times N$, pe care o vom nota cu W . O linie i acestei matrici conține ponderile asociate intrărilor într-un neuron de indice i , deci pentru calculul unei intrări ponderate h_i este suficient să cunoaștem doar linia respectivă din această matrice, și ieșirile din toți neuronii S_j ($j=1..N$). Acestea le vom grupa într-un vector S pe care îl vom numi *vector de stare* al rețelei Hopfield.

La rețelele feedforward, semnalul emis de un strat de neuroni (un vector) este transmis stratului următor, pornind de la primul strat și ajungând la ultimul, în ordinea activării straturilor. Această secvență de operații determină întotdeauna aceeași funcționare a rețelei. La o rețea Hopfield, la fel ca și la modelul Ising, apare și factorul timp. Funcționarea rețelei recurente poate fi formalizată dacă vom considera timpul ca o secvență discretă $t=0,1,2,\dots,t_{stop}$ și aplicăm funcția de transfer pentru fiecare neuron astfel încât ieșirea să fie obținută întârziată cu o cuantă de timp:

$$S_i(t+1) = \text{sgn} \left[\sum_{j=1}^N w_{ij} \cdot S_j(t) \right] \quad (2-18)$$

Această relație se numește *legea de evoluție dinamică a rețelei Hopfield discrete*. Putem rescrie relația vectorial:

$$\bar{S}(t+1) = \text{SGN} \times W \times \bar{S}(t) \quad (2-19)$$

Unde SGN este operatorul:

$$\text{SGN} \triangleq \begin{pmatrix} \text{sgn}(\cdot) & 0 & \dots & 0 \\ 0 & \text{sgn}(\cdot) & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \dots & \text{sgn}(\cdot) \end{pmatrix} \quad (2-20)$$

Din punct de vedere computațional, o astfel de rețea are o altă utilitate decât rețelele feedforward. O rețea feedforward poate să învețe un set de tipare duble. Un astfel de tipar este o pereche de vectori $\mathbf{i}=(i_1,\dots,i_n)$ și $\mathbf{d}=(d_1,\dots,d_m)$ unde n este numărul de intrări și m este numărul de ieșiri (vezi capitolul anterior). După învățare, prezentând la intrarea rețelei o componentă i dintr-un tipar vom obține la ieșire un vector \mathbf{o} foarte apropiat de vectorul \mathbf{d} care este pereche cu vectorul de intrare din tipar. Deoarece intrarea poate fi diferită de ieșire, rețelele feedforward se mai numesc și sisteme heteroasociative.

Rețelele Hopfield sunt rețele autoasociative. Ele pot învăța un set de tipare singulare. Un astfel de tipar poate fi un vector N dimensional cu componente 1 și -1 . Notăm baza de antrenare cu $L=\{X^1,\dots,X^p\}$. Funcționarea rețelei Hopfield este următoarea: se ia un vector N dimensional oarecare, se transformă componentele sale

în valori binare 1 și -1 și se setează neuronii rețelei cu aceste valori. Considerăm această stare ca fiind asociată momentului $t=0$ și o notăm cu $S(t_0)$. Aplicând legea dinamică de evoluție, rețeaua va trece printr-o serie de stări până când va ajunge într-o stare din care nu mai poate să evolueze. Notăm această stare cu $S(t_{stop})$. De obicei se poate observa că această stare este echivalentă cu unul din tiparele din baza de învățare. Se poate demonstra că tiparul atins la momentul t_{stop} este tiparul cel mai apropiat (considerând distanța Hamming ca metrică pentru spații binare N dimensionale) de starea inițială $S(t_0)$. Numărul de pași ($t=1,2,\dots,t_{stop}$) diferă de la caz la caz.

Pentru a obține această evoluție, Hopfield a inițializat ponderile rețelei (matricea W) aplicând o variantă a regulii de învățare Hebbiene [Hecht90]. De exemplu dacă am avea un singur tipar de învățare $X=(X_1,\dots,X_N)$ formula de calcul pentru ponderi ar fi:

$$w_{ij} = \frac{1}{N} X_i \cdot X_j \quad (2-21)$$

care semnifică următorul lucru:

- a. Dacă doi neuroni au starea 1, ponderea dintre ei va fi $1/N$ (pozitivă sau excitatorie).
- b. Dacă doi neuroni au stări diferite, ponderea este $-1/N$ (negativă sau inhibitorie).
- c. Dacă doi neuroni au starea -1, ponderea este tot excitatorie.

Dacă am avea mai multe tipare de învățare X^u ($u=1,p$) ponderile sunt obținute prin însumarea (superpoziția) pentru fiecare tipar a ponderilor obținute cu (2-21):

$$w_{ij} = \frac{1}{N} \sum_{u=1}^p X_i^u \cdot X_j^u \quad (2-22)$$

Mecanismul care asigură convergența rețelei către un tipar învățat poate fi explicat în felul următor. Se definește o funcție de energie a rețelei [Hertz&al91]:

$$H(t) = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N w_{ij} \cdot S_i(t) \cdot S_j(t) \quad (2-23)$$

Se demonstrează că legea de evoluție dinamică minimizează această funcție. Variația de la un moment discret de timp la altul este:

$$\Delta H = H(t+1) - H(t) = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N w_{ij} \cdot S_i(t+1) \cdot S_j(t+1) + \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N w_{ij} \cdot S_i(t) \cdot S_j(t) \quad (2-24)$$

Demonstrația se poate face în felul următor: considerăm că se modifică doar starea unui singur neuron i , $S_i(t+1) = -S_i(t)$. În acest caz (2-24) devine:

$$\Delta H = -\sum_i S_i(t+1) \sum_{j \neq i} w_{ij} S_j(t) + \sum_i S_i(t) \sum_{j \neq i} w_{ij} S_j(t) \quad (2-25)$$

Înlocuind $S_i(t+1)$ cu $-S_i(t)$ în (2-25), vom obține:

$$\Delta H = -\sum_i [-S_i(t)] \sum_{j \neq i} w_{ij} S_j(t) + \sum_i S_i(t) \sum_{j \neq i} w_{ij} S_j(t) = 2 \sum_i S_i(t) \sum_{j \neq i} w_{ij} S_j(t) \quad (2-26)$$

Dacă vom considera matricea \mathbf{W} ca fiind o matrice simetrică (și acest lucru este asigurat prin legea de învățare dată de (2-22)) și în plus având toate elementele de pe diagonala principală egale cu 0 (lucru care nu afectează nici învățarea, nici funcționarea) putem observa că:

$$h_i(t+1) = \sum_{j \neq i} w_{ij} S_j(t) \quad (2-27)$$

Formula finală pentru variația energiei va fi:

$$\Delta H = 2 \sum_i S_i(t) \cdot h_i(t+1) \quad (2-28)$$

Semnul intrării ponderate $h_i(t+1)$ va fi întotdeauna același cu semnul ieșirii neuronului $S_i(t+1)$ datorită funcției *signum*. Deci semnul termenului $S_i(t)h_i(t+1)$ va fi întotdeauna negativ pentru că $S_i(t+1) = -S_i(t)$. Energia rețelei descrește la orice modificare de stare în rețea, care este conformă cu legea dinamică de evoluție. Privind rețeaua dintr-o perspectivă fizică, ea reprezintă un sistem care se îndreaptă întotdeauna spre o stare de minim energetic.

Dar cel mai important fapt este următorul: *tiparele de învățare reprezintă minime pentru funcția de energie*. Aceasta poate fi reprezentată ca o suprafață în raport cu spațiul 2ⁿ dimensional al tuturor stărilor posibile și este simplu de demonstrat (pornind de la formulele energiei și ale legii dinamice de evoluție) că punctele \mathbf{X}^u ($u=1,p$) corespunzătoare pe această suprafață sunt puncte de minim, care mai sunt denumite și atractori.

O stare inițială oarecare $S(t_0)$ va duce la obținerea celui mai apropiat minim de pe suprafața de energie. Din păcate punctele de minim corespunzătoare tiparelor de învățare nu sunt singurele puncte de minim de pe această suprafață. În primul rând, se observă că pentru fiecare tipar de învățare \mathbf{X}^u ($u=1,p$) apare și un minim suplimentar numit și *minim opus* care este corespondent punctului $-\mathbf{X}^u$, care este imaginea inversată binar a vectorului de învățare. Atingerea unui astfel de minim are următoarea semnificație funcțională: starea inițială este cea mai diferită posibil de tiparul peste a cărui negată am ajuns. Un astfel de răspuns este de obicei destul de util pentru că ne indică foarte clar că starea inițială nu seamănă deloc cu tiparele învățate.

Slăbiciunea fundamentală a acestui model o reprezintă *stările mixate* (spurious states [Hertz&al91]). Ele apar atunci când avem mai mult de două tipare de învățare.

De exemplu pentru trei tipare X^1, X^2, X^3 vom avea:

$$X_i^{mix} = \text{sgn}(\pm X_i^1, \pm X_i^2, \pm X_i^3) \quad (2-29)$$

Avem 8 posibilități pentru toate combinațiile de semne, care vor duce la apariția a 8 puncte de minim local, care nu corespund nici unei stări identice cu unul dintre cele trei tipare de învățare. Deoarece prin legea de evoluție dinamică nu putem să mergem decât în sensul descreșterii energiei (nu putem "urca" pe suprafața de energie) orice apropiere în funcționarea rețelei către un minim local ne va conduce inevitabil la o fixare în această stare nedorită din care nu vom mai putea ieși.

Concluzia pe care o putem trage este ca rețeaua Hopfield este un model computațional care nu funcționează întotdeauna [Hecht90] [Haykin94]. Avem nevoie de modele mai puternice care să asigure o funcționalitate mai robustă.

2.2.3. Modelul Hopfield continuu

Sugestia lui Hopfield de a folosi în scopuri computaționale sisteme inspirate din modele fizice, a fost studiată de specialiștii în modelarea proceselor fizice. Una din cele mai interesante abordări, este [Pineda87]. Alternativa propusă este de a studia o rețea neuronală (cu topologie Hopfield - toți neuronii sunt legați între ei) ca *sistem dinamic*. Modelul matematic al sistemelor dinamice a fost folosit cu succes în magnetism, supraconductibilitate, mecanica fluidelor etc. și este bine pus la punct din punct de vedere teoretic [Wiggins90]. O rețea neuronală privită prin prisma sistemelor dinamice neliniare (numim în acest caz rețeaua *sistem neurodinamic*) are trei proprietăți principale:

- i. *Este un sistem cu foarte multe grade de libertate.* Dimensiunea spațiului în care evoluează starea unui sistem neurodinamic este dată de numărul de neuroni ai sistemului. Vectorul de stare al sistemului are dimensiunea N , fiecare componentă reprezentând starea momentană a unui neuron. În cazul analizat vom avea starea reprezentată de un număr real, spre deosebire de modelul Hopfield unde aveam doar două posibilități discrete 1 și -1. Starea sistemului va fi un punct în spațiul euclidian N dimensional. Concluzia interesantă ce se poate trage relativ la această proprietate este că puterea de calcul a unui astfel de sistem este dată de numărul de grade de libertate. Simplificând, am putea spune că creierul uman este o mașină computațională deosebit de puternică pentru că are aproximativ $2 \cdot 10^{11}$ grade de libertate.
- ii. *Este un sistem neliniar.* Ecuațiile care descriu starea sistemului în timp sunt obligatoriu ecuații neliniare. Conform demonstrațiilor lui Minsky [Minsky&Papert69] și Smolensky [Smolensky86] un sistem neurodinamic liniar nu este o mașină computațională adecvată. Ea poate rezolva doar un număr limitat de probleme simple.
- iii. *Este un sistem disipativ.* Sistemele dinamice disipative au proprietatea de stabilitate

asimptotică globală. Ceea ce înseamnă că sistemul se va opri într-o stare stabilă indiferent de starea inițială din care este pornit. Modelul Hopfield care prin legea sa de evoluție dinamică minimizează funcția de energie a sistemului, este un sistem disipativ (dar discret).

Un sistem dinamic disipativ, neliniar, continuu, cu N grade de libertate, este de obicei descris printr-un set de N ecuații diferențiale cuplate:

$$\frac{dx_i}{dt} = G_i(\bar{\mathbf{x}}) \quad (i=1, N) \quad (2-30)$$

Unde x_i sunt componentele unui vector de stare N dimensional \mathbf{x} . Fiecare stare este cuplată de celelalte stări prin funcția neliniară G_i . Traiectoria vectorului \mathbf{x} în spațiul euclidian N dimensional se numește flux N dimensional. Setul de ecuații asigură convergența acestui flux într-o stare stabilă [Wiggins90].

Sistemele dinamice disipative pot avea o comportare relativ complicată în raport cu stările stabile pe care le pot atinge. O stare stabilă poate fi:

- O orbită într-un spațiu M dimensional ($M \ll N$) numită și atractor oscilant.
- O orbită într-un spațiu cu dimensionalitate fracționară (Hausdorff-Besicovitch) numită și atractor straniu [Wiggins90].
- Un punct în spațiul N dimensional.

Din punct de vedere computațional ne interesează doar ultimul caz și vom fi obligați să alegem doar modele matematice care asigură acest tip de convergență.

Dacă alegem setul de ecuații:

$$\frac{dx_i}{dt} = G_i(\bar{\mathbf{W}}_i, \mathbf{I}, \bar{\mathbf{x}}) \quad (i=1, N) \quad (2-31)$$

Unde $\bar{\mathbf{W}}_i$ este un vector ($w_{i1}, w_{i2}, \dots, w_{iN}$) care reprezintă constantele de cuplare între mărimea x_i și restul de mărimi de stare x_j ($j=1, N$), \mathbf{I} reprezintă un set de parametrii care influențează din exterior sistemul, va fi asigurată convergența sistemului într-un punct. Aceste puncte sunt soluții ale sistemului neliniar de ecuații în x_i ($i=1, N$):

$$G_i(\bar{\mathbf{W}}_i, \mathbf{I}, \bar{\mathbf{x}}) = 0 \quad (i=1, N) \quad (2-32)$$

Vom nota aceste soluții cu \mathbf{x}^{inf} .

Un astfel de sistem de ecuații se rezolvă iterativ, plecând de la un vector de stare inițial \mathbf{x}^0 , simulând practic evoluția sistemului în timp până într-o stare stabilă \mathbf{x}^{inf} . Comportamentul unui astfel de sistem neurodinamic poate fi exploatat pentru a executa două tipuri de activități:

- autoasociere (avem atunci o rețea Hopfield continuă)
- heteroasociere temporală

Ce este heteroasocierea temporală? O rețea feedforward obișnuită este capabilă de heteroasociere, dar aceasta se petrece independent de timp. Asocierea intrare-ieșire se face practic instantaneu. Prezentăm o intrare rețelei și după propagarea înainte (teoretic instantanee) a semnalelor, obținem ieșirea rețelei de pe ultimul strat. Dacă am avea o secvență de intrări ordonate în timp, în cazul cel mai general o funcție vectorială de timp care generează un vector n dimensional ($n \ll N$ - numărul de grade de libertate), $f_m(t)=x'$, am putea selecta dintre cele N stări ale sistemului, n stări pe care să le forțăm pe valorile (x'_1, \dots, x'_n) . Neuronii asociați unor astfel de stări vor reprezenta neuroni de intrare în rețea. La fel, am putea considera o submulțime O a componentelor vectorului de stare al sistemului $x=(x_1, \dots, x_N)$, $O=\{x''_1, x''_2, \dots, x''_m\}$ ($m \ll N$) ca reprezentând un vector de ieșire x'' . Mai simplu spus, am selectat dintre neuronii rețelei o mică parte, pentru a avea ca la rețelele feedforward o intrare și o ieșire. Restul de neuroni îi considerăm neuroni ascunși și starea lor nu ne interesează în timpul funcționării rețelei.

Discutând în termeni de sisteme dinamice, dacă funcția $f_m(t)$ este o funcție continuă (descrie un flux n dimensional fără discontinuități) atunci secvența continuă de valori obținută la ieșire până la atingerea unei stări stabile va reprezenta graficul unei funcții vectoriale $f_{out}(t)=x''$. Se pune problema dacă ieșirea, care reprezintă un flux m dimensional, este continuă. În cazul în care cele n stări pe care le considerăm de intrare sunt forțate pe niște valori impuse se pot întâmpla două evoluții nedorite.

În primul rând, forțarea unor stări poate duce la imposibilitatea atingerii stărilor de echilibru de tip punctiform. Am spus că un astfel de fenomen este exclus la clasa de sisteme descrisă, dar stabilitatea punctiformă este respectată doar dacă sistemul este lăsat să evolueze liber, fără a modifica mărimile de stare.

În al doilea rând, chiar dacă am obține la ieșire doar stări stabile punctiforme, succesiunea lor în timp nu ar reprezenta un flux m dimensional continuu. S-a demonstrat că un astfel de sistem are în mod cert o comportare discontinuă [Pineda87]. Sistemul poate fi influențat din afară și prin alte căi decât alterarea forțată a stărilor care sunt considerate de intrare. În ecuația (2-31) s-a prevăzut posibilitatea introducerii unei influențe exterioare prin setul de mărimi I . Când Hopfield a încercat să construiască un model continuu plecând de la modelul său discret a gândit această mărime ca fiind echivalentă cu valorile de prag asociate neuronilor (threshold). Fiecare neuron (stare a sistemului) va avea asociată o astfel de mărime care este constantă pe timpul funcționării rețelei și care trebuie calculată la învățare. Pentru funcționarea rețelei continue, Hopfield a propus următorul model descris prin setul de ecuații diferențiale:

$$\frac{du_i}{dt} = -u_i + g(u_i + I_i) \quad (i=1, N) \quad (2-33)$$

Unde u_i sunt intrările asociate neuronilor. O astfel de ecuație poate fi interpretată în felul următor:

Variația infinitesimală a mărimii de stare este egală cu diferența dintre starea curentă u_i și starea următoare $g(h_i)$, considerând că timpul necesar aplicării funcției de transfer este egal cu cuanta infinitesimală dt .

Valoarea I_i va fi interpretată ca o valoare de prag, care va decide prin depășirea ei, în cazul fiecărui neuron dacă acesta este aprins sau nu. Această decizie este necesară pentru că de obicei la o rețea neuronală care are funcții de transfer reale și continue (de exemplu sigmoidă) nu ne interesează valorile reale pe care le au la ieșire neuronii, ci faptul dacă ei sunt aprinși sau nu.

Pineda a dat acestor mărimi sensul de influență exterioară. Putem considera că intrarea într-un astfel de sistem este dată de setul I_i . Alegând o stare inițială x_0 și aplicând un vector de intrare I de dimensiune n (numai unui număr limitat de stări) sistemul, neavând alterate stările din afară, va putea ajunge într-o stare stabilă punctiformă. Putem ilustra (figura 2-8) diferența care apare între o rețea Hopfield continuă (numită uneori și asociatorul continuu) și un astfel de sistem imaginat de Pineda (numit de el filtru adaptiv).

Dacă în figura 2-8a se observă că valorile I și W sunt fixate, intrarea este o stare x_0 , iar ieșirea este o stare x^{mf} , în cazul ilustrat în 2-8b se va fixa starea x_0 și W , $I(t)$ va fi intrarea, ieșirea fiind o fracțiune m -dimensională din fluxul $x^{mf}(t)$. Dar pentru că scopul pe care ni l-am propus este heteroasocierea temporală, trebuie analizat cum se comportă un astfel de sistem în cazul în care intrarea I se modifică în timp. Înainte de a studia filtrul adaptiv al lui Pineda din punctul de vedere al rețelelor neuronale (implementare, reguli de învățare) este important de făcut o analiză calitativă mai profundă a comportării acestui sistem. Pentru început trebuie explicate în detaliu câteva din proprietățile atractorilor (punctelor de stabilitate punctiformă). Dacă sistemul descris de ecuațiile (2-31) este lăsat să ruleze liber cu $I=0$ (fără perturbații externe) el se va opri într-un punct pe care îl notăm cu x^0 . De obicei pentru un sistem cu multe grade de libertate vom avea și un număr mare de atractori (nr. atractori $> N$). Atractorul în care sistemul ajunge la atingerea stării de echilibru este întotdeauna determinat de starea inițială x_m (ca și la rețeaua Hopfield discretă). Fiecare atractor are asociată o zonă din spațiul euclidian N dimensional, numită *bazin de atracție*.

Definiție

Toate punctele x care fiind considerate ca stări inițiale x_m duc la evoluția

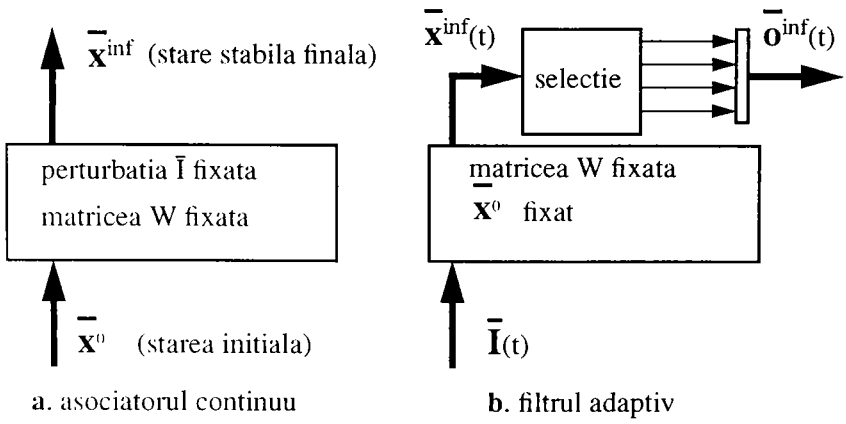


figura 2-8.

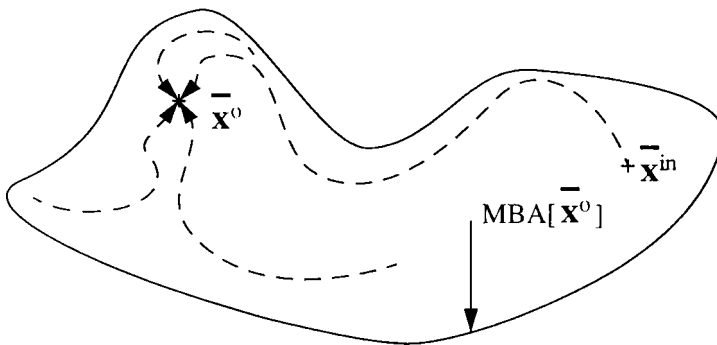


figura 2-9.

sistemului într-un atractor x^0 , formează bazinul de atracție al punctului x^0 .

Aceste bazine pot avea discontinuități (pot fi formate din zone disjuncte ale spațiului euclidian N dimensional) sau pot avea margini de tip fractal. Vom neglija aceste complicații posibile pentru că ele nu sunt relevante pentru analiza noastră. Figura 2-9 reprezintă un astfel de bazin de atracție. Notăm cu $MBA[x^0]$ marginea bazinului de atracție. Toate punctele din interiorul conturului aparțin bazinului de atracție, cele din exterior reprezentând un alt bazin de atracție. Dacă vom altera starea staționară x^0 a sistemului din acest bazin, aplicând o perturbație exterioară I , va apărea un alt atractor notat cu $x^{mf}(I)$ și un alt bazin de atracție de o altă formă, ca în figura 2-10. Noul atractor poate conține (ca și în figura 2-10) pe x^0 în bazinul său de atracție. Acest fenomen depinde de x^0 , de forma bazinului său de atracție și de I . Este foarte indicat să se întâmple aceasta, pentru că o perturbație I care "pierde" starea x^0 din bazinul de atracție al stării stabile $x^{mf}(I)$ va duce la evoluții nedorite, ceea ce se va vedea și în continuare. Dacă perturbația exterioară I variază continuu și lent (adiabatic) în timp, ne putem imagina următoarea evoluție (timpul este în intervalul continuu $[t_0, t_{stop}]$):

i. Sistemul este stabil și se află într-un atractor x^0 și este perturbat cu valoarea $I(t_0)$ pe care o notăm cu I_0 .

ii. Într-un interval de timp foarte scurt (practic instantaneu) sistemul se va relaxa și va ajunge în atractorul $x^{mf}(I_0)$. Putem considera că în acest timp perturbația $I(t)$ nu își va modifica valoarea (va rămâne tot I_0).

iii. Sistemul este adus din nou în starea x^0 .

iv. Se aplică o nouă valoare a perturbației I_1 , care este foarte puțin diferită de I_0 .

v. Sistemul va evolua într-o stare stabilă $x^{mf}(I_1)$.

vi. Se va repeta secvența iii-iv-v pentru o serie de perturbații $I(t)$ cu $t=t_1, t_2, \dots, t_{stop}$ unde I_1 și I_{j+1} sunt două puncte foarte apropiate în spațiul euclidian N dimensional. La fel, $t_j - t_{j-1}$ este un interval de timp foarte mic, dar mult mai mare decât timpul necesar relaxării sistemului într-o stare stabilă.

De obicei se va observa că punctele $x^{mf}[I(t)]$ vor reprezenta un flux N dimensional continuu [Wiggins90]. Cu o singură observație. Considerăm că $MBA[x^{mf}(I_j)]$ își va schimba forma în timp ca în figura 2-11. Dacă $MBA[x^{mf}(I_j)]$ nu traversează punctul x^0 , în mod cert fluxul N dimensional $x^{mf}(I_j)$ este continuu. Dacă datorită perturbației $I(t)$ marginea bazinului de atracție al stării stabile curente traversează punctul x^0 în momentul t_k , se va produce în mod cert o discontinuitate, ca în figura 2-12. O soluție pentru evitarea acestei situații este ca pasul iii din algoritmul care generează evoluția sistemului să fie înlocuit cu următorul:

iii'. Sistemul este lăsat în starea stabilă $x^{mf}[I(t_j)]$ când se aplică noua perturbație I_{j+1} .

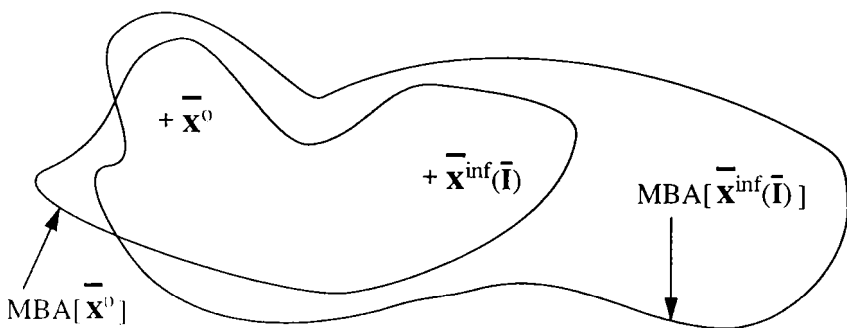


figura 2-10.

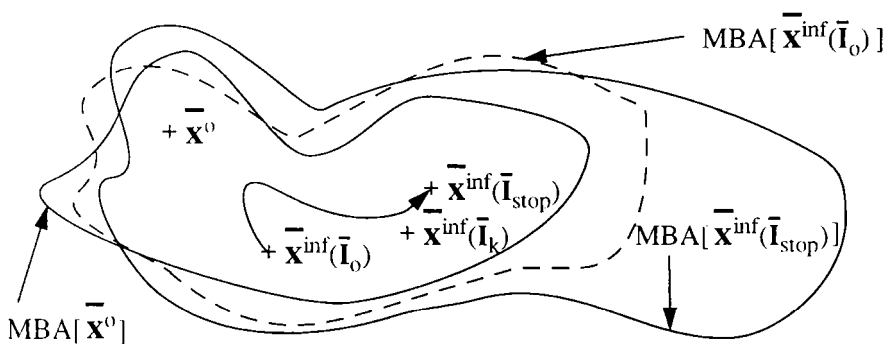


figura 2-11.

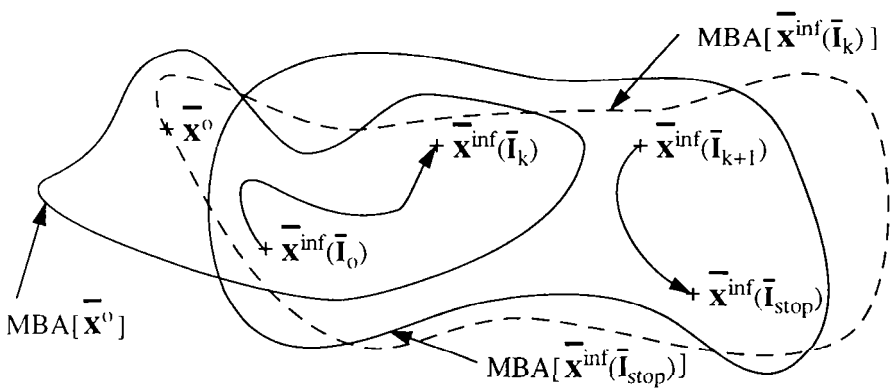


figura 2-12.

Din păcate această îmbunătățire nu garantează că fluxul $x^{inf}[I(t)]$ va fi întotdeauna continuu. Dar, plecând de la această problemă se poate face o observație mult mai relevantă. Se observă din figura 2-12 că pentru perturbația I_k avem doi atractori diferiți. Dar dacă lucrăm într-o zonă apropiată de o stare x^0 care este un attractor al sistemului liber numai unul dintre aceștia poate fi atins. Numim această zonă *regiune de operare*. Alegerea acesteia este foarte importantă pentru că de ea depinde funcționarea normală a sistemului.

Scopul inițial pe care mi l-am propus este de a mapa peste o rețea neuronală (modelată și implementată de exemplu de un astfel de sistem neurodinamic), un spațiu semantic (conceptual) unde starea unor neuroni individuali sau a unor grupe de neuroni sau configurații ale stărilor acestora să reprezinte concepte (valori semantice). Este normal ca o regiune de operare să fie într-o relație de similaritate cu o regiune dintr-un spațiu conceptual mapat peste spațiul euclidian N dimensional.

2.2.4. Algoritm de învățare pentru un sistem neurodinamic

Modelul prezentat anterior se pretează cel mai bine la o implementare analogică și nu una discretă. Este posibilă și o implementare discretă care să simuleze evoluția sistemului continuu descris de ecuațiile (2-31). Ea va fi mult mai lentă dar poate da rezultate mult mai bune din punct de vedere funcțional decât alte tipuri de rețele. Deoarece factorul timp nu ne interesează în contextul acestei discuții, vom prezenta algoritmul de învățare fără a intra în detalii de implementare. Algoritm de învățare ce va fi prezentat este continuu și o versiune discretizată a acestuia va fi prezentată în subcapitolul următor. Vom considera sistemul neurodinamic descris de ecuațiile:

$$\frac{dx_i}{dt} = -x_i + g_i \left(\sum_{j=1}^N w_{ij} \cdot x_j + I_i \right) \quad (i=1, N) \quad (2-34)$$

Fiecare mărime de stare x_i ($i=1, N$) este asociată cu starea momentană a unui neuron. Am convenit să privim această rețea recurentă ca având o intrare și o ieșire. O mică parte din neuronii rețelei vor reprezenta neuronii de intrare. Vom nota mulțimea lor cu A ($\text{card}(A)=n$). Mulțimea neuronilor de ieșire o vom nota cu O ($\text{card}(O)=m$). Un neuron poate fi în același timp și de intrare și de ieșire. Mulțimea H a neuronilor "ascunși" reprezintă factorul care permite heteroascierea temporală. Pentru o funcționare cât mai bună este indicat ca numărul neuronilor din H să fie mult mai mare decât numărul neuronilor din A și O . Intrarea o vom nota ca și la rețelele Hopfield cu X_i ($i=1, N$). Este util să definim funcția:

$$\Theta_{i\bullet} = \begin{cases} 1 & \text{daca neuronul } i \text{ apartine multimii } \Phi \\ 0 & \text{altfel} \end{cases} \quad (2-35)$$

Intrarea în rețea este dată în (2-34) de vectorul I . Dacă unitatea i ($i=1, N$) este de intrare atunci $I_i=X_i$, altfel $I_i=0$. Acest lucru se poate ilustra prin relația:

$$I_i = X_i \cdot \Theta_{iA} \quad (i=1, N) \quad (2-36)$$

Se observă că doar n astfel de relații contează.

Vom presupune inițial că intrările X_i sunt constante în timp. La ieșire vom dori să obținem un set de valori notate d_i (valori dorite). Funcționarea dorită a rețelei este dată de relația:

$$d_i = x_i^m \cdot \Theta_{iO} \quad (i=1, N) \quad (2-37)$$

La fel, doar m astfel de relații sunt nule. Pentru a obține o astfel de comportare a sistemului, trebuie ajustate ponderile w_{rs} ($r,s=1,N$) în mod corespunzător. Vom aplica metoda utilizată în algoritmul de propagare înapoi a erorilor de la rețelele feedforward. Se va minimiza o funcție E care măsoară distanța euclidiană între poziția atractorului dorit și cea a atractorului obținut. Funcția este definită în modul următor:

$$E = \frac{1}{2} \sum_{i=1}^N J_i^2 \quad (2-38)$$

unde:

$$J_i = (d_i - x_i^m) \cdot \Theta_{iO} \quad (i=1, N) \quad (2-39)$$

Vom considera că de data un set de ecuații similare cu (2-34) ale cărui variabile sunt elementele matricii W . Putem spune că matricea W definește un punct în spațiul $N \times N$ dimensional și trebuie deplasat acest punct astfel încât să minimizăm funcția E . Traiectoria din spațiul N^2 dimensional pe care trebuie să ne deplasăm este paralelă cu gradientul lui E , dar în direcția opusă. Pentru un singur element w_{rs} al matricii W ecuația acestei mișcări este:

$$\tau_w \frac{dw_{rs}}{dt} = - \frac{\partial E}{\partial w_{rs}} \quad (r, s=1, N) \quad (2-40)$$

Unde τ_w este o constantă de timp (foarte mică) care definește deplasarea (obligatoriu mică) cu care w_{rs} se modifică.

Ecuațiile (2-40) descriu regula de învățare a unui astfel de sistem. Pentru a obține performanțe mai bune am putea folosi metoda momentului de inerție [Rumelhart&al86] ori altele care implică și calculul derivatei de ordinul doi [Pearlmutter90], dar aceste îmbunătățiri țin mai mult de implementare. Componentele gradientului se calculează prin derivatele parțiale:

$$\tau_w \frac{dw_{rs}}{dt} = - \frac{1}{2} \frac{\partial}{\partial w_{rs}} \sum_{k=1}^N J_k^2 = \sum_{k=1}^N J_k \frac{\partial x_k^m}{\partial w_{rs}} \quad (r, s=1, N) \quad (2-41)$$

Derivata lui x_k^{inf} în raport cu w_{rs} se poate calcula ținând cont de faptul că \mathbf{x}^{inf} este o soluție a sistemului de ecuații (2-32) deci satisface ecuația algebrică neliniară:

$$\mathbf{x}_i^{\infty} = \mathbf{g}_i \left(\sum_{j=1}^N w_{ij} \cdot \mathbf{x}_j^{\infty} + I_i \right) \quad (i=1, N) \quad (2-42)$$

Derivând în raport cu w_{rs} ambele părți ale ecuației vom obține:

$$\frac{\partial \mathbf{x}_i^{\infty}}{\partial w_{rs}} = \mathbf{g}'_i \left(\sum_{j=1}^N w_{rj} \cdot \mathbf{x}_j^{\infty} + I_r \right) \cdot \sum_{j=1}^N \left\{ \frac{\partial w_{ij}}{\partial w_{rs}} \mathbf{x}_j^{\infty} + w_{ij} \frac{\partial \mathbf{x}_j^{\infty}}{\partial w_{rs}} \right\} \quad (i=1, N) \quad (2-43)$$

și notăm cu:

$$\mathbf{u}_r^{\infty} = \sum_{j=1}^N w_{rj} \cdot \mathbf{x}_j^{\infty} + I_r \quad (r=1, N) \quad (2-44)$$

Pentru simplificarea notației vom folosi funcția delta a lui Kronecker:

$$K_{ij}^{\Delta} = \begin{cases} 1 & \text{daca } i=j \\ 0 & \text{daca } i \neq j \end{cases} \quad (2-45)$$

Se observă că:

$$\frac{\partial w_{ij}}{\partial w_{rs}} = K_{ir}^{\Delta} \cdot K_{js}^{\Delta} \quad (2-46)$$

Cu ajutorul acestei funcții putem rescrie (2-43) în felul următor:

$$\sum_{j=1}^N K_{ij}^{\Delta} \frac{\partial \mathbf{x}_j^{\infty}}{\partial w_{rs}} = \mathbf{g}'_i (\mathbf{u}_r^{\infty}) \left\{ K_{ir}^{\Delta} \cdot \mathbf{x}_s^{\infty} + \sum_{j=1}^N w_{ij} \frac{\partial \mathbf{x}_j^{\infty}}{\partial w_{rs}} \right\} \quad (i=1, N) \quad (2-47)$$

Aducând toate derivatele în partea stângă obținem:

$$\sum_{j=1}^N \left[K_{ij}^{\Delta} - \mathbf{g}'_i (\mathbf{u}_r^{\infty}) \cdot w_{ij} \right] \frac{\partial \mathbf{x}_j^{\infty}}{\partial w_{rs}} = K_{ir}^{\Delta} \cdot \mathbf{g}'_i (\mathbf{u}_r^{\infty}) \cdot \mathbf{x}_s^{\infty} \quad (i=1, N) \quad (2-48)$$

Notăm cu:

$$\mathbf{L}_{ij} = K_{ij}^{\Delta} - \mathbf{g}'_i (\mathbf{u}_r^{\infty}) \cdot w_{ij} \quad (2-49)$$

Și vom obține ecuația matriceală:

$$\mathbf{L} \cdot \frac{\partial \bar{\mathbf{x}}}{\partial \mathbf{w}_{rs}} = K_{ir}^{\Delta} \cdot g'_r(u_r^{\bar{\mathbf{w}}}) \cdot \mathbf{x}_s^{\bar{\mathbf{w}}} \quad (2-50)$$

Dacă vom înmulți această ecuație cu (\mathbf{L}^{-1}) obținem:

$$\frac{\partial \bar{\mathbf{x}}}{\partial \mathbf{w}_{rs}} = (\mathbf{L}^{-1}) \cdot K_{ir}^{\Delta} \cdot g'_r(u_r^{\bar{\mathbf{w}}}) \cdot \mathbf{x}_s^{\bar{\mathbf{w}}} \quad (2-51)$$

Putem rescrie această formulă ca și set de ecuații și pentru a avea aceeași notație ca și în setul de ecuații (2-41) vom înlocui indicele i cu k :

$$\frac{\partial x_k^{\bar{\mathbf{w}}}}{\partial w_{rs}} = (\mathbf{L}^{-1})_{kr} \cdot g'_r(u_r^{\bar{\mathbf{w}}}) \cdot \mathbf{x}_s^{\bar{\mathbf{w}}} \quad (k=1, N) \quad (2-52)$$

Înlocuind \mathbf{L} vom obține:

$$\tau_w \frac{dw_{rs}}{dt} = g'_r(u_r^{\bar{\mathbf{w}}}) \cdot \mathbf{x}_s^{\bar{\mathbf{w}}} \cdot \sum_{k=1}^N J_k \cdot (\mathbf{L}^{-1})_{kr} \quad (2-53)$$

Care este tocmai rezultatul care ne interesează. Dacă vom privi mai atent această formulă vom observa asemănarea cu formula care reglează modificarea ponderilor la rețelele feedforward prin algoritmul de învățare prin backpropagation.

Din păcate, pentru a efectua ajustările tuturor ponderilor va trebui să calculăm inversa unei matrici, lucru care este extrem de intensiv din punct de vedere computațional. O metodă care evită acest calcul este următoarea. Notăm cu:

$$y_r^{\bar{\mathbf{w}}} = g'_r(u_r^{\bar{\mathbf{w}}}) \cdot \sum_{k=1}^N J_k \cdot (\mathbf{L}^{-1})_{kr} \quad (r=1, N) \quad (2-54)$$

Și vom obține:

$$\tau_w \frac{dw_{rs}}{dt} = y_r^{\bar{\mathbf{w}}} \cdot \mathbf{x}_s^{\bar{\mathbf{w}}} \quad (r, s=1, N) \quad (2-55)$$

Care pot fi interpretate ca ecuațiile unui sistem dinamic. Rescriem (2-54) în felul următor:

$$\sum_{r=1}^N L_{rk} \frac{u_r^{\bar{\mathbf{w}}}}{g'_r(u_r^{\bar{\mathbf{w}}})} = J_k \quad (k=1, N) \quad (2-56)$$

Înmulțind cu $g_k(u_k^{\bar{\mathbf{w}}})$ și înlocuind cu forma explicită a lui \mathbf{L} (2-49) obținem:

$$0 = -y_k + g_k(u_k) \left\{ \sum_{r=1}^N w_{rk} \cdot y_r + J_k \right\} \quad (k=1, N) \quad (2-57)$$

Se poate observa că soluțiile y_k^{inf} ale acestui sistem liniar de ecuații sunt atractori ai sistemului dinamic:

$$\frac{dy_k}{dt} = -y_k + g_k(u_k) \left\{ \sum_{r=1}^N w_{rk} \cdot y_r + J_k \right\} \quad (k=1, N) \quad (2-58)$$

Se poate observa în metoda de mai sus, calculul unei matrici inverse prin metoda relaxării. Acest algoritm de învățare poate fi implementat avântos și relativ simplu într-o modalitate analogică (metoda digitală e mai complicată și mult mai lentă [Williams&Zipser90]).

2.2.5. Algoritm de propagare înapoi a erorilor pentru rețele recurente discrete

Deoarece o implementare analogică este o soluție doar pentru rețele dedicate, utilizate cu precădere în controlul automat, care folosesc un număr limitat de neuroni, trebuie găsită o modalitate de a implementa metoda de învățare prezentată anterior pe un calculator cu arhitectură clasică, mergând pe un alt principiu, pentru a evita viteza foarte scăzută de procesare. Primul pas pentru realizarea unei astfel de metode de antrenare este discretizarea factorului timp. Vom considera timpul ca o succesiune $t=1,2,\dots,stop$. Substituirea unei mărimi continue cu una discretă va modifica evident și modelul matematic prezentat anterior, pe care va trebui să-l reevaluăm în această nouă interpretare. Ronald Williams și David Zipser au dezvoltat un model matematic discret plecând de la modelul lui Pineda, reușind să stabilească și pașii unui algoritm de învățare foarte eficient [Williams&Zipser89] [Williams&Zipser90] [Williams&Peng90].

Al doilea pas este descrierea unei arhitecturi adecvate pentru rețeaua neuronală, care să permită o mai ușoară interpretare a secvențelor de timp discrete. Pentru aceasta, vom considera că rețeaua are două straturi. Primul, un strat de calcul, format din N neuroni care au o funcție de transfer neliniară $g(h_i)$, unde h_i este suma ponderată a intrărilor. Pentru simplificare vom considera că toți neuronii au aceeași funcție de transfer (de exemplu sigmoidă) deși acest lucru nu este o restricție, fiecare neuron sau grup de neuroni putând avea funcția sa particulară de transfer $g_i(h_i)$. Al doilea strat, este format doar din neuroni de distribuție (fan-out) care trimite fiecare către toți cei N neuroni din stratul de calcul semnalul pe care îl primește la intrare. Acest semnal provine de la unul din neuronii de calcul. De exemplu, putem considera că neuronul i din stratul de calcul își va trimite semnalul său de ieșire către neuronul i' din stratul de distribuție.

Pentru a defini intrarea și ieșirea din această rețea vom proceda în felul următor. Vom considera un număr de m ($m \ll N$) neuroni de pe stratul de calcul ca fiind neuroni de ieșire, construind vectorul de ieșire o , semnalul lor de ieșire fiind transmis și neuronilor corespunzători din stratul de distribuție. La fel, putem considera și un număr de n neuroni din stratul de distribuție ca fiind neuroni de intrare, dar pentru o mai bună claritate este preferabil să adăugăm n neuroni suplimentari stratului de distribuție. Aceștia vor avea la intrare vectorul x de dimensiune n care va reprezenta influența externă asupra acestui sistem. Pentru o mai bună funcționare se recomandă adăugarea unui neuron pe stratul de distribuție care va avea intrarea și ieșirile întotdeauna forțate pe -1. Existența acestui neuron va duce la crearea factorului de prag pentru neuronii de calcul. Arhitectura generală a unei astfel de rețele este prezentată în figura 2-13.

Se observă că va rezulta o matrice de ponderi $N \times (N+n)$ dimensională. Vom nota cu $L=N+n$ numărul total de neuroni din stratul de distribuție. Pentru a avea o notație cât mai clară vom considera că toate ieșirile din stratul de calcul sunt V_i ($i=1, N$), iar intrările în acești neuroni sunt z_j ($j=1, L$). Pentru a introduce și factorul timp, considerăm că pentru un neuron oarecare i din stratul de calcul este necesară o cantitate de timp de la r la $r+1$ pentru a genera ieșirea V_i . Această ieșire este instantaneu preluată de cei N neuroni corespunzători din stratul de distribuție, care tot în mod instantaneu, o transmit către toți neuronii de calcul. Se observă că la trecerea fiecărei cuante discrete de timp, neuronii rețelei execută în mod sincron un pas de calcul. Relația matematică pe care o folosim pentru a ilustra funcționarea rețelei este următoarea:

$$V_i(t) = g(h_i(t)) = g\left(\sum_{j=0}^L w_{ij} \cdot z_j(t-1)\right) \quad (i=1, N) \quad (2-59)$$

Recurența rețelei este dată de faptul că intrările z_j de pe stratul de distribuție provin în parte de la neuronii de pe stratul de calcul. Numai n neuroni de pe stratul de distribuție primesc intrarea din exterior (vectorul x). Evident, această influență exterioară trebuie să fie sincronă cu funcționarea rețelei, deci mărimea x , care în mod normal este o funcție de timp și trebuie să fie eșantionată la momentele $t=1, 2, \dots, stop$.

O observație importantă este că discontinuități prea mari ale lui x pot duce la fenomenul de discontinuitate a fluxului de ieșire prezentat la rețelele continue. Mărimile z_j pot fi definite în felul următor:

$$z_j(t) = \begin{cases} -1 & \text{dacă } j=0 \\ x_j(t) & \text{dacă } j=1, n \\ V_{j-n}(t) & \text{dacă } j=n+1, L \end{cases} \quad (j=0, L) \quad (2-60)$$

iar ieșirea rețelei poate fi dată prin relația:

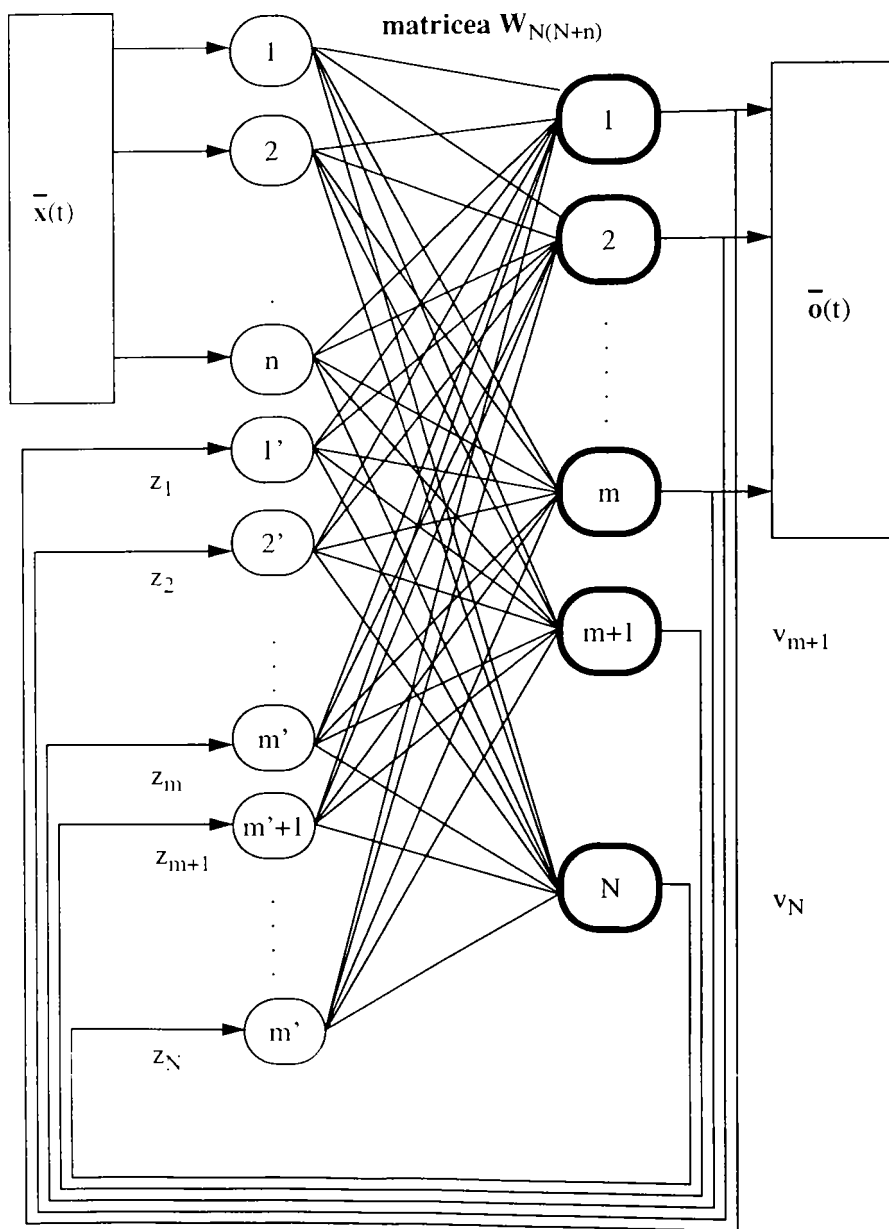


figura 2-13.

$$o_i(t) = V_i(t) \quad (i=1, m) \quad (2-61)$$

Aceste trei relații (2-59, 2-60, 2-61) sunt suficiente pentru a construi un algoritm care să simuleze funcționarea acestui tip de rețea. Suplimentar, mai putem scrie:

$$h_i(t) = \sum_{j=0}^L w_{ij} \cdot z_j(t-1) \quad (i=1, N) \quad (2-62)$$

Deoarece rețeaua are o întârziere de o cantă discretă de timp de la momentul aplicării intrării până la generarea ieșirii, apare următoarea problemă. Dacă la momentul 1 prezentăm rețelei intrarea $x(1)$, abia la momentul 2 vom obține o ieșire utilă $o(2)$. Pentru o pornire normală este indicată următoarea procedură. Deplasăm intrarea cu o cantă de timp înapoi ($t=0, 1, 2, \dots, stop-1$) și vom obține ieșirile $o(1), o(2), \dots, o(stop)$. Se recomandă ca la momentul $t=0$ toate ieșirile V să fie forțate pe zero.

În concluzie, se poate observa că dacă vom prezenta rețelei o secvență de vectori ordonați în timp $x(0), x(1), x(2), \dots, x(stop-1)$ vom obține la ieșire secvența de vectori $o(t), t=1, \dots, stop$. Vom nota aceste secvențe cu \mathbf{X} și \mathbf{O} . O pereche (\mathbf{X}, \mathbf{D}) pe care dorim ca rețeaua să o poată reproduce cât mai fidel va reprezenta un tipar de învățare. O astfel de rețea poate fi învățată mai multe astfel de tipare $\{(\mathbf{X}^s, \mathbf{D}^s) s=1, p\}$, adică o bază de învățare.

Învățarea trebuie să aducă elementele matricii \mathbf{W} într-o stare în care \mathbf{O}^s trebuie să fie cât mai aproape de \mathbf{D}^s pentru intrarea \mathbf{X}^s , pentru oricare $s=1, p$. Regula de învățare se deduce într-un mod asemănător cu regula backpropagation și regula de la rețelele recurente continue. Considerăm eroarea pătratică medie ca fiind o funcție de ponderi:

$$E(\mathbf{W}) = \frac{1}{2} \sum_{s=1}^p \sum_{t=1}^{stop} \sum_{i=1}^m [d_i^s(t) - o_i^s(t)]^2 \quad (2-63)$$

Pentru un singur tipar vom avea:

$$E^s(\mathbf{W}) = \frac{1}{2} \sum_{t=1}^{stop} \sum_{i=1}^m [d_i^s(t) - o_i^s(t)]^2 \quad (s=1, p) \quad (2-64)$$

Pentru fiecare tipar s vom putea calcula o ajustare a ponderii w_{ij} care se aplică în felul următor:

$$w_{ij}^{nou} = w_{ij}^{vechi} + \Delta w_{ij}^s \quad (s=1, p) \quad (2-65)$$

Pentru a calcula ajustarea, aplicând din nou metoda gradientului, obținem:

$$\Delta w_{ij}^s = -\eta \frac{\partial E^s}{\partial w_{ij}} \quad (2-66)$$

Derivata parțială a erorii se calculează în felul următor:

$$\begin{aligned} \frac{\partial E^s}{\partial w_{ij}} &= \frac{1}{2} \sum_{t=1}^{stop} \sum_{k=1}^m \frac{\partial [d_k^s(t) - o_k^s(t)]^2}{\partial w_{ij}} = \\ &= -\frac{1}{2} \sum_{t=1}^{stop} \sum_{k=1}^m 2 [d_k^s(t) - o_k^s(t)] \frac{\partial o_k^s(t)}{\partial w_{ij}} = \\ &= -\sum_{t=1}^{stop} \sum_{k=1}^m \delta_k^s(t) \frac{\partial o_k^s(t)}{\partial w_{ij}} \quad (i=1, N; j=1, L) \end{aligned} \quad (2-67)$$

unde:

$$\delta_k^s = [d_k^s(t) - o_k^s(t)] \quad (k=1, m) \quad (2-68)$$

Notăm cu:

$$r_{kij}^s(t) = \frac{\partial o_k^s(t)}{\partial w_{ij}} \quad (2-69)$$

Această mărime reprezintă influența pe care o are o pondere asupra unei ieșiri o_k pentru un tipar s , la momentul t , și se calculează în felul următor (considerăm pe i și j ca fiind fixați):

$$\begin{aligned} r_{kij}^s(t) &= \frac{\partial}{\partial w_{ij}} g(h_k^s(t)) = \frac{\partial}{\partial w_{ij}} g\left(\sum_{l=0}^L w_{kl} \cdot z_l^s(t-1)\right) = \\ &= g'(h_k^s(t)) \sum_{l=0}^L \frac{\partial}{\partial w_{ij}} (w_{kl} \cdot z_l^s(t-1)) = \\ &= g'(h_k^s(t)) \left[\frac{\partial w_{k0}}{\partial w_{ij}} \cdot z_0^s(t-1) + \dots + \frac{\partial w_{kj}}{\partial w_{ij}} \cdot z_j^s(t-1) + \dots + \sum_{l=0}^L w_{kl} \frac{\partial z_l^s(t-1)}{\partial w_{ij}} \right] = \end{aligned}$$

$$= g' (h_k^s(t)) \left[K_{ik}^\Delta \cdot z_j^s(t-1) + \sum_{l=0}^L w_{kl} \frac{\partial z_l^s(t-1)}{\partial w_{ij}} \right] \quad (k=1, m) \quad (2-70)$$

Al doilea termen al sumei se calculează ținând cont de relația (2-60) care desparte neuronii de distribuție alocați în parte pentru recurență și în parte pentru aplicarea intrării $x(t)$:

$$\sum_{l=0}^L w_{kl} \frac{\partial z_l^s(t-1)}{\partial w_{ij}} = \sum_{l=0}^n w_{kl} \frac{\partial x_l^s(t-1)}{\partial w_{ij}} + \sum_{l=n+1}^L w_{kl} \frac{\partial z_l^s(t-1)}{\partial w_{ij}} \quad (2-71)$$

Primul termen al sumei este egal cu zero pentru că vectorul x nu depinde de w_{ij} . În al doilea termen facem o schimbare a indicilor:

$$\sum_{p=1}^N w_{k(p-n)} \frac{\partial z_{p-n}^s(t-1)}{\partial w_{ij}} = \sum_{p=1}^m w_{k(p-n)} \frac{\partial o_p^s(t-1)}{\partial w_{ij}} + \sum_{p=m+1}^N w_{k(p-n)} \frac{\partial v_p^s(t-1)}{\partial w_{ij}} \quad (2-72)$$

Dacă vom considera că în acest caz (la învățare), toate ieșirile sunt importante și vom nota cu:

$$y_p^s(t) = \begin{cases} o_p^s(t) & \text{dacă } p=1, m \\ v_p^s(t) & \text{dacă } p=m+1, N \end{cases} \quad (2-73)$$

în acest caz membrul drept al relației (2-72) va deveni egal cu:

$$\sum_{p=1}^N w_{k(p-n)} \frac{\partial y_p^s(t-1)}{\partial w_{ij}} \quad (2-74)$$

Putem considera că:

$$r_{pij}^s(t) = \frac{\partial y_p^s(t)}{\partial w_{ij}} \quad (p=1, N) \quad (2-75)$$

Dacă vom rescrie relația (2-70) cu noua notație va rezulta:

$$r_{kij}^s(t) = g' (h_k^s(t)) \left[K_{ik}^\Delta \cdot z_j^s(t-1) + \sum_{p=1}^N w_{k(p-n)} \cdot r_{pij}^s(t-1) \right] \quad (k=1, N) \quad (2-76)$$

Relație care se observă că este recursivă relativ la timp. Pentru a calcula mărimea r la momentul t , pentru o ieșire k , trebuie să cunoaștem r pentru toate ieșirile la momentul $t-1$. Înlocuind în (2-70) vom obține regula de învățare:

$$\Delta w_{ij}^s = \eta \sum_{t=1}^{stop} \sum_{k=1}^m \delta_k^s(t) \cdot r_{kij}(t) \quad (2-77)$$

Această formulă aplicată la antrenare se numește *regula de propagare înapoi în timp a erorilor* și este metoda folosită în experimentele care sunt descrise în continuare în teză.

3. Realizarea legăturii între simbolic și conexiunist

Un sistem cognitiv acceptă la intrare un set de cunoștințe, îl procesează și generează la ieșire un alt set de cunoștințe. Sistemul gândit de mine în prima fază, prelua un grup de propoziții aflate într-o relație de ordine (reprezentate deja în formă de grafuri conceptuale) și a le furniza în ordine unei rețele neuronale recurente. Aceasta genera la ieșire un alt set ordonat de grafuri conceptuale, care reprezentau rezultatele dorite. Perechile de seturi de propoziții ($SET_{intrare}$, $SET_{ieșire}$) reprezentau baza de antrenare a rețelei. Două paradigme au fost enunțate în această fază:

I. Procesul de generare al ieșirii este analog cu traiectoria stării rețelei văzută ca sistem dinamic, iar intrarea este privită ca perturbația exterioară aplicată sistemului.

II. O proiecție a unui punct din spațiul stărilor rețelei poate fi interpretată ca și un graf conceptual (și invers).

3.1. Procesul cognitiv privit ca și traiectorie

Paradigma care stă la baza metodologiei propuse în teză este că: un proces cognitiv poate fi realizat ca o succesiune de stări apropiate, într-un spațiu cu foarte multe grade de libertate. Mai exact, o traiectorie, mai mult sau mai puțin continuă. Se pleacă dintr-un punct de start, dat pe de o parte de o regiune de lucru din spațiul euclidian și pe de altă parte de un punct inițial care reprezintă intrarea sistemului (unul sau mai multe grafuri conceptuale). Sistemul poate evolua liber sau perturbat (de alte intrări - grafuri conceptuale), ajungând într-un punct final. Rezultatul poate fi acest punct final ori acesta plus o serie de puncte intermediare de pe traiectoria parcursă.

Ideea este cunoscută ([Rumelhart&al86] - și fiind investigată încă din 1982) și a dus la tentativa de realizare a raționamentului ca traiectorie a unui sistem dinamic

[Smolensky86]. Trebuie însă precizată diferența între *raționament* și *proces cognitiv*. Un raționament dă rezultate care sunt acceptate ca fiind adevărate din punct de vedere logic, sau măcar este indicat un grad de confidență (factor de certitudine, valoare de probabilitate). Termenul de proces cognitiv se referă în acest context la acea clasă de procesări simbolice care dau rezultate cu sens (ca și operația de "join" de la grafurile conceptuale). Termenul este uzitat în comunitatea științifică anglo-saxonă sub numele de "common-sense reasoning" și implică o robustețe sporită a procesului la informații nerelevante sau contradictorii.

Dacă vom considera că o rețea Towell are o stare, dată de vectorul complet al ieșirilor neuronilor săi, iar rețeaua efectuează un raționament (logic corect) prin mecanismul de "aprinde" succesiva a straturilor sale (vezi capitolul 1), iar vectorul este reprezentat într-un spațiu euclidian cu un număr de dimensiuni egal cu numărul de unități din rețea, atunci și aici putem vorbi de o traiectorie. Spațiul folosit de rețea, se reduce la un hiper cub cu latura 2, centrat în origine. Stările posibile sunt plasate aproximativ pe colțurile acestui hiper cub. Atunci când un neuron își schimbă valoare de ieșire din aproximativ -1 în aproximativ 1 (și putem presupune că își schimbă continuu starea), "starea" rețelei se mută (pe o traiectorie) dintr-un colț al hiper cubului în alt colț apropiat.

Smolensky a demonstrat formal [Smolensky86] că în cazul în care avem un spațiu simbolic și un spațiu euclidian, se poate defini un izomorfism între un raționament simbolic și o traiectorie a unui sistem dinamic compus din unități conexiuniste, cu condiția ca unitățile să fie liniare. Tot în respectiva lucrare, el demonstra că izomorfismul nu e adevărat pentru unități neliniare. *Ceea ce doream eu să experimentez, era să antrenez o rețea neurodinamică cu unități neliniare (conform teoriei prezentată în capitolul 2) cu perechi de texte aflate în relație semantică și să află dacă la prezentarea la intrare a unui text nou (format din cuvinte care au fost întâlnite toate la antrenare în alte texte), la ieșire se va obține un text format din propoziții care au sens, și mai mult, dacă textul întreg formează un discurs unitar. Mai mult, doream să văd dacă elementele generice ale unui text cu semantică (introducere, dezvoltare, deznodământ) se păstrează și în cazul unui text nou generat.*

O posibilă aplicație a unei astfel de rețele este prezentată în capitolul 4, și poate fi utilizată pentru planificarea robustă a unor acțiuni. Pentru început, trebuia dovedit că o astfel de rețea poate fi antrenată, iar învățarea va conduce la capacitatea de generalizare. Primul pas important în stabilirea metodologiei propuse a fost găsirea unei metode de translație a grafurilor conceptuale în vectori de activare conexiuniști (puncte din spațiul euclidian multidimensional).

3.2. Translatarea informației simbolice în format conexiunist

Privind o rețea neuronală ca și sistem cognitiv, se observă imediat necesitatea definirii intrării și ieșirii din acest sistem. Intrarea este de regulă o secvență de structuri

simbolice eșalonate în timp. Pentru că sistemul are o funcționare impusă de intrare, putem considera că ieșirea este tot o secvență de forme abstracte. Rezultă din această structură funcțională că un astfel de sistem cognitiv este un sistem de dialog care poate fi antrenat astfel încât ca la un stimul (aparținând unei mulțimi finite de stimuli) să dea un răspuns (aparținând unei mulțimi finite de răspunsuri) [Wang&Waibel92]. Dacă aceste două mulțimi sunt suficient de largi și acoperitoare pentru o anumită temă a dialogului, putem să ne folosim de capacitatea de generalizare a rețelelor neuronale pentru a furniza stimuli din același domeniu, dar care nu au fost învățați și a observa dacă răspunsurile sunt viabile. Apariția după învățare a unui astfel de comportament poate fi un argument experimental pentru a demonstra că metodologia folosită este potrivită.

3.2.1. Paradigma sub-simbolică.

Vom numi paradigma sub-simbolică orice structură teoretică formală care își propune transformarea unei forme abstracte (structură simbolică cu o formalizare exactă și valoarea semantică asociată) într-o altă formă abstractă printr-un proces computațional aparținând unui model conexiunist sau care are ca finalitate utilizarea informației rezultate într-un sistem conexiunist.

Fiind o meta-teorie, Paradigma Sub-simbolică poate fi privită ca o compunere armonioasă a mai multor aparate teoretice. Putem identifica un număr de cinci teorii necesare pentru a avea o acoperire completă a tuturor problemelor legate de această metateorie [Smolensky&al94]:

I. Reprezentarea: Este necesar un formalism matematic exact pentru a descrie o mapare a unei forme abstracte peste tiparul de activitate al unei rețele neuronale.

II. Procesarea: Trebuie descris într-un mod matematic detaliat cum procesele cognitive de nivel înalt se pot desfășura într-un mod distribuit printr-o procesare conexiunistă.

III. Formarea: Procesarea conexiunistă duce în final la generarea unor noi tipare de activitate care trebuie interpretate printr-un aparat matematic astfel încât să putem reconstitui din acestea noi structuri simbolice.

IV. Gramatica: Orice structură simbolică trebuie descrisă în termeni preciși cu ajutorul limbajelor formale.

V. Semantica: Tot printr-un formalism unitar, structurilor simbolice (care au o informație semantică intrinsecă) trebuie să atașăm informație semantică, prin definiții și asocieri (eventual ponderate).

Pentru a fi cât mai explicit, în continuare voi prezenta un exemplu de paradigmă sub-simbolică - rețeaua Elman - [Elman90] care s-ar preta *aparent* la scopul practic

propus. Voi descompune teoria care stă la baza rețelelor Elman relativ la cele cinci componente prezentate anterior.

I'. Reprezentarea în rețelele Elman: Intrarea într-o rețea Elman se construiește în felul următor:

- Se stabilește un vocabular V finit de cuvinte care vor putea fi acceptate de rețea.
- Se construiește un strat de intrare format din $\text{card}(V)$ neuroni de "fan-out", unde fiecare neuron va avea asociat un cuvânt din V .
- Se baleiază cu ajutorul unei ferestre de lungime mică (7 cuvinte de obicei) textul care trebuie prezentat sistemului. Lungimea ferestrei se alege astfel încât să nu poată apărea același cuvânt de două ori în fereastră. Fereastra va avansa câte o poziție către sfârșitul textului la fiecare moment de timp t_i (timpul este o succesiune de cuante discrete). Dacă avem de exemplu textul:

"Lucrarea mea de doctorat este realizată în mare parte în timpul meu liber și în special duminica".

și ne alegem o fereastră de lungime de 5 cuvinte, la momentul t_0 vom prezenta rețelei următoarea mulțime de simboluri:

$$v(t_0) = \{ \text{lucrarea , mea , de , doctorat , este} \}$$

Vom alege din această mulțime ordonată un element pivot (de obicei cuvântul aflat în poziție centrală). Pentru mulțimea $v(t_0)$ pivotul va fi cuvântul "de". Forțăm pe valoarea 1 neuronii de intrare asociați cuvintelor din mulțimea $v(t_0)$, restul rămânând setați pe valoarea 0. La momentul t_1 vom avea:

$$v(t_1) = \{ \text{mea , de , doctorat , este , realizată} \}$$

cu pivotul "doctorat" și vom seta pe valoarea 1 neuronii asociați cuvintelor din $v(t_1)$. Se repetă procedeul de baleiere până când se epuizează textul. Ieșirea se formează în felul următor:

- Se construiește un strat de ieșire format tot din $\text{card}(V)$ neuroni, având ca și în stratul de intrare, câte un neuron asociat câte unui cuvânt din V . Dorim ca la momentul t_i , atunci când prezentăm rețelei mulțimea $v(t_i)$ cu pivotul corespunzător, la ieșire să se aprindă (să aibă valoarea cât mai aproape de 1) întotdeauna un singur neuron, cel asociat cuvântului imediat următor pivotului din $v(t_i)$, iar restul să aibă valori cât mai apropiate de 0. Această sarcină pare destul de ciudată pentru că la momentul t_i noi deja prezentăm rețelei cuvântul pe care îl dorim în acel moment la ieșire (împreună bineînțeles cu restul cuvintelor din fereastră), dar putem demonstra capacitatea de predicție a sistemului.

II'. Procesarea în rețelele Elman: Pentru realizarea task-ului descris anterior ne alegem un set cât mai mare de texte formate pe cât este posibil din cuvinte aparținând vocabularului V . Vom antrena rețeaua cu ajutorul acestor texte (ignorând cuvintele care

nu aparțin vocabularului ales). Elman a observat că scopul pe care și l-a propus putea fi implementat doar cu o rețea semi-recurentă (și nu una pur feed-forward) ca și cea prezentată în figura 3-1.

Pentru învățarea tiparelor care provineau din textele alese (un tipar este format dintr-o mulțime de cuvinte $v(t)$ pentru intrare și un cuvânt reprezentând poziția imediat următoare pivotului din $v(t)$), Elman a folosit o variantă a algoritmului backpropagation, adaptată pentru a suporta recurența dintre stratul ascuns și stratul de intrare. Învățarea a fost rapid convergentă (mai ales datorită faptului că vectorii de ieșire sunt liniar independenți) și eroarea a putut fi scăzută sub 2% [Omlin&Giles92]. Evident că și rețeaua Elman (la fel ca și alte modele conexiuniste) are un moment de saturație atunci când încercăm să o antrenăm cu prea multe tipare. Dar o alegere judicioasă a dimensiunii bazei de învățare duce la o funcționare normală. La acest punct se poate observa că aparatul matematic folosit nu explică de ce rețeaua este capabilă să predicteze cuvântul care apare după pivot [Smolensky92].

III'. Formarea: Metoda de a extrage informațiile din acest sistem cognitiv este următoarea: se prezintă sistemului un singur cuvânt din vocabular și se memorează într-o variabilă (de tip vector m -dimensional, cu componente reale) asociată acestui cuvânt starea indusă în stratul ascuns (care are m neuroni). Pentru toate cuvintele din vocabular vom obține astfel câte un punct asociat în spațiul euclidian m dimensional.

IV'. Gramatica: Este un aspect ignorat în această paradigmă sub-simbolică. Singura informație gramaticală care este furnizată sistemului este topica (ordinea) cuvintelor în text.

V'. Semantica: La fel ca și aspectul sintactic, cel semantic este ignorat cu desăvârșire, o informație semantică posibilă fiind cea dată de asocierile posibile între cuvinte apropiate.

Concluzia pe care o putem trage după analiza metodologiei Elman, interpretată ca și paradigmă sub-simbolică, este că trebuie insistat mai mult pe aspectele gramaticale și semantice care nu au fost încă suficient explorate de cercetarea în domeniul integrării procesării simbolico-conexiuniste.

3.2.2. Reprezentarea intrării în rețeaua neuronală aleasă

Rețeaua neuronală pe care am ales-o are un număr de n neuroni de intrare (vezi figura 2-13) care au rol de "fan-out" și pot primi și distribui valori reale. Intrarea momentană în această rețea trebuie în mod necesar să apară ca și un vector n dimensional cu componente reale. Problema este cum transformăm un text într-o succesiune de asemenea vectori. Pentru că în această lucrare dorim să construim o paradigmă sub-simbolică nouă, abordarea reprezentării intrării trebuie făcută într-un mod cât se poate de formalizat, evitând soluțiile implementaționiste (vezi critica

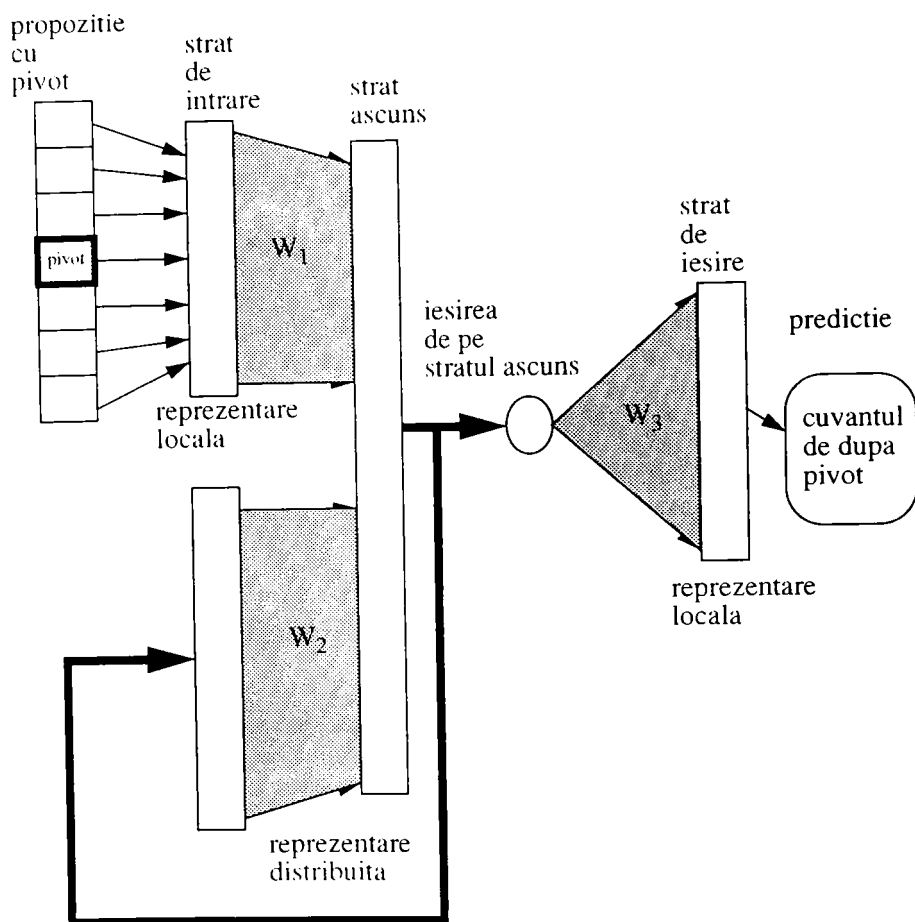


figura 3-1.

acestora în [Smolensky92]).

Reprezentările (transformarea intrării într-o formă vectorială) pot fi clasificate în două categorii:

- locale
- distribuite.

Într-o reprezentare locală fiecărui neuron i se asociază o entitate din mulțimea entităților de intrare (cuvinte, eșantioane de semnal, porțiuni de imagine). De exemplu în rețeaua Elman fiecărui neuron de pe stratul de intrare i se asociază un singur cuvânt. Pentru procesarea simbolică, acest mod de reprezentare este tipic local. La fel, în rețelele Towell (un exemplu va apărea în continuare) fiecărui neuron din rețea (indiferent pe care strat se află) i se asociază un concept [Towell91].

• Majoritatea cercetătorilor (vezi analiza direcțiilor neurosimbolice din [Alexandre96]) apreciază că reprezentările locale reprezintă doar o fază de infantilism în acest domeniu și nici o paradigmă sub-simbolică serioasă nu poate să se bazeze pe această metodă (care rămâne totuși viabilă pentru aplicațiile foarte simple). Se susține (argumentat) ideea folosirii reprezentărilor *distribuite* ca fiind cele mai adecvate pentru paradigmele sub-simbolice [Ajjanagadde94].

• În acest caz o anumită entitate care trebuie reprezentată într-un grup de neuroni nu va fi mapată peste acel grup prin aprinderea unui singur neuron asociat, ci printr-un *tipar de activitate* al aceluși grup. Ca un foarte sugestiv exemplu, maparea dintre mulțimea cuvintelor din vocabularul V și stările induse de acestea în stratul ascuns din rețeaua Elman, determină o reprezentare distribuită. Formal spus:

$$\Psi : Voc \rightarrow \mathbf{R}^n \quad (3-1)$$

În cazul rețelei Elman această reprezentare se folosește exclusiv pentru obținerea ieșirii finale. Pentru a descrie formal o reprezentare distribuită pentru intrarea în rețeaua aleasă voi prezenta următoarele definiții (o varietate inițială a lor poate fi găsită în [Smolensky&al94]):

Definiție

Un tipar de activitate al unei rețele neuronale (conexiuniste) cu n neuroni (unități) este un vector n dimensional cu componente reale:

$$\vec{v} \triangleq (v_1, v_2, \dots, v_n) \quad v_i \in \mathbf{R} \quad (i=1, n) \quad (3-2)$$

Definiție

Spațiul stărilor V asociat unei rețele neuronale este zona spațiului real n dimensional unde pot apărea tiparele de activitate ale rețelei.

Aceasta poate fi infinită sau finită, depinzând de funcțiile de transfer ale neuronilor componenți. De exemplu dacă toți neuronii posedă ca și funcție de transfer sigmoida, această zonă va fi un hiper cub n dimensional cu latura 1 care va avea un colț în origine și laturile de-a lungul axelor sau ortogonale pe axe. Pentru tangenta hiperbolică, spațiul V va fi un hiper cub cu latura 2, centrat în origine.

Putem folosi mai multe tipuri de structuri simbolice pentru reprezentarea informației (cunoștințelor) simbolice, ca de exemplu secvența (string) de simboluri, arborele cu noduri conținând simboluri etc.

Limbajul natural este format dintr-o mulțime finită de cuvinte. Dacă dorim să construim un sistem cognitiv care să manipuleze concepte denotate prin simboluri provenite din limbajul natural, acesta va putea lucra doar cu o mulțime finită de simboluri (pe care o vom nota cu F). Pentru o anumită aplicație particulară, F va fi o mulțime limitată ca dimensiune și grupată din punct de vedere semantic [Winograd80]. Nu toate structurile simbolice formate cu ajutorul cuvintelor din F sunt utilizabile. După [Smolensky90] ele trebuie să îndeplinească trei condiții:

- i. Să fie corecte gramatical ("well-formedness").
- ii. Să aibă sens (valoare semantica validă).
- iii. Să fie adevărate (din punct de vedere empiric și chiar logic acolo unde este posibil).

Definiție

Se numește viabilă, orice structură simbolică, formată pe baza unei gramatici, având sens și care este general acceptată ca fiind adevărată.

Definiție

Se numește S spațiul structurilor simbolice viabile folosite pentru o anumită aplicație.

În acest moment avem toate componentele necesare pentru a defini legătura dintre un text și o rețea neuronală.

Definiție

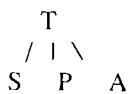
Se numește reprezentare conexiunistă distribuită aplicația:

$$\Psi : S -> V \quad (3-3)$$

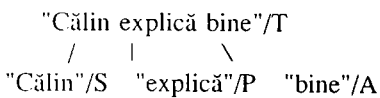
dacă aplicația este o bijecție și nici unei structuri simbolice din S nu îi corespunde vectorul nul.

Observație: Pentru o anumită alegere particulară a vectorilor din V reprezentarea poate fi locală.

Pentru a îndeplini cele două condiții puse de definiția anterioară trebuie construită o schemă de mapare între S și V care să se bazeze pe o *regulă de decompoziție* peste S . Pentru a defini acest concept trebuie explicat mai întâi ce înseamnă conceptul de *rol*. De exemplu într-o secvență care are o structură definită printr-o regulă de ordonare, fiecare poziție poate fi privită ca o variabilă care poate lua anumite valori din mulțimea F . Deci putem interpreta structura ca o succesiune fixă de variabile numite roluri care pot fi completate cu simboluri. Pentru o secvență de exact trei simboluri putem construi macheta (r_1, r_2, r_3) . În acest caz pentru structura simbolică particulară "a b a" se poate scrie forma $(a/r_1, b/r_2, a/r_3)$ în care rolurile sunt completate (legate) cu valori din F [Smolenksy90]. La fel, putem construi arborii sintactici. Având de exemplu arborele:



Unde rolurile T (propoziție), S (subiect), P (predicat), A (atribut) pot fi completate cu valori simbolice ca de exemplu:



Definiție

Se numește regulă de decompoziție aplicația:

$$\mu_{F \times R} : F \times R \rightarrow f/r(s) \quad (s \in S, f \in F, r \in R) \quad (3-4)$$

unde:

- F este mulțimea simbolurilor (vocabularul)
- R este mulțimea rolurilor
- S este spațiul structurilor viabile
- $f/r(s)$ este o construcție care denotă că legarea f/r este corectă pentru o anumită structură simbolică viabilă s (face parte din această structură sau poate face parte când structura este incompletă).

Pentru o structură simbolică s particulară putem avea mai multe perechi (f,r) care se potrivesc în această structură. Din acest motiv o regulă de decompoziție poate induce aplicația (70), care construiește o mulțime a tuturor perechilor (f,r) care sunt corecte pentru o anumită structură simbolică viabilă s .

$$\beta : S \rightarrow 2^{F \times R} \quad (3-5)$$

Deci:

$$s \mapsto \{(f, r) \mid f/r(s) \text{ corectă}\} \quad (3-6)$$

Aplicația β se numește *reprezentare simbol-rol*. Pentru o anumită structură s_0 vom avea mulțimea finită:

$$\beta(s_0) = \{(f_1, r_1), (f_2, r_2), \dots, (f_k, r_k)\} \quad (3-7)$$

Construcția $f/r_i(s_0)$ poate fi interpretată ca o construcție clauzală de tip predicat(argument) unde legătura f/r_i poate fi interpretată ca și un predicat logic având argumentul s_0 .

Structura simbolică viabilă s_0 este din punct de vedere logic considerată adevărată, deci toate clauzele $f_i/r_i(s_0)$ le putem considera logic adevărate. În acest caz predicatul:

$$\Pi(s_0) = \bigwedge_{(f_i, r_i) \in \beta(s_0)} f_i/r_i(s_0) \quad (3-8)$$

este adevărat [Szirbik93]. Dacă mulțimea F are cardinalitatea n putem construi foarte simplu o aplicație:

$$\Psi_F : F \rightarrow V_F \quad (V_F = R^n) \quad (3-9)$$

care mapează un simbol peste spațiul stărilor unui grup de n neuroni. În mod necesar aplicația trebuie să fie o bijecție și nici un simbol să nu aibă asociat vectorul nul. Deci pentru fiecare simbol vom avea:

$$\Psi_F(f) = \bar{v}_f \quad (3-10)$$

Este ideal ca mulțimea tuturor vectorilor de reprezentare $\{v_f\}$ să formeze o bază în spațiul vectorial n dimensional (să fie liniar independenți). În practică această condiție este foarte greu de îndeplinit, dar putem construi această mapare astfel încât gradul de acoperire al vectorilor să fie cât mai mic.

Cel mai simplu ar fi să construim această mapare ca pe o reprezentare locală (ca și la intrarea în rețelele Elman). Fiecărui simbol din F îi corespunde câte un neuron și mulțimea vectorilor $\{v_f\}$ va reprezenta o bază. Două considerente ne împiedică să procedăm în acest mod:

i. Metodologic. S-a observat (și în unele cazuri demonstrat teoretic) că o reprezentare distribuită este mai performantă [Smolensky90][Smolensky92].

ii. Practic.

ii-a) Numarul de neuroni n trebuie să fie în acest caz egal cu $\text{card}(F)$. Se va observa în continuare că numărul de neuroni alocați acestei reprezentări influențează stabilirea numărului de neuroni din întreaga rețea, număr care poate ajunge exagerat de mare [Maclin&Shavlik94].

ii-b) În cazul modificării lui F trebuie să adăugăm/eliminăm neuroni din rețea.

În concluzie, este mai productiv să ne stabilim din start un număr fix de neuroni (mai mic decât $\text{card}(F)$) și să încercăm să creem iterativ $\{v_f\}$ astfel încât orice nou simbol adăugat să fie asociat cu un vector v_f care diferă pe cât posibil de vectorii deja stabiliți și de combinațiile lor liniare.

În același mod vom construi aplicația:

$$\Psi_R : R \rightarrow V_R \quad (V_R = R^m) \quad (3-11)$$

care trebuie să îndeplinească aceleași proprietăți ca și Ψ_F . Pentru rolurile care aparțin unui arbore binar se poate construi maparea într-un mod foarte elegant în felul următor [Smolensky92]: *în funcție de adâncimea maximă a arborelui ne stabilim valoarea m a dimensionalității spațiului vectorial V_R (sau altfel spus numărul de neuroni alocați pentru reprezentarea rolurilor).*

Stabilim calea de la rădăcină până în poziția în care se află rolul respectiv. Parcurgem această cale construind vectorul v_f din componente 1 și -1 în funcție de direcția (stânga sau dreapta) pe care o avem de la un nod la altul. Dacă nu am ajuns la adâncimea maximă completăm restul vectorului cu 0. Un exemplu este dat în figura 2- care reprezintă un arbore binar cu adâncimea 4. Pentru r_1 , r_2 și r_3 vom avea:

$$\Psi_R(r_1) = (-1, -1, 0, 0) \quad (m=4) \quad (3-12)$$

$$\Psi_R(r_2) = (1, -1, 1, -1) \quad (6-13)$$

$$\Psi_R(r_3) = (1, 1, 1, 0) \quad (3-14)$$

Smolensky a dat o soluție și mai elegantă, recursivă, bazată pe calculul tensorial, pentru arborii binari de adâncime teoretic infinită, mapare care se face peste un spațiu V_R de dimensiune finită [Smolensky92].

Legătura simbol/rol o putem de asemenea reprezenta ca un tipar de activitate peste un grup de neuroni. Metoda matematică pe care o folosim este *produsul tensorial*. Dacă avem:

$$\Psi_F(f_i) = \bar{F}_i \quad (3-15)$$

și:

$$\Psi_R(r_i) = \bar{R}_i \quad (3-16)$$

Reprezentarea legăturii f_i/r_i este:

$$\overline{f_i/r_i} = \overline{f_i} \otimes \overline{r_i} = \overline{b_i} \quad (3-17)$$

ceea ce înseamnă următorul calcul la nivel de element. Dacă:

$$\overline{f_i} = (f_{i1}, f_{i2}, \dots, f_{in}) \quad (3-18)$$

$$\overline{r_i} = (r_{i1}, r_{i2}, \dots, r_{im}) \quad (3-19)$$

atunci tensorul b_i va fi format din toate produsele posibile dintre elementele vectorilor f_i și r_i .

$$\overline{b_i} = (f_{i1} \cdot r_{i1}, f_{i1} \cdot r_{i2}, \dots, f_{i1} \cdot r_{im}, f_{i2} \cdot r_{i1}, \dots, f_{in} \cdot r_{im}) \quad (3-20)$$

sau altfel spus:

$$\overline{b} = \overline{f} \otimes \overline{r} \quad \text{unde } b_{\rho\phi} = f_{\rho} \cdot r_{\phi} \quad \rho=1, n; \phi=1, m \quad (3-21)$$

Dacă Ψ_F și Ψ_R sunt bijecții, putem defini aplicația:

$$\Psi_b : F \times R \rightarrow V_F \otimes V_R \quad (3-22)$$

$$f/r \rightarrow \Psi_F(f) \otimes \Psi_R(r) \quad (3-23)$$

care mapează legătura simbol/rol peste un grup de $m \times n$ neuroni pentru că:

$$V_R \otimes V_R = V_B \quad (3-24)$$

unde V_B este un spațiu vectorial cu $n \times m$ dimensiuni. Demonstrația faptului că Ψ_b este tot o bijecție este imediată.

Urmează să analizăm cum am putea combina tensorii rezultați din aplicația Ψ_b pentru a construi o reprezentare pentru predicatul $\Pi(s_0)$. Deoarece acesta este construit pe baza operatorului logic și, întrebarea pe care ne-o punem în acest moment este:

- Cum se poate reprezenta în mod conexiionist conjuncția logică?

Într-o rețea conexiionistă, care susține o reprezentare locală a conceptelor, conjuncția se poate face cu ajutorul funcției de transfer unipolare sau chiar a unei sigmoide foarte abrupte [Towell91]. Vom da în continuare un exemplu ilustrativ. Dacă avem:

$$p_3 = p_1 \wedge p_2 \quad (3-25)$$

$$p_1 = a_1 \wedge a_2 \wedge a_3 \quad (3-26)$$

$$p_2 = a_4 \wedge a_5 \quad (3-27)$$

Aceste relații se pot transforma într-o rețea feedforward cu mai multe straturi (ca și rețeaua Towell prezentată în primul capitol). Dacă vom nota funcția de transfer cu f , putem calcula ieșirea prin formulele:

$$o(p_1) = f(o(a_1) + o(a_2) + o(a_3)) \quad (3-28)$$

$$o(p_2) = f(o(a_4) + o(a_5)) \quad (2-28)$$

$$o(p_3) = f(o(p_1) + o(p_2)) = f\left(\sum_{i=1}^5 o(a_i)\right) \quad (3-30)$$

Rezultat care se datorează funcției unipolare de transfer. Observația pe care o putem face relativ la acest exemplu este că rezultatul se obține printr-o superpoziție (însușire).

Putem extinde această metodă de reprezentare conexiunistă a conjuncției și pentru reprezentările distribuite. Dacă pentru o mulțime de propoziții $\{p_i\}$ cu valoare de adevăr avem reprezentările distribuite:

$$\bar{x}_i = \Psi(p_i) \quad (3-31)$$

Pentru conjuncția lor vom avea:

$$\Psi\left(\bigwedge_i p_i\right) = \sum_i \Psi(p_i) \quad (3-32)$$

Se observă că se păstrează proprietățile de asociativitate și comutativitate. Dacă în (85) înlocuim p_i cu clauzele predicative $f_i/r_i(s_0)$ vom obține:

$$\begin{aligned} \Psi\left(\bigwedge_i (s_0)\right) &= \Psi\left(\bigwedge_i f_i/r_i(s_0)\right) = \\ &= \sum_i \Psi(f_i/r_i(s_0)) = \end{aligned} \quad (3-33)$$

$$= \sum_i \Psi_F(f_i) \otimes \Psi_R(r_i) \quad (f_i, r_i) \in \beta(s_0) \quad (3-34)$$

Și astfel, în final putem obține rezultatul dorit, adică reprezentarea:

$$\Psi_T : S \rightarrow V \quad (V = V_F \otimes V_R) \quad (3-35)$$

prin descompunerea în două aplicații:

$$\Psi_T(s) = \Psi \left(\bigwedge_{(f_i, r_i) \in \mathbf{p}(s)} \Psi_b(f_i/r_i) \right) = \quad (3-36)$$

$$= \sum_{(f_i, r_i) \in \mathbf{p}(s)} \Psi_F(f_i) \otimes \Psi_R(r_i) = \bar{v}_s \quad (3-37)$$

Această traducere din spațiul simbolic (conceptual) în spațiul conexionist (euclidian) se numește **reprezentare tensorială**.

Observații:

- Dacă Ψ_F și Ψ_R sunt bijecții atunci și Ψ_T este o bijecție. Dacă $\{v_f\}$ și $\{v_r\}$ sunt baze ale spațiilor V_F și V_R atunci și $\{v_f \otimes v_r\}$ vor fi o bază în spațiul V . Acesta este din păcate un ideal greu de atins dar o alegere a reprezentărilor Ψ_F și Ψ_R astfel încât să avem un grad de acoperire cât mai mic va induce aceeași proprietate în Ψ_T . Se poate demonstra că Ψ_T nu duce niciodată la o reprezentare prin vectorul nul, cu condiția ca Ψ_F și Ψ_R să nu aibă mapări de genul $v_f=0$ și $v_r=0$.
- Dacă Ψ_F și Ψ_R sunt reprezentări locale, Ψ_T va fi o reprezentare conexionistă tensorială locală.
- Dacă una din Ψ_F și Ψ_R este locală și cealaltă distribuită, Ψ_T va fi o reprezentare conexionistă tensorială semilocală.
- Dacă Ψ_F și Ψ_R sunt distribuite, Ψ_T va fi o reprezentare conexionistă tensorială distribuită.

3.2.3. Implementarea conexionistă a reprezentării tensoriale.

Trecerea de la o structură simbolică la vectorul $v(s)$ pare a fi deosebit de complicată, relativ la simplitatea modelului computațional conexionist și ar necesita o interfață masivă cu rețeaua aleasă, care din punctul de vedere al implementării ar fi mai complicată decât rețeaua însăși. Problema poate fi destul de simplu rezolvată prin schemă de implementare dată în figura 3-2. Se observă că în acest caz vom avea m ($m = \text{card}(F)$) neuroni de "fan-out" care vor fi legați fiecare la p neuroni din zona I (de intrare) și p ($p = \text{card}(R)$) neuroni de "fan-out" care vor fi legați fiecare la m neuroni din zona I, după schema de distribuție din figură.

Dacă prezentăm rețelei un vector b_j , rezultat al reprezentării Ψ_b pentru o legătură simbol/rol f/r_j , vom aplica vectorul f_j peste neuronii f_j ($j=1, m$) și vectorul r_j peste

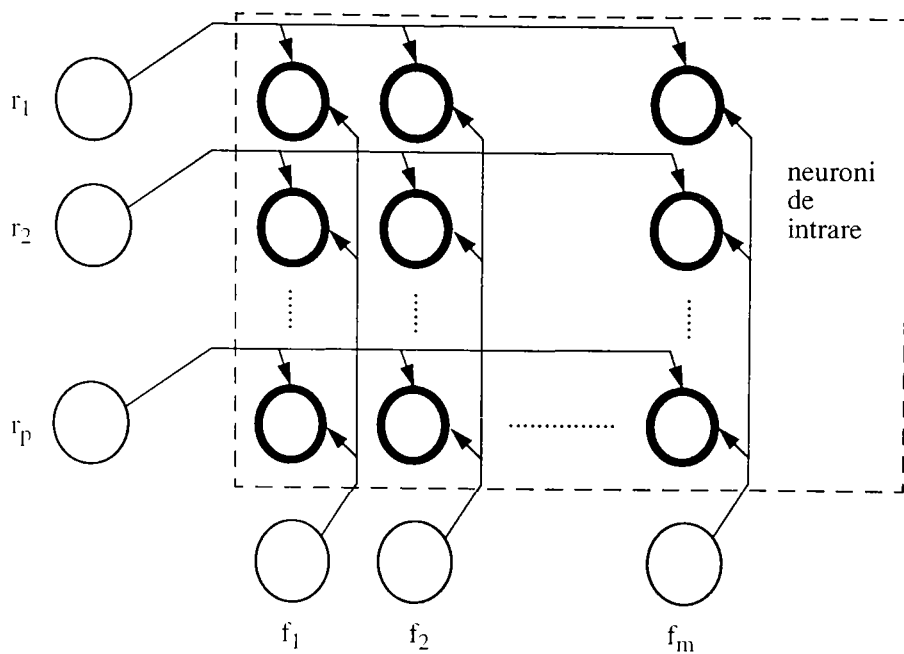


figura 3-2.

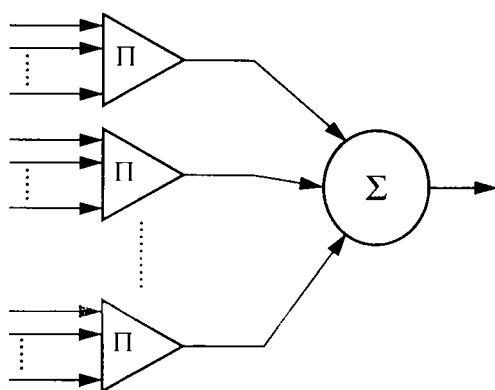


figura 3-3.

neuronii r_k ($k=1,p$). Pentru a putea implementa produsul tensorial, neuronii din zona I trebuie să aibă următoarea schemă (relativ simplă) de funcționare:

- Cele două intrări pe care le are fiecare neuron din zona I (una de la neuronii f , și cealaltă de la neuronii r), trebuie înmulțite, rezultatul fiind starea de activare a neuronului respectiv.

De exemplu, pentru neuronul din poziția $j=s$ și $k=t$ din zona I activarea va fi dată de formula:

$$v_{st} = f_s \cdot r_t \quad (3-38)$$

Prin această foarte simplă schemă am realizat primul pas al reprezentării tensoriale. Ea poate fi aplicată foarte simplu ca implementare software sau ca implementare hardware.

Pentru a implementa complet reprezentarea tensorială a unei structuri simbolice s mai avem nevoie de superpoziția tuturor vectorilor b_i rezultați din reprezentarea Ψ_b a elementelor mulțimii $\beta(s)$. Implementarea se poate face în felul următor: după ce s-a calculat mulțimea $\beta(s)$ se calculează iterativ $\Psi_b(f_i/r_i)$ pentru fiecare pereche (f_i, r_i) . După fiecare iterație rezultatul este memorat aditiv de neuronii din zona I. Dacă avem $\text{card}(\beta(s))=u$ atunci vom obține în final:

$$v_{st} = \sum_{i=1}^u f_s \cdot r_t \quad (s=1, m; t=1, p) \quad (3-39)$$

Formulă care reprezintă metoda practică de introducere a unei structuri simbolice într-o rețea neuronală recurentă prin intermediul a n neuroni ($n = m \times p$) de intrare.

Dacă u ar fi o valoare constantă (de obicei ea diferă de la o structură simbolică la alta, deși s-ar putea normaliza la o valoare acoperitoare), am putea folosi pentru implementare (inclusiv hard) clasa de neuroni numită "sigma-pi". Aceasta a fost intens studiată și există metode verificate de accelerare computațională prin procesoare vectorizate și/sau paralele, pentru implementările software și VLSI analogic pentru implementările hardware [Hecht90][Beiu96]. Un astfel de neuron are funcția de transfer:

$$v_i = \sum_{\sigma=1}^u w_{i\sigma} \prod_{j=1}^z w_{i\sigma j} \cdot I_{\sigma j} \quad (i=1, n) \quad (3-40)$$

unde u este numărul de vectori care trebuie însumați și z numărul de intrări care trebuie înmulțite între ele. Se observă că se poate atașa fiecărei intrări σ câte o pondere și fiecărei intrări j câte o pondere. Schema unui astfel de neuron este dată în figura 3-2.

O rețea formată din astfel de neuroni poate implementa superpoziția de produse tensoriale ponderate între tensori de rang superior. În cazul nostru produsul este format întotdeauna doar din doi factori ($z=2$) f_s și r_t . Ponderile sunt toate 1 iar u poate varia

între diferite valori. Privind schema de funcționare a neuronului "sigma-pi" putem face observația că prin reprezentarea noastră folosim doar o mică parte a potențialului computațional oferit de astfel de neuroni și ar fi o idee bună încercarea de a căuta noi metode care să folosească din plin resursele unor astfel de neuroni.

3.3. Metoda propusă pentru translatarea grafurilor conceptuale

Vom considera că avem pentru problema noastră, un univers al discursului structurat într-un mod prezentat în capitolul 2. Relativ la grafurile conceptuale din universul discursului vom defini trei mulțimi:

- **F** este mulțimea conceptelor individuale (fillers)
- **R** este mulțimea conceptelor generice - clasele (roles)
- **L** este mulțimea relațiilor (links)

Mulțimea $S = F \cup R$ este vocabularul de bază al universului discursului. Putem mapa aceste trei spații peste trei spații vectoriale, definind aplicațiile:

$$\Psi_F: F \rightarrow V_F \quad (3-40)$$

$$\Psi_R: R \rightarrow V_R \quad (3-42)$$

$$\Psi_L: L \rightarrow V_L \quad (3-43)$$

Cea mai bună reprezentare pentru aceste aplicații este una locală, care generează mulțimi de vectori reprezentând baze în spațiile vectoriale V_F , V_R și V_L . În acest caz se pot construi foarte ușor și aplicațiile inverse Ψ_X^{-1} .

Textul specificației, care este stimulul de start, poate fi transformat într-o secvență de grafuri conceptuale $GC_1, GC_2, \dots, GC_{stop}$. O problemă foarte importantă este că între aceste propoziții există *corelații de context*. Dacă vom translata fiecare graf separat într-un vector, aceste corelații vor fi pierdute. Soluția este de a combina "temporal" grafurile între ele. Fiecarui graf i se va aplica operația canonică "join" cu toate grafurile care sunt *înaintea* lui în secvența temporală GC_i ($t=1, t-1$). La extrem, se poate întâmpla ca ultimul graf (GC_{stop}) să fie combinat cu toate grafurile anterioare (de la $t=1$ la $t=stop-1$). În [Thomason&Touretzky90] se apreciază ca redundanța informației semantice oferite unei rețele conexiuniste este benefică și cu cât această redundanță este mai mare, cu atât rețeaua va funcționa mai aproape de comportamentul dorit. Această operație de "join" are un efect benefic și asupra traiectoriei parcurse în spațiul stărilor [Pineda88]. Dacă pe parcursul introducerii punctelor corespondente grafurilor conceptuale, nu facem "salturi" prea mari în spațiul vectorial al intrărilor, atunci nu vom avea "salturi" nici în spațiul stărilor sistemului, respectiv în spațiul

vectorial de ieșire. Dacă propozițiile introduse ca puncte sunt apropiate ca semantică, se cere ca distanța euclidiană între punctele corespunzătoare să fie mică. În acest spirit trebuie imaginată aplicația Ψ .

Vom considera o propoziție oarecare (după ce asupra ei s-a aplicat "join") GC_i . Avem o mulțime a conceptelor, a claselor și a relațiilor pentru acest graf. Pentru fiecare nod concept cu indicele i , vom calcula vectorul b_i ($i=1, N$ - numărul de concepte din graful respectiv) cu formulele:

$$\overline{x}_i = \overline{F}_i \otimes \overline{r}_i \quad (3-44)$$

unde:

$$\overline{F}_i = \Psi_F(\text{concept}) \quad \text{și} \quad \overline{r}_i = \Psi_R(\text{clasă}) \quad (3-45)$$

Vectorul b_i va reprezenta asocierea concept/clasă (filler-role binding). Pentru a reprezenta legăturile dintre un nod concept și vecinii săi (presupunem că are M_i concepte vecine în graf) vom calcula:

$$\overline{rel}_k = [(\overline{F}_k \otimes \overline{r}_k) \otimes I_k] \cdot cf_k \quad (3-46)$$

unde $f_k \otimes r_k$ este asocierea concept/clasă a vecinului k iar $I_k = \Psi$ (relația dintre conceptul i și conceptul k). Valoarea cf_k este un factor de conexiune care depinde doar de direcția arcelor care trec prin relația conceptuală respectivă. Dacă avem situația $i \rightarrow \text{relație} \rightarrow k$ atunci $cf=1.0$ și dacă avem $i \leftarrow \text{relație} \leftarrow k$ atunci $cf=0.5$. Mai există și situații de tipul $i \rightarrow \text{relație} \leftarrow k$ și $i \leftarrow \text{relație} \rightarrow k$, dar în acest caz se va alege $cf=0.75$. După ce am calculat toți vectorii care codifică relațiile, le vom combina prin superpoziție într-un singur vector:

$$\overline{srel}_i = \frac{1}{M_i} \sum_{k=1}^{M_i} \overline{rel}_k \quad (3-47)$$

În final, pentru fiecare concept i , vom calcula un vector care conține informație codificată despre semantica respectivului concept și a clasei din care face parte precum și a vecinilor săi și măsura în care este legat de aceștia. Vom nota:

$$\overline{B}_i = \overline{x}_i \otimes \overline{srel}_i \quad (3-48)$$

Pentru toate conceptele dintr-un graf conceptual GC_i , prin superpoziție vom calcula un singur vector:

$$\overline{v}_t = \sum_{i=1}^N \overline{B}_i \quad (3-49)$$

care este un vector obținut prin desfășurarea unui tensor de gradul 5.

Se poate observa că aplicația tensorială propusă nu generează o bază în spațiul vectorilor de intrare. Din acest motiv, va fi mai greu să definim o aplicație inversă Ψ^{-1} . Puteam defini aplicația Ψ în felul următor:

$$\overline{v}_t = \bigotimes_{i=1}^N \overline{b}_i \quad , \quad \overline{b}_i = (\text{cf. 25}) \quad (3-50)$$

unde:

$$\overline{srel}_i = \bigotimes_{k=1}^{M_i} \overline{rel}_k \quad , \quad \overline{rel}_k = (\text{cf. 23}) \quad (3-51)$$

dar această variantă ar avea rangul $(M_i(\overline{rel}_k)+2(x_i))*N(v_i)$. Dacă avem doar 5 concepte și în medie 2 vecini/concept, rangul unui astfel de tensor este 20 (mult prea mare pentru a fi fezabil). În plus, rangul tensorului (care dă și dimensiunea vectorului de intrare) este variabil, în funcție de numărul de concepte și vecini. Acest lucru ar impune o rețea neuronală cu o dimensiune **maximă** pentru spațiul vectorilor de intrare, maxim dat de cel mai voluminos graf conceptual care poate fi întâlnit în cursul raționamentului. Mai mult, la ieșire, nu am putea ști care este dimensiunea curentă a vectorului de ieșire. Deci, se poate trage concluzia că avem nevoie de o dimensiune fixă și prima variantă este mai bună.

Dacă avem dimensiunile $V_i=V_R=V_L=5$, pentru formula cu superpoziție (25), dimensiunea spațiului de intrare (și prin consecință, și a celui de ieșire) este $5^7=3125$, care este enormă. Dacă pentru Ψ_x vom folosi reprezentări locale, înseamnă că vocabularul nostru e format din 5 concepte, 5 clase și 5 relații, care duce la un univers al discursului atât de sărac, încât termenul de univers "este" superflu. Din acest motiv, pentru a mări vocabularul, fără a afecta dimensiunile spațiilor de intrare/ieșire, trebuie să renunțăm la reprezentările locale pentru Ψ_x , folosind următoarele reprezentări distribuite [Szirbik&al95]:

- pentru Ψ_f vom codifica (binar de exemplu) numărul de ordine din tabela de individualități, pe un vector cu lungime fixă (pentru lungimea 5 vom avea 32 de individualități posibile).

- pentru Ψ_r , având o ierarhie de clase de adâncime fixă putem construi un vector binar de aceeași lungime, prin metoda "left-right" [Smolensky92] (pentru adâncimea 5 vom avea 31 de clase reprezentabile).

- pentru Ψ_l putem păstra reprezentarea locală, pentru că în general, numărul relațiilor este mic în orice univers al discursului și se obișnuiește să se facă o standardizare a lor în contextul operațiilor canonice cu grafuri conceptuale. Pentru că am descris în capitolul și un arbore ierarhic al relațiilor, putem aplica și aici metoda "left-right", în caz că sunt prea multe relații. Se poate observa că vom avea în total, un vocabular de 68 de cuvinte.

3.4. Concluzii după prima etapă

Plecând de la cele două paradigme enunțate, putem să schițăm o posibilă metodologie pentru realizarea unui sistem cognitiv bazat pe rețele neurodinamice și translatate tensorială. Un astfel de sistem va primi la intrare un set de propoziții, care sunt transformate în puncte în spațiul intrărilor rețelei, formând o traiectorie. Această intrare va fi percepută ca o perturbație, rețeaua fiind inițializată într-un atractor cu o regiune de lucru care conține proiecția intrării la nivel de spațiu total al stărilor. Rețeaua va evolua pe o traiectorie, spațiul de ieșire conținând o proiecție a traiectoriei complete, care este compusă din puncte care vor fi traduse în grafuri conceptuale.

Se poate observa că procesarea efectivă se face exclusiv prin rețeaua neuronală. Translatarea grafurilor conceptuale în vectori de activare neuronală (și invers) este o sarcină auxiliară, care nu afectează absolut deloc evoluția intrinsecă a sistemului. Din acest motiv, un astfel de sistem s-ar încadra din punctul de vedere al clasificării propuse în introducere, în clasa sistemelor neurosimbolice unificate. La prima vedere ar părea un sistem hibrid translațional, dar în aceste sisteme, translația cunoștințelor are și rol de procesare suplimentară, pe când aici nu apare așa ceva. Mai mult, după cum s-a prezentat în subcapitolul 3.2.3 translatarea tensorială poate fi implementată (în sensul graf→vector) prin o serie de arhitecturi conexiuniste, care deși nu sunt uniforme cu rețeaua recurentă folosită, se bazează pe aceeași paradigmă computațională.

Se va vedea însă că un astfel de model este simplist și are o serie de manifestări care nu permit funcționarea sa corectă. *Soluția este una hibridă, în care pe lângă procesarea pur neuronală, au loc și procesări pur simbolice.* Capitolele următoare prezintă argumentele pentru susținerea aceste afirmații.

Partea II-a
Dezvoltarea modelului propus

4. Modelul hibrid translațional

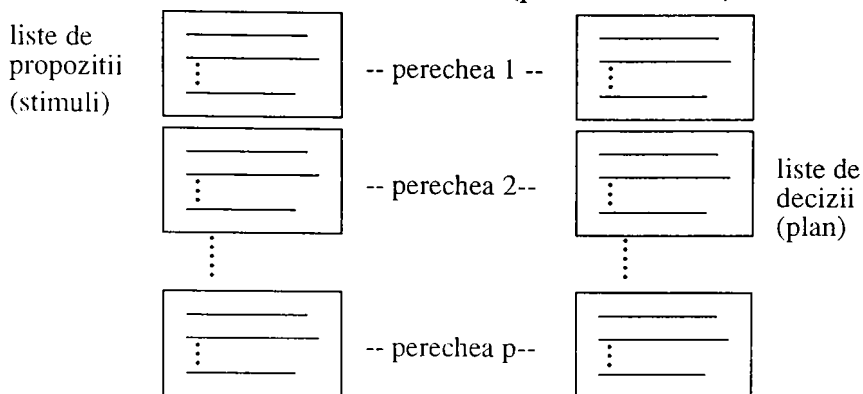
Metodologia care folosește pentru realizarea unui proces cognitiv simbolic o rețea neurodinamică și o schemă de translație tensorială a avut nevoie de o confirmare practică prin implementarea unei aplicații de test. Găsirea aceste aplicații s-a dorit a fi un proces de raționament, în genul în care o rețea Towell poate efectua raționamente. Aplicația pe care am intenționat să o dezvolt, era bazată pe un *model în buclă închisă* (vezi mai jos), care presupunea dezvoltarea a două rețele neurodinamice și a interfețelor simbolice aferente, cu o aplicație imediată în planificare. Sarcina s-a dovedit a fi mult prea complexă, dezvoltându-se doar o singură rețea, dar o serie de rezultate și observații s-au dovedit valoroase pentru dezvoltarea modelului pe alte căi (prezentate în capitolele 5 și 6).

4.1. Modelul cognitiv conexionist în buclă închisă

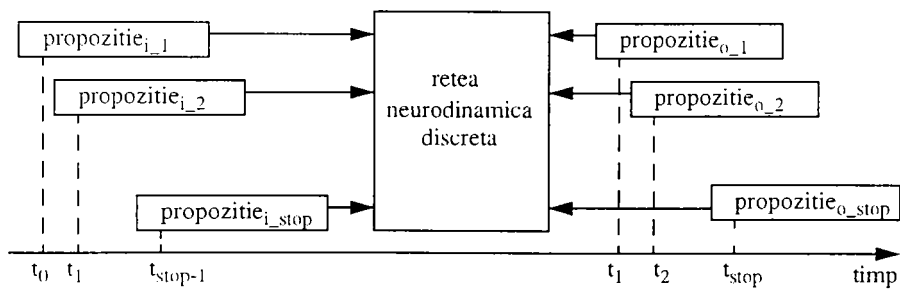
Acest model cognitiv a fost propus de Rumelhart, Smolensky, McClelland și Hinton [Rumelhart&al86]. El este constituit din două rețele neuronale (a căror topologie nu este precizată - dar din punctul meu de vedere sunt rețele neurodinamice) legate în buclă închisă. Se consideră că un sistem cognitiv, primește stimuli din mediul înconjurător, pe care îi analizează și sintetizează, luând apoi decizii care sunt aplicate asupra mediului. O serie secvențializată de stimuli va duce la generarea unei secvențe de decizii. Dacă aplicăm asupra sistemului toți stimuli pe rând, culegând din ieșire deciziile generate, fără să luăm în considerare influența lor imediată asupra mediului, vom comite o eroare. Sistemul este lipsit de "feed-back", deciziile fiind generate fără să se țină cont de modificările pe care acestea le aduc asupra mediului. Necunoscând starea curentă a mediului, putem lua decizii greșite. (vezi figura 4-1)

Soluția este de a modifica stimuli în raport cu mediul. Modificările pe care le facem asupra mediului prin aplicarea deciziilor, trebuie să modifice în mod dinamic

a. baza de antrenare (pur simbolica)



b. antrenarea cu o pereche



c. utilizarea

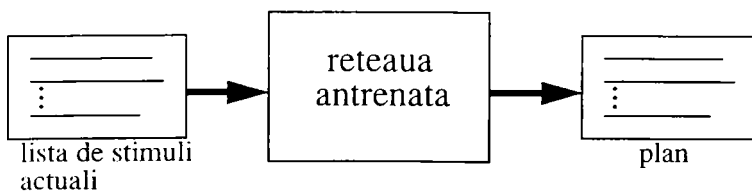


figura 4-1.

natura și conținutul stimulilor. Pentru aplicarea acestei strategii este însă necesar ca să construim sistemul cognitiv în cadrul mediului în care își va desfășura activitatea (vezi figura 4-2). Experimental, acest lucru este foarte costisitor și uneori imposibil, pentru că trebuie să creăm un mediu de laborator și să "antrenăm" sistemul în funcție de acesta. Acest mod de lucru este eficient într-o fază de testare, dar nu în una de construire (considerând construirea ca fiind de natură conexionistă, prin învățare).

O soluție mai bună este construirea unei alte rețele conexioniste, pe baza cunoștințelor acumulate despre mediu, care să-l "simuleze". Practic, o astfel de rețea, ar avea la intrare, deciziile luate de prima rețea și ar genera chiar stimulii pe care mediul l-ar oferi sistemului, ca în figura 4-3. În acest caz, dacă oferim sistemului doar un prim stimul, de "start" și îl deconectăm de la mediul exterior, acesta va genera o primă decizie, care la rândul ei va genera prin rețeaua de modelare, un nou stimul, ș.a.m.d., până când nu se mai poate lua nici o decizie, sau nu se mai modifică mediul virtual. Evident, se poate întâmpla și fenomenul unei ciclări infinite, care nu este de dorit. Acest model are o plauzibilitate psihologică, putând fi explicația pentru modul în care oamenii își crează scenarii (de exemplu, cum ne imaginăm că se va desfășura o discuție cu o persoană cunoscută și ale cărei reacții le putem anticipa) [Mäkeläinen94].

Există două considerente legate de designul acestui sistem: sincronizarea și comunicarea. Sincronizarea celor două rețele este legată de caracteristica temporală a rețelelor folosite. Dacă ambele rețele sunt continue, nu se pune explicit problema sincronizării, doar cel mult la nivel de realizare pseudo-continuă a traiectoriilor. Cazul în care o rețea este continuă și una discretă nu merită a fi luat în considerare. Dacă ambele rețele sunt discrete, sincronizarea se poate face la nivel de pași de evoluție discretă (vezi capitolul 2), iar cea mai naturală cale este echivalarea "unu-la-unu" a pașilor (existând și posibilitatea de "unu-la-doi", etc). O implementare cu două rețele continue este foarte costisitoare, pentru că impune fie o simulare digitală a unor sisteme dinamice neliniare, fie o implementare dedicată în hardware. Soluția pe care am folosit-o în continuare a fost cea de sincronizare discretă "unu-la-unu".

Comunicarea se poate face pe două căi care par foarte asemănătoare ca efect, dar diferențiază net sistemul rezultat pe două clase diferite la nivelul taxonomiei neurosimbolice (vezi Introducerea). Cea mai naturală pare a fi de a furniza ieșirea uneia dintre rețele, în formă vectorială, la intrarea celeilalte rețele. Pentru a obține rezultatele în forma accesibilă de grafuri conceptuale, se vor traduce doar acei vectori de ieșire care sunt considerați necesari. Se poate observa că în acest caz sistemul își păstrează caracterul de uniformitate, procesarea făcându-se doar la nivel neuronal, chiar dacă este separată pe două rețele. Cealaltă alternativă este de a traduce fiecare vector generat de cele două rețele în graf conceptual, iar înainte de a-l furniza ca intrare, de a-l traduce din nou în formă conexionistă (prin translatarea tensorială prezentată în subcapitolul 3.3). Deși o astfel de schemă pare redundantă și costisitoare din punct de vedere al translatării, se va vedea în continuare că este singura posibilă.

Translatarea vectorilor de ieșire în grafuri conceptuale duce la discretizarea

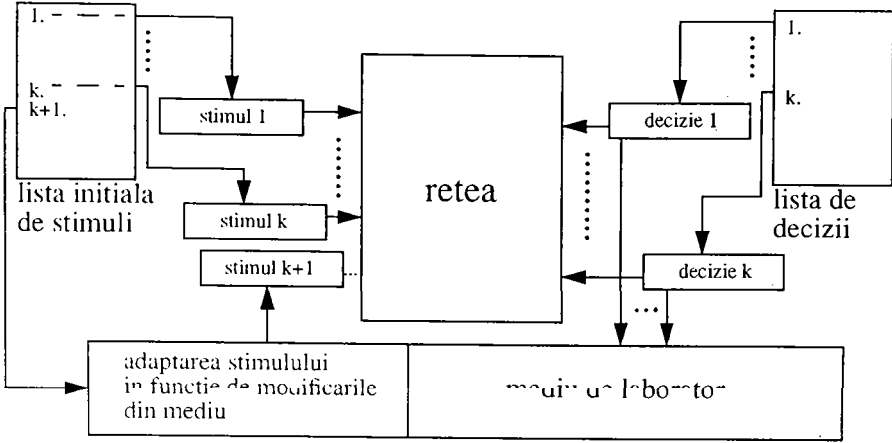


figura 4-2.

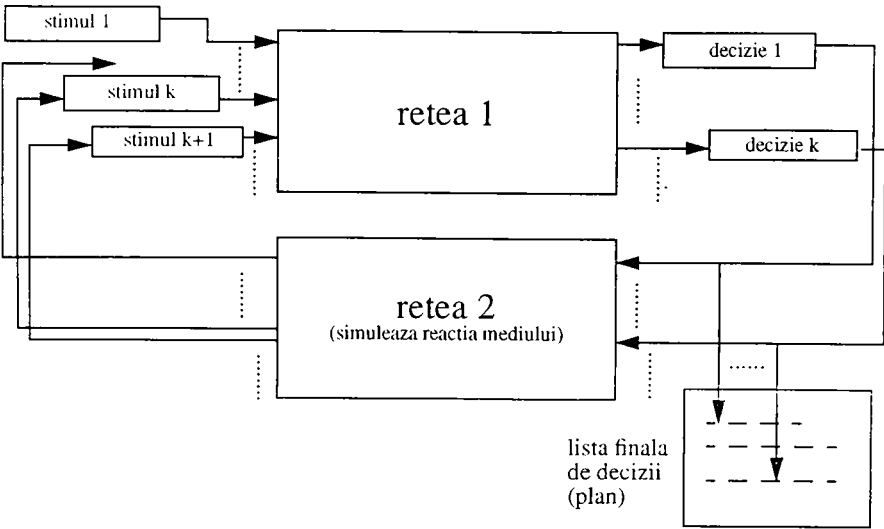


figura 4-3.

informației, pentru că vectorii rezultați la ieșirea rețelelor au valori continue în intervalul $[-1.0, 1.0]$. Mai mult, în cazul în care valorile sunt afectate de mici erori inerente legate de funcționarea aproximativă a rețelelor, translația *vector* \rightarrow *graf-conceptual* \rightarrow *vector* duce la transformarea valorilor ușor eronate în valori corecte, printr-un proces de restabilire a informației simbolice reprezentate de vectorii conexioniști de intrare. După cum se va vedea în continuare, această restabilire este necesară. Am denumit această schemă de comunicare, *translație-retranslație*.

4.2. Principul divergenței

Datorită faptului că procesarea neuronală introduce inerent un factor de aproximare (sau mai bine spus, o eroare), desfășurarea unui raționament în buclă închisă va duce la o deteriorare progresivă a calității propozițiilor generate. În loc să avem după un număr de iterații de raționament o planificare bună, vom obține una total necorespunzătoare. Acest lucru se petrece din două motive:

i. *Modelul lui Rumelhart este incomplet.* Plauzibilitatea lui este evidentă doar dacă există un mecanism de control, care supervizează cele două rețele conectate. În termeni psihologici, acest supervisor are ca și corespondent gândirea serială, care analizează fiecare pas în dezvoltarea listei de decizii. În termeni computaționali este vorba de un proces, care trebuie să verifice "calitatea" informației vehiculată între cele două rețele.

ii. *Erorile din sistem se acumulează progresiv,* pe măsură ce se trece de la o iterație la următoarea. Acest fenomen poate duce la imposibilitatea translației vectorilor în grafuri conceptuale, aceștia fiind prea afectați de "zgomotul" introdus de erorile acumulate.

Următoarea demonstrație este foarte simplă și imediată:

Propoziția 1

Rețeaua 1 și rețeaua 2 au fost antrenate cu bazele de învățare $B_1 = \{S_1, S_2, \dots, S_p\}$, respectiv $B_2 = \{T_1, T_2, \dots, T_r\}$, unde $S_i = (S_i^{in}, S_i^{out})$ și $T_j = (T_j^{in}, T_j^{out})$, cu condiția $\{T_j^{out}\}_{j=1,r} \subseteq \{S_i^{in}\}_{i=1,p}$.

Propoziția 2

Dacă rețeaua 1 primește la intrare $X_i \approx S_i^{in}$ ($i \in \{w | w=1,p\}$) ea va genera la ieșire $Y_j \approx S_i^{out}$. Se va introduce o eroare suplimentară nenulă, conform:

$$|X_i - S_i^{in}| - |Y_j - S_i^{out}| = \epsilon_1 \quad (13)$$

Propoziția 3

Aplicând Y_j la intrarea rețelei 2, și considerând că $Y_j = T_j^{in}$ ($j \in \{w | w=1,r\}$) aceasta va

genera $X_2 \approx T_j^{out}$. Se va introduce o eroare suplimentară nenulă:

$$|Y_1 - T_j^{in}| - |X_2 - T_j^{out}| = \epsilon_2 \quad (14)$$

Propoziția 4

Aplicând iterativ propozițiile 2 și 3 vom obține un șir $\epsilon_1, \epsilon_2, \dots, \epsilon_p, \dots$

Orice sistem conexiunist respectă regula propagării erorii. Dacă erorile medii de la intrare sunt în relația $\epsilon_1^{in} < \epsilon_2^{in}$ atunci și erorile de la ieșire sunt în aceeași relație $\epsilon_1^{out} < \epsilon_2^{out}$. Se poate observa din propozițiile 2 și 3 că:

$$|X_1 - S_i^{in}| < |X_2 - S_k^{in}| \quad (15)$$

unde $S_k^{in} = T_j^{out}$ (conform propoziției 1 trebuie să găsim o astfel de de identitate). Putem afirma:

Consecință (regula divergenței)

Sirul ϵ_i ($i=1, f$) nu este convergent.

Eroarea crescând în mod progresiv, la un moment dat t vom ajunge în situația că nici una din valorile X, Y vehiculate, nu mai seamănă cu valorile din bazele de învățare B_1, B_2 , rezultatele generate fiind nefolositoare.

Această regulă a divergenței (pe care am ridicat-o la rang de principiu), l-am observat experimental, în toate cazurile în care am lucrat cu rețele aflate în buclă închisă. Din acest motiv, folosirea comunicării bazate pe translație-retranslație este obligatorie. Acest fapt va schimba natura sistemului, din sistem uniform în unul hibrid, pentru că procesul de translație-retranslație este unul simbolic, în special prin faptul că apelează la structurile simbolice ierarhice atașate grafurilor conceptuale utilizate. Pentru că sistemele hibride sunt puternic criticate de către susținătorii modelului uniform, am considerat necesar să clarific într-un subcapitol separat această opțiunea a mea de a utiliza modelul hibrid.

4.3. Renunțarea la uniformitate și folosirea unei structuri hibride

Punctul de plecare în argumentare poate fi *principiul uniformității*. Acesta este definit în felul următor:

Principiul uniformității

O problemă computațională poate fi rezolvată și implementată eficient dacă se folosește doar aparatul formal al unei singure arhitecturi computaționale.

Principiul a fost introdus de Newell (1980), reluat și rafinat de Langley (1987) (vezi [Newell90]), ambele abordări analizând cazul modelelor (arhitecturilor) cognitive clasice, iar de pe pozițiile conexionismului a fost argumentat de Smolensky [Smolensky1988] și Jordan (1986). O mulțime de alți autori ([Clark89], [Shavlik91], chiar [Smolensky92], în plus [VanGelder95], [Sharkey95]), au combătut în perioada 89-95 acest principiu.

Aplicarea principiului uniformității presupune definirea unui set de operații de bază, foarte simple, care să definească o arhitectură computațională. Extremele care apar prin această aplicare duc la două direcții (mai exact mentalități) de cercetare în modelarea proceselor cognitive (procesarea simbolică cu atribute de inteligență). Partizanii principiului prin aplicarea metodelor sintactice și logice, susțin următoarea aserțiune:

Orice rezultat al unei activități cognitive poate fi explicat printr-un model matematic derivat din logică și teoria limbajelor formale, sau probabilități.

Cei din tabăra conexionistă (aderenți ai principiului uniformității) susțin că:

Orice rezultat al unei activități cognitive poate fi obținut folosind modelul computațional conexionist.

Argumentele pe care le vom dezvolta în continuare, sugerează că ne putem opune principiului uniformității, combătând avantajele sale pe toate planurile (claritate, implementabilitate, robustețe, consistență matematică, plasticitate, transparență). Alternativa este de a propune o imagine a unui sistem cognitiv, ca o "combinare" a unei multitudini de arhitecturi computaționale, fiecare fiind adaptată pentru anumite probleme specifice, fiecare fiind univoc definită de un set de operații de bază. Aceste operații sunt distincte de la o arhitectură la alta.

Pentru a ilustra dificultățile care apar relativ la capacitatea de a realiza artificial procese cognitive cu cele două soluții extreme prezentate anterior, vom folosi ca instrument ajutător, noțiunea de funcție. Aparent, rezultatul cognitiv care se obține prin aplicarea unei anumite funcții asupra unor date inițiale, se referă doar la CE funcție este invocată și nu ne interesează CUM este calculat rezultatul funcției. Dacă reducem modelul conexionist și pe cel sintactico-logic la mașina Turing (și este posibil [Giles&al95]) am putea demonstra că putem calcula la fel de bine aceleași funcții cu ambele arhitecturi, evident cu rezultate diferite în ceea ce privește viteza de calcul, și în acest caz soluția de implementare este dată doar de indicii de performanță.

Se poate considera că problema (mai exact greșeala) este dată de aserțiunea (probabil insuficient analizată) dată de propoziția "calculând aceleași funcții". Oare o mașină care calculează 7×8 prin asocierea dintr-o memorie asociativă a rezultatului $7 \times 7 = 49$ cu adunarea finală a lui 8, face "aceiași funcție" cu o mașină care adună pe 7 de 8 ori într-un acumulator? În sensul că prezentând două numere și o operație obținem un rezultat; același în ambele cazuri, răspunsul ar fi cu certitudine **da**, e vorba de

aceiași funcție. Modul în care se ajunge la rezultat va fi diferit. Din punct de vedere al vitezei și al robusteții, prima variantă (deși neconvențională) este mai bună. A doua variantă este mai simplă de implementat dar consumă timp, iar în cazul defectării sistemului chiar în timpul evaluării funcției, rezultatul primei versiuni va fi mai apropiat de rezultatul final în majoritatea cazurilor decât rezultatul celei de-a doua (paradoxal, modelul conexionist este mai bun într-un astfel de caz particular, de implementare a tablei înmulțirii).

Exemplul dat poate fi puțin derutant, pentru că se discută aici despre procese cognitive, nu de operații aritmetice, deci funcția de mai sus trebuie privită prin prisma cercetătorului din domeniul Inteligenței Artificiale și nu a elevului de școală primară. Mai ales pentru că performanța umană în domeniul procesării cognitive se măsoară în primul rând prin viteză, minimizarea erorii și robustețea raționamentului. Deci, în acest context, este important să analizăm nu atât CE anume face funcția, ci CUM este calea prin care face. Ajungând din nou la cele două soluții extreme de implementare, memorarea asociativă a datelor este o caracteristică de bază a modelului conexionist. Funcții din clasa celei prezentate pot fi realizate prin rețele conexioniste. Este însă exagerat a afirma că oferind o astfel de implementare (apropiată de modul de calcul uman, care are memorată asociativ tabla înmulțirii) putem explica microstructura tuturor proceselor cognitive. Există chiar o analogie favorită a partizanilor modelelor conexioniste uniforme, cea dintre modelul sintactico-logic cu fizica Newtoniană și dintre modelul conexionist cu fizica cuantică. Putem cita:

"Se poate argumenta că procesarea simbolică convențională oferă un model macroscopic, analog mecanicii Newtoniene, pe când modele noastre [conexioniste] oferă o analiză microscopică, analoagă teoriei Cuantice. Observați că, în domeniul lor de acoperire, aceste două teorii pot face predicții precise asupra evoluției și comportamentului obiectelor din lumea reală... Totuși, în anumite situații, teoria Newtoniană nu este corectă. În aceste situații, trebuie să ne bazăm pe teoria Cuantică."

(Rumelhart & McClelland, PDP, 1986, p.125, sublinierile îmi aparțin)

Pare a fi o analogie rezonabilă, dar e corectă numai în cazul în care teoria Newtoniană ar fi corectă pentru o anumită clasă de fenomene. Dar nu este așa. Mai exact, putem afirma, că în anumite situații, teoria Newtoniană (deși incorectă la nivel de detaliu), este aplicabilă, iar teoria Cuantică explică și de ce este așa. Doar în acest caz am respecta și principiul uniformității aplicat la fizică. Dar se poate observa că în această situație nu mai există aceeași analogie între modelul convențional și cel conexionist. Revenind la ideea paradigmei subsimbolice [Smolensky88], ar însemna că descrierile seriale, simbolice ale macroproceselor cognitive ar fi aproximări ale unei descrieri de detaliu, care este oferită de nivelul de microprocesare conexionistă. Putem cita din nou:

"Noi vedem macroteoriile ca aproximații ale unei structuri cu detalii de distribuire microscopică..."

(Rumelhart & McClelland, PDP, 1986, p.126)

Ceea ce e evident că nu corespunde realității. Dimpotrivă, macroprocesele nu sunt deloc aproximații, ele sunt fenomene deterministe, modelate printr-un aparat formal, care este considerat chiar prea rigid. În contrast, modelul conexionist este tocmai cel care introduce aproximarea.

Această analogie greșită poate duce cercetătorul, într-o primă fază de entuziasm "conexionist" către o subestimare a modelelor convenționale. Aplicând în mod necondiționat principiul uniformității în realizarea proceselor cognitive conexioniste putem face o mare greșeală, folosind instrumente deloc adecvate pentru diverse clase de procese cognitive, care pot fi relativ ușor rezolvate cu arhitecturi convenționale. Modelul conexionist NU ESTE uniform, relativ la domeniul cognitiv. Mai mult, sunt situații în care modelele convenționale sunt complete, într-un mod în care fizica Newtoniană nu este niciodată completă.

Exemplele sunt ușor de dat. Dacă un subiect uman simulează efectiv un proces simbolic serial, în același mod în care acesta ar fi realizat de o mașină Van Neumann, procesul va fi corect. Punct. Nu "aproximativ corect". Evident că suportul fizic al unei astfel de simulări este tot un sistem conexionist. Dar nu este relevant din punct de vedere procesului simbolic ce se petrece la acest nivel, pentru că nivelul de definire (prin formalism) al procesului nu ia în considerare absolut deloc un model conexionist. Deci, în mod clar, în această situație principiul uniformității este dăunător. Practic, pentru acest tip de procese (să le numim seriale), putem introduce conceptul de mașină virtuală, care poate rula pe orice fel de mașină (Van Neumann, conexionistă, sau chiar mintea umană). Rezultatul este cel care contează și el este mereu același (corect și fără aproximări).

În acest caz, întrebarea care se pune, este cum poate mintea umană să aibă capacitatea de a rezolva probleme tipice pentru procesarea simbolică, serială, sintactico-logică, dacă ea este fundamental o arhitectură conexionistă, deci aproximativă. Speculația pe care o putem face, este că mașina conexionistă pe care o implementează creierul poate ocazional să simuleze o mașină Von Neumann. Dacă sistemul tipic uman de procesare al informației își desfășoară activitatea prin asociere și "relaxarea" într-o soluție, fără a folosi operațiile logice, de ce putem efectua raționamente bazate pe matematică, sau deducții științifice solide? Cum putem raționa noi pur logic, dacă operațiile de bază nu au nimic în comun cu logica? Răspunsul este foarte interesant: *Capacitatea noastră de a executa procese formalizabile, guvernate de reguli precise este dată de abilitatea de a crea instrumente artificiale, exterioare sistemului cognitiv, care sunt reprezentări fizice ale problemelor, pe care le putem manipula prin scheme foarte simple pentru a obține în final rezultatele unor probleme abstracte foarte dificile.*

Sunt două clase distincte de task-uri computaționale de acest fel (clase pentru care modelele conexioniste nu se pretează):

- i. procese de raționament serial, în care ordinea operațiilor este vitală.
- ii. procese de raționament "generativ", în care structuri simbolice întregi sunt produse prin aplicare unor reguli (de inferență simbolică - unificare și pattern-

matching) asupra unei baze de cunoștințe reprezentabile la nivel simbolic.

Ca exemple putem da la punctul 1: planificarea conștientă a acțiunilor, raționamentul logic, demonstrația matematică și gândirea abstractă. La punctul ii., cel mai important task, pentru mintea umană, este generarea discursului.

În mod cert, genul de procese cognitive prezentat mai sus este "conștient" pentru persoana care le execută. Asocierile se execută însă mult mai des, iar în majoritatea cazurilor subiectul nu este conștient că un astfel de proces cognitiv (de asociere) este în curs de desfășurare. Aceste procese de asociere pot fi mai multe active simultan, conștientizarea putându-se face doar în momentele de "stabilitate" conexionistă, când spunem că asocierea dă rezultate. Procesele seriale se desfășoară însă într-un ritm mult mai lent, fiind necesară întreaga lor conștientizare și fixarea în context pentru fiecare pas al procesului. Mai mult, pentru a nu pierde "cursul" unui astfel de proces, avem nevoie de memorie externă și reprezentări semiotice adecvate, cu semantică asociată, astfel că în majoritatea cazurilor suntem nevoiți să folosim hârtia și creionul. Atunci când facem acest lucru, folosim instrumentele formale externe menționate mai sus, pe care le-am învățat să le manipulăm. Ele sunt inventate de noi, cel mai adesea făcând parte din fondul general de instrumente formale ale științei. Dacă de exemplu, avem o serie de numere mari și trebuie să adăugăm 2 la fiecare și să obținem suma totală, în mod cert, o persoană obișnuită va avea nevoie de hârtie și creion. Și mai avem nevoie de acele reprezentări exterioare pentru numere și operații aritmetice. Mediul exterior pe care ni l-am creat devine o extensie esențială a minții noastre. Chiar dacă vom învăța să executăm aceste operații "în gând", renunțând la instrumentele exterioare, acest lucru este posibil pentru că ne construim un "model mental" al structurilor exterioare folosite, iar manipularea acestora nu se face exterior, ci interior, dar serial și conștient. Dar nu trebuie să pierdem din vedere că funcții de tipul $7 \times 7 = 49$ se realizează prin simplă asociere.

Făcând o sumă a acestor abilități cognitive, putem specula că sunt trei capacități combinate, care permit minții umane să execute funcții complexe, cu un grad mare de abstractizare:

- a. capacitatea asociativă fundamentală (simplă în esență)
- b. capacitatea de a modela mental mediul înconjurător
- c. capacitatea de a manipula fizic instrumentele exterioare și de a percepe efectele acestor manipulări.

Practic, putem reduce aceste capacități, la una de bază, care este pur neuronală ca operații de bază (relaxarea termodinamică perturbată). Dar acest lucru nu este relevant pentru că nivelele b. și c. sunt specifice pentru mașina Von Neumann, care separă datele procesate de funcțiile care le procesează (lucru care nu se întâmplă în modelele conexioniste). Slăbiciunea arhitecturii Von Neumann aplicată în modelele cognitive este dată de faptul că suntem constrânși să tratăm toate procesele cognitive ca depinzând de manipularea instrumentelor simbolice externe. Asta face ca modelele cognitive derivate prin această arhitectură să fie limitative și rigide.

Noul start în acest domeniu, oferit de modelul conexiunist, ne oferă o serie de mecanisme care pot completa neajunsurile arhitecturilor convenționale. Greșeala este de a crede că aceste noi modele le pot înlocui complet. Problema este cum vom reuși să completăm lacunele sistemelor convenționale cu ajutorul modelului conexiunist. O alternativă atractivă este de a construi sisteme conexioniste, care prin învățare sau construcție, să fie capabile de a executa procese seriale, abstracte. Din păcate, acest lucru nu este de loc simplu. Instrumentele de reprezentare a realității fizice s-au format în cursul istoriei, fiind în majoritatea cazurilor, rafinate prin mai multe generații, și aduse în forma finală și păstrate, prin transmitere culturală. În plus, numărul și complexitatea lor este foarte mare. Implementarea acestor instrumente prin rețele conexioniste ar fi plauzibilă din punct de vedere al paralelismului model uman - mașină, dar nu ar fi eficientă tehnologic.

Concluzia pe care o putem extrage din această discuție, este că avem nevoie de sisteme cognitive care au ambele modele încorporate. Un astfel de sistem ar avea două nivele: unul rapid, bazat pe modelul conexiunist, care să genereze rezultate prin asociere, și unul convențional, care să obțină rezultatele prin procese seriale, pur simbolice. Implementarea celor două nivele trebuie să se facă renunțând la principiul uniformității, utilizând modelele convenționale când este cazul și modelul conexiunist pentru calculul de tip "intuitiv". Legătura dintre cele două arhitecturi este o problemă de comunicație și translație. Am arătat în capitolul dedicat paradigmei subsimbolice (3.2) că există posibilitatea de a folosi structuri simbolice în contextul modelelor conexioniste.

4.4. Antrenarea unei rețele neurodinamice cu grafuri conceptuale

4.4.1. Dimensiunea rețelei folosite

Se obișnuiește să se aleagă în cazul utilizării sistemelor neurodinamice un număr de unități de calcul (neuroni) pentru care dimensionalitatea spațiului de intrare/ieșire (însurate) să reprezinte 10-20%. De exemplu, pentru o aplicație de robotică cu trei grade de libertate la intrare și 6 la ieșire este absolut necesar să avem 50-100 de neuroni în total (numărul depinzând și de dimensiunea bazei de învățare). Pentru cazul procesării simbolice propuse s-ar deduce că am avea nevoie de $3 \cdot 10^4 - 6 \cdot 10^4$ neuroni și ar rezulta o matrice de ponderi cu $1 - 2 \cdot 10^8$ elemente, adică sute de milioane de numere în virgulă flotantă, respectiv peste un GByte de memorie necesară. Din acest punct de vedere, la nivelul anului 1997, tehnologic și practic problema nu reprezintă nimic deosebit, dar algoritmul de învățare folosit [Williams90] are complexitatea $O(n^4)$, unde n este numărul de neuroni. Am avea deci nevoie de un număr de $4 \cdot 10^{16}$ iterații pentru o singură epocă (și de obicei avem nevoie de mii de epoci) și se observă imediat că aceste ordine de mărime pentru timpul de rulare necesar pentru o antrenare completă și care sunt în jur de 10^{20} , sunt impracticabile.

O propunere interesantă, făcută de Smolensky este de a folosi pentru astfel de probleme (unde vectorii nu reprezintă mapări ale realității fizice - mărimi cu dimensiune

euclidiană - ci sunt combinații binare rezultate din produse și superpoziții tensoriale) aceiași neuroni și pentru intrare și pentru ieșire, renunțând complet la neuronii ascunși. Explicația este simplă: majoritatea intrărilor în rețea la un moment t sunt 0, iar în pasul următor, la ieșire se vor obține puține intrări setate pe 1. Astfel, perturbația totală care se aplică asupra rețelei și modificarea în starea ei (ieșirea) de la un moment de timp la altul, sunt ambele foarte mici. Nu trebuie uitat faptul că starea unităților este o stare cuprinsă în intervalul $[0,1]$, dar la ieșire, vom considera că orice stare care are valoarea mai mare decât 0.5 este echivalentă cu 1. La intrare, perturbațiile de tip 1 vor fi ponderate de valorile calculate la învățare pentru respectivele intrări.

Se poate compara această rețea cu una care are doar 10% din neuroni afectați pentru intrări și ieșiri, dar în care modificările aparute în starea acestora între două momente succesive de timp sunt mult mai mari. Altă explicație este următoarea: capacitatea de învățare a unei rețele este dată în primul rând de dimensiunea matricii de ponderi. Este clar că pentru un număr relativ mic de propoziții care generează vectori binari (cu puține elemente setate pe 1) matricea de ponderi este arhisuficientă dacă are $3125 \cdot 3125 = 9,765,625$ elemente (aproximativ 80 MBytes pentru o reprezentare în virgulă flotantă). Însă din motive de timp de rulare, o rețea cu mii de unități, este prea mare. A trebuit deci să reduc dimensiunea rețelei la ordinul sutelor pentru a putea rula antrenări fezabile și care se conțină un număr suficient de epoci.

4.4.2. Timpul necesar antrenării și paralelizarea algoritmului

Am observat că rulând algoritmul de învățare pentru un număr de 3125 unități și o matrice corespunzătoare de 80 Mbytes pe o stație SUN Sparc într-o rețea cu serverul de fișiere la distanță, timpul de rulare pentru o singură epocă ar atinge aproximativ 900 de zile!. Timpul scade la aproximativ 300 de zile dacă matricea este memorată pe discul local, iar stația este folosită exclusiv pentru acest scop. Dacă am avea un sistem cu peste 80 MBytes memorie RAM (nu am avut la dispoziție astfel de resurse în prima fază), timpul s-ar reduce cu un ordin de mărime (sau chiar mai mult). Totuși, chiar și în această situație, pentru o singură epocă, timpul de rulare rămâne de ordinul zilelor. Singura soluție pentru obținerea unor timpi de rulare acceptabili este procesarea paralelă. Implementarea algoritmului de învățare în PVM [Szirbik&Somlo 95] și rulare pe 4 stații în paralel, a relevat posibilitatea unei accelerări foarte apropiate de factorul maxim de "speed-up" (4 în acest caz). Spre deosebire de restul algoritmilor de învățare (un exemplu clasic este backpropagation-ul aplicat la rețele feedforward), algoritmul lui Williams și Zipser (a cărui variantă modificată prezentată în capitolul 2, am folosit-o și eu) este unul care se pretează perfect la procesarea paralel-distribuită.

Mai multe rezultate experimentale relativ la paralelizare (care nu este subiectul acestei teze) pot fi găsite în [Szirbik95b]. Aceste rezultate au fost obținute prin rulare programului pe un sistem de calcul cu 80 de procesoare Sparc, fiecare cu 32 de MBytes memorie RAM locală și o memorie "comună" RAM de 128 MBytes (accesată ceva mai lent decât memoria locală). Procesoare comunicau prin mesaje, printr-o rețea

Ethernet specială, cu o rată de 100Mbs. Acest sistem de calcul, numit ZOO (al Facultății de Matematică și Informatică de la Universitatea Liberă din Amsterdam), dispunea de un sistem de operare artizanal, AMOEBA [Tannenbaum&Rennese89] care permitea multitasking-ul, dar și rezervarea întregului sistem pentru un singur program, în perioadele în care sistemul nu era folosit (8 p.m.-8 a.m. sau vineri 8 p.m. - luni 8 a.m.). Cele mai bune rezultate erau obținute în acest mod de rulare "rezervat". Partea de programare era foarte mult simplificată prin existența unui compilator (dedicat) pentru limbajul de programare cu facilități paralel-distribuite Orca [Bal89]. Limbajul fiind foarte elegant și eficient, testarea de aplicații ca aceasta era relativ facilă, rămânând în sarcina utilizatorului numai problemele aplicației cu care acesta se confruntă.

Utilizând aceste resurse de calcul, pentru rețeaua de 3125 de unități și un set restrâns de tipare (3-5 tipare) formate din secvențe de 10 ($t=0, \dots, 9$) grafuri conceptuale de dimensiuni reduse care erau traduse în vectori, am obținut că durata unei epoci era de 3 ore și 20 de minute. Având în vedere că ulterior am remarcat că sunt necesare un număr de aproximativ 500-1500 epoci pentru astfel de task-uri de învățare, acest timp de rulare nu era corespunzător. Singura soluție pe care am găsit-o a fost reducerea dimensiunii rețelei, astfel încât să putem încadra aproximativ 500 de epoci într-un interval de 12 ore, deci aproximativ 1 minut/epocă. Aceste concluzii au fost obținute prin experimente succesive, ajungând la următoarele rezultate:

nr. unități	timp pentru 1 tipar (sec)	timp pentru 2 tipare (sec)	timp pentru 3 tipare (sec)	timp pentru 4 tipare (sec)
3125	3210	6100	9050	12020
2000	1120	2080	3050	4030
1500	550	1004	1537	1940
1000	101	191	288	375
800	30	57	85	112
700	15	28	43	54
600	9	17	25	31

în condițiile în care nimeni nu mai folosește sistemul. Deci, experimentele se puteau desfășura doar într-un anumit interval de mărime al rețelei, de [700,800] unități de calcul. Utilizarea unui număr atât de mic de unități, face dificilă aplicarea reprezentării tensoriale. Fiind vorba de un tensor de gradul 5, dimensionalitățile spațiilor vectoriale V_F , V_R și V_L trebuie să fie de maxim 3 sau 4 (pentru că $4^5=1024$ și $3^5=243$). Acest lucru înseamnă că putem avea 8-16 individualități, 8-16 clase și 8-16 relații posibile. Am ales pentru experimentare o variantă cu 16 individualități, 14 clase (reprezentarea "left-right" generează $2^4 \cdot 2$ combinații de codificare posibile) și 8 relații codificate binar complet (000, 001, ..., 111). Rezultă exact 768 de unități și ne încadrăm foarte bine în

intervalul dorit:

nr. unități	timp pentru 1 tipar (sec)	timp pentru 2 tipare (sec)	timp pentru 3 tipare (sec)	timp pentru 4 tipare (sec)
768	22	41	60	79

Dar având un număr atât de mic de cuvinte în universul discursului, este imposibil de construit un set de tipare care să depășească un nivel semantic acceptabil. Din acest motiv am folosit doar simboluri literale:

- individualități: $F = \{ A, B, C, D, E, F, G, H, I, J, K, L, M, N, O, P \}$
- clase $R = \{ \alpha, \beta, \gamma, \delta, \epsilon, \zeta, \eta, \theta, \sigma, \chi, \lambda, \mu, \nu, \xi \}$
- relații $L = \{ r, s, t, u, w, x, y, z \}$

cu care am construit grafurile conceptuale și ierarhia de clase. Într-o fază în care încercam să testez doar modalitățile de antrenare, nu era necesar ca aceste simboluri să aibă o semantică asociată.

Un tipar era format din 10 grafuri conceptuale la intrare și 10 la ieșire. Au fost făcute experimente de antrenare a unei rețele cu 786 unități pentru 1, 2, 3 și 4 tipare. Scopul urmărit a fost nu de a învăța tiparele ci de a cerceta puterea de generalizare a rețelei. Aceasta capacitate poate fi determinată prin următorul experiment:

- se învață rețeaua câteva tipare, care au cam același "pattern" în timp, dar la "nivele" diferite
- după ce eroarea a scăzut sub 3% (în astfel de aplicații 5-1% este intervalul cel mai folosit [Hilario96]), se dă rețelei o componentă care respectă "pattern"-ul temporal al intrărilor învățate (indiferent de "nivel")
- dacă vom obține la ieșire o componentă care respectă "pattern"-ul temporal și semantic al ieșirilor învățate, atunci scopul a fost atins.

Nu era necesar în acest moment testarea unui model cu două rețele în buclă închisă, un sistem cu o singură rețea fiind suficient pentru testarea puterii de generalizare a unei rețele recurente neurosimbolice.

4.4.3. Testarea generalizării

Primele tentative de construire a unei baze de învățare pentru un model cu o singură rețea a plecat de la ideea unui posibil *dialog* om-sistem cognitiv artificial. Utilizatorul uman furnizează propoziții, într-o secvență care implică la fiecare pas, receptarea și interpretarea răspunsului dat de sistemul automat. Utilizatorul introduce

o propoziție $i(t)$ și primește imediat răspunsul $o(t+1)$. Pe baza acestui răspuns și a propozițiilor vehiculate între om și sistem în momentele anterioare, el va continua dialogul cu o nouă propoziție $i(t+1)$, care se află într-o relație de asociere semantică cu cele furnizate de utilizator sau generate de sistem. Din punct de vedere al plauzibilității unui astfel de sistem, putem consemna apropierea în spațiul euclidian a punctelor care reprezintă propoziții cu sens apropiat și corespunzător, crearea unei traiectorii în spațiul stărilor rețelei; un alt aspect este că o rețea recurentă are memorie, în sensul că păstrează informație despre punctele prin care a trecut.

Asocierea semantică între două sau mai multe propoziții se referă la vocabularul folosit, la clasele care apar în grafurile conceptuale rezultate și un anumit pattern temporal al "acțiunilor" (relațiilor între concepte) care rezultă din dialog, în sensul că anumite propoziții apar în mod obligatoriu după alte propoziții. Folosirea unui vocabular cu un număr foarte mic de cuvinte (vezi subcapitolul anterior) a făcut imposibilă construirea unor dialoguri în limbaj natural care să aibă sens. *Dar putem considera ca un limbaj inteligent și inteligibil de comunicare poate fi format și dintr-un număr foarte mic de simboluri.* O serie de regularități pot fi observate în dialogurile în limbaj natural, iar pe baza acestora se pot construi dialoguri care aparent nu au sens, dar pot testa un astfel de sistem minimal. Am construit o bază de învățare formată din patru dialoguri.

• primul tipar:

$i1(0) : (\alpha(A), r, \beta(B))$	\rightarrow	$o1(1) : (\gamma(C), s, \delta(D))$
$i1(1) : (\beta(B), t, \varepsilon(E))$	\rightarrow	$o1(2) : (\delta(D), u, \zeta(F))$
$i1(2) : (\varepsilon(E), w, \eta(G))$	\rightarrow	$o1(3) : (\zeta(F), x, \theta(H))$
$i1(3) : (\eta(G), y, \sigma(I))$	\rightarrow	$o1(4) : (\theta(H), z, \chi(J))$
$i1(4) : (\sigma(I), r, \lambda(K))$	\rightarrow	$o1(5) : (\chi(J), s, \mu(M))$
$i1(5) : (\lambda(K), t, v(N))$	\rightarrow	$o1(6) : (\mu(M), u, \xi(O))$
$i1(6) : (v(N), w, \alpha(P))$	\rightarrow	$o1(7) : (\xi(O), x, \beta(A))$
$i1(7) : (\alpha(P), y, \gamma(B))$	\rightarrow	$o1(8) : (\beta(A), z, \delta(C))$
$i1(8) : (\gamma(B), r, \varepsilon(D))$	\rightarrow	$o1(9) : (\delta(C), s, \zeta(E))$
$i1(9) : (\varepsilon(D), t, \eta(F))$	\rightarrow	$o1(10) : (\zeta(E), u, \eta(F))$

• al doile tipar:

$i2(0) : (\alpha(P), y, \gamma(B))$	\rightarrow	$o2(1) : (\gamma(C), z, \varepsilon(D))$
$i2(1) : (\gamma(B), r, \varepsilon(D))$	\rightarrow	$o2(2) : (\gamma(D), s, \delta(B))$
$i2(2) : (\varepsilon(D), t, \eta(F))$	\rightarrow	$o2(3) : (\delta(B), u, \zeta(E))$
$i2(3) : (\eta(F), w, \sigma(I))$	\rightarrow	$o2(4) : (\zeta(E), x, \theta(I))$
$i2(4) : (\sigma(I), y, \xi(A))$	\rightarrow	$o2(5) : (\beta(A), z, \delta(B))$
$i2(5) : (\alpha(A), r, \beta(B))$	\rightarrow	$o2(6) : (\delta(B), s, \zeta(F))$
$i2(6) : (\beta(B), t, \varepsilon(E))$	\rightarrow	$o2(7) : (\zeta(F), u, \eta(H))$
$i2(7) : (\varepsilon(E), w, \eta(H))$	\rightarrow	$o2(8) : (\beta(B), x, \delta(D))$
$i2(8) : (\eta(H), y, \sigma(J))$	\rightarrow	$o2(9) : (\theta(I), z, \chi(M))$
$i2(9) : (\sigma(J), r, \lambda(K))$	\rightarrow	$o2(10) : (\chi(M), s, \mu(O))$

• al treilea tipar:

$i_3(0) : (v(M), w, \alpha(O))$	\rightarrow	$o_3(1) : (\xi(B), x, \beta(P))$
$i_3(1) : (\alpha(O), y, \gamma(D))$	\rightarrow	$o_3(2) : (\beta(P), z, \delta(D))$
$i_3(2) : (\gamma(D), r, \epsilon(F))$	\rightarrow	$o_3(3) : (\delta(D), s, \zeta(F))$
$i_3(3) : (\epsilon(F), t, \eta(G))$	\rightarrow	$o_3(4) : (\zeta(F), u, \eta(G))$
$i_3(4) : (\delta(G), w, \sigma(K))$	\rightarrow	$o_3(5) : (\theta(H), x, \gamma(C))$
$i_3(5) : (\sigma(K), y, \xi(O))$	\rightarrow	$o_3(6) : (\gamma(C), z, \delta(G))$
$i_3(6) : (\alpha(A), r, \beta(B))$	\rightarrow	$o_3(7) : (\gamma(D), s, \delta(G))$
$i_3(7) : (\beta(B), t, \epsilon(D))$	\rightarrow	$o_3(8) : (\delta(G), u, \zeta(E))$
$i_3(8) : (\epsilon(D), w, \eta(F))$	\rightarrow	$o_3(9) : (\zeta(E), x, \theta(I))$
$i_3(9) : (\eta(F), y, \sigma(I))$	\rightarrow	$o_3(10) : (\theta(I), z, \chi(J))$

• al patrulea tipar:

$i_4(0) : (\lambda(L), t, v(M))$	\rightarrow	$o_4(1) : (\mu(M), u, \xi(O))$
$i_4(1) : (v(M), w, \alpha(O))$	\rightarrow	$o_4(2) : (\xi(O), x, \beta(P))$
$i_4(2) : (\alpha(O), y, \gamma(B))$	\rightarrow	$o_4(3) : (\beta(P), z, \delta(B))$
$i_4(3) : (\alpha(P), r, \beta(A))$	\rightarrow	$o_4(4) : (\delta(B), s, \zeta(F))$
$i_4(4) : (\beta(A), t, \epsilon(E))$	\rightarrow	$o_4(5) : (\delta(C), u, \zeta(G))$
$i_4(5) : (\epsilon(E), w, \eta(F))$	\rightarrow	$o_4(6) : (\zeta(G), x, \theta(I))$
$i_4(6) : (\eta(F), y, \sigma(I))$	\rightarrow	$o_4(7) : (\theta(I), z, \chi(J))$
$i_4(7) : (\sigma(I), r, \lambda(L))$	\rightarrow	$o_4(8) : (\chi(J), s, \mu(M))$
$i_4(8) : (\epsilon(E), t, \eta(G))$	\rightarrow	$o_4(9) : (\zeta(F), u, \mu(O))$
$i_4(9) : (\delta(D), w, \sigma(K))$	\rightarrow	$o_4(10) : (\theta(H), x, \gamma(B))$

Se pot observa o serie de regularități (care nu sunt respectate întotdeauna, tocmai pentru a sublinia aspectul de "common-sense reasoning"), o parte impuse de structura modelului și o parte definite de structura unor dialoguri în limbaj natural care au fost studiate pentru găsirea unor "pattern"-uri. Structura modelului a impus folosirea unor propoziții simple, cu structuri sintactice identice, formate din două concepte și o relație, pentru a avea o schemă cât mai simplă de decodificare din vector în graf conceptual (vezi capitolul 3). Lungimea dialogurilor trebuie să fie aceeași, adică 10 pași ceea ce este impus de natura algoritmului Williams&Zipser. Individualitățile trebuie să aparțină unor clase, care au fost definite după cum urmează:

$\alpha \supset \{A, O, P\}$	$\beta \supset \{A, B, P\}$	
$\delta \supset \{C, D, E\}$	$\gamma \supset \{B, C, D\}$	$\epsilon \supset \{E, F, D\}$
$\zeta \supset \{E, F, G\}$	$\eta \supset \{F, G, H\}$	$\theta \supset \{H, I, K\}$
$\sigma \supset \{I, J, K\}$	$\chi \supset \{J, M, N\}$	$\lambda \supset \{K, L, M\}$
$\mu \supset \{M, O, P\}$	$v \supset \{N, M, O\}$	$\xi \supset \{A, B, O\}$

Regularitățile induse de mine (nu cele impuse de structura modelului) au fost:

- la intrare pot apare jumătate din relații, iar la ieșire cealaltă jumătate.
 - relațiile apar în aceeași secvență, și la intrare și la ieșire:
 intrare : (r,t,w,y) ieșire : (s,u,x,z)
- și sunt sincronizate unele cu celelalte, chiar dacă prima relație diferă de la un dialog

la celălalt.

- conceptul al doilea dintr-o propoziție $i/o(t)$, apare ca prim concept în propoziția următoare $i/o(t+1)$.
 - anumite concepte nu pot apărea decât împreună cu anumite clase.
- (Observație: aceste reguli nu sunt respectate în toate circumstanțele).

Aceste regularități impun ca dialogurile "artificiale" construite să aibă un aspect structural foarte apropiat de dialogurile din limbaj natural.

Antrenarea rețelei s-a făcut prima oară cu un singur tipar. Cu toate că această antrenare a fost dificilă (vezi subcapitolul următor) ea a putut fi adusă în marja de eroare de 3%. Dar cu toate că furnizând la intrare secvențe de propoziții asemănătoare cu $\{i_1(k); k=0,9\}$, obțineam rezultate foarte apropiate (în majoritatea cazurilor, identice) cu secvența $\{o_1(k); k=1,10\}$, este prea puțin pentru a argumenta o generalizare. Am antrenat rețeaua cu 2, 3 și 4 tipare, observând că învățarea devine foarte dificilă cu cât numărul de tipare crește și cu cât tiparele sunt mai diferite de "pattern"-ul comun, după cum se vede și din graficul 4-4 (a,b,c,d.). Dificultatea nu s-a datorat numai creșterii numărului de epoci necesare, ci mai ales unui fenomen numit bifurcație, care apărea cu atât mai des cu cât numărul de tipare era mai mare, iar diversitatea lor creștea. Trebuie accentuat mai ales faptul că efortul de antrenare al celor patru tipare mai sus prezentate până la o eroare de sub 5% a fost enorm, de aproximativ 125 de ore de rulare pe sistemul ZOO în mod "unic-utilizator" (luând în considerare și reprizele de antrenare nereușită).

Sistemul cognitiv rezultat a fost testat în felul următor: utilizatorul furniza o propoziție, sistemul răspundea cu altă propoziție, iar utilizatorul continua dialogul, ținând cont de propozițiile sale anterioare și de propoziția dată de sistem. Este indicat ca utilizatorul să respecte (mai mult sau mai puțin) configurația structurală a tiparelor învățate. Două tipuri de rezultate au ieșit în evidență:

- "conforme"
- cu bifurcații

Voi da în continuare câte un exemplu concret pentru ambele situații.

Dialogul "conform" respectă "pattern"-ul tiparelor învățate, chiar dacă există mici abateri de la acestea. Dacă am început un dialog cu:

$$i_test1(0) : (v(N), w, \alpha(P))$$

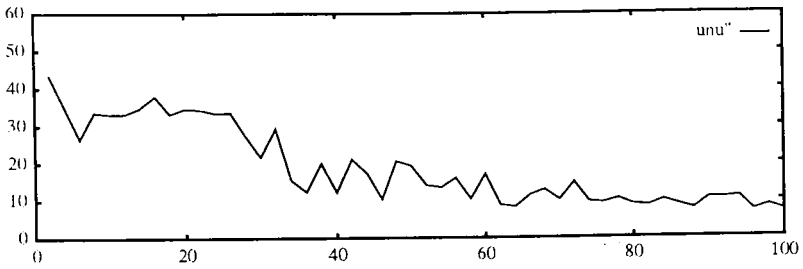
am respectat structura primei propoziții din tiparul 3, modificând însă individualitățile M și O cu N și P (respectând apartenența la clase). Răspunsul a fost:

$$o_test1(1) : (\xi(B), x, \beta(P))$$

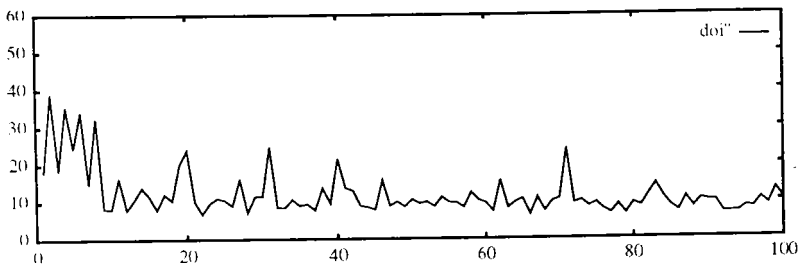
iar continuarea:

$$i_test1(1) : (\beta(B), y, \gamma(C)) \rightarrow o_test1(2) : (\beta(P), z, \delta(C))$$

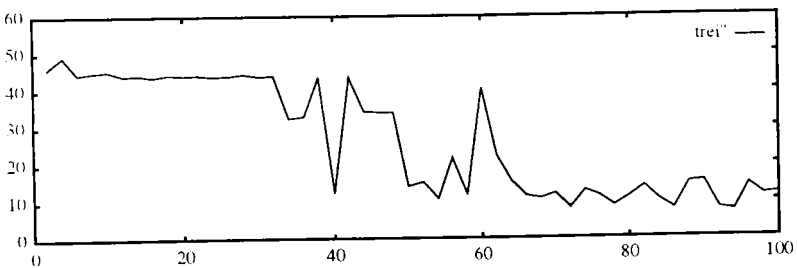
$$i_test1(2) : (\gamma(C), r, \varepsilon(D)) \rightarrow o_test1(3) : (\delta(C), s, \zeta(E))$$



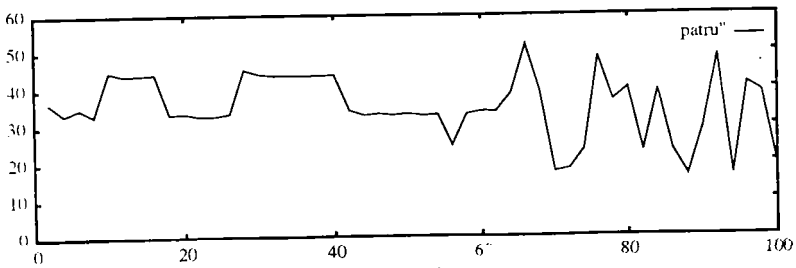
a. un singur tipar



b. doua tipare



c. trei tipare



d. patru tipare

figura 4-4.

$$i_test1(3) : (\varepsilon(D), t, \eta(F)) \rightarrow o_test1(4) : (\zeta(E), u, \eta(F))$$

Pentru a avea o traiectorie totuși ușor schimbată, am modificat tripletul (clasă, relație, clasă), dar doar în măsura în care exista altul cu aceeași relație în baza de antrenare (în cazul următor, din tiparul 2):

$$i_test1(4) : (\eta(F), w, \delta(I)) \rightarrow o_test1(5) : (\zeta(F), x, \theta(I))$$

Conform tiparului învățat, sistemul ar fi trebuit să răspundă cu tripletul $(\theta(H), x, \gamma(C))$. Interesant este că a răspuns tot cu relația x , care urma în secvența impusă (s, u, x, z) și cu individualitățile F și I care apăreau în momentul anterior și în intrare. Dacă am întrerupt regularitatea dată de succesiunea circulară a relațiilor din intrare (r, t, w, y) , sărind peste y :

$$i_test1(5) : (\sigma(I), r, \lambda(L)) \rightarrow o_test1(6) : (\theta(I), z, \chi(J))$$

sistemul rămâne la ieșire pe traiectoria știută din tiparul 3 și răspunzând cu un triplet care îl conține pe z (din tiparul 2) și cu individualitatea I , care apare și în intrare. La pasul următor se produce resincronizarea între relații:

$$i_test1(6) : (\lambda(L), t, v(N)) \rightarrow o_test1(7) : (\zeta(E), u, \mu(M))$$

Dacă am continuat să revin în "pattern"-ul tiparului 3:

$$i_test1(7) : (\varepsilon(D), w, \eta(F)) \rightarrow o_test1(8) : (\zeta(E), x, \theta(I))$$

$$i_test1(8) : (\eta(F), y, \sigma(I)) \rightarrow o_test1(9) : (\theta(I), z, \chi(J))$$

iar răspunsurile au fost cele așteptate. Pentru că se poate continua și pe traiectorii de dialog de mai mult de 10 pași, am introdus propoziții din tiparul 1:

$$i_test1(9) : (\sigma(I), r, \lambda(K)) \rightarrow o_test1(10) : (\chi(M), s, \mu(O))$$

$$i_test1(10) : (\lambda(K), t, v(N)) \rightarrow o_test1(11) : (\mu(M), u, \xi(O))$$

$$i_test1(11) : (v(N), w, \alpha(P)) \rightarrow o_test1(12) : (\xi(O), x, \beta(A))$$

$$i_test1(12) : (\alpha(P), y, \gamma(B)) \rightarrow o_test1(13) : (\beta(A), z, \delta(C))$$

obținând rezultatele așteptate. Sistemul avea tendința să se înscrie pe o traiectorie care urma traiectoria anterioară și era influențată de traiectoria dată de utilizator la intrare, dar numai în cazul în care acesta respecta cât de cât traiectoriile antrenate. Sistemul era robust, reușind în majoritatea cazurilor să revină la traiectorii învățate, chiar dacă apăreau pași care erau "perturbatori".

Dialogurile cu bifurcații apăreau atunci când utilizatorul depășea o anumită limită în ceea ce privește abaterile de la traiectoriile de învățare antrenate. Salturi neașteptate în spațiul de intrare, duceau la salturi și ciclări la ieșire. Există posibilitatea ca într-un univers al discursului cu semantică, un dialog de acest tip să aibă un posibil sens, dar relativ la structurile traiectoriilor învățate, aceste ieșiri nu pot fi asociate cu ele sub nici o formă. Unul din exemplele cele mai tipice este următorul:

$i_test2(0) : (\alpha(P), r, \beta(A)) \rightarrow o_test2(1) : (\delta(B), s, \zeta(F))$
 $i_test2(1) : (\delta(C), s, \zeta(E)) \rightarrow o_test2(2) : (\delta(C), s, \zeta(D))$
 $i_test2(2) : (\varepsilon(E), t, \eta(G)) \rightarrow o_test2(3) : (\gamma(C), s, \delta(D))$
 $i_test2(3) : (\zeta(F), u, \eta(H)) \rightarrow o_test2(4) : (\zeta(E), x, \theta(F))$
 $i_test2(4) : (\eta(G), w, \sigma(I)) \rightarrow o_test2(5) : (\alpha(A), w, \gamma(B))$
 $i_test2(5) : (\zeta(G), x, \theta(H)) \rightarrow o_test2(6) : (\alpha(A), w, \gamma(B))$
 $i_test2(6) : (\alpha(A), y, \gamma(B)) \rightarrow o_test2(7) : (\delta(E), s, \zeta(F))$
 $i_test2(7) : (\beta(B), z, \delta(C)) \rightarrow o_test2(8) : (\zeta(E), u, \mu(O))$
 $i_test2(8) : (\delta(D), w, \sigma(J)) \rightarrow o_test2(9) : (\zeta(E), u, \mu(O))$
 $i_test2(9) : (\lambda(L), t, v(N)) \rightarrow o_test2(10) : (\delta(E), u, \zeta(F))$
 $i_test2(10) : (\sigma(I), r, \lambda(L)) \rightarrow o_test2(11) : (\delta(E), u, \zeta(F))$

și așa mai departe, fără ca ieșirea să se fie asemănătoare cu cea învățată sau să se sincronizeze cu intrarea, chiar dacă reveneam la o traiectorie cunoscută (rețeaua are memorie, și odată trimisă pe o traiectorie perturbată și menținută mai mulți pași pe aceasta, nu mai revine la traiectoriile știute).

Dialogurile conforme arată că sistemul poate să reproducă în mod robust ceea ce a învățat. Chiar și dialogurile cu bifurcații sunt utile, pentru că ele sunt un indicator că sistemul nu acceptă intrări care nu au fost antrenate. Se poate argumenta că demonstrarea experimentală a puterii de generalizare a rețelei are nevoie de un număr mai mare de teste, recte un număr mai mare de baze de antrenare învățate (de o certă varietate). Faptul că antrenarea este atât de costisitoare a făcut imposibilă reluarea învățării pentru un nou set de tipare. În capitolele 5 și 6, se va arăta cum se pot realiza sisteme hibride cu un efort mai redus, modificând arhitectura modelului.

4.5. Fenomenul de bifurcație la antrenare - căi de rezolvare

Spre deosebire de rețelele feedforward, cele recurente prezintă atât la funcționare, cât și la antrenare, un comportament nedorit, dar caracteristic sistemelor dinamice neliniare, datorat fenomenului de bifurcație. În termeni laici, aceasta este o modificare bruscă a stării unui sistem dinamic, datorată modificării lente și continue a unui parametru (de regulă, un parametru de control, sau chiar o perturbație externă cu gradient redus). Subiectul bifurcațiilor la antrenare este tratat în [Pineda88], iar cel al bifurcațiilor la antrenare apare pentru prima oară în [Pearlmutter90]. Există o serie de metode propuse pentru evitarea lor la antrenare [Doya93], dar concluzia generală din literatura la zi, este că această problemă a comportamentului haotic la antrenare (un alt termen preferat pentru această problemă) este pe departe cea mai dificilă pentru cei care dezvoltă rețele recurente [Fahlman97].

4.5.1. Analiza fenomenului și corecțiile posibile

Pentru a explica ce reprezintă bifurcațiile, în continuare, se prezintă un exemplu. Dacă vom considera că funcționarea unui sistem dinamic este descrisă de ecuația:

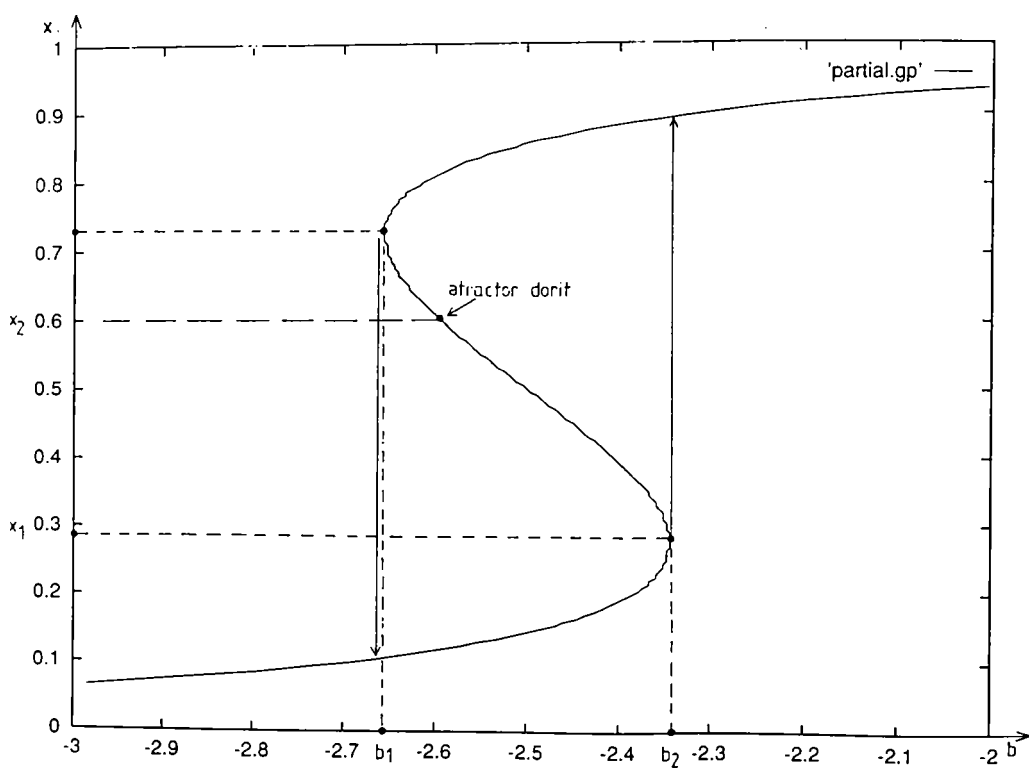


figura 4-5.

$$\sigma \cdot \frac{dx(t)}{dt} = -x(t) + \Gamma(Wx(t)) \quad (4-4)$$

unde:

$$\Gamma(x)_i = \frac{1}{1 + e^{-x_i}} \quad (4-5)$$

La limită se poate considera că sistemul poate fi și discret:

$$x(t + \Delta t) = \Gamma(Wx(t)) \quad (4-6)$$

Să presupunem că avem un sistem cu o singură unitate neuronală:

$$\sigma \cdot \frac{dx(t)}{dt} = -x(t) + \frac{1}{1 + e^{-(wx+b)}} \quad (4-7)$$

Leșirea x , în cazul în care b este o valoare constantă, va converge către un punct care satisface ecuația:

$$x = \frac{1}{1 + e^{-(wx+b)}} \quad (4-8)$$

Dacă vom interpreta această ecuație ca pe o funcție $f(b)=x$, pentru valoarea $w=5.0$, graficul funcției va arăta ca în figura 4-5 (pentru $b \in [-3, -2]$). Toate punctele de pe curba din figură reprezintă atractori pentru sistemul lăsat liber cu o perturbație b constantă ($db/dt=0$). Dacă în timpul relaxării sistemului vom modifica valoarea b , printr-o funcție de preferință liniară $f(t)=b$, ($f'(t)/dt^2=0$), crescătoare sau descrescătoare, vom obține următorul fenomen. Să presupunem că atractorul țintă (final) este în acest caz, $x_f=0.6$, iar valoarea inițială a perturbației este $b=0$. Descrscând linear valoarea perturbației, starea sistemului se va deplasa către valori mai mici, apropiate de starea dorită. Dar când vom atinge valoarea b_1 din figura 4-5, atractorul imediat următor care este pe ramura de jos a funcției și se produce o modificare abruptă a valorii de stare x . Dacă prin mărirea valorii de intrare b , dorim să ajungem la atractorul țintă, vom avea o situație similară când $b=b_2$ pentru că se va produce un salt pe ramura superioară. Se poate observa că în intervalul (x_1, x_2) este o zonă a spațiului stărilor care nu poate fi niciodată atinsă pornind din afara ei, în cazul în care $w=5.0$. Este de remarcat faptul că pentru un sistem discret, avem un comportament aproximativ identic.

Din experimentele efectuate cu secvențele de propoziții prezentate în subcapitolul anterior, am observat că fenomene care pot fi interpretate ca bifurcații apar atunci când intrarea este ca și "pattern" temporal față de intrările antrenate. Dacă la un moment dat se furnizează sistemului o propoziție la care acesta nu se aștepta, iar această propoziție este foarte "depărtată" ca și clase, individualități și mai ales ca relație (care trebuie să fie grupate în secvențele repetitive din baza de antrenare), atunci se va

produce cu siguranță un salt brusc în spațiul de ieșire. Acest salt poate să nu fie o bifurcație, dar comportamentul anormal (mai ales ciclic) care urmează, indică clar un astfel de fenomen. Pentru că antrenarea unui sistem neurodinamic (continuu sau discret) presupune evoluția altui sistem dinamic (vezi ecuația 2-55) cu $N \times N$ ecuații, către un atractor din spațiul ponderilor unde eroarea să fie minimă, fenomenul de bifurcație apare și în această situație. El se manifestă prin salturi bruște în spațiul definit de ponderile w_{ij} și mai ales prin creșteri bruște ale erorii de la o epocă la alta. Aceste salturi sunt deosebit de dezavantajoase din punct de vedere al unei antrenări rapide, pentru că ele în majoritatea covârșitoare a cazurilor, plasează matricea ponderilor într-o zonă a spațiului $R^{N \times N}$ care este foarte depărtată de atractorii doriți, iar revenirea presupune un număr exagerat de epoci suplimentare. De obicei, am observat că în cazul în care apărea o bifurcație, task-ul de antrenare respectiv trebuia abandonat. Au fost și situații în care saltul era atât de mare pe unele axe ale spațiului ponderilor, încât s-a ajuns ca ponderile să fie numere reale nereprezentabile prin magnitudinea lor. Și în acest caz, antrenarea era compromisă.

Există și posibilități de a detecta dacă s-a produs o bifurcație, de a corecta traiectoria compromisă și de a evita zonele în care s-au produs anterior bifurcații. Cea mai simplă metodă de a observa dacă s-a produs o bifurcație la antrenare este de a măsura variația erorii de la o epocă la alta. Dacă eroarea globală crește brusc cu 10%, înseamnă că aproape sigur s-a produs un salt nedorit al valorii ponderilor. Dar pot să apară și situații în care creșterea erorii să fie relativ mică (sub 5%), dar să se fi produs o bifurcație, care să modifice substanțial matricea ponderilor, fără ca acest lucru să afecteze eroarea. Un alt aspect relevant și în [Pearlmutter95], este că după ce s-a produs o bifurcație care să fie un salt relativ scurt (fără a afecta prea mult eroarea), zona în care am ajuns *după* bifurcație, "oferă" posibilitatea imediată a unei alte bifurcații, care însă de această dată să modifice catastrofal ponderile și într-un mod ireversibil.

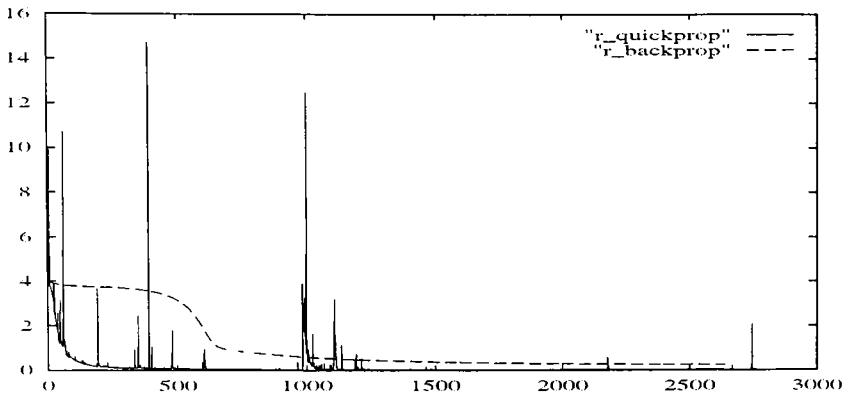
O metodă mai sigură de detecție este calculul distanței euclidiene între punctele din spațiul ponderilor, atinse după două epoci succesive. Dacă vom considera indicele k pe post de contor al epocilor:

$$d = \sqrt{\sum_{i=1}^N \sum_{j=1}^N (w_{ij}(k+1) - w_{ij}(k))^2} \quad (4-9)$$

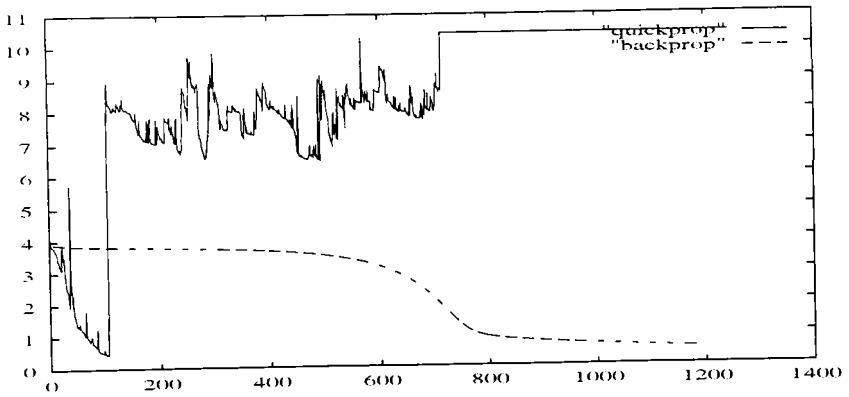
Dacă d este mai mare decât rata de învățare înmulțită cu media aritmetică a ponderilor în pasul k , atunci în mod sigur s-a produs o bifurcație. Însă reciproca acestei aserțiuni nu este adevărată. Saltul poate să se producă pe foarte puține axe (la limită chiar una), iar distanța d să rămână mică, pentru că celelalte ponderi nu se modifică decât neglijabil. În consecință, trebuie verificat pentru fiecare pondere:

$$(w_{ij}(k+1) - w_{ij}(k)) < \eta \cdot w_{ij}(k) \quad (4-10)$$

Cel mai acut dezavantaj al apariției bifurcației, este necesitatea de a relua antrenarea (de preferință din alt punct inițial setat aleator). Dacă punctul care a fost atins corespundea unei erori relativ mici, reluarea din start este deosebit de costisitoare.



a. fara bifurcatii



b. aparitia unei bifurcatii

figura 4-6

Unii autori, ca de exemplu, [Fahlman97] și [Pearlmutter95] recomandă reluarea antrenării dintr-un punct apropiat de locul unde s-a produs bifurcația, modificând rata de antrenare. Pentru a realiza acest stil de reluare, este necesară menținerea unei matrici cu ponderile anterioare, numită pivot (poate fi cu 1, 2, 3 sau mai mulți pași de antrenare înainte de momentul bifurcației). Repornind din punctul pivot, cu o rată de antrenare schimbată (de obicei micșorată) se observă [Doya93] [Szirbik95a] [Szirbik96] că mai devreme sau mai târziu bifurcația se produce oricum.

O metodă folosită pe larg în căutarea euristică a soluțiilor într-un spațiu de mari dimensiuni, este metoda de a genera aleator vecini apropiați ai punctului pivot. Acesta este memorat, iar traiectoria se reia dintr-o astfel de vecinătate. Dacă după un număr relativ mic de pași de antrenare, se produce o bifurcație, se revine în punctul pivot și se generează un nou vecin din care se reia antrenarea. Esențial este aici de stabilit ce înseamnă "număr relativ mic de pași". Dacă traiectoria s-a depărtat de punctul unde se putea produce o bifurcație și se produce o nouă bifurcație dar în alt punct, este inutil să revenim în pivotul stabilit înaintea primei bifurcații, fiind necesară menținerea unui pivot "glisant", care urmează traiectoria curentă cu un număr de pași în întârziere, iar în cazul în care se produce bifurcația, se consideră ca pivot, punctul în care este pivotul "glisant" în acel moment. Necesară este în acest context și memorarea punctului de bifurcație anterior, pentru a vedea dacă ne bifurcația următoare se produce tot în același punct. Bifurcații repetate în același punct, înseamnă necesitatea reluării traiectoriei dintr-un punct aflat la distanță de acest punct "critic".

O altă metodă de folosi saltul aleator este de a genera în același episod de bifurcație, mai multe puncte în jurul pivotului. Pentru toate aceste puncte, se calculează eroarea globală a sistemului cu matricea de ponderi asociată punctului respectiv și se alege ca punct de repornire a traiectoriei cel care are eroarea minimă. Pentru a evita în traiectoriile de antrenare punctele de bifurcații care au fost atinse deja, este recomandat ca ele să fie memorate (deși acest lucru este foarte costisitor, în cazul meu, memoram doar ultimul punct detectat). Dacă traiectoria punctului curent intră într-o hipersferă de rază r (stabilirea valorii se face prin încercări, în experimentele mele, valoarea preferată a fost $r=0.3$), atunci se va declanșa un proces de generare aleatoare a unor puncte de reluare a traiectoriei, care să fie în afara sferei.

O metodă de accelerare a convergenței către atractorii cu eroare globală mică a fost utilizarea regulii de antrenare quickpropagation [Fahlman89], care poate fi aplicată atât la rețelele feedforward cât și la cele recurente. Utilizarea acesteia a dus la o scădere a erorii mult mai rapidă (vezi figura 4-6a), în raport cu metoda clasică, reprezentată cu linie întreruptă), dar și la o creștere a frecvenței bifurcațiilor, care făcea antrenarea foarte dificilă (vezi figura 4-6b). O metodă de accelerare distribuită eficientă a quickpropagationului recurent, care pornea în paralel pe mai multe traiectorii și le abandona pe cele care ajungeau la o bifurcație, balansând continuu efortul de calcul între procesoare, este prezentată pe larg în [Somlo96].

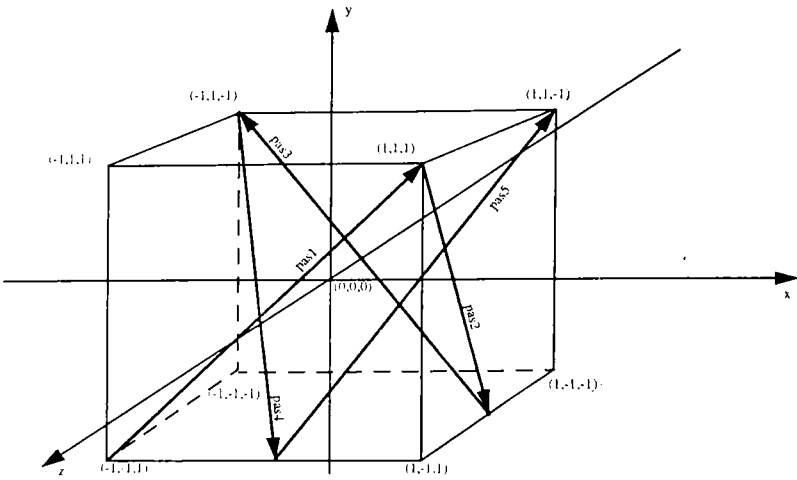
4.5.2. Influența formei tiparelor de antrenare

Studiind lucrările care tratau în special incidența fenomenului de bifurcație la antrenare [Tsong&Cottrell93] [Nagayama&Akamatsu94], am observat că în cazul experimentelor mele, fenomenul apărea mult mai des decât în alte genuri de experimente. Pentru că algoritmiile de antrenare erau aceiași, cauza putea fi numai baza de antrenare. Pentru a fi sigur de aceasta, am construit o nouă bază de antrenare, care trebuia să fie cât mai diferită de de prima. Oarecum întâmplător, am ales o metodă simplă și rapidă prin trasarea automată a unor curbe discrete într-un hiper-cub cu latura 2, centrat în origine. Pur și simplu, îmi alegeam aleator un punct de start și generam aleator alte 9 puncte succesive, ca vecinătăți ale punctului anterior generat. Pentru ca traiectoria generată aleator să aibă totuși o formă, am folosit regula empirică: "direcția traiectoriei să fie oarecum aceeași". Aceasta se poate obține, modificând substanțial coordonatele doar pe un număr mic de axe, pe celelalte modificarea fiind mică (dar nu și neglijabilă). Pentru calibrarea parametrilor generatorului, am folosit un instrument de lucru în 3D, care genera rezultate de tipul celor ilustrate în figura 4-7b. Caracteristic pentru traiectoriile generate este sunt pseudo-continue (o interpolare spline de exemplu, duce la o curbă netedă). Pentru că alegeam puncte inițiale cu valori mici, curbele erau aproape de origine, adică de centrul cubului. Pentru a evita aglomerarea tiparelor într-o singură zonă, am construit intenționat și curbe care erau apropiate de frontiera hiper-cubului.

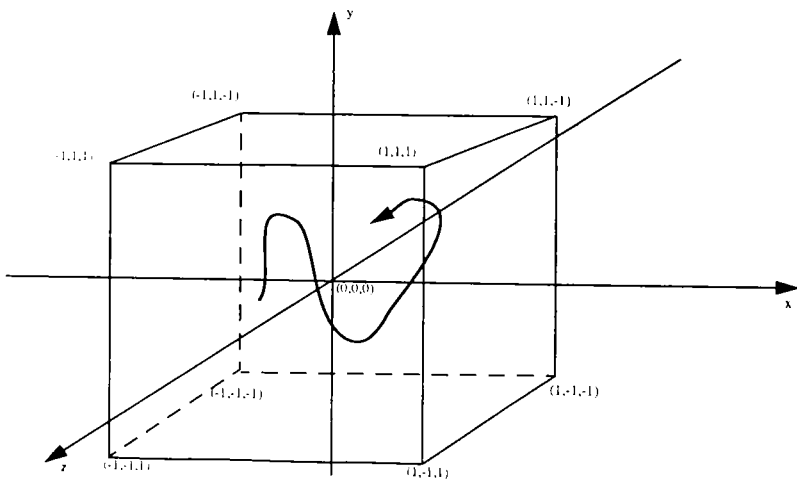
Am comparat rezultatele antrenării, relativ la numărul de bifurcații rezultate. Experimentele au fost făcute cu cu grupe diferite de tipare, pentru un număr fixat de epoci (450), fără ca eroarea să fie monitorizată, doar numărul de bifurcații, care era notat după 150 de epoci. Tiparele provenite din grafuri conceptuale sunt notate GC, iar cele generate aleator AL. Pentru că observasem din experimentele de până atunci că la o rată de antrenare mică, bifurcațiile apar mai des, în toate experimentele am folosit o rată de învățare constantă de valoare foarte mică (0.05). Rezultatele au fost următoarele:

TIPARE	număr de bifurcații la 150 de epoci	număr de bifurcații la 300 de epoci	număr de bifurcații la 450 de epoci
2 tipare GC	8	11	25
2 tipare AL apropiate de origine	0	0	0
2 tipare AL depărtate de origine	1	1	2
4 tipare GC (toate)	12	22	##
4 tipare AL	1	1	2
4 tipare AL apropiate de origine	0	0	1
4 tipare AL depărtate de origine	2	2	5

(## - matricea de ponderi a ajuns în mod repetat să conțină valori nereprezentabile)



a. o traiectorie CG



b. o traiectorie NC

figura 4-7.

Chiar și din acest unic set de date experimentale se poate deduce că rata de apariție a bifurcațiilor pentru tiparele GC este foarte mare, pe când tiparele AL au performanțe excelente relativ la stabilitatea antrenării.

Explicația este dată și în [Pearlmutter95], ca rezultat al unor experimente asemănătoare. Dacă tiparele de antrenare sunt curbe netede, și pseudo-continue (fără salturi semnificative între două puncte succesive) plasate într-o zonă centrală a hipercubului de stare al rețelei, antrenarea se desfășoară de regulă fără bifurcații. Matematic, explicația este dată de liniaritatea funcției de transfer sigmoidale în jurul originii. Astfel, sistemul dinamic reprezentat de mecanismul de antrenare poate fi considerat cvasi-liniar, iar un astfel de sistem are un comportament instabil doar în situații bine determinate.

În contrast cu tiparele netede și centrate, cele rezultate din grafuri conceptuale (pe baza translației tensoriale), sunt formate din secvențe de tensori formați în majoritate din valori binare (în cazul tangentei hiperbolice 1 și -1), care mai pot conține doar valorile 0.5, 0.75, -0.75 și -0.75, induse de formula (3-23). Pentru o schemă tensorială simplificată (cu rezultate de gradul 3), o ilustrare a unei astfel de traiectorii este dată în figura 4-7a. Tiparele de tip GC sunt succesiuni de puncte care formează traiectorii extrem de discontinue în spațiul stărilor rețelei, având chiar și salturi între "colțurile", hipercubului, ambele caracteristici (discontinuitatea și depărtarea de zona liniară din jurul originii) favorizând într-un mod exagerat bifurcațiile la antrenare.

Observația pe care am făcut-o după aceste experimente, a fost că dificultatea antrenării unor astfel de sisteme rezidă nu atât în duratele necesare pentru atingerea erorii impuse, cât în evitarea fenomenului de bifurcație, care mărește cu un ordin de mărime aceste durate.

4.5.3. O schemă extinsă de translație

O soluție pentru evitarea bifurcațiilor, care pare cea mai simplă și directă, este de a folosi un alt gen de tipare decât tiparele GC și anume tipare care prin interpolare spline reprezintă curbe netede și centrate în hipercubul stărilor rețelei (care sunt similare cu tiparele AL). Dificultatea este că aceste noi tipare trebuie să conțină informația din grafurile conceptuale inițiale. Am numit această clasă de tipare, NC (netede și centrate).

Forma tiparelor GC este datorată metodei de translație tensorială. Pentru a obține tipare NC, ar fi necesar să modificăm sau să completăm metoda de translație. Deși am căutat o metodă automată de transformare, nu am reușit să o dezvolt. Am folosit o metodă semi-manuală, în care alocam puncte în centrul hipercubului pentru propozițiile din baza de antrenare, astfel încât propozițiile cu semantică apropiată să fie apropiate, iar cele cu semantică diferită, să fie depărtate. Pentru a avea tipare care să reprezinte curbe netede, am alocat punctele astfel încât secvențele de propoziții să genereze astfel

de traiectorii. Dacă două secvențe conțineau propoziții comune, cele două curbe de antrenare se intersectau.

La început am folosit instrumentul 3D, pe care l-am folosit pentru generarea tiparelor AL. Am păstrat dimensiunea rețelei tot la 768 unități, dar am redus drastic dimensiunea intrării și a ieșirii, din motivul simplu ca erau dificil de generat puncte într-un spațiu cu dimensiune atât de mare. Experimental, am ajuns la concluzia că sistemul este cel mai simplu de antrenat, utilizat (iar tiparele NC de generat) atunci când dimensiunea intrării și ieșirii sunt de 80 adică aproximativ 10% din totalul unităților, valoare ideală dată și în [Pineda88]. O primă metodă ad-hoc pe care am folosit-o a fost divizarea componentelor vectorilor de natură binară cu 2.0, 2.5 și 3.0, pentru a aduce valorile acestora în zona centrală a hipercubului. Deși numărul de bifurcații a scăzut (la 15 bifurcații pe primele 300 de epoci), acest rezultat nu era cel scontat, dorind o reducere cu un ordin de mărime. Obținusem tipare relativ centrale dar nu și netede, iar dimensiunea nu a fost redusă de 10 ori, așa cum mi-am propus.

Pentru transformarea vectorilor obținuți prin translatore tensorială din grafuri conceptuale în vectori ai tiparelor NC dorite, este necesară o schemă de mapare, care să transforme vectorii de natură binară în vectori reali, cu componente cuprinse între -0.5 și 0.5 și de o dimensiune 10 ori mai mică. Din moment ce maparea se face propoziție cu propoziție, se poate constitui un tabel de mapare, unde pentru fiecare propoziție există o intrare conținând propoziția în forma sa simbolică și vectorul real ales pentru a face parte din tiparele NC. Nu ar mai fi necesară translatore tensorială, pentru că pur și simplu prin "pattern-matching", propoziția ar putea fi identificată în tabel și s-ar furniza rețelei vectorul real asociat. La ieșire, s-ar putea calcula distanța între punctul reprezentat de vectorul de ieșire în spațiul 80 dimensional și toate punctele descrise de vectorii din tabel, iar cel mai apropiat este cel asociat propoziției care va fi selectată ca și ieșire. La funcționarea sistemului, apar însă propoziții care nu fac parte din baza de antrenare și care nu aparțin tabloului static descris mai sus. Schema de mapare propusă este în consecință inutilizabilă.

Soluția găsită de mine a fost extinderea schemei de translație tensoriale prin două rețele neuronale suplimentare, menită să transforme vectorii eminent binari în vectori reali și invers. La intrare, după aplicarea translației tensoriale și obținerea vectorului cu componente binare, acesta din urmă este furnizat unei rețele neuronale feedforward cu un număr de intrări egal cu dimensiunea de ieșire a formulei tensoriale de translație și cu 80 de ieșiri. Ieșirile se normalizează pe intervalul (-1.0,1.0) de pe intervalul specific al tiparelor NC, (-0.5,0.5). Rețeaua feedforward este antrenată cu perechi (propoziție, vector-real) din tabloul static prezentat anterior. Pentru translația inversă, se folosește o altă rețea feedforward cu 80 de intrări și un număr de ieșiri egal cu intrările primei rețele feedforward. Această rețea se antrenează cu perechi (vector-real, propoziție) construite din același tablou static. Dacă o nouă propoziție, care nu apare în tabel, este furnizată sistemului în faza lui de funcționare, aceste două rețele suplimentare vor genera punctele potrivite în spațiul redus dimensional al tiparelor NC și în spațiul tiparelor GC.

Apare un avantaj legat de generarea tiparelor NC. Pentru că aceste două rețele sunt capabile de generalizare, nu este necesar să le antrenăm cu toate perechile de date din tabloul static. Corespunzător, nu este necesar să generăm toate punctele care formează tiparele NC. Am generat doar jumătate din acestea, restul fiind obținut prin generalizare, obținându-se chiar o netezime mai bună în acest caz (pentru curbele NC), decât în cazul în care erau generate semi-manual. Alegerea jumătății potrivite din totalul de propoziții am făcut-o astfel încât să acopere toate simbolurile din vocabularul folosit. Alt avantaj major este posibilitatea folosirii unei scheme de translație care să lucreze cu reprezentări cât mai puțin "înghesuite", pentru simbolurile din V_F , V_R și V_L (vezi capitolul 3). Aceasta va îmbunătăți considerabil performanța translației inverse, din vector de reprezentare cu componente binare în graf conceptual. Dimensiunea vectorilor din tiparele GC va crește foarte mult, dar nu va afecta dimensiunea rețelei recurente. Relativ la vocabularul folosit și la noile reprezentări binare, dimensiunea acestor vectori devine egală cu 12634. Schema sistemului îmbunătățit este dată în figura 4-8.

Structura rețelei feedforward FFI am stabilit-o la 12634-100-80 unități, cu o bază de antrenare de 20 de tipare. Resursele necesare pentru o astfel de rețea sunt relativ mari (aproximativ 14 MBytes pentru matricile de ponderi), antrenarea necesitând peste 8000 de epoci, adică aproximativ 59 de ore de rulare pe singur procesor SPARC - algoritmul fiind neparalelizat - folosind un program de antrenare clasic "quickpropagation", o variantă a algoritmului backprop (importat de pe "site"-ul Neurosoft Repository de la Carnegie-Mellon University, Pittsburgh). Rezultatul a fost foarte bun, eroarea putând fi redusă sub 3% din prima încercare. Stabilirea dimensiunii stratului ascuns s-a făcut anterior, la antrenarea rețelei FFO (vezi figura 4-8), care a fost dezvoltată înaintea rețelei FFI.

Rețeaua FFO are o particularitate care mi-a atras atenția și din acest motiv m-am gândit la utilizarea unui alt tip de rețea. La antrenarea rețelelor feedforward, pierderea de timp se datorează în principal găsirii unui număr potrivit de neuroni pentru stratul ascuns. De obicei se începe cu un număr mic de unități, se testează dacă eroarea impusă poate fi atinsă într-un număr rezonabil de epoci, iar dacă nu se poate, se mărește numărul de unități de pe stratul ascuns cu o fracțiune și se pornește procesul de antrenare din start. Aceste cicluri de reglare a dimensiunii rețelei pot fi deosebit de costisitoare ca și timp de calcul irosit. Ele sunt necesare pentru a stabili un necesar de memorie optim pentru rețea. Relativ la dimensiunea stratului ascuns, necesarul de memorie pentru ponderi poate să difere foarte mult. De exemplu, regula empirică folosită de obicei pentru stabilirea primei valori alese este:

$$\lfloor 1/2 * (nr_neuroni_intrare + nr_neuroni_ieșire) + nr_tipare \rfloor / 2$$

Pentru rețeaua mea specifică (FFO) dimensiunea stratului ascuns dată de formula de mai sus, ar fi fost de 3189 de unități. Ceea ce înseamnă o primă matrice de ponderi de 12634*3189 de valori în virgulă flotantă, adică un necesar de aproximativ 330 de MBytes, iar pentru a doua matrice 3189*80 de ponderi (aprox. 2 MBytes). Comparete cu valorile pe care le-am utilizat (aprox. 14 MBytes pentru ambele matrici), ac

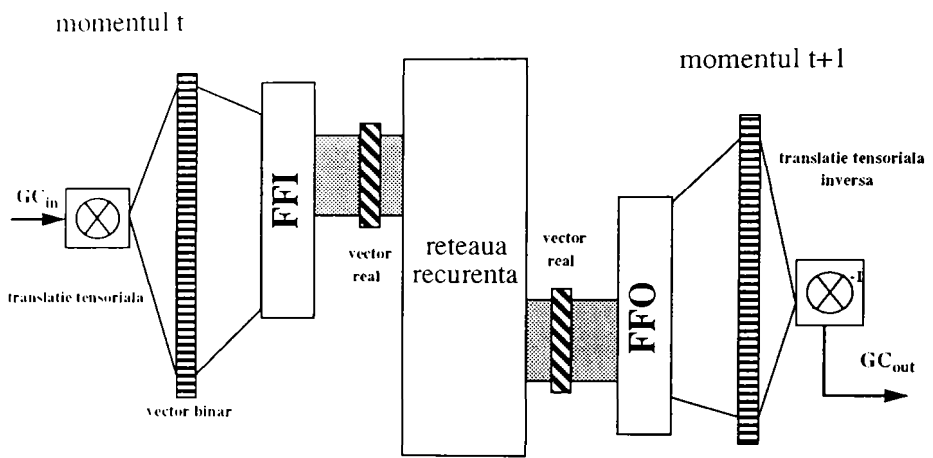


figura 4-8.

este valori sunt mult prea mari. Dimensiunea optimă (aproximativ 100 de unități pe stratul ascuns) a fost detectată cu o variantă a algoritmului de antrenare "cascade-correlation" [Fahlman&Lebiere91], care își ajustează din mers numărul de neuroni de pe stratul ascuns. Dezavantajul este că o arhitectură în cascadă funcționează mult mai încet (ca și când ar fi o rețea cu un număr de straturi egal cu numărul de neuroni ascunși), dar antrenarea este relativ rapidă (m-a costat aproximativ 80 de ore pe o mașină SPARC). Ideea de a folosi "cascade-correlation" mi-a venit datorită faptului că acest algoritm se comportă excelent în situația când intrările sunt reale, iar ieșirile binare, vezi [Fahlman&Lebiere91].

4.5.4. Concluzii

Repetând antrenarea rețelei recurente cu tiparele NC (pe mașina ZOO), timpul de antrenare s-a redus cu aproximativ un ordin de mărime. Acest fapt s-a datorat reducerii la jumătate a numărului de epoci necesar pentru atingerea erorii impuse (5%), dar mai ales reducerii drastice a numărului de bifurcații la antrenare. În 735 de epoci s-au produs doar de trei ori bifurcații, corectarea traiectoriei fiind eficientă, fără a se mai obține valori exagerat de mari în matricea de ponderi și fără a fi necesară reluarea din start a procesului de antrenare. Verificând capacitatea de generalizare a rețelei, cu cele două tipare de test prezentate anterior am obținut aproximativ același gen de rezultate. Ceea ce validează această structură îmbunătățită a sistemului.

Bineînțeles, există un preț pentru această evitare forțată a fenomenului de bifurcație la antrenare. Este necesară introducerea a două noi rețele, care ocupă un spațiu important de memorie și care cer și ele un timp destul de lung de antrenare, chiar mai lung decât al rețelei recurente - datorită variantei neparalelizate. Însă în termeni absoluți acest timp este mult mai mic decât timpul necesar pentru antrenarea rețelei recurente (aproximativ cu două ordine de mărime). Posibilitatea de a folosi un spațiu cu foarte multe dimensiuni pentru tiparele GC duce la tentația de a folosi un vocabular mult mai bogat, dar aceasta ar conduce la o reprezentare "înghesuită" a simbolurilor din vocabular și din o serie de experimente, am observat că antrenarea rețelelor FFI și FFO (mai ales), devine extrem de costisitoare.

Putem rezuma metodologia propusă la următorii pași:

- alegerea vocabularului (clase, individualități și relații)
- alegerea a două reprezentări pentru V_F , V_R și V_L (una "largă" și alta "înghesuită").
- stabilirea propozițiilor folosite și a tiparelor de tip secvențe de propoziții.
- determinarea dimensiunii vectorilor GC obținuți prin translație tensorială pentru ambele reprezentări, notate D pentru reprezentarea "largă" (dimensiune 3125 poziții) și d pentru cea "înghesuită" (dimensiune 768).
- stabilirea dimensiunii rețelei recurente egală ca număr de unități cu d (intrarea și ieșirea sunt $n=d/10$).
- jumătate din propoziții (în mod obligatoriu trebuie să acopere tot vocabularul folosit) sunt asociate cu puncte din spațiul n dimensional, astfel încât tiparele să devină în acest

spațiu netede și centrate în zona originii.

- stabilirea dimensiunii rețelelor FFI (D intrări și n ieșiri) și FFO (n intrări și D ieșiri).
- stabilirea prin "cascade-correlation" a numărului de unități ascunse din FFO (notat h), prin antrenarea până la atingerea unui nivel propus de eroare (recomand 3%).
- antrenarea rețelei FFI cu "quickpropagation", cu aceleași perechi ca și pentru FFO, pe o structură fixă cu $D-h-n$ unități până la atingerea nivelului propus de eroare (tot 3%).
- antrenarea rețelei recurente, până la un nivel de eroare de 5%.

Cele mai costisitoare activități sunt antrenarea rețelei recurente, care cere un enorm efort computațional (96% din total), și asocierea propozițiilor la vectori în spațiul n dimensional, care cere un efort enorm din partea celui care dezvoltă sistemul (90% din total), în situația când această operațiune se face semi-manual. Automatizarea asocierii structurilor simbolice la un spațiu euclidian, astfel încât să se păstreze cerințele tiparelor NC este o problemă care ar putea reprezenta o direcție de cercetare separată. Există rezultate în acest sens (vezi de exemplu [Crucianu&Memmi92]), dar un excelent tutorial recent publicat pe această temă [Pfeifer96] nu indică nici o metodă pe care aș fi putut-o folosi în contextul specific al modelului dezvoltat de mine.

5. Modelul multimodular omogen

Sistemul prezentat în capitolele precedente, se bazează pe o singură rețea recurentă de mari dimensiuni, un vocabular de o dimensiune relativ scăzută și o schemă de translatăre bidirecțională. Se folosesc două tipuri distincte de date pentru reprezentarea cunoștințelor: vectorii de activare neuronali și grafurile conceptuale. Este un sistem omogen, pentru că procesarea propriu-zisă - raționamentul - se desfășoară prin calculele specific conexioniste. Natura parțial hibridă este dată de restaurarea de tip "translație-retranslație", care este aplicabil doar modelul în buclă închisă și de verificările posibile de consistență logică (dar acestea nu sunt neaparat necesare).

Complexitatea antrenării acestei rețele este cu aproximație N^4P , iar în cazul în care lungimea tiparelor temporale discrete este aproximativ egală cu N , complexitatea este de gradul 5 înmulțită cu P - numărul total de tipare de antrenare. Timpul necesar pentru antrenarea celor două rețele feedforward este nesemnificativ în raport cu timpul necesar rețelei recurente. Pentru că nu există soluții prea bune pentru reducerea complexității algoritmilor de antrenare recurenți [Pearlmutter 95], singura cale de a îmbunătăți sistemul este să reducem N , numărul total de neuroni ai rețelei recurente. Pare o soluție paradoxală, pentru că scade capacitatea de antrenare, odată cu dimensiunea vocabularului. Dar există și alternativa *distribuirii*.

Idea mea a fost să folosesc mai multe rețele recurente mici [Szirbik 95b], care au fiecare, un vocabular propriu, foarte redus. Dar dacă fiecare rețea mică lucrează cu câteva cuvinte pe care numai ea le poate procesa, există și cuvinte care pot fi acceptate doar de câteva rețele și există un fond comun de simboluri utilizabile de toate rețelele, atunci numărul total de simboluri poate fi foarte mare. Fiecare rețea recurentă este însoțită de cele două rețele feedforward (FFI și FFO) necesare pentru completarea translației simbolic-neuronale și neuronal-simbolice, astfel încât tiparele de antrenare să nu faciliteze apariția bifurcațiilor la antrenare. În prezentarea care urmează, am făcut abstracție de aceste două rețele, fără să intru în amănunte legate de implementarea și performanțele lor. Aceasta pentru că antrenarea lor este banală, singura problemă fiind

asocierea unor puncte din spațiul N dimensional la propozițiile folosite. Dimensiunea rețelei recurente care efectuează procesarea cognitiv-neuronală este întotdeauna în funcție de vocabularul asociat ei, mai exact rezultă din formula de translație tensorială care folosește cea mai compactă reprezentare pentru simboluri (vezi capitolul 3). Pentru că toate rețelele folosite sunt identice ca structură, am numit această arhitectură, *modelul multimodular omogen*. Este important de menționat că acesta este un contra-exemplu relativ la metodologia propusă în teză, iar dezavantajele sale sunt analizate în profunzime, pentru a constitui argumentarea necesară pentru modelul prezentat în capitolul următor.

5.1. Considerații dimensionale

Cel mai evident avantaj al unei arhitecturi cu mai multe rețele recurente este posibilitatea utilizării unui vocabular de mari dimensiuni. Vom numi în continuare rețelele folosite (împreună cu partea aferentă de translație) *module conexiuniste - MC*. Putem considera că fiecare modul poate avea la intrare sau genera la ieșire grafuri conceptuale care sunt compuse dintr-un vocabular limitat la maximum 20 de cuvinte. Am observat că în general, pentru un set relativ mare de propoziții (314), distribuția morfologică optimă, pentru a putea construi propoziții cu sens normal (ne-metaforic sau prea factual), este următoarea: 30% verbe, 35% substantive, 35% adjective. Avem în cazul vocabularului de 20 de cuvinte: 6 verbe, 7 substantive și 7 adjective. Aplicând formula tensorială prezentată în cap.3.1. vom obține o rețea cu aproximativ 200 de unități (40000 ponderi = 320kBytes), mult inferioară ca dimensiuni decât una cu 1000 de unități (1 milion de ponderi = 8MBytes) care poate utiliza un vocabular de maximum 50 de cuvinte. În același spațiu de 8MBytes pot fi memorate ponderile a aproximativ 25 de module "mici". Să presupunem următoarea situație: Avem 3 MC, iar distribuția vocabularului este următoarea:

- 5 cuvinte "locale" fiecărui modul
- 10 cuvinte "universale", care pot fi utilizate de toate modulele
- 5 cuvinte "comune" două câte două (vezi figura 5-1)

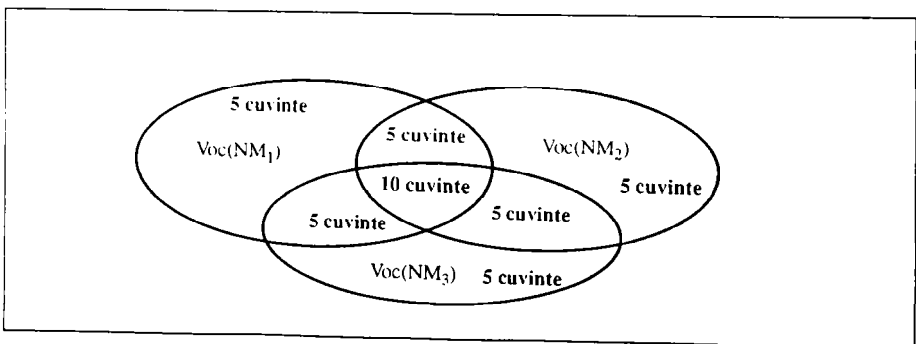


figura 5-1

Se poate observa că avem un total de 40 de cuvinte. Dacă numărul de module este crescut la 25, și folosim același stil de distribuție, vom avea un total de:

$$25 \cdot 5 + 12 \cdot 5 + 10 = 200 \text{ de cuvinte}$$

Folosind 400 de module (necesar de memorie = 128 MBytes) și considerând că cuvintele comune (nu cele "universale") sunt distribuite pe grupe mai mari (aproximativ de 50 de module), atunci vom avea un total de:

$$400 \cdot 5 + 50 \cdot 5 + 10 = 2260 \text{ de cuvinte}$$

care este un vocabular foarte mare. Însă vom vedea în continuare, din analiza modului de funcționare că nu este avantajos să avem un număr atât de mare de cuvinte locale (individuale unui singur modul) - acestea reprezentând majoritatea din cele 2260 de cuvinte.

5.2. Comunicarea între module

Comunicarea se desfășoară prin intermediul unei zone comune, numită "blackboard" (BB). Toate modulele au acces la această zonă putând citi și scrie propoziții din/pe BB. Din punct de vedere temporal, se pot alege două variante constructive:

- varianta sincronă
- varianta asincronă.

Din motive de simplitate am ales modelul sincron, cel asincron neoferind nici un avantaj din moment ce informația era transmisă în pachete delimitate de conținutul propoziției. Schema de funcționare a unui astfel de set de module este dată în figura 5-2. Intuitiv, funcționarea este următoarea: se scriu pe BB informațiile (propozițiile) de start, apoi modulele care pot citi aceste propoziții (adică toate cuvintele din respectiva propoziție fac parte din vocabularul modului respectiv) le acceptă la intrare și evoluează din starea inițială într-o altă stare staționară. Se poate întâmpla ca starea finală să fie identică cu starea inițială. În acest caz se consideră că modulul nu a generat nici o ieșire (vom numi acest modul, *staționar*). Propozițiile care au fost citite de pe BB de către module nestaționare, se consideră "consumate" și sunt șterse.

Modulele care sunt nestaționare și au consumat câte o propoziție din BB vor genera o propoziție de ieșire. Setul de propoziții astfel obținut va fi scris pe BB și va reprezenta pachetul de informație necesar unui nou ciclu ca cel descris, pe care îl vom numi pas de raționament. În momentul citirii propozițiilor de pe blackboard, apar o serie de probleme legate de prioritate (ordinea în care sunt citite propozițiile).

i. o propoziție poate conține un set de cuvinte care poate fi acceptat de mai multe module conexionate. În acest caz, am decis ca toate modulele care pot accepta propoziția să o și citească.

ii. pot apărea la un moment dat pe BB mai multe propoziții care pot fi acceptate de un singur modul. În acest caz vom avea o coadă de așteptare, pentru că un modul nu poate

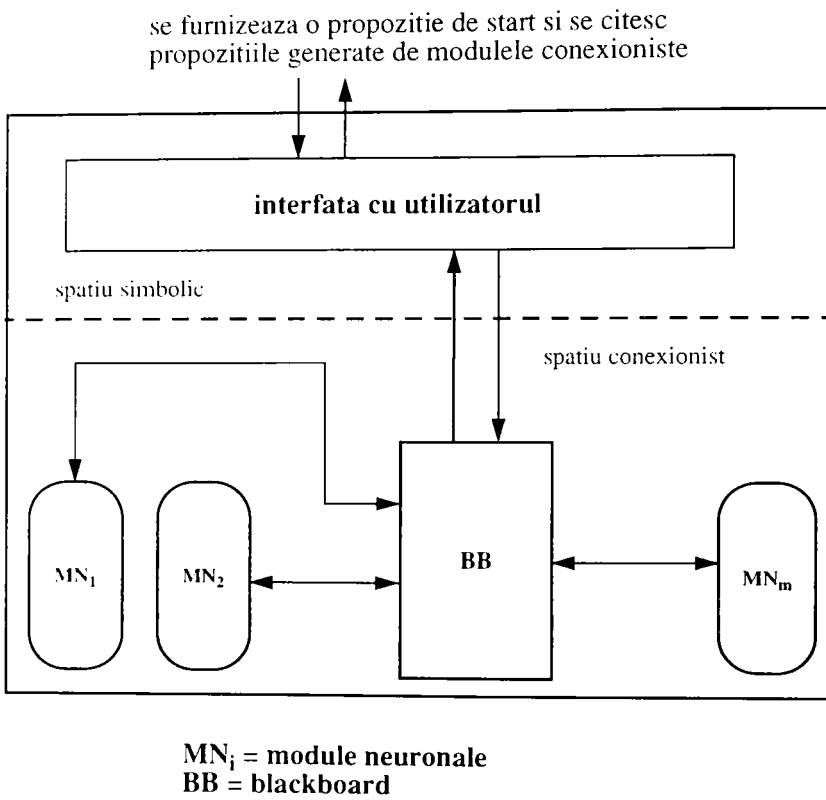


figura 5-2.

accepta decât o singură propoziție într-un pas de raționament. Se poate imagina și o schemă în care fiecare modul își consumă într-un singur pas toate propozițiile din coada de așteptare, dar am observat că în acest caz, modulele tind să se distanțeze la nivel semantic, propozițiile de pe BB fiind foarte diferite și neducând la un fir de raționament unitar, ci la generarea unor propoziții foarte diferite între ele. Ca raționamentul să poată evolua, e necesar ca toate modulele să fie grupate la un moment dat în aceeași zonă semantică.

Dacă folosim schema cu o_propoziție/un_pas_de_raționament, se pot imagina mai multe scheme de rezolvare a priorității propozițiilor. O soluție "brută" este aplicarea unei operații de "join" grafurilor conceptuale din coada de așteptare. Există însă cazuri în care această operație să nu ducă la un singur graf final, pentru că între anumite grafuri din coada de așteptare, operația de "join" nu este posibilă. Evident, această metodă este posibilă numai dacă reprezentarea propozițiilor se face pe BB în formă de grafuri conceptuale. Când operația de "join" ducea la un singur graf final, se producea o "aglomerare" a informației (majoritatea elementelor vectorului GC erau 1) la intrarea rețelei neuronale FFI, și de obicei rețeaua recurentă evolua într-o stare stabilă de unde nu mai putea fi scoasă de nici o altă intrare (ca un fel de saturare), adică devenea tot timpul staționară.

Altă soluție ar fi superpoziția vectorilor conexioniști care reprezintă propoziții. În acest caz, nu mai e necesară reprezentarea la nivel de graf a propozițiilor pe BB. Superpoziția se face prin însumarea vectorilor de intrare. Dar comparând rezultatul obținut prin translatarea după o operație de "join", cu vectorul obținut prin superpoziție (plecând de la aceleași propoziții inițiale), se poate observa o mare diferență. Acest lucru este normal, pentru că metoda superpoziției directe este o simplificare a metodei de translatare prezentată în capitolul 3. Efectuând experimentele cu ambele variante, am ajuns la concluzia că nici una din ele nu este potrivită. Prima este aplicabilă doar când propozițiile sunt în formă de graf conceptual și operația de "join" este posibilă între toate propozițiile din șirul de așteptare al unui modul (rezultă un singur graf). A doua metodă simplifică în mod grosolan procesul de translatare al unui text într-un vector, pierzând foarte multă informație.

Alternativa era ordonarea propozițiilor care pot fi acceptate de un modul (sau mai multe module simultan) în coada de așteptare, ele fiind consumate una câte una. A fost deci necesară găsirea unor criterii de ordonare. Primul criteriu de ordonare este cât se poate de evident, ordinea în coada de așteptare fiind dată de ordinea venirii (FIFO sau LIFO). Am ales FIFO (first-in, first-out). Dar în acest caz, se poate întâmpla ca mai multe propoziții care aparțin aceleiași cozi de așteptare să apară în același pas de raționament (ca ieșire simultană a mai multor module). În acest caz, am adoptat o stabilire conform unei distribuții aleatoare uniforme a ordinii de intrare în coada de așteptare, pentru că nu am găsit nici un criteriu viabil pentru stabilirea ordinii.

Mai apare un fenomen de care a trebuit să țin cont. De multe ori, o propoziție generată într-un pas de raționament este deja pe BB (undeva într-o coadă de așteptare). Cea mai simplă metodă de evitare a includerii multiple a unei propoziții din coada de

așteptare era verificarea dacă propoziția există deja și în acest caz se ignora noua apariție. Am ajuns la concluzia că o astfel de abordare pierde din informație, pentru că apariția succesivă a aceleiași propoziții (informații) are și ea semnificația ei. Am decis ca în astfel de cazuri, propoziția să fie "avansată" în coada de așteptare, până pe primul loc. Și în această situație poate să apară o situație de conflict. De exemplu, se poate întâmpla ca într-o coadă de așteptare să avem trei propoziții pe poziții diferite și acestea apar din nou într-un pas de raționament:

$$\dots P_j, \dots P_j, \dots P_k, \dots, P_{first} \quad (\rightarrow \text{sensul cozii})$$

În algoritmul de căutare al cozii potrivite pentru o propoziție nou generată și constatarea existenței anterioare a propoziției respective, aceste operații se vor face secvențial (de regulă în ordinea de numerotarea a modulelor care au generat propoziții). Să presupunem că ordinea în care tratăm propozițiile este P_k, P_j, P_i . În acest caz, aplicând metoda descrisă mai sus, coada va avea următorul aspect:

$$\dots P_{first}, P_k, P_j, P_i$$

pentru că tratând propoziția P_k prima dată, am pus-o pe prima poziție, care apoi a devenit a doua și apoi a treia, ajungând la o ordine total diferită de cea dată de pașii de raționament. Acest lucru se datorează numai și numai faptului că ordinea de tratarea fost cea descrisă. În alta ordine s-ar fi obținut o prioritate dată de această nouă ordine. Ca să păstrez totuși ordinea inițială, am grupat propozițiile deja existente într-o nouă coadă provizorie, care păstra ordinea inițială:

$$P_i, P_j, P_k$$

indiferent de ordinea tratării (în același pas) acestor propoziții și am concatenat această coadă provizorie cu coada inițială, din care s-au șters propozițiile existente din coada provizorie:

$$\dots P_{first}, P_i, P_j, P_k \quad (\rightarrow)$$

Se poate observa că rezolvând toate aceste probleme, comunicarea între module și suportul ei au devenit o parte importantă a sistemului.

În ceea ce privește nivelul de reprezentare, putem alege între două strategii fundamentale diferite pentru a realiza comunicarea între module, și anume:

- comunicarea la nivel conexiunist
- comunicarea la nivel simbolic (pe care am folosit-o).

În primul caz, informație care circulă între module este reprezentată prin vectori de activare neuronali. Apar două probleme: propagarea erorilor (conform regulii divergenței - vezi subcapitolul 4.2) dar mai ales devine foarte dificilă *identificarea modulelor potrivite care pot accepta o anumită propoziție*. La nivel simbolic, am reprezentat propozițiile ca și grafuri conceptuale (sau într-o variantă mai simplă, ca și

structuri predicative).

În cazul în care restaurăm ieșirea de tip vectorial a unui modul în forma de graf conceptual (și chiar facem și o verificare aplicând principiul rezoluției dacă propoziția nou creată nu duce la o inconsistență), eliminăm micile erori care pot apare în structura unui vector de ieșire. Bineînțeles, această structură simbolică trebuie transformată din nou în vector pentru a fi aplicată ca intrare la o altă rețea neuronală. Sistemul acesta este mult mai complicat, dar elimină posibilitatea apariției fenomenului de acumulare a erorilor și a eșuării complete a procesului de raționament. Dezavantajul este mai ales din cauza funcționării mai lente, și nu atât din cauza complicării implementării. Practic, fiecare modul va apela la *aceiași* schema de translație, directă și inversă, dar cu vocabularul aferent.

Problema cea mai importantă care apare în cazul comunicării la nivel vectorial este identificarea modulelor potrivite care pot accepta la intrare vectorii respectivi. Schema de funcționare pe care mi-am propus-o se bazează pe ideea că propoziția nou creată va fi acceptată de unul sau mai multe module care au un vocabular care include toate cuvintele din propoziție. Acest lucru nu este cunoscut numai dacă se descifrează vectorul care reprezintă propoziția (și timpul necesar pentru această operație este cu un ordin de mărime mai mare decât timpul necesar pentru operația inversă). Aceasta a fost prima metodă pe care am aplicat-o, furnizând vectorul inițial modulelor care puteau să accepte propoziția, fără să mai translatez grafurile conceptuale refăcut într-un nou vector din care am eliminat practic micile erori introduse de procesarea conexiunilor. Numai că aplicând această idee, am observat acumularea unor erori, până la momentul în care translația din vectorii de ieșire nu se mai putea produce și sistemul se bloca. Abia atunci am definit principiul divergenței.

Ca și concluzie putem spune că avem o colecție de module, care pot fi privite ca niște subsisteme de sine stătătoare, de același tip ca și întregul sistem prezentat în capitolul 4. Fiecare modul este format dintr-un translator din graf în vector și un translator invers, din vector în graf, și o rețea neuronală recurentă, care primește ca intrare vectori și generează tot vectori (prin perturbare dintr-un punct stabil și relaxare în alt punct stabil). Comunicarea se realizează prin intermediul unei structuri de date de tip blackboard (BB) care poate conține numai grafuri conceptuale, care sunt menținute în mai multe șiruri de așteptare aferente fiecărui modul. Trebuie să mai existe și un subsistem de administrare a BB, care să aleagă șirul (sau șirurile) potrivit pentru fiecare propoziție (graf generat) și să stabilească ordinea în șirurile de așteptare.

Înainte de a face primele experimente, mi-am dat seama că ideea de a folosi cuvinte individuale pentru un modul era restrictivă. Este destul de simplu de observat (deja din stadiul experimentelor facute cu creionul pe hartie, ca o pre-simulare a sistemului) că aceste cuvinte individuale pot crea probleme în funcționarea sistemului. Ideea principală a sistemului propus mai sus este ca raționamentul să fie efectuat mai ales prin comunicarea mai multor module care pot efectua pași limitați de raționament. Existența unor cuvinte individuale, nu va permite transmiterea informației la alte

module, din simplul motiv că acestea nu pot accepta propoziții cu cuvinte pe care nu le cunosc. În mod practic, ceea ce se întâmplă este că aceleași module care încep procesul de raționament îl și continuă pentru că numai ele cunosc vocabularul cumulat al propozițiilor care sunt generate. Restul modulelor ar putea să fie folosite în raționament, doar dacă se întâmplă ca o propoziție generată să fie compusă din cuvinte care nu sunt individuale aceluși modul, dar sunt cunoscute și de alte module, care vor putea accepta respectiva propoziție.

Această problemă rămâne și pentru cazul că nu mai folosim cuvinte individuale, dar gradul de acoperire al vocabularelor modulelor este mic. Cu cât vom crește gradul de acoperire, cu atât va crește și rata de implicare a modulelor în raționament. La limită, în cazul (nedorit) în care toate modulele au același vocabular (adică ele cunosc toate cuvintele din sistem), toate modulele vor fi implicate în fiecare pas de raționament, fapt care nu este productiv și mărește dimensiunea rețelelor constituente ale modulelor. La cealaltă extremă avem situația în care un singur modul efectuează toți pașii de raționament, restul rămânând nefolosii. Se poate observa că problema alegerii vocabularului pentru fiecare modul și a gradului de acoperire al vocabularelor între module este foarte importantă în funcționarea sistemului.

5.3. Alegerea vocabularelor modulelor, experimente

5.3.1. Scopul unui astfel de sistem

Dacă scopul sistemului monomodular era demonstrarea generalizării unor texte și a capacității rețelelor recurente de a fi folosite în astfel de aplicații, scopul pe care mi l-am propus pentru utilizarea sistemului multimodular a fost realizarea unor raționamente automate. Practic, un astfel de sistem expert are același rol ca și un sistem expert implementat prin metode pur simbolice, și din punct de vedere al comportamentului văzut din exterior, sistemul este același cu un sistem pur simbolic. Nu am dorit să realizez un sistem care să funcționeze mai bine sau să aibă capacități pe care un sistem simbolic nu le are. Dorișta mea a fost să realizez un sistem care să funcționeze la fel ca și un sistem expert, să demonstrez că se pot implementa sisteme bazate pe rețele neuronale recurente și să explorez problemele care apar la dezvoltarea și implementarea sistemului și eventual să identific avantaje pe care sistemele simbolice nu le au.

Este evident că un astfel de sistem este extrem de costisitor în comparație cu sistemul simbolic care face același lucru. Necesarul de memorie este cu două ordine de mărime mai mare (la același volum de cunoștințe), iar viteza de lucru este cu un ordin de mărime mai mic. Deși există posibilitatea de a realiza aceleași funcții (raționamente) și într-un caz și în celălalt, pentru sistemul multimodular conexiionst, avem nevoie de resurse de calcul care nu sunt justificate. În mod normal, un sistem de raționament automat care are aceste capacități, va fi întotdeauna implementat prin metodele "clasice". Singurul avantaj al unei implementări conexiioniste este de a arăta că o versiune bazată pe rețele recurente care comunică la nivel simbolic este

funcțională.

Aceste dezavantaje care țin de intensitatea computațională a sistemului conexionist, pot fi evitate prin utilizarea unui hardware adecvat (pentru că implementarea se bazează pe ideea modularității extensibile) și anume a unor procesoare neuronale [Ienne97] care pot implementa rețele recurente și pot comunica prin intermediul unui procesor obișnuit, care ar realiza secvențierea operațiilor de raționament. Dar marele dezavantaj rezidă în dezvoltarea sistemului. Dacă în cazul sistemului pur simbolic, extensia cunoștințelor s-ar face doar prin simpla adăugare la baza de cunoștințe a unor noi propoziții, în cazul sistemului conexionist multimodular, adăugarea unui nou set de propoziții implică adăugarea unor noi module în sistem. Aceasta implică antrenarea rețelelor recurente și acest lucru, făcut de fiecare dată, este extrem de costisitor. Repet însă, demonstrarea, chiar și experimentală, a funcționalității unui sistem expert neurodinamic este un succes în sine, pentru că nimeni nu realizat așa ceva până la nivelul anului 1997.

5.3.2. Stabilirea vocabularului

Am construit vocabularul plecând de la un sistem de raționament automat bazat pe reguli, funcționând la nivel predicativ, și care folosea *modus ponens* pentru găsirea unor propoziții adevărate, plecând de la un set de propoziții de start. Propozițiile puteau fi exprimate fie ca și clauze Horn, fie ca și grafuri conceptuale foarte simple. Toate cuvintele care apăreau în baza de cunoștințe a acestui sistem expert, au format vocabularul viitorului sistem conexionist. Cuvintele din vocabular au fost împărțite ca și în cazul sistemului monomodular în trei categorii: relații (predicate), individualități (constante) și clase (variabile). Acest lucru este necesar pentru a aplica formulele de translație tensorială prezentate în capitolele anterioare.

Nu am folosit cuvinte din limbajul natural, ci am folosit simboluri fără semnificație semantică, pentru a avea un mai bun control asupra proceselor de raționament și a nu fi influențat în primă fază de sensul propozițiilor generate în sistemul conexionist. Pentru că vocabularul era destul de mare la nivelul întregului sistem, l-am distribuit în așa fel, încât toate modulele să "cunoască" un set minim de cuvinte comune (cele care sunt cele mai des întâlnite) din toate cele trei categorii de cuvinte posibile iar anumite cuvinte să apară la un număr mai mic de module (la jumătate, la un sfert, ș.a.m.d.). Nu am folosit cuvinte individuale, chiar și cele mai rar folosite cuvinte fiind cunoscute de cel puțin două module. Un alt motiv pentru care am evitat cuvinte explicite din limbajul natural, cu o semantică bine definită, a fost același pe care l-am mai invocat în capitolul 4 și anume imposibilitatea definirii unui univers al discursului cu un vocabular relativ mic (deși în cazul de față, vocabularul a fost mai mare).

Pentru a împărți cuvintele care nu sunt general cunoscute de toate modulele (să le numim cuvinte individuale) am încercat două tehnici: una bazată pe structura propozițiilor din sistem (în acest caz cel care dezvoltă sistemul trebuie să aleagă el

modulul potrivit pentru fiecare cuvânt) sau metoda pur aleatoare. Metoda pur aleatoare, nu este însă potrivită, pentru că în faza ulterioară de antrenare a rețelelor recurente, se poate întâmpla ca o propoziție care este necesară pentru un anumit modul, să nu fie posibil de reprezentat pentru modulul respectiv, pentru că un cuvânt din această propoziție nu aparține vocabularului aferent modulului. Singura metodă sigură este identificarea propozițiilor de intrare și de ieșire pentru fiecare modul și stabilirea vocabularului care formează propozițiile necesare antrenării unui modul.

Dar cum stabilim care sunt aceste propoziții? Aici am identificat cea mai mare slăbiciune a unui astfel de sistem, și anume că *trebuie prefigurate liniile posibile de raționament* și pe baza acestora, făcute seturi de propoziții care vor conduce la configurația modulelor. Din păcate, această metodă este bazată pe prezumția că noi trebuie să cunoaștem toate cazurile în care sistemul va ajunge de la premise la concluziile finale. Într-un sistem clasic, nu trebuie decât să adăugăm noi propoziții în baza de cunoștințe, iar în funcție de acestea, putem atinge diverse concluzii finale. În cazul sistemului multimodular conexiionist, din cauza necesității distribuirii vocabularului, trebuie să cunoaștem apriori modul în care raționamentul va fi făcut în timpul utilizării sistemului.

Pentru că într-un sistem bazat pe *modus ponens*, propozițiile pot fi grupate pe nivele într-un graf (de tip arbore) de raționament, am împărțit propozițiile care apar în raționamentele sistemului în funcție de poziția lor în acest graf. Un exemplu posibil de mod de construire a sistemului este dat în continuare. Considerăm următoarea colecție de cunoștințe:

- if $s(A, \alpha(B))$ and $t(D, \alpha(B))$ then $r(A, \alpha(B))$.
- if $s(A, \alpha(C))$ and $t(D, \alpha(C))$ then $r(A, \alpha(C))$.
- if $t(D, \alpha(B))$ and $u(D, \beta(E))$ then $r(D, \beta(E))$.
- if $t(D, \alpha(C))$ and $u(D, \beta(E))$ then $r(D, \beta(E))$.
- if $t(D, \alpha(B))$ and $u(D, \beta(F))$ then $r(D, \beta(F))$.
- if $t(D, \alpha(C))$ and $u(D, \beta(F))$ then $r(D, \beta(F))$.
- if $u(D, \beta(E))$ and $v(G, \gamma(H))$ then $r(G, \gamma(H))$.
- if $u(D, \beta(F))$ and $v(G, \gamma(H))$ then $r(G, \gamma(H))$.
- if $u(D, \beta(E))$ and $v(G, \gamma(I))$ then $r(G, \gamma(I))$.
- if $u(D, \beta(F))$ and $v(G, \gamma(I))$ then $r(G, \gamma(I))$.
- if $w(A, B)$ and $x(A, C)$ then $s(A, \alpha(C))$.
- if $w(A, B)$ and $x(D, B)$ then $s(A, \alpha(B))$.
- if $x(A, C)$ and $y(D, B)$ then $t(D, \alpha(B))$.

- $\underline{\text{if}} \ x(A, C) \ \underline{\text{and}} \ y(D, E) \ \underline{\text{then}} \ t(D, \alpha(C)).$
 $\underline{\text{if}} \ x(D, B) \ \underline{\text{and}} \ y(D, B) \ \underline{\text{then}} \ t(D, \alpha(B)).$
 $\underline{\text{if}} \ x(D, B) \ \underline{\text{and}} \ y(D, E) \ \underline{\text{then}} \ t(D, \alpha(C)).$
 $\underline{\text{if}} \ y(D, B) \ \underline{\text{and}} \ z(D, F) \ \underline{\text{then}} \ u(D, \beta(F)).$
 $\underline{\text{if}} \ y(D, E) \ \underline{\text{and}} \ z(D, E) \ \underline{\text{then}} \ u(D, \beta(E)).$
 $\underline{\text{if}} \ y(D, B) \ \underline{\text{and}} \ z(D, F) \ \underline{\text{then}} \ u(D, \beta(F)).$
 $\underline{\text{if}} \ y(D, E) \ \underline{\text{and}} \ z(D, E) \ \underline{\text{then}} \ u(D, \beta(E)).$
 $\underline{\text{if}} \ z(D, F) \ \underline{\text{and}} \ w(G, I) \ \underline{\text{then}} \ v(G, \gamma(I)).$
 $\underline{\text{if}} \ z(G, H) \ \underline{\text{and}} \ w(G, I) \ \underline{\text{then}} \ v(G, \gamma(H)).$

Se poate observa că aceste propoziții sunt formate cu ajutorul unui vocabular de 21 de cuvinte, împărțite în felul următor:

- 9 relații: r, s, t, u, v, w, x, y, z
- 3 clase: α , β , γ
- 9 individualități: A, B, C, D, E, F, G, H, I

Putem figura acest sistem expert bazat pe modus ponens ca pe un graf cu 12 noduri (vezi figura 5-3), în care nodurile 8, 9, 10, 11, 12 sunt noduri de intrare, iar 1, 2, 3 sunt noduri de ieșire. Putem scrie că nodurile au următoarele propoziții asociate:

1. $r(A, \alpha(B \text{ or } C))$
2. $r(D, \beta(E \text{ or } F))$
3. $r(G, \gamma(H \text{ or } I))$
4. $s(A, \alpha(B \text{ or } C))$
5. $t(A, \alpha(B \text{ or } C))$
6. $u(D, \beta(E \text{ or } F))$
7. $v(G, \gamma(H \text{ or } I))$
8. $w(A, B)$
9. $x(A, C) \text{ or } x(D, B)$
10. $y(D, B) \text{ or } y(D, E)$
11. $z(D, F) \text{ or } z(G, H)$
12. $w(G, I)$

Vom aloca noduri la module conexiuniste, se poate observa cum graful a fost partiționat în patru zone (I, II, III și IV), corespunzătoare funcționalității celor patru module care vor implementa sistemul multimodular. Fiecare modul va fi antrenat să accepte propozițiile posibile din intrare, generând conform regulilor sistemului expert aferente zonei respective de graf, propozițiile corespunzătoare de ieșire. Din punct de vedere al vocabularului se poate observa că avem următoarea distribuție:

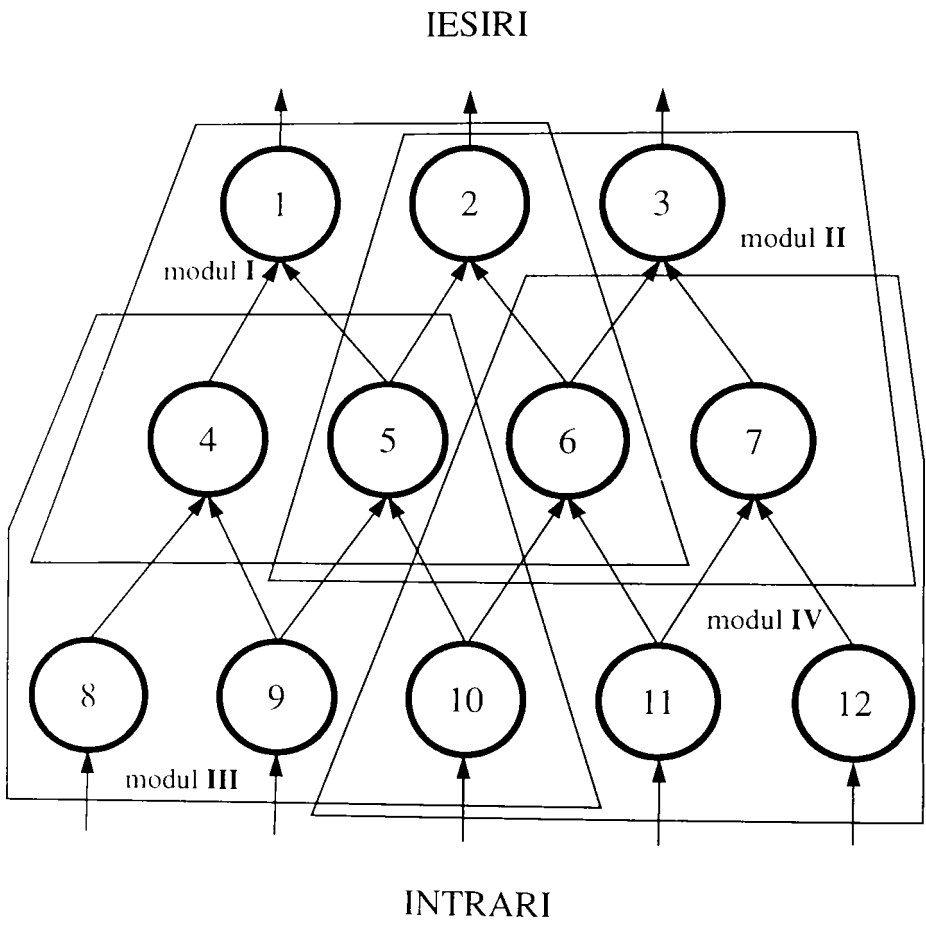


figura 4-3.

nr. modul	relații	roluri	individualități	total cuvinte/modul
I	r,s,t,u	α, β	A,B,C,D,E,F	12
II	r,t,u,v	α, β, γ	B,C,D,E,F,G,H,I	15
III	s,t,w,x,y	α	A,B,C,D,E	11
IV	u,v,w,y,z	β, γ	B,D,E,F,G,H,I	14

Putem optimiza numărul de cuvinte/modul, pentru că acesta ne va da în final numărul de unități al rețelei neuronale, astfel încât numărul de cuvinte/modul să fie cât mai mic în raport cu numărul total de cuvinte din sistem. Aici, numărul de cuvinte a scăzut de la 21 la 15 (pentru cel mai mare modul), ceea ce pare o scădere destul de mică, dar trebuie ținut cont că dimensiunea rețelei neuronale este dată de formula de translatare tensorială. Pentru cazul cel mai avantajos, cel al modulului III, vom avea un vector de lungime 125. Pentru cazul cel mai dezavantajos (II) vom avea un vector de lungime 1252 (destul de mult, luând în considerație rezultatele prezentate în capitolul 4.) Bineînțeles, numai în cazul în care folosim o reprezentare locală pentru clase, relații și individualități (fillers). Am arătat însă în capitolul anterior că folosirea rețelelor suplimentare FFI și FFO reduce la cât dorim noi dimensiunea rețelei recurente (în limita scopului propus - cel de a putea învăța tiparele necesare).

Se poate observa că alegerea acestor vocabulare este un proces deosebit de laborios și care este destul de dificil de a fi făcut în mod automat. Pentru fiecare mecanism de raționament (principiul rezoluției, mașină abstractă Warren, abducție, etc) trebuie ca acest proces de distribuire al vocabularului general să fie făcut de către proiectantul sistemului, care în plus, trebuie să identifice toate perechile de propoziții care sunt intrări/ieșiri ale unui modul, astfel încât să pregătească modulele pentru antrenarea prin algoritmul de backpropagation recurrent.

5.3.3. Experimentele efectuate

Din cauza faptului că ar fi fost o investiție foarte mare de timp să realizez un sistem cu propoziții cu semantică bine definită, am ales o cale de experimentare mai simplă, deși mai puțin spectaculoasă, fără a asocia simbolurilor folosite o anumită semantică (cuvinte din limbajul natural). Și calea pe care am urmat-o pentru realizarea modulelor a fost diferită de cea prezentată în subcapitolul anterior, pentru că am plecat în primul pas de la vocabular și nu de la reguli. Apoi am construit regulile (fiind constrâns de vocabularul strict definit) în așa fel încât modulele care rezultau să aibă cam același număr de cuvinte și aceeași distribuție a claselor, relațiilor și individualităților. Chiar am forțat acest lucru, astfel încât în final am obținut ca toate modulele să aibă același număr de cuvinte în vocabularul propriu și aceeași distribuție pe categorii de cuvinte. Mecanismul de raționament simbolic pe care am încercat să-l simulez cu modulele conexioniste era cel de modus-ponens, acesta fiind după părerea

mea foarte potrivit pentru implementarea neuronală. Vocabularul de la care am plecat a fost următorul:

$F = \{ A, B, C, D, E, F, G, H, I, J, K, L, M, N, O, P, A1, B1, C1, D1, E1, F1, G1, H1, I1, J1, K1, L1, M1, N1, O1, P1 \}$

$\text{card}(F) = 32$

$R = \{ \alpha, \beta, \gamma, \delta, \epsilon, \zeta, \eta, \theta, \iota, \kappa, \lambda, \mu, \nu, \xi \}$

$\text{card}(R) = 14$

$L = \{ q, r, s, t, u, w, x, y, z, q1, r1, s1, t1, u1, w1, x1, y1, z1 \}$

$\text{card}(L) = 18$

adică un total de 64 de cuvinte. Distribuția pe individualități (32), clase (14) și relații-predicate (18), este conformă cu distribuția acestor categorii sintactice în propoziții de tip "if-then" folosite uzual în limba engleză [Towell91].

Împărțind vocabularul pe module, am ales o distribuție uniformă, în care toate modulele aveau același număr de cuvinte și același număr de cuvinte în fiecare categorie. În urma mai multor încercări de distribuire cea mai echilibrată (din punct de vedere al resurselor necesare) a reieșit cea care avea 9 individualități, 5 clase și 5 predicate. O restricție majoră poate să apară în cazul în care atașăm individualitățile la clase și construim o ierarhie de clase din mulțimea R (o lattice). În acest caz, individualitățile care sunt atașate unui modul, trebuie ca în mod necesar să poată aparține doare claselor atașate aceluși modul. Dar din moment ce nu lucrăm la nivel semantic, acest lucru nu a creat nici o problemă, pentru că puteam să atașez din start oricare individualitate oricărei clase.

Construirea unei ierarhii de clase a simplificat problema codificării binare a cuvintelor folosite. Pentru că cele mai multe cuvinte erau din categoria individualităților, am optat pentru codificarea acestora la o schemă left-right, care a redus lungimea vectorului rezultat la 5. Aplicând formula tensorială dezvoltată în capitolul 4, lungimea unui vector de reprezentare a unei propoziții de tip intrare/ieșire este de numai 125, pentru că într-o astfel de propoziție poate fi doar un singur predicat, iar produsele tensoriale clasa/individualitate se suprapun prin adunare. Se poate observa o reducere substanțială a dimensiunii vectorilor cu care lucrăm, de la 12634 în cazul prezentat în capitolul 4, la 125 (doua ordine de mărime). Am ales să folosesc un număr total de 250 de unități_recurente/modul, cu un număr de intrări și ieșiri egal cu 125. Același număr de intrări și ieșiri l-am folosit pentru rețelele FFI și FFO (pentru fiecare modul - a rezultat o structură fixă de 125-125-125 unități). Rezulta un necesar de memorie numai pentru rețeaua recurentă de:

$$250[\text{ieșiri}] \times (250[\text{intrări recurente}] + 125[\text{intrări}]) = 93750[\text{ponderi}]$$

adică un necesar de memorie pentru fiecare modul de aproximativ de 1MByte (dat fiind și un surplus nesemnificativ pentru FFI și FFO).

Numărul de module este astfel restricționat de memoria ocupată de un modul. Pentru a nu fi blocat de această restricție și de memoria limitată a sistemului folosit de mine (8MBytes) am decis ca numai un singur modul să fie la un moment dat activ în memorie, restul de module (mai exact, matricile de ponderi) să fie menținute pe disc (unde dispuneam de în medie de 130 până la 300MBytes) iar atunci când era nevoie de un anumit modul, să încarc ponderile acestuia de pe disc. Evident că acest stil de lucru este foarte lent, dar nu are nici o importanță dacă un proces de raționament durează câteva minute. Numarul de maxim de module care puteau fi folosite într-o astfel de configurație era în jur de 100, fără ca să fiu în pericol să nu am spațiu disponibil. La limită, în condiții ideale, puteam să folosesc chiar și 300 de module.

Configurația finală a modulelor a fost stabilită în felul următor: am împărțit vocabularul de 64 de cuvinte, în mulțimi de 19 cuvinte, conform schemei prezentate, astfel încât să asigur un grad de acoperire cât mai mare, și să am un număr cât mai mare de module. În final, am obținut un număr de 19 module, pe care le-am grupat într-o structură ca în figura 5-4. Suprapunerea trapezelor care reprezintă module în figura de mai sus, indică și suprapunerea vocabularelor. Fiecare modul poate primi o propoziție ca și intrare, și poate genera la ieșire o propoziție prin simularea inferenței prin modus-ponens. Relativ la schema de funcționare prezentată în subcapitolul dedicat comunicării între module (vezi 5.1.2) aici apare o problemă legată de conectorii logici. În exemplul dat în subcapitolul anterior, se putea observa că la intrarea unui modul puteau apărea mai multe propoziții deodată, iar acest lucru era dat de folosirea conectivei "and" (conjuncție logică) care apărea în regulile "if-then" ale bazei de cunoștințe inițiale. În acest caz, ar fi fost necesară superpoziția celor doi vectori care reprezentau cele două propoziții conectate prin "and". Pentru a nu complica schema de funcționare a sistemului multimodular, am preferat să folosesc schema de secvențiere a propozițiilor rezultate și să admit la intrarea unui modul doar o singură propoziție la un moment dat. Acest lucru face ca să dispară conectivele din graful de decizie al unui sistem pur simbolic izomorf cu cel conexiunist, pentru că nu putem asocia nici un fel de conective logice cu operațiile de secvențiere din șirurile de așteptare de pe BB descrise în 5.2.1.

Ceea ce am dorit să fac era o verificare a funcționalității, în sensul că atunci când furnizăm sistemului un set de propoziții care sunt premise de lansare a unui raționament, să se obțină propoziții pe nivelul final de ieșire. Nu era important dacă rezultatul era corect din punct de vedere al logicii predicatelor de gradul 1, ci dacă pur și simplu apărea un rezultat (o propoziție) care fusese antrenat anterior în cadrul modulelor care asigurau ieșirea finală a sistemului. Dacă apăreau propoziții care nu fuseseră antrenate, însemna că sistemul nu funcționa corespunzător. Fiecare modul a fost antrenat cu câteva perechi de propoziții generate cu ajutorul vocabularului redus al modulului, fiecare propoziție reprezentând un fel de regulă "if-then". Pentru că rețelele recurente au o componentă temporală, era foarte importantă ordonarea tiparelor de antrenarea formate de perechile de propoziții, astfel încât punctele în spațiul stărilor rețelelor să fie aștezate pe o traiectorie formate din puncte apropiate. Acest lucru se putea ușor realiza, studiind vectorii asociați propozițiilor (prin translația tensorială și transformarea vectorului cu FFI). Se definește în acest caz o normă a vectorilor de

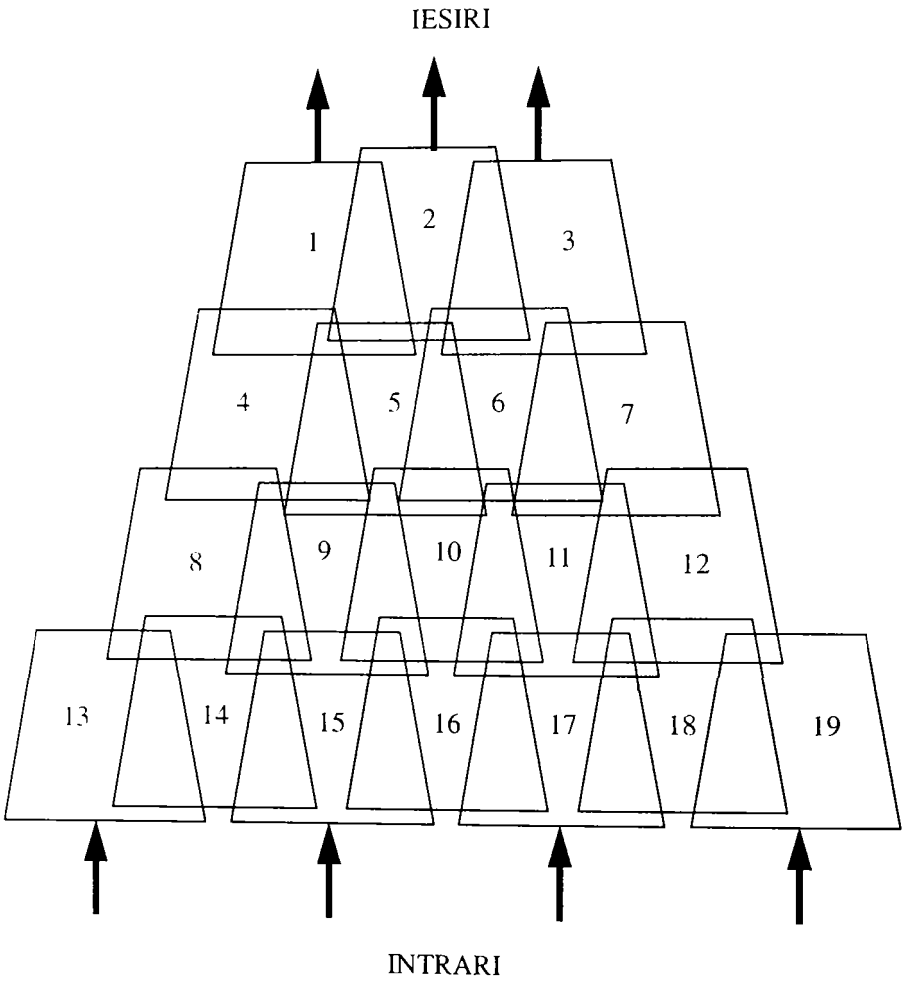


figura 5-4.

intrare/ieșire, dată de numărul de poziții din vector setate pe valoarea 1, și se ordonează crescător vectorii, astfel încât să formeze o traiectorie discretă de o formă cat mai puțin sinuoasă în hipercubul de dimensiune 125.

Experimentul care testa funcționarea sistemului, a decurs în felul următor:

- i. se poziționează rețele din toate modulele (starea modulelor - un vector real 250 dimensional - este menținută în RAM) pe un punct din spațiul 250 dimensional dat de prima propoziție de intrare și prima propoziție de ieșire de pe traiectoria descrisă de toate propozițiile învățate de modulul respectiv. Mai exact, se iau cele 125 de valori din vectorul de intrare și se setează cu aceste valori cele 125 de unități ascunse ale rețelei recurente. La fel, se iau cele 125 de valori din vectorul de ieșire și se setează cu aceste valori unitățile de ieșire ale rețelei. Această setare a stării tuturor rețelelor este o operațiune deosebit de "artificială", dar nici un raționament nu poate fi pornit dacă starea inițială este setată aleator (sau este pusă pe 0). Practic, sistemul este pus într-o stare care deja generează un răspuns plecând de la un set de premise (care sunt propozițiile de intrare de pe prima poziție a traiectoriilor de antrenare a modulelor 13-19) și ajungând la un set de concluzii finale.
- ii. Se furnizează sistemului un set de premise (care în mod obligatoriu, trebuie să fie propoziții învățate de modulele 13-19) și se pornește schema de comunicare între module, prin intermediul blackboard-ului BB, până se obțin propoziții de ieșire din modulele 1-3. În acest moment, procesarea se consideră încheiată. Se poate observa asemănarea între acest sistem și un sistem bazat pe forward-chaining, bazat pe reguli predicative.

O serie de aspecte funcționale neașteptate (unele nedorite) au fost observate în decursul experimentării acestui sistem. Primul efect l-am denumit "ciclare simplă". Pentru că există un singur blackboard care deservește toate modulele, propozițiile care sunt generate de un modul, vor fi "consumate" imediat de același modul, pentru că vocabularul propozițiilor de intrare, este identic cu vocabularul propozițiilor de ieșire. Pentru a elimina acest fenomen, am introdus pentru fiecare modul, un pas de acceptare suplimentar, care verifica dacă propoziția de intrare nu face parte din setul de propoziții de ieșire învățate. În caz afirmativ propoziția de intrare respectivă era eliminată din șirul de așteptare al modului respectiv. Efectul nedorit al ciclării simple era acela că starea rețelei era dusă de obicei într-un punct unde rămânea staționară, indiferent de propoziția pe care o primea în intrare. Alt fenomen, mai rar întâlnit, dar care ducea tot la stări staționare nedorite ale unor module, era ciclarea prin două module. Se putea întâmpla în unele situații, ca propoziția generată de un modul (de exemplu 10) care primea intrare o ieșire a unui alt modul (de exemplu 15) să genereze o ieșire, care din punct de vedere al vocabularului să fie acceptată din nou modulul care a generat prima ieșire (în cazul exemplului, tot 15). Evident că propoziția nou acceptată nu face parte din setul de propoziții de intrare învățate, dar nici din setul de propoziții de ieșire învățate (conform noii reguli, ar fi ștearsă). De obicei, acceptarea unei astfel de propoziții duce tot la o stare staționară. Se poate imagina o schemă de evitare a acestui fenomen, dar nu am mai implementat așa ceva, pentru că fenomenul era relativ rar.

Există un caz ideal de funcționare a sistemului, care dă și rezultatele dorite.

Dacă se iau propozițiile de intrare care sunt pe a doua poziție a traiectoriilor de antrenare ale modulelor 13-19 (după ce sistemul a fost pus în așteptare conform pasului 1 de lucru descris anterior) și se aplică aceste propoziții ca și premise de intrare, vom obține la ieșire propozițiile învățate de pe a doua poziție a traiectoriilor de ieșire a modulelor 1-3. Lasând modulele în starea în care au fost aduse de acest raționament (fiecare a făcut un pas și am observat că nici unul nu a rămas staționar - adică a făcut pasul dar s-a întors în starea anterioară) și aplicând ca premise propozițiile de intrare care urmează pe traiectoriile de antrenare ale modulelor 13-19, vom obține la ieșire tot propozițiile de ieșire care urmează pe cele de la procesul anterior de raționament (pe traiectoriile de antrenare ale modulelor 1-3). Procesul continuă în același mod până se epuizează toate propozițiile de pe traiectoriile de intrare care sunt premise.

Acest caz ideal presupune ca toate traiectoriile de antrenare să fie formate din același număr de puncte (propoziții), iar la intrare (ca premise) să fie date tuturor modulelor de intrare (13-19) câte o propoziție și numai una. Se poate observa că această funcționare este foarte restrictivă, fiind din punct de vedere al performanțelor cognitive (capacitate de inferență) sub nivelul unui sistem similar pur simbolic. Dar partea cea mai interesantă a funcționării acestui sistem, cea care mi-a dat noi idei pentru continuarea aplicării rețelelor recurente în raționamentul automat au fost experimentele care nu au dus la o funcționare ideală, ba chiar au dus la blocare sistemului sau la obținerea de rezultate care nu au fost prevăzute prin antrenare.

Cel mai plauzibil scenariu de funcționare al unui astfel de sistem este că el primește câteva premise la intrare, care sunt aplicabile doar unui număr mic de module care au rol de a prelua premise (în experimentul meu, modulele 13-19) și acest lucru duce în final la obținerea unei propoziții (sau mai multe) de la modulele care au rol de furniza ieșirea utilă sistemului (în cazul prezentat, 1-3). În cel mai bun caz, dacă numărul de propoziții de intrare (premise inițiale) este mai mic decât numărul de module de intrare, atunci raționamentul va implica în continuare doar propoziții care sunt generate de aceste module activate. Dacă aceste propoziții de start au fost cele de pe a doua poziție a traiectoriilor de antrenare, raționamentul se va desfășura normal, până se va genera o răspuns la la modulele care asigură ieșirea. Chiar și în cazul în care am dat o singură propoziție ca și premisă inițială.

În cazul în care în procesul următor de raționament apar premise care sunt și ele pe poziția a doua a traiectoriilor (și nu au mai fost prezentate anterior sistemului) funcționarea dorită nu mai este asigurată în continuare pentru că acestea pot să ducă la activarea unor module care au fost deja activate și au făcut un pas pe traiectoriile de antrenare. Aceste module, care sunt în "avans" față de propozițiile pe care le primesc la intrare, datorită faptului că sunt activate la de către o propoziție care este "trecută" pentru ele, vor fi module staționare. Acest fenomen duce la o "blocare prin lipsă de informație" (starvation) cunoscut din sistemele multi-agent, și la imposibilitatea obținerii unui rezultat.

Se poate observa că pentru funcționare normală a sistemului este nevoie de "sincronicitate", adică toate modulele să fie din punctul de vedere al traiectoriei discrete

al stării interne în același punct. Orice întârziere a unui modul, prin nefolosirea sa într-un proces de raționament, duce la fenomenul de starvation prin atingerea de către aceste module a unor stări staționare. Spre deosebire de acest model de sistem, un sistem pur simbolic, bazat pe reguli propoziționale de tip "if-then" și pe inferență de tip modus ponens, funcționează și fără ca această sincronicitate să fie necesară.

5.4. Concluzii

Deși rezultatele nu sunt deosebit de încurajatoare, sistemul experimentat are importanța lui practică. El demonstrează că se poate (cu o serie de restricții - destul de serioase) realiza un raționament automat cu ajutorul unor rețele recurente. Dacă acest lucru nu era posibil cu ajutorul unei rețele de mari dimensiuni (din considerente computaționale), prin divizarea acestei rețele în mai multe rețele de mici dimensiuni, care necesită fiecare în jur de 1 - 4 de ore pentru antrenarea cu 4-5 perechi de propoziții (pe o stație obișnuită), se poate implementa un astfel de mecanism de raționament.

Modelul prezentat ar fi trebuit să genereze prin natura lui conexiunistă și distribuită, trei însușiri care l-ar fi făcut superior unui sistem pur simbolic și unuia bazat pe o structură cu un singur bloc. Aceste însușiri sunt:

- i. generalizarea;
- ii. robustețea la intrări perturbate (și într-o mai mică măsură generalizarea tipică conexiunistă, deși acest deziderat era mai greu de atins, dată fiind natura duală, conexiunist simbolică a sistemului);
- iii. plasticitatea: posibilitatea de a adăuga noi module, fără a le mai modifica pe cele vechi, și acest lucru să ducă la o mai bună funcționare a sistemului, în sensul că "baza de cunoștințe" este mai mare.

Din păcate, aceste deziderate pe care le-am căutat în dezvoltarea modelului inițial nu au putut fi atinse. Generalizarea nu se poate realiza pentru că ar furniza propoziții care nu au fost învățate de către modulele dedicate pentru intrare (chiar dacă au un vocabular adecvat) și aceasta conduce în mod inevitabil la fenomenul de modul staționar și la starvation.

Robustețea ar putea fi obținută în cazul în care am avea vectori de lungime mai mare, rezultați prin translația tensorială a unor propoziții bazate pe un vocabular mare și am avea situații în care propozițiile ar avea mici deviații de la traiectoria antrenată. În acest caz, rețeaua recurentă ar fi suficient de robustă ca să nu se abată de la traiectoria învățată. Dar în cazul nostru, vocabularul este foarte mic, și modificarea unui cuvânt duce la modificarea destul de puternică a vectorilor de reprezentare. Lucrul care duce la abaterea de pe traiectoria stării interne a rețelei suficient de mult ca să se ajungă în stări care nu mai pot fi modificate prin propozițiile învățate. Concluzia este în acest caz, că cu cât mai mică este rețeaua și mai ales vocabularul aferent, cu atât modulul este mai puțin robust la intrări perturbate sau incomplete.

Plasticitatea implică cunoașterea modulelor existente și a vocabulelor lor. Dacă este nevoie de adăugarea de noi cunoștințe, trebuie ca toate cuvintele adăugate în sistem să fie deja cunoscute măcar de o parte a modulelor deja existente. Adăugarea unui cuvânt necunoscut implică construirea unui modul suplimentar, ca cel puțin două module să cunoască acel cuvânt. Mai mult, modulele nou adăugate trebuie să aibă un rol foarte clar în lanțurile de raționament, și să respecte principiile de sincronicitate ca sistemul să poată funcționa normal. Toate aceste restricții fac ca adăugarea de noi module să fie un proces foarte laborios, mult mai complicat decât în cazul unui sistem pur simbolic.

Există însă pentru această arhitectură un avantaj foarte important. Resursele necesare (timp, memorie) sunt mult mai mici decât pentru o variantă echivalentă monomodulară. Dacă n ar fi numărul total de unități necesare, iar k gradul complexității polinomiale care dă necesarul pentru o anumită resursă, este evident că dacă:

$$\sum_{i=1}^m c_i = n \quad (5-1)$$

(m este numărul de module, c_i numărul de unități pentru fiecare modul) atunci:

$$\sum_{i=1}^m (c_i)^k < n^k \quad (5-2)$$

Am mai remarcat încă un fenomen: **cu cât modulele sunt mai echilibrate ca număr de unități, cu atât resursele necesare sunt mai mici**. Dacă, la limită, $c_i=c$ ($i=1,m$), acest fenomen este susținut de relația (demonstrația ei este imediată):

$$\sum_{i=1}^m c^k \leq \sum_{i=1}^m (c+x_i)^k \quad (5-3)$$

unde $x_i \in \mathbb{Z}$, $|x_i| < c$ și:

$$\sum_{i=1}^m x_i = 0 \quad ; \quad m \cdot c = n \quad (5-4)$$

Acste observații au influențat modul în care am partajat modulele în modelul prezentat în continuare.

6. Modelul multimodular eterogen

Cea mai importantă concluzie preliminară la care am ajuns după experimentele și analiza modelului multimodular a fost că un sistem de raționament automat complex (cu un vocabular și o complexitate a informației mari) nu poate fi implementat cu ajutorul instrumentelor conexioniste actuale, utilizând mai ales rețele recurente. Motivul nu este neapărat necesarul de resurse computaționale (și aceasta este o problemă, dar nu insurmontabilă). O rețea neuronală, sau un grup care cooperează pentru rezolvarea unei probleme au nevoie de *supervizare*, pentru control și deblocare.

În sprijinul afirmației că procesarea pur neuronală nu este o soluție practică, putem face o analogie cu modelul uman de rezolvare a problemelor. Dacă problema este simplă, iar omul care o rezolvă are experiență în domeniu, răspunsul "îi vine" imediat în "minte", prin analogie cu alte probleme mai mult sau mai puțin asemănătoare. Dacă problema este dificilă și presupune utilizarea unui formalism și a unor pași concreți de rezolvare, aflați într-o secvență suficient de lungă, omul se vede în situația de a folosi o reprezentare grafică exterioară. El manipulează simboluri, aplicând reguli stricte de combinare / procesare / extragere. Doar la un nivel al pașilor intermediari între diversele operații, intervin decizii "imEDIATE", rezultate din experiență și "intuiție". Este adevărat că în mintea umană, formalismele legate de reprezentările simbolice sunt memorate și utilizate la nivel pur neuronal, dar analizând această schemă de lucru, putem spune că într-un proces de raționament, efectuat pentru rezolvarea unei probleme, exprimarea și procesarea simbolică, nu numai că sunt necesare, dar ca și importantă, ele se află la un nivel superior celui analogic-neuronal. Mai ales în sensul că ele sunt cele care "conduc" raționamentul.

La nivelul tehnicii și ingineriei actuale, din punct de vedere practic, realizarea unui sistem inteligent pur neuronal este inutilă și extrem de costisitoare. Majoritatea task-urilor intelectuale care presupun cunoștințe aprofundate și un volum mare de lucru, pot fi rezolvate prin programarea obișnuită, pe suportul computațional existent (firește există și excepții - generarea discursului de exemplu). Însa toate aceste task-uri

presupun fiecare, realizarea unor operații care sunt foarte asemănătoare (omogenitate). Problema apare când taskul propus pentru rezolvare presupune apelarea la operații care sunt foarte diferite, sau reprezentarea informației este extrem de eterogenă (sau cu accente de incertitudine). În astfel de cazuri, se folosește strategia "divide et impera", task-ul fiind împărțit în mai multe sub-task-uri specializate, care sunt rezolvate într-o anumită ordine, utilizând la nevoie, arhitecturi computaționale foarte diferite. Supervizarea realizării task-urilor este o problemă care presupune "inteligentă", pentru că agentul (uman sau automat) care se desfășoară o astfel de activitate, trebuie să aibă cunoștințe extrem de eterogene.

6.1. Modelul compozițional aplicat în inteligența artificială

În domeniul inteligenței artificiale clasice, răspunsul la problema eterogenității task-urilor a fost **modelul compozițional** [Chang&Keisler73] [Cohen83&85] [Treur88]. El este aplicat și în alte domenii, cum ar fi sistemele de operare distribuite sau bazele de date. După cercetări și experimentări care implicau sisteme urișe (vezi ca exemple [Giunchiglia&Weyhrauch88] sau [Vilain&McAllester89]), înglobând un volum foarte mare de cunoștințe, asupra cărora se aplicau mecanisme foarte complicate de inferență, s-a observat că o specializare și modularizare în cadrul acestor sisteme este necesară. Cunoștințele au fost grupate pe diverse grade de similaritate, iar mecanismele de inferență au fost simplificate, ele fiind aplicabile în acest caz asupra unei singure clase de cunoștințe. Practic, un sistem mare, a fost împărțit în mai multe module mici, mai ușor de construit, administrat și adăugat. Noi module, conținând cunoștințe într-un format nou, pot fi adăugate la sistem. Din acest motiv modelul se numește compozițional.

Ideea sistemului multimodular omogen, mi-a venit după încercările nereușite cu sistemul monomodular, pentru că la Vrije Universiteit exista un framework pentru sisteme compoziționale, numit **DESIRE** (**DES**ign and **SI**pecification of **RE**asoning modules) [Treur90] [Treur92]. Prima idee pe care am împrumutat-o din acest framework a fost aceea de a comunica prin blackboard, într-un limbaj universal (grafuri conceptuale, în cazul meu) cunoscut de toate modulele. Spre deosebire de modelul meu, DESIRE este eterogen, modulele sale nefiind identice ca funcționalitate interioară. Doar limbajul de comunicare (un gen de propoziții asemănătoare logicii predicatelor de gradul întâi, trivalorică și cu operatori modali, fără cuantificatorul existențial) era același pentru toate modulele, fiecare modul având în componența sa, translațoare din/în acest limbaj. De exemplu, pentru implementarea unui joc, se pot construi mai multe module care rezolvă problema "mutării" următoare (prin metode diferite), altele care supervizează "strategia" pe termen lung, astfel încât, în final, o variantă de optim (de obicei local) să fie atinsă într-un pas de joc. Stilul acesta de lucru, pe lângă avantajul simplificării prin modularizare, mai are și avantajul execuției paralele, prin alocarea diverselor module la procesoare mai mult sau mai puțin specializate.

Ideea de bază a colectivului de la VU a fost implementarea unor module foarte simple, capabile să facă operații banale de raționament, care să fie utile în cât mai multe posibile aplicații. Dezvoltând aceste aplicații, sistemul este lărgit treptat, module care servesc numai unor aplicații, urmând să devină utile și noilor aplicații (dacă este cazul). Sistemul se dezvoltă gradual, devenind din ce în ce mai mare (nici un modul construit nu este abandonat), pentru rezolvarea unor probleme noi, utilizându-se cât mai mult din vechile module. Sub "acoperișul" inteligenței artificiale, se refolosea o idee mai veche, din ingineria programării, cea a reutilizării codului. Însă aici, aplicarea acestui principiu are o semnificație mai subtilă. Refolosirea unui modul, înseamnă refolosirea unor cunoștințe "dobândite" (prin programare) anterior pentru rezolvarea unei probleme diferite dar care are o parte oarecum analoagă cu noua aplicație.

Dezideratul "creșterii" inteligenței unui sistem prin simpla adăugare a unor module nu a fost încă atins. În acest model, creșteri cantitative nu duc la salturi calitative. Cercetătorii de la VU cred că punctul important este protocolul de comunicare între module. Diverse tehnici de comunicare au fost folosite, evitându-se construirea unui modul de supervizare, care să preia problemele, să le dividă, să le secvențeze și să le aloce modulelor. Din păcate, un protocol "natural" de comunicare nu a fost încă găsit, fiind necesară pentru fiecare tip de aplicație, construirea unui modul dedicat de secvențiere.

Din discuțiile săptămânale avute cu colectivul de AI de la VU (în timpul șederii mele la Amsterdam), am observat o reticiență în a folosi module bazate pe calcule analogice (mai ales datorită tradiției pur simbolice a grupului respectiv). Pornind de la această observație și de la concluzia că un sistem de raționament automat **trebuie** să fie dirijat de o componentă simbolică, am încercat inserarea unor module conexiuniste în acest model. Prima soluție, era pur și simplu adăugarea unor module care funcționau la același nivel ca și cele simbolice, comunicând prin blackboard (transformând grafulurile conceptuale în propoziții de gen DESIRE și invers). Piedica principală a fost vocabularul foarte bogat al părții simbolice a sistemului. Un alt motiv pentru care nu am folosit această metodă, a fost că tot ce puteau face modulele conexiuniste recurente dezvoltate de mine, puteau face și module simbolice cu același rol, pentru că la nivelul problemelor care puteau fi rezolvate de DESIRE nu era nevoie de generalizare și de robustețe.

Analizând funcționarea unei aplicații de generare a discursului utilizând un sistem bazat pe framework-ul DESIRE, s-a constatat [vanLangevelde92] că problema principală care apărea la funcționare și care apărea și în alte aplicații, era *starvation* (problemă apărută și în cazul modelului multimodular omogen - vezi capitolul 5). Modulele nu își puteau continua linia de raționament din cauza unei propoziții lipsă, care trebuia în acest caz, furnizată pe baza contextului în care se afla sistemul, context dat de informația de pe blackboard și eventual starea internă a unor module active. Interesant că pentru probleme oarecum asemănătoare, în speță, generarea unor discursuri cu semantică și temă aproape identică, dar cu amănunte diferite, era nevoie de *aceleași* propoziții. Metoda de lucru prin care utilizatorul furniza propozițiile

necesare continuării raționamentului era denumită la VU "offering hints". Proiectanții sistemului puteau face o analiză experimentală a unei anumite aplicații și puteau adăuga un modul care furniza în mod automat aceste propoziții în funcție de context. Acest stil de lucru nu asigură însă o facilitate automată de Machine Learning. Soluția la care m-am gândit a fost adăugarea și antrenarea unor module conexioniste pentru "offering hints".

6.2. Schema pe doua nivele

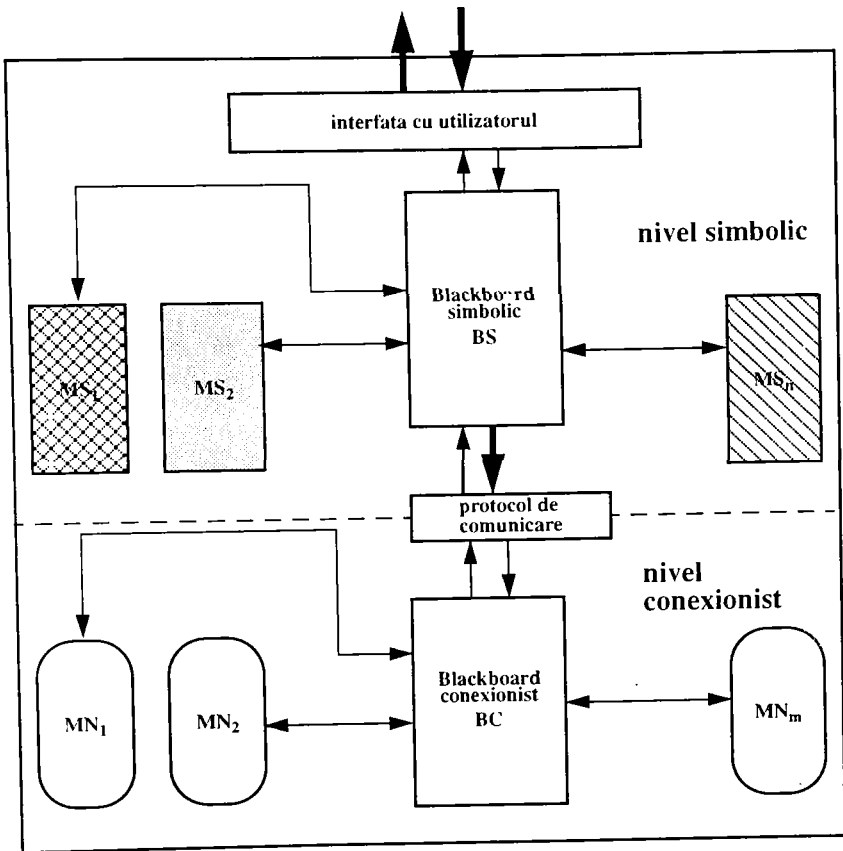
Din discuția de mai sus, rezultă că partea simbolică a sistemului este componenta de bază, cea care preia informația de start, execută linia de raționament, și furnizează informația. Partea neuronală este folosită doar pentru deblocarea sistemului, atunci când acesta întâmpină evenimente de tip "starvation-points". Am considerat că cel mai potrivit mod de a realiza hibridizarea este utilizarea unei "gâtuirii" (bottleneck) a informației între două nivele distincte. Un nivel va fi pur simbolic, compus din module care cooperează prin intermediul unui blackboard, și un alt nivel compus din module conexioniste, care cooperează prin intermediul unui alt blackboard. Singura legătură între cele două nivele este asigurată prin scrierea unor informații de pe un blackboard pe celălalt. Schema unui sistem de acest tip este dată în figura 6-1. Analizând viabilitatea unei astfel de scheme bazate pe bottleneck, comparând-o cu alte idei care erau aplicabile în această situație și cu alte metode propuse în literatură (de exemplu: [VanGelder95], [Hilario96], [Khosla&Dillon96]), am concluzionat că este cea mai plauzibilă arhitectură. Deși nu era explicit enunțat, se poate observa și în arhitectura sistemelor propuse acest principiu al gătuirii informației (precum și al separației funcționale între cele două stiluri diferite de procesare).

6.2.1. Funcționarea părții simbolice a sistemului

Pentru că sistemul propus era unul pur experimental, menit să demonstreze utilitatea unui nivel conexionist, am folosit un număr relativ mic de module pur simbolice, toate bazate pe mecanisme de raționament cu forward-chaining și backward-chaining, iar limbajul de comunicare a fost redus la structuri predicative cu maximum două predicate. Am încercat două tipuri diferite de protocoale de comunicare. În principiu, aceste module simbolice funcționează sincron, citind propoziții de pe blackboard la un moment dat, și generând în pasul următor propoziții care se scriu pe blackboard. Trei strategii diferite de citire/scriere pot fi folosite:

- i.- după ce s-au citit propozițiile, cele care au fost citite se șterg.
- ii.- după ce s-au citit propozițiile, toate sunt șterse.
- iii.- toate propozițiile rămân.

Primele două strategii sunt utile în cazul în care modulele nu au o stare internă, care



MS_i = module simbolice (eterogene)
 MN_i = module conexiuniste (omogene)

figura 6-1.

se modifică dinamic, iar alocarea propozițiilor la module se face static (doar pe baza vocabularului). În cazul în care alocarea unei propoziții la un modul depinde și de starea internă a acestuia (punctul unde a ajuns el în raționamentul local pe care îl efectuează) [Treur92], atunci toate propozițiile pot să rămână pe blackboard.

Pentru simplitate, am ales a doua strategie de citire/scriere. Aparent rudimentară, ea este însă mult mai potrivită pentru funcționarea sistemului fără ca acesta să aibă module dedicate pentru supervizarea executării anumitor task-uri particulare. Strategiile care implică lăsarea tuturor propozițiilor pe blackboard, implică o puternică supervizare, atât din punctul de vedere al alocării propozițiilor (sub-task-urilor) la module cât și din punct de vedere al secvențierii operațiilor (unele module trebuie să execute neapărat operații înaintea altora, deși acestea pot fi "pornite", în caz contrar, apărând situații de deadlock).

Din punct de vedere al modului de lucru, am ales varianta nesupervizată (numită și "naturală"), în care modulele recunosc care sunt propozițiile de care au nevoie la un moment dat (acestea fiind prezente pe blackboard) și doar prin simpla comunicare/citire/generare de propoziții, task-ul general este rezolvat. La nivel teoretic, funcționarea unui astfel de sistem este demonstrată [Treur90], în felul următor: se construiește o bază de cunoștințe formată din reguli de aceeași formă (omogenă) și care poate fi folosită într-un raționament prin un singur mecanism de inferență, iar aceasta bază KB, este divizată pe module care pot executa raționamente parțiale pe grupuri mici de reguli (kb). Formalizând blackboard-ul sub forma unui modul dinamic, care își schimbă setul de reguli în timp, și alocând fiecărui modul o stare, care se schimbă în funcție de starea curentă a blackboard-ului și a operațiilor de raționament locale unui modul, se poate efectua cu sistemul compozițional orice raționament care se putea executa și cu sistemul inițial (monomodular - KB). Există mult mai multe variante de sisteme compoziționale [Treur92], dar acesta este singurul care funcționează fără modul de supervizare. Limitarea majoră a acestui sistem este firește omogenitatea, toate modulele fiind de fapt la fel. Singurul avantaj este ușoara mentenanță, noi cunoștințe putând fi adăugate prin construirea unor noi module (mici și în special lizibile) și prin funcționarea eficientă (doar strictul necesar de module participă în raționament, reducând timpul de procesare).

Problema apare atunci când adăugăm module care nu sunt identice ca și formă de reprezentare internă a cunoștințelor și/sau mecanism de raționament. Formalismul care definește starea internă a modulelor trebuie să fie singular (plecând de la modulul dinamic reprezentat de blackboard), și el se bazează pe unicitatea acestor structuri per întregul sistem. Dacă două module sunt diferite nu se mai poate descrie o stare a modulelor (numai empiric). Până la nivelul anului 1996 (Decembrie) nu s-a realizat nici măcar unificarea pentru două module bazate pe backward-chaining, respectiv forward-chaining în logica bivalentă, care ar fi primul pas [Treur97]. Mai mult, eforturile s-au concentrat în direcția cercetării unui formalism adecvat pentru modulele de supervizare (secvențiere prin grafuri de tranziție, alocarea nodurilor, găsirea unui model cât mai general pentru aceste module).

Adăugarea unui modul cu structură internă diferită, transformă sistemul într-un sistem eterogen. Atâta timp cât modulul nou adăugat se supune regulilor compoziționalității (adică poate citi/scrie propoziții de pe blackboard) nu există nici un impediment ca sistemul să fie viabil din punct de vedere funcțional. Experimentele cu astfel de sisteme eterogene nesupervizate, au relevat apariția unor puncte de starvation. Se pare că problema este insuficiența informației care asigură legătura între module (adică informația din întregul sistem nu mai poate fi refăcută ca un set de reguli KB, care să ducă la construirea unui graf de decizie conectat). Problema poate să apară și în cazul în care sistemul este omogen, și se adaugă un nou modul. În acest caz însă, se poate face o verificare a faptului că graful derivat din KB-ul general este conectat sau nu. În cazul în care sistemul este eterogen, această verificare automată nu se mai poate face [Treur90]. Singura soluție este o analiză funcțională a sistemului pe un număr cât mai mare de exemple de procesare, și adăugarea informației necesare pentru evitarea punctelor de starvation. Metoda este evident empirică, laborioasă, și nu exclude lăsarea unor potențiale puncte de starvation nerezolvate. Din acest motiv, sistemele compoziționale eterogene (chiar și cele cu module de supervizare) permit interacțiunea utilizatorului pentru furnizarea unor informații necesare pentru deblocarea sistemului.

Observație: starea internă a modulelor, într-un sistem omogen, nu este relevantă pentru utilizatorul care urmărește linia de raționament. Această stare este folosită pentru alocarea propozițiilor la module, și pentru demonstrarea teoretică a viabilității sistemului. În cazul sistemului eterogen fără module de supervizare, utilizatorul trebuie să vadă o "stare" a sistemului. Această pseudo-stare este de obicei dată de configurația blackboard-ului, și propozițiile care au fost citite anterior.

6.2.2. Funcționarea nivelului conexionist al sistemului eterogen

Dezavantajul principal al metodei de lucru prin "offering hints" este că sistemul nu "reține" (învață) din experiență, astfel încât, în aceleași situații, utilizatorul trebuie să ofere aceeași informație de fiecare dată. Metoda de corectare este doar studiul acestor cazuri, și adăugarea de noi reguli unor module, sau chiar a unor noi module. Ideea mea a fost să automatizez acest proces de eliminare al punctelor de intervenție utilizator, prin adăugarea unui nivel conexionist, compus dintr-o populație de module neuronale, ca cel prezentat în capitolul anterior. Am denumit noua arhitectură *modelul eterogen multimodular hibrid pe două nivele* [Szirbik97] (vezi figura 6-1). Natura eterogenă este dată de nivelul simbolic, care are în componență module de tipologie internă diferită, iar natura hibridă este dată de existența a două nivele, unul simbolic și unul conexionist. Spre deosebire de sistemul prezentat în capitolul 5, în care modulele rezultau ca subcomponente ale unei graf de decizie inițial, și în care distribuirea vocabularului general trebuia făcută cu foarte multă strictețe, altfel sistemul neavând o funcționalitate utilă, aici lucrurile sunt mai simple. Pur și simplu, nivelul conexionist trebuie să "urmărească" linia de raționament care se desfășoară la nivel simbolic, și în momentul în care apare un punct de starvation, să poată oferi informația

necesară.

Primul pas pentru realizarea nivelului conexiunist este stabilirea configurației modulelor. Aceasta se face ca și în cazul precedent (capitolul 5) prin împărțirea vocabularului general în subseturi care se acoperă (intersecția lor nu este vidă). Este natural ca modulele să aibă un vocabular care să poată genera propozițiile necesare pentru rezolvarea de puncte de starvation. Se memorează aceste propoziții, se grupează în seturi care au același vocabular, iar aceste seturi vor fi alocate fiecare unui modul conexiunist. În funcție de dimensiunea fiecărui set de cuvinte, se poate stabili dimensiunea rețelei recurente interne. Modul în care modulele sunt antrenate, va fi prezentat în subcapitolul 6.3. Ideal ar fi ca nivelul conexiunist să fie compus dintr-un singur modul care să poată oferi "hints", pentru că în acest caz antrenarea ar fi mult mai simplă, iar problema dificilă a distribuției dicționarului ar dispărea. Pentru că vocabularul necesar ar fi foarte mare, din considerente prezentate în capitolul 4, este imposibil de realizat un astfel de modul. Se poate argumenta că vocabularul necesar pentru "hints" este mai redus decât vocabularul întregului sistem, dar vom vedea că este nevoie de toate cuvintele care sunt cunoscute în cadrul părții simbolice a sistemului.

Modulele conexioniste execută și ele "pași" de lucru, care sunt deplasări pe traiectoriile interne posibile (traiectorii discrete). Ideea de bază a funcționării sistemului este că modulele simbolice și modulele conexioniste execută pași în mod *sincron*. Dacă un număr de module de pe nivelul simbolic execută un pas în linia de raționament, modulele conexioniste (cele care sunt activate conform regulii vocabularului - vezi capitolul 5) execută și ele un pas pe traiectoriile din spațiul stărilor rețelelor recurente. Problemele care sunt imediat următoare sunt stabilirea unui protocol de lucru între modulele conexioniste și stabilirea unui protocol de comunicare între cele două nivele. Deși în principiu, schema de funcționare a nivelului conexiunist era aceeași ca și la sistemul multimodular omogen, aceasta era puternic afectată de modul în care se efectua comunicarea între cele două nivele. Pentru că raționamentul era dirijat de nivelul simbolic, nivelul conexiunist fiind oarecum "auxiliar", nu mai este nevoie de o schemă foarte complicată, care să poată efectua un raționament de sine stătător, ci este suficientă și o schemă mai simplă, care să poată oferi informația necesară deblocării sistemului la momentul potrivit.

6.2.3. Posibile moduri de comunicare între cele două nivele

Principiul pe baza căruia am încercat să stabilesc un stil de comunicare între cele două nivele a fost principiul bottleneck-ului informațional. Modulele simbolice nu pot comunica direct cu cele conexioniste (și invers). Acest lucru are și o motivație psihologică, în cursul unui raționament, omul își poate concentra atenția numai asupra unor aspecte de moment, un set restrâns de atomi de informație (propoziții, formule, părți de desen) care sunt în fața lui, reprezentate simbolic. Minteia lui este concentrată în acel moment, doar la câteva concepte (experimentele cu subiecți umani [Quillian&

Collins69] au relevat ca nu pot fi mai mult de 7, cu o medie de 3-5), acestea fiind în general ultimele asupra cărora și-a concentrat atenția. În modelul meu, comunicarea între nivele se face exclusiv între cele două blackboard-uri. Am notat blackboard-ul nivelului simbolic cu BS, iar blackboard-ul nivelului conexionist cu BC.

Cele două nivele funcționează sincron, la un pas pe linia de raționament simbolic îi corespunde un pas pe traiectoriile neuronale. Între aceste două evenimente există un interval, ele neavând loc simultan. Prima dată se execută pasul pe nivelul simbolic, ceea ce duce la scrierea unor propoziții pe BS. Întotdeauna, aceste propoziții sunt copiate pe BC. După aceasta copiere, se execută pasul la nivel conexionist, care duce la generarea unor propoziții noi, care sunt scrise pe BC. Am identificat patru metode posibile (mai sunt și altele, dar nu am considerat că au sens) de gestionare a informației de pe BC:

I. Dacă nu există în pasul următor de pe nivelul simbolic un starvation point, atunci, înainte de citirea noilor propoziții de pe BS pe BC, toate propozițiile de pe BC sunt șterse. Dacă apare un punct de întrerupere la nivel simbolic, propozițiile nou generate pe BC sunt scrise pe BS și se încearcă un nou pas la nivel simbolic. Dacă propoziția lipsă a fost generată corect la nivel conexionist, noul pas la nivel simbolic va reuși și conform cu această metodă, înainte de a scrie noile propoziții de pe BS pe BC, BC este ștersă în totalitate.

II. Doar propozițiile consumate de pe BC sunt șterse (indiferent dacă este vorba de starvation point sau nu). În acest caz, pe BC pot rămâne propozițiile neconsumate care au fost citite de pe BS, sau propoziții care au fost generate de nivelul conexionist, toate fiind grupate în cozi de așteptare, în modul descris în capitolul anterior. Se poate întâmpla ca propoziții generate la nivelul simbolic să nu poată fi citite la nivel conexionist, din cauza vocabularului limitat al modulelor (acestea cunosc toate cuvintele, dar pot exista propoziții care să fie formate din cuvinte care aparțin unor module diferite). Aceste propoziții nu trebuie scrise pe BC, iar pentru aceasta am folosit un "filtru" care analizează compoziția propoziției din punct de vedere al vocabularelor posibile.

III. Propozițiile generate la nivel conexionist, sunt scrise pe BS, indiferent dacă a apărut un punct de întrerupere sau nu, înaintea pasului executat la nivel simbolic (și sunt folosite în acest pas de raționament, dacă apare starvation point). Conform strategiei inițiale, după un pas de raționament simbolic, toate propozițiile de start sunt șterse de pe BS, înainte de a scrie noile propoziții generate. Relativ la ce se întâmplă pe BC, putem avea două subcazuri:

IIIa. Toate propozițiile de pe BC sunt șterse.

IIIb. Numai propozițiile consumate (la nivel simbolic, pentru deblocare, sau la nivel conexionist) sunt șterse. Această ultimă variantă, poate fi rafinată, în sensul că doar propozițiile consumate la nivel simbolic să fie șterse, cele consumate la nivel

conexionist, urmând să fie păstrate - în ideea că eventual, o astfel de propoziție generată la nivel conexionist, va fi utilă în unul din pașii următori.

Aceste patru metode au fost prezentate în ordinea complexității date de implementarea lor. Din experimente, am constatat că cea mai eficientă metodă era prima, care este și cel mai simplu de implementat. Asta nu înseamnă că celelalte nu sunt bune. În special IIIb, este o metodă foarte interesantă, care duce la rezolvarea unor probleme care apar aplicând metodele I și II. Se mai poate observa că aceste metode implică gradual (de la prima la a patra) o creștere a schimbului de informație între cele două nivele (referindu-ne numai la direcția BC către BS, cealaltă direcție fiind constantă). În primul caz comunicarea este limitată doar la momentele de starvation, în al doilea caz se observă o aglomerare a informației pe BC, iar în al treilea caz, fiecare pas conexionist contribuie la pașii de pe nivelul simbolic.

6.3. Antrenarea modulelor nivelului conexionist

Primul lucru care trebuie cunoscut la antrenarea nivelului conexionist este numărul de module. Acest număr poate fi constant sau variabil. Am ales varianta mai simplă (statică) în care numărul de module rămâne neschimbat. Mai mulți factori influențează numărul de module create. Primul este numărul total de cuvinte folosit la nivelul simbolic. Din experiențele efectuate cu sistemul multimodular omogen, am observat că numărul optim de cuvinte/modul (lucrând cu stații de lucru obișnuite este de 15-22 cuvinte. Am concluzionat că numărul maxim de module este dat de formula empirică:

$$\text{nr_module} = (\text{nr_total_cuvinte} / 22) * \text{grad_de_acoperire_global}$$

Tot din experiențele anterioare, am constatat ca valoarea uzuală pentru *grad-de-acoperire-global* este în jur 1.5 (exprimarea matematică a gradului de acoperire global este dată în 6.4.3).

6.2.1. Vocabularul folosit la nivel conexionist

Factorul cel mai important în stabilirea numărului de module este însă numărul de starvation points, adică numărul total de propoziții care trebuie furnizate nivelului simbolic. Acestea pot fi identificate manual, prin analiza funcțională a sistemului (compozițional-eterogen-pur simbolic), sau automat, în timp ce sistemul este folosit de utilizator, care în momentele de blocare, oferă "hints". S-a observat (vezi [Treur90], [Treu92]), că pentru mai multe linii de raționament asemănătoare, punctele de întrerupere apar în general în același moment (măsurat prin numărul de pași efectuat de la începutul raționamentului). Mai mult, propozițiile care sunt necesare pentru deblocarea raționamentului sunt foarte asemănătoare din punct de vedere al

vocabularului folosit. Această observație m-a condus la decizia de a construi un modul separat pentru fiecare grup de propoziții asemănătoare care apar aproximativ în același moment.

Definiție

Un Punct General de Blocare (PGB) este un moment t (variabil) din linia de raționament a mai multor linii de raționament $\{L_i\}_{i=1,n}$ ($n = \min(m)$) unde apare un punct de întrerupere datorită fenomenului de starvation. Valoarea t nu este fixă, ea fiind calculată față de o valoare fixă t_j , ca $t_j + l - k$, definindu-se astfel un interval de apariție.

Observație: Valorile m și k se stabilesc de către proiectantul sistemului, în funcție de dimensiunile vocabulelor modulelor care rezultă.

Am lucrat cu valorile $m = 4$ și $k = 1$. Pentru detectarea automată a acestor PGB-uri, sistemul trebuie încercat de utilizator de un număr cât mai mare de ori, acesta oferind informația lipsă în momentele de blocare. Doar în cazul în care informația oferită este utilă (duce la continuarea și terminarea raționamentului), propoziția este reținută ca făcând parte din grupul de propoziții legat de PGB. Acest grup poate avea un număr mai mare de propoziții decât n și în acest caz ar trebui împărțit în grupe mai mici de propoziții. Împărțirea se face pe baza vocabularului, astfel încât pentru subgrupele de propoziții să avem un număr cât mai mic de cuvinte. Pentru că am făcut această împărțire manual (deși se poate construi un program care să facă această aplicație automat), nu am fost foarte riguros în păstrarea regulii ca exact patru propoziții să facă parte dintr-un subgrup (uneori erau și 5 sau 6, uneori doar 2 sau 3). Relativ la exemplele pe care le-am folosit, împărțirea unui grup în subgrupe s-a întâmplat destul de rar, de obicei un PGB având asociat un singur modul conexiunist. Mai mult, în experimentele pe care le-am făcut, am observat ca o tendință că sunt puține starvation points care să poată fi grupate conform definiției PGB-urilor. Pot să apară și starvation points care să nu îndeplinească condițiile impuse de valoarea n . Am luat decizia ca și pentru acestea să construiesc module, chiar dacă a fost cazul de a forma un modul pentru o singură propoziție (deși astfel de situații individuale sunt ușor de rezolvat, prin intervenția proiectantului la nivelul simbolic - printr-o simplă adăugare a unor reguli într-unul din modulele simbolice). Avantajul acestei metode de rezolvare conexiunistă s-a observat mai târziu.

Pentru a putea furniza o propoziție, un modul conexiunist bazat pe o rețea recurentă, trebuie să parcurgă o traiectorie în spațiul stărilor acestei rețele. Traiectoria este discretă, iar fiecare punct reprezintă o propoziție translatată din spațiul structurilor simbolice în spațiul vectorial real N -dimensional. Se poate observa de aici, că vocabularul folosit de un modul, care primește la intrare propoziții provenite de la nivelul simbolic, prin scrierea de pe BS pe BD, nu poate fi limitat doar la vocabularul propozițiilor asociate unui PGB. Pentru a putea urmări raționamentul desfășurat la nivel simbolic, este bine ca modulele conexiuniste, să poată recunoaște TOATE

cuvintele folosite în sistem. După cum se va vedea, datorita schemei de antrenare pe care am folosit-o, metoda mea necesită acest lucru.

Ceea ce am încercat să evit la stabilirea configurației modulelor conexioniste a fost o legătură directă inter-nivele, între anumite module simbolice și anumite module conexioniste. Am fost tentat să fac la început acest lucru, pentru ca PGB-urile puteau fi asociate și cu perechi de module simbolice (starvation points-urile apăreau cu frecvență ridicată "între" două module simbolice care puteau fi identificate). Astfel s-ar fi putut crea o legătură directă între trei module, două simbolice și unul conexionist, care s-ar fi susținut reciproc, prin schimb de informație, pentru evitarea blocării. Însă folosind o astfel de comunicare, aș fi renunțat la principiul de "bottleneck" informațional și chiar la naturala compozițională a sistemului (modulele simbolice nu ar mai fi comunicat doar prin BS, ci și prin intermediul modulelor conexioniste). S-ar fi renunțat și la BC, de care nu era nevoie în această schemă (modulele simbolice comunicând direct cu modulele conexioniste), astfel încât ar fi dispărut comunicarea între modulele conexioniste. Această comunicare s-a dovedit a fi foarte importantă, pentru că ea a fost motivul pentru un comportament avantajos, neașteptat în timpul experimentelor (pe larg prezentat în subcapitolul 6.4). O tentație asemănătoare poate fi și asocierea intenționată (la nivel de proiectare inițială) a vocabularelor unor module simbolice cu vocabulare ale modulelor conexioniste, mapări ale unor cuvinte din zona simbolică cu reprezentarea lor în zona neuronală fiind întâlnite în literatură (vezi [Sun92]). Această legare am considerat-o ca nefiind necesară, pentru că oricum, în timpul construirii modulelor prin schema MILD (prezentată în subcapitolele următoare), o astfel de mapare se realizează automat.

6.3.2. Traiectoriile de antrenare și alocarea lor la module

Datorita comunicării permanente între BS și BC (la fiecare pas, conținutul BS este copiat pe BC), modulele conexioniste pot avea acces la informația care este folosită pentru efectuarea raționamentului simbolic. Pentru a putea oferi informație utilă la momentul oportun, ele trebuie să decidă între posibilele linii de raționament simbolic care pot să existe. Numărul de linii de raționament posibile este foarte mare (dat de o formulă combinatorială în raport cu numărul de propoziții care apar în sistem). Dar "hint"-ul corespunzător este întotdeauna asociabil cu linia de raționament din care face parte. Dacă cunoaștem ce propoziții s-au scris pe BS până la un anumit moment, putem prevedea cu ușurință ce propoziție ar urma (și lipsește). Ideea mea a fost de a lega linia de raționament care duce la un starvation point cu o traiectorie neuronală care se termină într-un punct care reprezintă propoziția căutăată. Pentru că o rețea neuronală recurentă poate fi antrenată cu mai multe traiectorii diferite, un modul conexionist poate rezolva mai multe starvation points (grupate într-un PGB). Dacă o rețea neuronală primește ca intrare o serie de propoziții (translatate în vectori) pe care le-a învățat anterior, ea va evolua pe o traiectorie predictibilă, la capătul căreia putem plasa propoziția care rezolvă punctul de întrerupere care a blocat raționamentul simbolic.

Aceste propoziții care formează traiectoriile de antrenare, pot fi identificate, urmărind interacțiunea utilizatorului cu sistemul simbolic. De la startul raționamentului, toate propozițiile care au fost plasate pe BS, sunt memorate separat pentru fiecare pas executat, și în ordinea apariției lor. La un moment dat, pot să apară mai multe propoziții pe BS, dar va trebui ca numai una per pas să fie reținută (modalitatea de eliminare a restului de propoziții se face pe bază de vocabular). Numai atunci când apare un punct de întrerupere, propozițiile memorate până în acel moment vor fi folosite pentru construirea bazei de antrenare pentru rețeaua neuronală. Mai mult, punctul de întrerupere *trebuie* rezolvat cu succes de utilizator; informația oferită de el trebuind să ducă la finalizarea raționamentului dorit. În acest caz favorabil, propozițiile liniei de raționament, de la început, până la punctul de întrerupere, trebuie memorate într-o listă ordonată pentru a fi utilizată în procesul de învățare a rețelelor recurente. Ultima propoziție din listă va fi tocmai cea oferită de utilizator. Am numit aceste propoziții terminale, PHINT (phrase for hint). Întregul șir de propoziții care duce la PHINT, l-am denumit WL-PHINT (Whole Line to PHINT).

Dacă extragem toate cuvintele din care sunt compuse propozițiile dintr-un WL-PHINT obținem un vocabular de o anumită mărime. Grupând vocabularele mai multor WL-PHINT-uri (prin reuniune), a căror PHINT-uri aparțin aceluiași PGB, obținem de obicei un vocabular cu o dimensiune foarte mare în raport cu numărul total de cuvinte (o treime sau chiar mai mult de jumătate). Pentru construirea modulelor este important (din punct de vedere al resurselor computaționale) ca vocabularul intern asociat rețelei recurente, să aibă dimensiuni rezonabile (20, maximum 30 de cuvinte - vezi capitolele 4 și 5). Dacă am asocia fiecărui PGB un modul conexiunist și dimensiunea relativă a vocabularului PGB ar fi mare în raport cu numărul total de cuvinte, am restrânge vocabularul total la un număr prea mic de cuvinte (maximum 60). Reuniunea vocabularelor asociate tuturor PGB-urilor impune vocabularul total al sistemului simbolic folosit experimental. În cazurile practice, vocabularul de la nivel simbolic impune vocabularul folosit la nivel conexiunist.

Pentru sisteme cu un vocabular restrâns (sisteme destul de rar întâlnite), soluția ar fi construirea unui modul pentru fiecare PGB, fiecare WL-PHINT asociat PGB fiind transformat prin translație tensorială într-o traiectorie de antrenare de intrare pentru rețeaua neurodinamică. Am numit această traiectorie T2P (Trajectory to PHINT). Ieșirea va fi de tip predictiv, propoziția care se generează, este următoarea propoziție (asociată următorului pas) de pe traiectoria de intrare. Acest mod de antrenare este cel mai natural și avantajos pentru rețelele recurente [Williams&Zipser90]; chiar dacă există propoziții lipsă, sau unele propoziții sunt incomplete, starea rețelei se menține pe traiectoria antrenată. Apar însă două probleme.

a. Adăugarea unor noi module simbolice, poate mări vocabularul total al sistemului, lucru care influențează aproape sigur PGB-urile existente, vocabularul asociat lor fiind mărit de apariția a noi WL-PHINT-uri. Acest lucru duce la reconstruirea din start a modului conexiunist, cu toate implicațiile legate de mărirea matricii de ponderi și a timpilor de antrenare. Desigur, se poate dezvolta sistemul prin

adăugarea de noi module simbolice și fără a mări vocabularul, dar această restricție este foarte severă. Din acest motiv nu am folosit idea de asocia fiecărui PGB un modul conexiunist, ci am gândit o schemă mai flexibilă, pe care am numit-o "Module Integration and Linkage for Distribution" (MILD).

b. A doua problemă este însă mult mai importantă. O rețea recurentă de tipul celei utilizate prezintă o restricție importantă, legată de lungimea traiectoriilor discrete care urmează să fie antrenate. **Acestea trebuie să fie formate toate din același număr de puncte.** Numărul de propoziții din WL-PHINT-uri este diferit pentru același PGB și trebuie ca T2P-urile să aibă același număr de puncte. Am găsit mai multe metode de a rezolva această problemă. Una ar fi executarea unor operații de "join" între propoziții consecutive, astfel încât numărul final să fie mai fix. Dar nu se poate garanta întotdeauna că aceste operații duc la rezultatul dorit (lungimi egale). O metodă mult mai simplă este de a găsi experimental un număr fix de propoziții pentru un modul (pe care l-am notat NP - Number of Phrases) și a elimina primele TN_m - NP propoziții din WL-PHINT_m, unde TN_m este numărul total de propoziții din WL-PHINT-ul de indice m . O altă problemă care decurge din această metodă, este că uneori lungimea traiectoriilor este prea mică pentru a putea îndeplini sarcina rețelei de a prezice următoarele propoziții de pe traiectorie. Și pentru această problemă am găsit o rezolvare, care se îmbină foarte bine cu schema MILD.

Schema MILD am aplicat-o manual, deși dacă s-ar defini un formalism mai riguros al ei, s-ar putea implementa și automat (am renunțat la aceasta pentru că era destul de complicat din punct de vedere al implementării). Schema este următoarea: - Se alege un număr (de WL-PHINT-uri) care să fie în jur de n , numărul minim de propoziții PHINT pentru construirea unui PGB, și se alocă unui modul. Vocabularul rezultat nu trebuie să depășească o valoare maxim impusă (în cazul experimentelor mele, 30). Dacă există și alte WL-PHINT-uri asociate aceluiași PGB rămase neintegrate, se grupează și se alocă altui modul nou (care va avea un vocabular ușor diferit). Pentru PGB-uri care au număr de WL-PHINT-uri mai mic decât n , se va face câte un modul pentru fiecare (chiar dacă e vorba de PGB-uri cu un singur WL-PHINT), dar încercând păstrarea unui număr cât mai mic de module. Voi argumenta mai târziu, în discuțiile legate de experimentele efectuate, de ce este bine ca numărul de module rezultate prin MILD să fie cât mai mic.

După această împărțire, se stabilesc vocabularele, iar în funcție de acestea structura și dimensiunea rețelelor, precum și numărul NP pentru fiecare modul și se transformă WL-PHINT-urile în T2P-uri. Există o relație între numărul de T2P-uri învățate și lungimea lor. Relația este de proporționalitate inversă, crescând lungimea și scăzând numărul de traiectorii cu o rație egală, timpul de antrenare este același (pentru o epocă - vezi capitolul 3). Pentru a avea un echilibru între efortul de antrenare al modulelor conexiuniste rezultate, este bine ca pentru module care au multe T2P-uri de învățat, lungimea lor să fie mai mică, iar pentru module cu puține traiectorii, lungimea să fie mai mare. Aceasta poate duce la traiectorii prea scurte pentru modulele cu multe T2P-uri în baza de antrenare, astfel încât rețeaua să nu poată fi antrenată până

la o eroare suficient de mică ca să poată discrimina PHINT-ul corect. Soluția ar fi împărțirea WL-PHINT-urilor alese pentru acest modul în două noi grupe, fiecare ducând la un modul mai mic (și ca număr de traiectorii, dar și ca vocabular) și ușor de antrenat, chiar dacă mărim valoarea NP pentru fiecare. Dar repet, este foarte important să păstrăm numărul de module conexiunite cât mai mic.

6.3.3. Echilibrarea efortului de antrenare al rețelelor neuronale

O soluție pentru asigurarea reușitei discriminării între WL-PHINT-urile antrenate, fără a le împărți excesiv între module este atașarea unui modul special, pe care l-am denumit **modul de prefață**. Acest lucru presupune adăugarea unui modul suplimentar, dar stilul de adăugare nu implică o divizare a vocabularului inițial al modului "mare". Voi ilustra ideea cu un exemplu. Să presupunem că avem un modul cu un număr mare de WL-PHINT-uri asociate (n traiectorii), care au lungimea traiectoriilor de antrenare P , obținută prin scurtarea traiectoriilor inițiale (să presupunem că acestea aveau aproximativ $2 * P$). Pentru a obținem discriminarea corectă, am putea să extindem traiectoriile cu câte un pas ($P+1$, $P+2$, ș.a.m.d.) și să repornim antrenarea din start. Din păcate numărul de pași suplimentari necesar este destul de mare atingând de obicei $2 * P$. Experimental am observat că dacă avem m traiectorii de lungime P discriminate corect, și adăugăm încă m traiectorii (în total $2 * m$) aveam nevoie de aproximativ $2 * P$ pași. Având un număr dublu de traiectorii și existând relația pătratică între lungimea lor și timpul de antrenare, efortul de antrenare crește de 8 ori. Acest calcul nu ia deloc în considerare creșterea absolut necesară a numărului de unități necesare pentru ca antrenarea să reușescă cu același succes ca la rețeaua cu m unități (creșterea capacității de învățare a rețelei). Având mai multă informație de introdus în rețea prin antrenare, creșterea dimensiunii rețelei este absolut necesară. Timpul de antrenare este în relație cubică cu numărul de unități (vezi capitolul 3), astfel încât numai dublarea numărului de unități duce la mărirea cu un factor de 8 a timpului de antrenare inițial. Deci, o dublare a numărului de traiectorii, o dublare a lungimii lor, și o dublare a numărului de unități duce la o mărire cu factorul 64 a timpului de antrenare (relativ la o singură epocă).

Este evident că o astfel de soluție nu este viabilă. Unul din scopurile mele minore în această fază a fost ca toate modulele conexiunite să poată fi antrenate aproximativ cu același necesar de resurse computaționale (care să fie relativ mic). Chiar dacă acest deziderat nu este atins, suma totală a timpilor de antrenare este bine să fie cât mai mică. Am arătat că acest timp total este minim atunci când avem același număr de unități distribuite fiecărui modul (vezi subcapitolul 5-4). Altă observație este că mărirea numărului de traiectorii care trebuie învățate nu duce neapărat la o creștere cu aceeași rată liniară a numărului de unități. O astfel de dependență există și între lungimea discretă a traiectoriilor și numărul de unități. Este normal ca în această situație să căutăm o soluție care reduce lungimea traiectoriilor și nu numărul de traiectorii. Reducerea numărului de traiectorii duce și la reducerea vocabularului (care ar duce și la reducerea numărului de unități) dar am afirmat că nu dorim acest lucru,

din motive care vor fi ilustrate ulterior. În plus, reducerea vocabularului în acest caz, poate fi destul de redusă, aproape ne semnificativă, dacă WL-PHINT-urile care se despart în două grupuri, au cam același vocabular (și de obicei așa este).

Dacă avem un grup de n traiectorii T2P (să luăm pentru claritate n număr par) de lungime $2 * P - 1$, putem împărți traiectoriile în două grupuri egale de $n/2$ T2P-uri. Aceasta ar fi o soluție "paralelă". O soluție "serială", este să împărțim toate traiectoriile în două segmente contigue, una de la momentul I la momentul P, și alta de la momentul P la momentul $2 * P - 1$. Important este ca punctul final unei traiectorii de la I la P să fie identic cu punctul inițial al unei traiectorii de la P la $2 * P - 1$. Vom obține două grupuri de n traiectorii, care pot fi conectate spațial două câte două prin perechile de puncte "ultimul-primul". Vom antrena o rețea recurentă cu primul grup de traiectorii cu pași de la I la P și altă rețea cu al doilea grup de traiectorii cu pași de la P la $2 * P - 1$. Punctele din momentele $2 * P - 1$ reprezintă PHINT-uri, iar pozițiile asociate momentului P le-am denumit LINK-uri. Modulul nou creat, l-am denumit, modul de prefață (cel de-al doilea l-am numit **principal**). Dacă un modul conexiunist recunoaște o linie de raționament care se desfășoară la nivel simbolic (a cărei propoziții sunt transmise pe BC) și este unul din modulele de prefață, ea va evolua pe traiectoria internă până la punctul asociat propoziției LINK. Această propoziție este scrisă pe BC, și va fi citită de modulul (conexiunist) principal, iar în cazul în care evoluția de la nivelul simbolic urmează linia WL-PHINT-ului care a dus la construirea traiectoriei urmată în prima sa jumătate de modulul de prefață, modulul principal va urma aceeași traiectorie, până la punctul asociat PHINT-ului. Aceasta este partea de "Linkage", din schema MILD, folosită exclusiv din motive computaționale, pentru a echilibra task-urile de antrenare a rețelelor. Ca avantaj al metodei, se observă că efortul computațional pentru antrenare este același pentru cele două rețele (și nu implică nici o modificare a numărului de unități ale unei rețele pentru n traiectorii de lungime P).

În aceste moduri, PGB-urile prea "aglomerate", pot fi împărțite pe mai multe module, fără a renunța la idea de a avea cât mai puține module conexiuniste. O regulă empirică este de a utiliza metoda împărțirii WL-PHINT-urilor (asociate unui PGB) în mod "paralel" dacă PGB-ul se află la un moment relativ "de început" (în comparație cu liniile de raționament cele mai lungi) sau de aplica metoda "Linkage", cu un modul suplimentar de prefață, dacă PGB-ul este "la sfârșit". Pentru că nu am lucrat decât cu sisteme de mică dimensiune, nu a fost necesară metoda "Linkage" desfășurată pe trei sau mai multe module, deși nu văd nici un impediment ca așa ceva să nu funcționeze la fel de bine ca și schema cu un singur modul de prefață. Sunt și PHINT-uri care nu pot fi grupate în PGB-uri; pentru acestea am construit module separate, care au fost mult mai ușor de antrenat decât modulele derivate din PGB-uri cu n traiectorii, dar numărul lor a fost în general mic și nu a influențat efortul de antrenare. Resursele computaționale necesare pentru toate aceste "lightweight modules", au fost sub nivelul resurselor cerute de pildă de un modul principal împreună cu modulul lui de prefață.

6.4. Experimente

Scopul acestei cercetări, după cum am mai afirmat anterior, este validarea unui model care propune o metodologie diferită pentru construirea unui sistem hibrid (prin integrarea rețelelor neurodinamice). O validare teoretică a acestui tip de model este extrem de dificilă și după unii autori chiar imposibilă [Touretzky91] [Smolensky92], dată fiind natura analogică, puternic distribuită a rețelelor neuronale, în contrast cu natura simbolică, bine determinată, locală și unitară a sistemelor logice cu care sunt puse în contact. Un excelent exemplu de model de raționament automat care nu are un suport demonstrațional teoretic este MYCIN [Shortliffe&Buchanan76], utilizat în foarte multe sisteme operaționale, care au certe avantaje. Demonstrația pe care am încercat-o pentru modelul meu, va fi utilitaristă, printr-o implementare de tip "toy system", și care a încercat pentru început doar să ilustreze că un astfel de sistem funcționează, iar în plus, să arate că funcționarea nu este posibilă în lipsa nivelului simbolic (lucru discutat și în capitolul anterior), prin ilustrarea rolului necesar al acestuia.

Ca pe un fel de demonstrație prin inducție, am conceput partea experimentală în două faze, similare demonstrării unei formule care are un parametru k variabil. Dacă presupunem că pentru n formula este corectă, iar pentru $k+1$ putem demonstra că ea rămâne corectă, atunci este corectă pentru orice k . Fazele experimentului meu erau:

- i. Construirea unui sistem minimal, care să arate că sistemul funcționează, după schema dorită.
- ii. Extinderea sistemului (cu un "ordin de mărime") și testarea lui în aceleași condiții ca și la pasul 1. Dacă sistemul funcționează, putem presupune că modelul este viabil. Acest tip de abordare se numește în literatura dedicată sistemelor de raționament automat cu incertitudini, "weak demonstration" [Todd91], și este în general acceptat în comunitatea cercetătorilor în domeniul sistemelor expert cu componente analogice.

Deși nu a fost unul din scopurile urmărite inițial, s-a relevat experimental și rolul util al proprietății de generalizare a componentelor conexiuniste folosite. La alt nivel și cu alte metode, aptitudinea de generalizare a unei rețele recurente a fost încercată și în capitolul 4. Dacă experimentul acela nu a reușit decât parțial, din motive de intensitate computațională, prin reducerea complexității problemei de rezolvat și adaptarea ei în cadrul unui sistem hibrid cu o puternică parte simbolică, am obținut un rezultat asemănător pe care îl consider mai valoros decât simpla funcționalitate a sistemului de probă. *O observație foarte importantă este că punerea în evidență a capacității de generalizare a întregului nivel conexiunist este proporțională cu gradul de acoperire al vocabularelor modulelor conexiuniste.* Pentru că am considerat fenomenul interesant am studiat această proprietate de generalizare în o a treia fază demonstrativă, care a dus la eliminarea cuvintelor individuale din aceste vocabulare și suprapunerea lor la un nivel mai mare decât cerințele funcționale ale modelului.

6.4.1. Faza sistemului minimal

Pentru a obține o simplificare esențială a experimentelor, am renunțat (ca și în experimentele descrise la capitolul 4) la o semantică definită prin cuvinte alese din limbajul natural. Am folosit simboluri fără o reprezentare în lumea reală sau abstractă, de tipul c1 pentru clasa 1 ș.a.m.d. Experiența acumulată în timpul lucrului cu sistemul monomodular mi-a accentuat convingerea că utilizarea unor propoziții certe din limbajul natural este o frână în de experimente care caută funcționalitatea sistemului. Folosirea unor propoziții cu sens bine determinat duce la necesitatea unei plauzibilități cu sens psihologic a rezultatelor obținute. Din cauza nuanțelor și a semanticii, vagi a limbajului natural este foarte greu ca să obținem rezultate edificatoare. Alt motiv foarte serios este limitarea vocabularului folosit. Este foarte greu, de construit o bază de cunoștințe cu sens, legată de o problematică reală, cu un vocabular constrâns la un număr mic de relații, clase și individualități. Realizarea unei aplicații bazate pe modelul propus de mine ar fi în primul rând o căutare a aplicației potrivite și a rezolvării problemei mapării între forma de exprimare a cunoștințelor și a structurilor simbolice folosite. Din păcate, distanța de la cunoștințele exprimate în limbaj natural până la grafulile conceptuale folosite este destul de mare.

O analiză critică a experimentelor efectuate în continuare (și aici intră și cele din capitolele 4 și 5) ar putea avea ca punct central *lipsa semanticii asociată simbolurilor folosite*. O argumentație în favoarea metodei mele ar fi următoarea: omul a inventat de-a lungul istoriei o mulțime de limbaje pentru comunicare și reprezentare a cunoștințelor sale, folosind simboluri și mijloace foarte diverse. Toate limbajele au ca și caracteristică comună, posibilitatea stabilirii de asocieri între simbolurile (cuvinte sau pictograme) folosite, și/sau stabilirea unor ierarhii și categorizări ale acestora. Nu văd de ce nu putem alege o cale inversă pentru definirea unor structuri de test, în sensul ca mai întâi să definim limbajul (vocabular, categorii, ierarhii, asocieri) și apoi să definim structuri, chiar dacă acestea nu au o interpretare într-un univers real (semantică). Am încercat să aloc simbolurilor cuvinte ale limbajului natural, dar rezultatele au fost neconcludente. Consider ca aceasta se datorează complexității foarte mari a limbajului natural și mai ales a lumii înconjurătoare, descrisă de cuvintele noastre. Putem să ne imaginăm însă o lume mai simplă, artificial creată, descrisă de cunoștințele folosite de mine în experimente.

Cel mai simplu sistem conform modelului multimodular hibrid eterogen putea fi obținut pornind de la un sistem compozițional eterogen cu două module. Punctul de plecare a fost construirea unei baze de cunoștințe care să poată fi folosită prin forward chaining. Vocabularul pe care l-am conceput a fost minimal, 15 cuvinte care erau:

- 3 clase (c1, c2, c3) {roles}
- 5 relații (p1 - p5) {predicates, links}
- 7 individualități (f1 - f7) {fillers}

Pentru a obține puncte de întrerupere, am adăugat un modul compus din reguli care puteau fi folosite prin backward chaining astfel încât cele două module puteau comunica prin blackboardul asociat nivelului simbolic pentru executarea unor linii de

raționament. Schema de lucru a sistemului era următoarea. Utilizatorul scria o propoziție pe blackboard-ul nivelului simbolic BC, aceasta operație având sensul epistemologic de "este această propoziție adevărată?; dacă da, atunci în ce condiții?". Modulul simbolic bazat pe backward-chaining (l-am notat MB) stabilea o propoziție (ar putea fi mai multe dar pentru simplitate am ales varianta cu lanțuri de raționament inițiale formate din noduri cu o singură propoziție) care avea sensul "propoziția de la intrare este adevărată, dacă și numai dacă propoziția nou generată este adevărată. Modulul bazat pe mecanismul de forward chaining (notat MF) prelua această propoziție și o interpreta conform regulilor sale, în sensul "dacă se pune întrebarea P atunci trebuie să se pună și întrebarea Q". Întrebarea Q este pusă pe blackboard-ul BS și este citită de modulul MB, ș.a.m.d.

Notând propozițiile cu literele mici ale alfabetului, se poate ilustra setul de reguli din modulele simbolice MF și MB. Pentru modulul MB o schemă inițială a regulilor poate fi:

$$\begin{aligned} e &\rightarrow a \\ f &\rightarrow b \\ g &\rightarrow c \end{aligned}$$

$$\begin{aligned} x &\rightarrow h \\ y &\rightarrow i \\ z &\rightarrow j \end{aligned}$$

iar pentru MF:

$$\begin{aligned} ?e &\rightarrow ?h \\ ?f &\rightarrow ?i \\ ?g &\rightarrow ?j \end{aligned}$$

Plecând de la întrebarea ?a, sau ?b, ?c, putem avea trei linii de raționament care se termină cu propozițiile x, y, sau respectiv z. Nu mai există altă informație care să poată continua lanțul de raționament. În mod intenționat, pentru a introduce puncte de întrerupere, am modificat regulile din MB în felul următor:

$$\begin{aligned} e &\rightarrow a \\ f &\rightarrow b \\ g &\rightarrow c \\ x &\rightarrow (h \wedge m) \\ y &\rightarrow (i \wedge n) \\ z &\rightarrow (j \wedge o) \end{aligned}$$

Se poate observa destul de ușor că în pasul al treilea al raționamentului compozițional este necesară o propoziție suplimentară (de la caz la caz, m, n sau o). Este acceptat

faptul că prezența mai multor propoziții pe blackboard înseamnă din punct de vedere logic o conjuncție. Dacă utilizatorul consideră că propozițiile **m**, **n** sau **o** sunt întotdeauna adevărate și le poate furniza ca și "hint", această situație conduce la posibilitatea rezolvării automate a acestui punct de blocare.

Apariția unui punct de starvation în același pas, pentru linii diferite de raționament, face ca să putem identifica un PGB (care conține doar trei WL-PHINT-uri, față de patru câte am impus, dar pentru acest sistem minimal, am acceptat și o definiție cu trei). Cele trei WL-PHINT-uri, de lungime 3, vor fi:

- (a, e, m)
- (b, f, n)
- (c, g, o)

Acestea se transformă în traiectorii T2P prin translatore tensoriale (atenție, propozițiile nu sunt atomice, au o structură predicativă, obținută prin vocabularul prezentat) și sunt folosite ca bază de învățare pentru un singur modulul conexiionist (notat MC). Acesta a avut un număr de 250 de unități din care 125 erau de ieșire și un număr de $250 \cdot (250 + 150)$ de ponderi (necesar de memorie - 0.8MBytes), o dimensiune destul de mică - antrenarea fiind obținută cu un efort de 1451 de epoci până la o limitare a erorii de 5%.

Un exemplu de funcționare este următorul: utilizatorul pune întrebarea **a**, aceasta este scrisă pe BS și pe BC. Modulul MB produce propoziția **e**, aceeași propoziție fiind predicată și de modulul MC. Propoziția **e** este consumată de MF și este ștearsă de pe BS și se produce **h** care este scrisă pe BS. Dar propoziția **e** este transmisă conform protocolului de comunicare între nivele și pe BC de unde este consumată de MC, acest modul producând la ieșire propoziția **m**, conform traiectoriei asociată WL-PHINT-ului (a, e, m). Dat fiind faptul că la nivelul simbolic apare un punct de întrerupere, informația de pe BC este scrisă pe BS, acesta conținând simultan propozițiile **h** și **m** care vor fi acceptate în același pas de MB care va putea genera răspunsul așteptat **x**. Pentru celelalte trei linii de raționament, procesul este identic.

Se poate afirma în acest moment că rolul modulului conexiionist ar putea fi cu ușurință jucat de un modul simbolic mult mai simplu. Se mai poate accentua și predictibilitatea comportării modulului conexiionist. Dar tocmai posibilitatea utilizării rețelei neuronale este probată, nu căutarea unui alte metode de "machine learning". Există metode care se pot aplica unor sisteme compoziționale [Peng96], dar sunt limitate datorită complexității combinatoriale. Acest sistem este funcțional, modulul conexiionist funcționând așa cum m-am așteptat, iar antrenarea rețelei a fost un task doar de ordinul orelor pe o stație de lucru obișnuită. Este adevărat că multe din operațiile de construire ale modulului conexiionist le-am făcut manual, dar automatizarea lor este o problemă destul de simplă (relativ la identificarea punctelor de starvation, a WL-PHINT-urilor, a stabilirii structurii rețelei și a antrenării). Practic, un sistem complet automat, ar constata punctele unde utilizatorul dă informație

sistemului, ar grupa un număr de linii de raționament care conduc în aceste puncte și ar antrena o rețea recurentă pentru fiecare din aceste grupe.

Problema cea mai dificilă ar fi distribuirea liniilor de raționament în funcție de vocabularele asociate lor, astfel încât grupurile să aibă un vocabular cât mai mic. Aceasta problemă de optimizare s-ar putea să fie foarte complicată, dat fiind faptul că apar multe cerințe (număr de traiectorii, lungimi, număr de unități minim), fiind principalul motiv pentru care nu am încercat să intru în detaliile de implementare ale unui subsistem automat de identificare / construire / antrenare a modulelor conexiuniste. Am considerat că efortul merită concentrat în altă direcție. Ca și observație la sistemul minimal, se impune foarte clar, ca modulul (în acest context, chiar nivelul) conexiunist să "cunoască" tot vocabularul sistemului simbolic.

6.4.2. Faza sistemului extins

În timpul experimentului anterior, din start, în mod intenționat, am introdus în sistem cauzele pentru apariția momentelor de starvation. În acest al doilea set de experimente, nu am mai prevăzut din etapa de proiectare care sunt punctele de întrerupere. Acesta este și abordarea naturală, un proiectant de baze de cunoștințe nefăcând altceva decât să adauge noi cunoștințe, verificând dacă acestea duc la inconsistențe, sau în cazul sistemelor compoziționale, la puncte de întrerupere. Folosind acest stil de construire a modulelor simbolice, este dificil de prevăzut toate punctele de întrerupere, mai ales dacă sistemul are un număr relativ mare de reguli. Din acest motiv, găsirea punctelor de întrerupere trebuie făcută exclusiv prin încercări, prin traversarea unor linii de raționament de către utilizator (care furnizează și "hint"-urile necesare).

Ar fi fost convenabilă crearea un subsistem automat de identificare a punctelor de întrerupere. Nu am folosit o astfel de strategie, pentru ca ar fi fost dificilă de implementat, ci pur și simplu am schițat un număr de linii de raționament, cunoscând pașii de început și de sfârșit, lăsând intenționat un număr de reguli incomplete la aproximativ mijlocul liniilor de raționament. În mod normal, identificare informației lipsă se face prin cunoașterea semanticii propozițiilor care formează linia de raționament. În cazul meu, unde propozițiile nu au o semantică (relativ la cuvinte din limbajul natural) definită, identificarea propozițiilor care sunt necesare ca "hint"-uri, se face studiind regulile din modulul care urmează să continue linia de raționament. Modul acesta de lucru este poate nenatural, dar repet, cele două soluții erau fie stabilirea intenționată la proiectare a tuturor punctelor de întrerupere (și am remarcat că este dificil de făcut pe măsură ce sistemul crește în dimensiune), fie prin asignarea unei semantici propozițiilor din sistem (dar m-aș fi îndepărtat de problema în sine).

Am dezvoltat o baza de cunoștințe formată din reguli predicative, care putea fi utilizată prin mecanismul de forward-chaining. Am folosit un vocabular de 43 de cuvinte, împărțite pe categorii în felul următor:

- 6 clase (c1 - c6)
- 15 relații (p1 - p15)
- 22 individualități (f1 - f22)

Regulile le-am separat în două module pur simbolice, criteriul de separare fiind vocabularul, astfel încât cele două vocabulare asociate modulelor să fie cât mai echilibrate ca număr de cuvinte. Numărul de reguli rezultate a fost de 30 pentru fiecare modul. Se poate afirma creșterea sistemului "cu un ordin de mărime", deși într-un context pur simbolic este dificil de definit "ordinul de mărime". Vocabularul a crescut de trei ori, iar numărul de reguli de aproximativ 10 ori. Numărul de module simbolice a crescut doar cu un modul, al treilea modul simbolic fiind un modul bazat pe mecanismul de backward-chaining (primele două fiind module bazate pe forward-chaining).

Ca și la sistemul minimal, linia de raționament era începută de modulul bazat pe backward-chaining (notat MB), care citea informația pusă pe BS de utilizator (și care era sub forma unei întrebări). Modulul MB genera o propoziție rezultat, care era scrisă pe MB. În acest moment, ar fi fost necesară existența unui secvențiator care să aleagă între cele două module bazate pe forward-chaining (notate MF1 și MF2). Pentru a simplifica gestiunea liniilor de raționament, am ales o schemă de alegere a modulului care preia informația de pe BS bazată pe vocabular, într-un stil asemănător cu metoda folosită la sistemul multimodular omogen (conexionist - vezi capitolul 5). Am asociat fiecărui modul (MF1 și MF2) câte un vocabular, iar propozițiile care puteau fi citite de un modul de pe blackboard, erau doar acelea care aveau cuvintele componente membre ale vocabularului modulului. În cazul în care propoziția putea fi redusă la ambele vocabulare, am acceptat ca ambele module să fie capabile să citească propoziția și să genereze câte un rezultat. Cazul favorabil era ca numai un modul să genereze răspuns, iar în cazul în care apăreau două propoziții, se punea problema unei conjuncții logice, care să modifice comportamentul modulului care trebuia să preia linia de raționament (considerat de mine a fi modulul MB). Important este din punctul meu de vedere că o astfel de construcție duce la un comportament nedeterminist.

Am adaptat și modificat regulile din cele trei module astfel încât să obțin aproximativ zece linii de raționament (de la o propoziție start la una țintă) de lungimi aproximativ egale (12-14 pași). Modulul MB putea să pornească și pe alte linii de raționament, dar acestea nu au fost prevăzute intenționat, ele fiind un potențial de lucru pe care l-am folosit ulterior (au fost detectate în final încă zece linii neintenționate gândite). Studiind cele zece linii intenționate pregătite, am detectat 7 puncte de întrerupere, care se apăreau întotdeauna după o scriere pe BS efectuată de modulul MB. Figura 6-2 ilustrează locul unde apăreau aceste puncte. Cele 7 "starvation points" se puteau grupa în două PGB-uri, datorită atât vocabularului, cât și lungimii diferite a liniilor de raționament care conduceau la aceste puncte. Primul PGB constă din:

WL-PHINT1 - de lungime 5

WL-PHINT2 - de lungime 6

WL-PHINT3 - de lungime 7

iar al doilea din următoarele patru WL-PHINT-uri (4, 5, 6 și 7) de lungimi 8, 8, 9 și

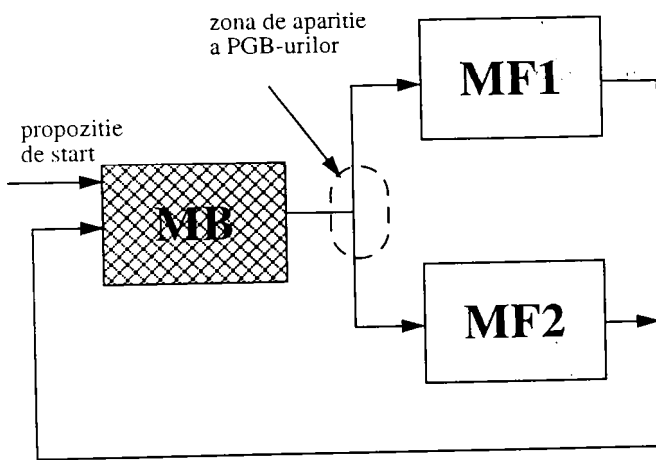


figura 6-2.

respectiv 10 pași. Conform schemei MILD de partiționare, am construit un modul conexiunist pentru învățarea celor trei traiectorii rezultate din WL-PHINT-urile 1, 2, și 3. Lungimea acestor traiectorii a fost limitată la 5. Acest modul conexiunist (notat MC1) a fost ușor de antrenat, dat fiind numărul mic de traiectorii. Pentru celelalte 4 traiectorii, de lungime limitată la 8, antrenarea devenea prea intensivă computațional, astfel încât am folosit un modul de prefață (MCP2), cu 4 traiectorii de lungime 4 și un modul principal (MC2), cu partea a doua a celor patru traiectorii, de lungime 4 (lungime totală - 8). Antrenarea celor trei module conexiuniste a fost echilibrată, dat fiind și echilibrul aproximativ între vocabularele folosite, 20 de cuvinte pentru MC1, 19 cuvinte pentru MC2 și 22 de cuvinte pentru MCP2. Dimensiunile rețelelor recurente folosite au fost:

- 300 cu 140 de unități de intrare și 140 de ieșire pentru MC1
- 280 cu 135 de unități de intrare și 135 de ieșire pentru MC2
- 300 cu 160 de unități de intrare și 160 de ieșire pentru MCP2

Am remarcat din nou că reuniunea vocabularelor modulelor conexiuniste reprezintă întregul vocabular de la nivelul simbolic.

Potrivirea acestor valori s-a dovedit o muncă destul de dificilă. Regula empirică care se poate aplica este: numărul de unități de intrare și de ieșire datorate vocabularului, înmulțit cu doi, dă numărul de unități al rețelei. În funcție de numărul de puncte care formează traiectoriile de antrenat, numărul de unități suplimentare (nededicat pentru ieșire) poate oscila. Dacă avem 4-5 pași regula aceasta duce la un număr de aproximativ 500 până la 1500 de epoci pentru fiecare rețea (de ordinul orelor pe o stație obișnuită). Dacă traiectoriile sunt mai lungi, numărul de unități suplimentare trebuie să crească, altfel numărul de epoci devenind mult mai mare (zeci de mii). Mai mult, creșterea numărului de unități și a lungimii traiectoriilor măresc simțitor (polinomial) timpul necesar pentru o epocă. Limita admisibilă consider că este maximum 30 de cuvinte/modul și lungimi de traiectorii de maximum 6 puncte (400 de unități). antrenarea pe o stație performantă putând să dureze zile și chiar săptămâni dacă numărul de epoci necesare atingerii unui minim de eroare de 5% este destul de mare.

Important a fost că am reușit să demonstrez că un astfel de sistem poate fi extins la un număr mai mare de module și pentru un vocabular mai consistent. Ceea ce a fost interesant și consider că este reușita majoră a acestui experiment este următorul fenomen. Descoperind noile linii de raționament, prin experimentarea funcțională a sistemului (ele ar fi putut fi descoperite și automat, funcționarea fiind nedeterministă doar din punctul de vedere al extinderii seturilor de reguli), am identificat și situații care puteau fi interpretate ca puncte de întrerupere. Mai exact, furnizând un "hint" aceste linii de raționament (relativ scurte, de 7-8 chiar 9 pași) puteau fi continuate prin pași posibili, până la o lungime aproximativ egală cu liniile construite intenționat (12-14 pași și chiar mai mult). Repet, aceste linii de raționament nu au o semantică precizată, relativ la limbajul natural, sau la realitate, iar sensul lor ca întreg nu poate fi precizat. Este însă plauzibil să afirm, că în cazul în care o astfel de linie de raționament este dezvoltată de un utilizator, ea să fie folositoare și dorită

(în cazul unui sistem cu semantică asociată simbolurilor utilizate). Evident că și afirmația contrară poate fi susținută, adică un "hint" nepotrivit să ducă la un rezultat nedorit. Este puțin credibil însă că două linii parțiale de raționament, din care una începe cu pasul de start și se termină la pasul t din motiv de starvation, putând fi continuată la pasul $t+1$ prin adăugarea unei singure propoziții lipsă, să nu fie de fapt o singură linie căreia îi lipsește o legătură.

Am identificat două astfel de linii incomplete (pe lângă cele 10 linii complete "neintenționate"). Nici nu m-aș fi gândit la astfel de linii incomplete dacă o comportare neașteptată a sistemului nu mi-ar fi atras atenția. Acest lucru s-a întâmplat în timpul în care căutam în mod empiric, "lucrând" cu sistemul, liniile de raționament introduse în mod nedeterminist. Această căutare se putea face în două moduri. Fie lucrând doar cu nivelul simbolic al sistemului, fie lucrând și cu nivelul conexiunist conectat, pentru a obține și liniile care au fost completate prin intervenția modulelor conexiuniste antrenate. *Lucrând cu întreg sistemul, am observat că una din cele două linii de raționament care nu fusese prevăzută inițial, era rezolvată prin intervenția unui modul conexiunist* (mai exact MC2) în pasul 6 (în pasul 2 al modului MC2). Mai exact, după pasul al șaselea al raționamentului, din motive de starvation, raționamentul era oprit. Pentru că am considerat că doar liniile cu minim 10 pași sunt complete, sistemul citea automat de pe BC pe BS, ultimele propoziții generate de modulele conexiuniste. Erau două astfel de propoziții, una generată de MC1, iar a doua de MC2. Dacă prima nu era "apropiată", nici ca vocabular, nici ca poziție în linia de raționament cu "hint"-ul necesar, a doua era chiar propoziția cerută pentru a continua raționamentul. A fost o întâmplare, dar acest lucru mi-a atras atenția asupra efectelor pe care le puteau induce modulele conexiuniste. Cu toate că am considerat că efectul constat era favorabil și îmbucurător, am căutat să găsesc și efectele nefavorabile ale intervenției prin "hint", oferită de nivelul conexiunist.

Efectul nefavorabil ar fi fost crearea unor linii de raționament "parazite" care să fie formate din mai multe linii parțiale (mai multe decât două), prin oferirea unui număr exagerat de mare de "hint"-uri care să se și potrivească. Căutând astfel de situații, nu am mai găsit decât o singură linie, care era incompletă doar într-un singur pas (pasul 8), deci putea fi considerată o linie favorabilă, care are nevoie de un singur hint. Mai mult, în cazul aceste linii incomplete, am descoperit că în pasul 6, pe BC apărea propoziția lipsă, care trebuia oferită ca hint în pasul 8. Din cauza protocolului de comunicare (cel mai simplu, notat cu unu roman în secțiunea 6.2.3) această propoziție era ștersă, nemaiputând fi utilizată în pașii ulteriori. Din acest motiv am încercat și alte protocoale de comunicare, cu riscul de a introduce linii parazite, prin aglomerarea de propoziții pe cele două blackboard-uri. Utilizând protocolul al doilea de comunicare, în care doar propozițiile consumate de pe BC erau șterse, propoziția respectivă rămânea pe BC și era folosită de sistem în momentul starvation point-ului. Dar nici acest lucru nu este suficient.

Al treilea protocol l-am folosit pentru a evita o situația în care propoziția care este utilă într-un pas ulterior, este consumată prematur de către modulele conexiuniste.

pentru că vocabularul acestora permite acest lucru. Nu acesta era cazul în exemplul de mai sus, pentru că doar MC2 putea consuma această propoziție, iar încercarea de a procesa această propoziție particulară în pasul respectiv, ducea la stagnare, adică rețeaua rămânea în același punct de pe traiectorie. Am considerat însă acest caz ca fiind foarte particular, probabilitatea ca propoziția utilă să fie consumată la nivel conexiunist și să dispară să fie foarte mare. Acest protocol, în ambele sale variante are însă dezavantajul că induce un potențial pentru apariția unor linii de raționament parazite prin păstrarea unui număr de propoziții mult prea mare pe cele două BB-uri. Însă la un sistem cu un număr atât de mic de module, relevanța acestui fenomen nu a putut fi detectată, nici o linie de raționament suplimentară nemaifiind detectată, chiar dacă se folosea cea mai complicată schemă de protocol (IIIb). Singura soluție pentru studiul acestui fenomen parazit, era construirea unui sistem și mai mare (mai ales ca număr de module). În plus, un astfel de sistem avea și avantajul că putea arăta dacă fenomenul favorabil de rezolvare automată a unor "hint"-uri nedeterminate inițial se păstrează și la o scară mai largă și reprezintă o posibilă caracteristică a acestui tip de model.

6.4.3. Faza de urmărire a generalizărilor

Pentru a studia efectul de generalizare constat anterior am apreciat că cea mai bună metodă este tot una experimentală. Metoda mea a fost de a utiliza un sistem mai mare, care dacă ar fi prezentat același comportament favorabil, într-o măsură egală cu creșterea lui relativă, ar fi reprezentat un argument necesar (dar nu și suficient - o abordare teoretică, prin formalizarea sistemului de mai sus este după părerea mea necesară, chiar dacă ar fi extrem de dificilă). Este însă destul de greu de apreciat ce înseamnă un sistem "mare", în raport cu cele două sisteme experimentale dezvoltate anterior. Pur și simplu am considerat că un sistem suficient de mare pentru studiul propus este sistemul cel mai voluminos care poate fi construit de către un singur experimentator, într-un timp limitat. Poate că ar fi necesare experimente cu un sistem și mai mare, dar aceasta presupune munca unui colectiv, timp, și mai multe resurse de calcul.

Prima dimensiune pe care am operat mărirea sistemului a fost cea a vocabularului. Numărul total de cuvinte a fost 120, împărțite pe categorii în felul următor:

9 clase,

43 de predicate,

68 de individualități.

Numărul de reguli care pot fi scrise cu un astfel de vocabular este foarte mare. A existat însă o limită la care să mă opresc, astfel încât baza completă de cunoștințe să fie lizibilă și coerentă (se pot genera reguli și aleator, dar probabilitatea de a avea linii de raționament este mult mai mică). Dacă aș fi păstrat aceiași reguli/modul ca la sistemul extins cu trei module simbolice, aș fi putut obține un număr mare de module, în jur de 20. Ar fi fost o creștere semnificativă a numărului de module, deci și o

creștere aparent semnificativă a sistemului. Deși informația înglobată sub formă de reguli era de aproximativ de zece ori mai mare, creșterea numărului de module, păstrând aproximativ constantă rata de reguli/modul, nu este o soluție prea bună, pentru că duce la fragmentarea exagerată a liniilor de raționament posibile, adică la apariția a prea multor puncte de întrerupere. Practic, am observat că folosind un număr mare de module, nici o linie de raționament nu este completă. Plauzibil este să considerăm că doar o mică parte din liniile de raționament posibile au linii de întrerupere (deși la sistemul minimal au fost trei linii și toate au fost întrerupte). Este greu de apreciat care este cea mai bună rată de **linii întregi/linii întrerupte** pentru sistemele "normale" (la sistemul extins, rata era de 11/9, dar era un sistem relativ mic).

Cu același număr de reguli, se puteau face diverse partiționări în module. Metoda de partiționare pe care am folosit-o a fost cea de a scrie toate regulile folosite prin forward-chaining într-un singur bloc (astfel încât să existe coerența necesară pentru realizarea unor linii de raționament) și apoi de a împărți aceste reguli între module care astfel încât regulile care erau înlănțuite să fie în module diferite (tocmai pentru a oferi posibilitatea de a produce linii de raționament compoziționale). La fel am procedat și cu regulile care erau evaluate prin backward-chaining, scriindu-le pe toate la un loc și apoi separându-le. Cu cât împărțeam regulile în mai multe module, cu atât obțineam mai puține linii complete de raționament. Acestea sunt necesare, pentru că există linii asemănătoare ca și structură a propozițiilor, care sunt incomplete, metoda de căutare prin similaritate facilitând identificarea lor. Această metodă nu va găsi toate liniile incomplete (considerând ca incompletă, orice linie din care lipsește doar o singură propoziție), dar este cel mai simplu mod de lucru. Dacă pentru un număr mare de module numărul de linii de raționament complete era zero, odata cu scăderea numărului de module, numărul de linii identificate complete era din ce în ce mai mare.

Cele mai bune rezultate le-am obținut pentru următoarea configurație: două module bazate pe backward-chaining și patru module bazate pe forward chaining; care au dus la identificarea a 55 de linii de raționament complete. Ce înseamnă în acest context o linie de raționament completă? Plecând de la sistemul inițial monobloc, putem identifica numărul de pași necesar pentru găsirea unui rezultat. Acest lucru se poate varia în funcție de "forma" grafului de decizie care poate fi construit pornind de la reguli. Cu același număr de reguli, putem obține linii de raționament lungi de 30-40 de pași (și atunci vor fi mai puține linii de raționament posibile) sau linii scurte (de 12-14 pași, ca și în cazul anterior - dar având multe linii de raționament). Pentru că posibilitatea de a identifica linii de raționament foarte lungi este mică, putând să existe prea multe puncte de întrerupere, am ales o variantă cu linii relativ scurte de raționament (14-16 pași). Ca schemă de funcționare, sistemul este prezentat în figura 6-3. Se poate argumenta că dimensiunea sistemului dată de lungimea liniilor de raționament nu a fost crescută proporțional cu celelate dimensiuni (vocabulary, număr de reguli, număr de module). Însă în acest caz sistemul ar fi fost foarte greu de controlat din punct de vedere al raționamentului însuși.

Pornind de la aceste linii de raționament am găsit un număr de 38 de linii de

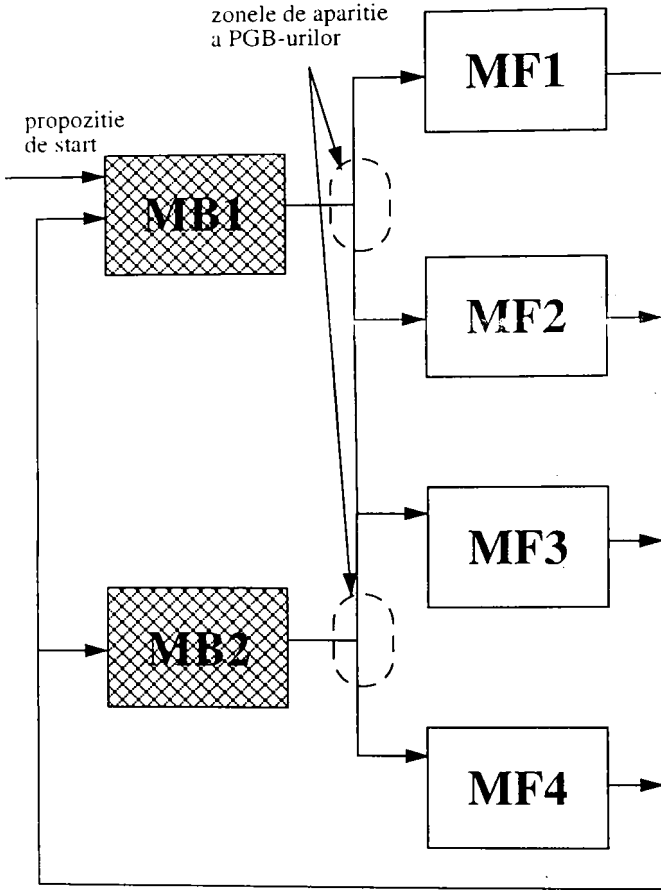


figura 6-3.

raționament incomplete (care completate cu o propoziție aveau lungimea de 14-16 pași). Se poate remarca că proporția de linii_întregi/linii_întrerupte este aproximativ aceeași ca și la sistemul anterior. Momentele predilecte în care apăreau întreruperile erau pașii 5, 6, 7, 8 și 9. Punctele în care apăreau situațiile de "starvation" erau înainte ca raționamentul să fie continuat de unul din modulele MF_i (vezi figura). Puteam să asociez punctele cu modulele "în fața" cărora apăreau, dar am decis gruparea lor să se facă doar pe baza vocabularului folosit. Pe baza acestor linii de raționament incomplete se puteau face traiectoriile de antrenare pentru nivelul conexiunist. Pentru că în această fază a experimentării am căutat să determin dacă generalizarea este o caracteristică a sistemului am folosit o metodă de "cross-validation". Am împărțit liniile de raționament întrerupte în două grupuri egale, urmărind ca și împărțirea după vocabular și lungime să fie proporțională. Pornind de la primul grup, am trasat 19 WL-PHINT-uri pe baza cărora am construit modulele nivelului conexiunist, folosind schema MILD.

Problema alocării acestor traiectorii la module conexiuniste s-a dovedit deosebit de dificilă, în acest punct fiind probabil necesară o metodă de optimizare, multicriterială, dar care ar reprezenta un subiect de cercetare separat. Totuși, până la acest nivel de complexitate, problema a mai putut fi rezolvată manual. Criteriile de partiționare ale celor 19 traiectorii erau vocabularele și lungimile lor. Prima împărțire am făcut-o pe baza vocabularelor, obținând patru grupe distincte, cu vocabulare destul de puțin apropiate, și care aveau un grad mic de acoperire.

i. primul grup era format din 8 WL-PHINT-uri, cu un vocabular total de 35 de cuvinte (foarte mare). Lungimea traiectoriilor era de 10 pași pentru 5 dintre traiectorii, una era de 8 pași, una de 7 și una de 4 pași. Aplicând schema MILD, am obținut pentru trei PGB-uri, centrate astfel:

- primul în pasul 10 (cu 5 traiectorii),
- al doilea în pasul 7 (cu 2 traiectorii),
- al treilea în pasul 4 (cu o traiectorie).

Am construit trei module conexiuniste MC1, MC2 și MCP1. MC1 și MCP1 (prefața lui MC1) erau antrenate cu cele 5 traiectorii de lungime 10, împărțite în semitraectorii de lungime 4 (pentru MCP1) și lungime 6 (pentru MC1). Pentru a nu face un modul separat pentru traiectoria cu lungimea 4 (oarecum izolată de restul), am antrenat această traiectorie în modulul MCP1. Cele două traiectorii de lungime 7, au fost baza de antrenare pentru MC2. Caracteristicile pentru aceste module sunt date în tabelul următor:

nume	număr traietorii	lungime traietorii	cardinalitate vocabulary	număr de unități (total + I/O)	necesar de memorie (MBytes)
MCP1	6	4	28	340 + 220	1.9
MC1	5	6	25	300 + 190	1.4
MC2	2	7	18	250 + 135	0.9

Efortul de antrenare a fost relativ dezechilibrat, pentru MCP1 fiind necesare câteva zile (aproximativ 34 de ore pentru 1200 de epoci) și trecerea prin câteva încercări nereușite din cauza unor bifurcații. Celelalte două module au putut fi antrenate cu prețul unei zile de rulare pentru fiecare. Se poate observa că apropierea de limita de 30 de cuvinte, duce la obținerea unor timpi de antrenare prohibitivi.

ii. al doilea grup era format din 7 WL-PHINT-uri, cu un vocabular foarte diferit de primul grup, format tot din 35 de cuvinte. Am împărțit aceste traiectorii în două PBG-uri. Primul conținea 4 traiectorii (de lungime 9, 9, 9 și 8) și am limitat lungimea T2P-urilor la 8 pași. Al doilea conținea 3 traiectorii (de lungime 5, 6 și 7) și am limitat lungimea T2P-urilor la 5 pași. Configurația modulelor construite pe baza acestui grup de WL-PHINT-uri este dată în tabelul următor:

nume	număr traietorii	lungime traietorii	cardinalitate vocabular	număr de unități (total + I/O)	necesar de memorie (MBytes)
MCP3	4	4	25	300 + 200	1.5
MC3	4	4	22	300 + 180	1.4
MC4	3	5	20	250 + 150	1.1

Efortul de antrenare a fost destul de mic, singurul modul care a avut nevoie de peste 24 de ore de rulare fiind MCP3. Am observat că păstrând un număr de unități suplimentare (diferența dintre numărul total și numărul de unități de I/O) mai mare decât în experimentele anterioare, numărul de epoci de antrenare necesar pentru scăderea erorii la 5% era destul de mic (340 pentru MC3 și 387 pentru MC4 față de valorile apropiate de 10^3 constatate anterior).

iii. mai rămâneau doua grupe, una formată din două WL-PHINT-uri de lungime 4 cu un vocabular foarte diferit de restul (22 de cuvinte) și cealaltă formată dintr-un WL-PHINT izolat de lungime 10 și deasemenea cu un vocabular foarte diferit de restul (15 cuvinte). Modulele rezultate au avut următoarele caracteristici:

nume	număr traietorii	lungime traietorii	cardinalitate vocabular	număr de unități (total + I/O)	necesar de memorie (MBytes)
MC5	2	4	22	300 + 180	1.4
MC6	1	10	15	300 + 125	1.2

Antrenarea a fost foarte simplă pentru aceste ultime module. Se poate observa că efortul este proporțional cu vocabularul, și mai puțin cu numărul de traiectorii. Lungimea conteaza destul de mult, dar doar în cazul în care numărul de traiectorii este

mai mare de 2.

Am obținut în final 8 module, a căror antrenare a fost totuși relativ echilibrată. Faptul că numărul de module conexioniste este aproape egal cu numărul de module simbolice este pură întâmplare și cred că nu există nici o formulă empirică care să definească vreo relație între aceste două numere. Doar vocabularul și factorul de acoperire, precum și structura traiectoriilor dă în final o schemă echilibrată, optimă pentru modulele conexioniste. Ceea ce trebuie remarcat însă este faptul că din vocabularul total de 120 de cuvinte, aceste module folosesc doar 88. Numărul însumat de cuvinte pe grupuri (35+35+22+15) este 107, dar există și un ușor factor de acoperire între vocabularele acestor grupuri. Factorul de acoperire este mult mai mare pentru vocabularele modulelor care aparțin aceluiași grup (MC5 și MC6 nu fac parte din același grup - au vocabulare diferite). Am afirmat anterior că este indicat ca nivelul conexionist să cunoască toate cuvintele din sistem. Și în acest caz, cunoașterea tuturor cuvintelor ar fi necesară, în cazul în care am fi antrenat mai multe module, pentru toate cele 38 de WL-PHINT-uri descoperite. Selectând numai jumătate din ele, am redus din păcate și vocabularul necesar (nu cu jumătate, dar semnificativ).

Căutând care din cele 19 traiectorii neantrenate (lăsate pentru "cross-validation") sunt completate cu "hint"-ul corespunzător la momentul potrivit, am aplicat toate cele patru tipuri de protocoale de comunicare între cele două nivele. Rezultatele prezentate au două aspecte: unul pozitiv, dat de numărul total de WL-PHINT-uri "rezolvate" (prezentat procentual, raportat la totalul de 19) și unul negativ, dat de numărul de linii de raționament care nu mai sunt terminate din cauza intervenției nedorite a modulelor conexioniste (un fenomen apropiat ca efect cu cel de "over-learning" de la rețelele feed-forward). Acest număr, este format atât din numărul de linii inițiale complete, și din numărul de WL-PHINT-uri "rezolvate" explicit prin antrenare. Am numit numărul de "rezolvări" favorabile, "**factor de generalizare**". Numărul de linii de raționament pur simbolic (în cazul curent, un total de 55 de linii descoperite) care sunt alterate de funcționarea nivelului conexionist este "**factorul de deteriorare la nivel simbolic**". Numărul de linii de raționament corectate prin nivelul conexionist (prin construirea explicită a modulelor - mă refer la cele 19 WL-PHINT-uri care au dus la antrenarea traiectoriilor), dar care nu sunt terminate din motivul funcționării defectuoase al acestuia l-am numit "**factor de deteriorare la nivel conexionist**". Există și un "**factor de deteriorare globală**" care le îmbină pe cele două anterioare, relativ la toate liniile de raționament care pot fi executate (mai puțin cele posibile prin WL-PHINT-urile neantrenate). Primele rezultate sunt ilustrate în tabelul următor, reflectând o trecere prin toate cele 55+19+19 traiectorii identificate:

protocol	generalizare (nr.traectorii, procent)	deteriorare la nivel simbolic	deteriorare la nivel conexiunist	deteriorare globală
I	2 (10.5%)	0	0	0
II	3 (15.7%)	0	4 (21.05%)	4 (5.4%)
IIIa	8 (42.1%)	12 (23.6%)	4 (21.05%)	16 (22.6%)
IIIb	6 (31.5%)	6 (10.9%)	3 (15.7%)	9 (12.1%)

Se poate observa că în cifre absolute cel mai bun protocol este primul (I). Aparent, cel mai bun este IIIa, care produce un factor de generalizare foarte mare, aproape jumătate din WL-PHINT-urile neantrenate ducând la traiectorii complete. Dar, o mare parte din traiectoriile inițiale (aproape un sfert) sunt pierdute. Câștigând 8 noi traiectorii, pierdem 16 (un minus de 8 traiectorii, față de cazul protocolul I, unde câștigam două traiectorii și nu pierdem nici una). Pentru protocolul II, pierdem în final o traiectorie, iar în cazul IIIb pierdem trei traiectorii.

Am căutat să gășesc metode care să îmbunătățească performanțele protocolului I și metode care să ducă la factori de deteriorare cât mai mici pentru protocoalele care dau generalizări bune. Au fost mai multe idei aplicate, dar *singura care a avut succes* a fost **mărirea vocabularului nivelului conexiunist**, până la numărul total de cuvinte din sistem. După mai multe încercări, am observat că așa se obține cea mai bună generalizare. O altă metodă, care a generat însă o creștere mică a fost cea a reducerii numărului de module și creșterea a gradului de acoperire. Aceasta a fost aplicată mai întâi și a condus apoi la ideea folosirii tuturor cuvintelor. Prima dată, am încercat să fac modificări relative numai la cele 19 WL-PHINT-uri pe baza cărora am construit primul set de module. Considerând că modulele MC5 și MC6 sunt prea dezechilibrate ca dimensiune și foarte depărtate ca vocabular de restul modulelor, cu prețul unor eforturi de antrenare mult mai mari (de ordinul sutelor de ore de rulare - adică săptămâni) am modificat structura modulelor în felul următor:

nume	număr traectorii	lungime traectorii	cardinalitate vocabulary	număr de unități (total + I/O)	necesar de memorie (MBytes)
MC1'	7	4	38	380 + 350	2.2
MC1	6	6	37	400 + 345	2.5
MC2	2	7	18	250 + 135	0.9
MCP3	4	4	25	300 + 200	1.5
MC3	4	4	22	300 + 180	1.4
MC4'	5	4	32	350 + 275	1.75

Cel mai greu de antrenat a fost modulul MCPI (care a primit prima jumătate din traiectoria dedicată în experimentele anterioare modulului MC6 și aproape întregul vocabular asociat), necesitând 225 de ore de rulare și mai multe ședințe din motive de bifurcații. Pot să spun că a fost cea mai frustrantă parte din toată munca desfășurată pe parcursul celor 4 ani de experimente. Totuși, în final, am obținut o eroare de sub 5%. Deși de o dimensiune mai mare și un timp_de_antrenare/epocă mai mare, MC1 (care avea acum în plus a doua jumătate a respectivei traiectorii) a putut fi antrenată mult mai ușor (109 ore și fără bifurcații). Cele două traiectorii de lungime 4 al modulului MC5 au fost alocate acum modulului MC4, prin trunchierea traiectoriilor acestuia de la 5 la 4 pași. Este important de observat aici că cele trei traiectorii inițiale au fost de lungimile 5, 6 și 7 pași. Voi nota cele trei traiectorii cu TFIVE, TSIX și TSEVEN, pentru că trunchierea lor este relevantă la nivelul experimentelor.

Am reluat experimentele cu noul nivel conexiunist format de data aceasta din cele 6 module, din care trei erau de dimensiuni mărite, ca vocabular și ca număr de traiectorii învățate (și ca unități și necesar de memorie, dar acest lucru nu este semnificativ la nivel funcțional). Am măsurat aceiași factori de performanță și deteriorare și am obținut următoarele rezultate:

protocol	generalizare (nr.traiectorii, procent)	deteriorare la nivel simbolic	deteriorare la nivel conexiunist	deteriorare globală
I	5 (26.3%)	0	0	0
II	3 (15.7%)	0	3 (15.7%)	3 (4.05%)
IIIa	10 (52.6%)	12 (23.6%)	3 (15.7%)	14 (20.2%)
IIIb	7 (36.8%)	7 (12.7%)	3 (15.7%)	10 (13.5%)

Rezultatele sunt mult mai bune pentru protocolul I (un plus de 5 linii de raționament - mai mult de un sfert din cele 19 lăsate pentru cross-validation), pentru protocolul II nu avem nici un raționament câștigat, iar pentru ultimele două avem o pierdere de 5, respectiv 3 linii. Este poate forțat să afirm că micșorarea numărului de module duce la o îmbunătățire pentru toate protocoalele. Totuși, pentru protocolul IIIb, se observă o mică creștere a generalizării (repet, la nivelul acesta de mărime al sistemului și la numărul de experimente diferite ca și configurație de module conexiuniste, este poate o întâmplare).

Significativ este însă că cele trei traiectorii care sunt deteriorate (sunt din cele 19 antrenate explicit) sunt TFIVE, TSIX și TSEVEN. TSIX și TSEVEN erau și în cazul anterior deteriorate, împreună cu două alte traiectorii de lungime inițială 9 (integrate în MCP3 și MC3), care au fost trunchiate la 8 pași. Se poate trage concluzia că trunchierea traiectoriilor, pentru a le putea antrena în module care învață și alte traiectorii mai scurte, este dăunătoare, ducând la deteriorarea aproape sigură a liniilor

de raționament asociate atunci când nu folosim protocolul I de comunicare. Același lucru l-am observat și făcând un alt experiment, integrând traiectoria modulului inițial MC6 nu în MCP1' și MC1' ci în MC2 (care avea două traiectorii de lungime 7). Am trunchiat traiectoria de lungime 10 la 7 pași, pentru a putea fi antrenată cu celelalte două și în plus a mai trebuit să măresc vocabularul MC3 la 32 de cuvinte. Antrenarea a fost mai dificilă (dar mai ușoară decât la MCP1' și MC1'), iar sistemul părea mult mai echilibrat. Rezultatul a fost că am obținut rezultate mai slabe decât în cazul prezentat (o generalizare de 15.7% pentru protocolul I), iar traiectoria trunchiată de la 10 la 7 apărea tot timpul la capitolul traiectorii deteriorate (pentru celelalte protocoale). Configurația și rezultatele acestui experiment sunt date în tabelele următoare:

nume	număr traiectorii	lungime traiectorii	cardinalitate vocabular	număr de unități (total + I/O)	necesar de memorie (MBytes)
MCP1	6	4	28	340 + 220	1.9
MC1	5	6	25	300 + 190	1.4
MC2'	3	7	32	350 + 240	2.0
MCP3	4	4	25	300 + 200	1.5
MC3	4	4	22	300 + 180	1.4
MC4'	5	4	32	350 + 275	1.75

(rezultate)

protocol	generalizare (nr.traiectorii, procent)	deteriorare la nivel simbolic	deteriorare la nivel conexiunist	deteriorare globală
I	3 (15.7%)	0	0	0
II	3 (15.7%)	0	4 (21.05%)	4 (5.4%)
IIIa	9 (47.3%)	12 (23.6%)	4 (21.05%)	16 (22.2%)
IIIb	6 (31.5%)	7 (12.7%)	4 (21.05%)	11 (14.8%)

Aceste rezultate indică totuși o ușoară creștere a factorului de generalizare pentru două din cele patru protocoale. Spre deosebire de setul de experimente anterior prezentat, în acest caz factorul de acoperire al vocabularelor este mai mic, pentru că în loc să distribuim vocabularul modulului MC6 în două module (MCP1' și MC1'), l-am reunit doar cu MC2, construind noul modul MC2'. Făcând aceeași operație de distribuire a vocabularului peste două module și nu unul singur, rezultatele sunt sensibil mai bune. Observația mea a fost că e mai bine să am module care să aibă vocabulare cât mai mari, iar intersecția lor (factorul de acoperire) să fie deasemenea cât mai mare.

Altă observație a fost relativă la numărul de module folosite. Un experiment care nu a căutat să obțină performanțe, a fost antrenarea fiecărei traiectorii cu un modul separat (au rezultat 19 module foarte simple, cu vocabulare izolate). Era normal ca generalizarea să fie zero, dar am constatat o creștere semnificativă a factorului de deteriorare global. Am grupat apoi traiectoriile două câte două (obținând 10 module), iar factorul de deteriorare a fost mai mic. Pentru experimentele cu 8 și apoi cu 6 module s-a observat o scădere a aproape liniară a factorului de deteriorare, după cum se observă și în tabelul următor, care ilustrează evoluția factorului de deteriorare globală odată cu scăderea numărului de module:

protocol	19 module	10 module	8 module	6 module
II	17.5%	8.1%	5.4%	4.05%
IIIa	44.5%	25.6%	22.6%	20.2%
IIIb	29.1%	20.2%	12.1%	13.5%

Am încercat să găsesc o explicație a acestei tendințe a factorului de deteriorare prin urmărirea liniilor de raționament în cele 4 cazuri de mai sus. Ca și la sistemul multimodular omogen, partea conexionistă a acestui sistem are nevoie de "sincronicitate", modulele conexioniste trebuind să urmărească niște linii de raționament pe care dacă le "pierd" și pornesc pe alte linii, conduc la o funcționare nedorită. Cu cât sistemul are mai multe module conexioniste, cu atât este mai greu de asigurat sincronicitatea, posibilitatea ca unele module să accepte propoziții la momente nedorite și să genereze propoziții nedorite fiind mai mare. Cu cât numărul de module este mai mic, cu atât e mai puțin posibil ca acestea să fie "deviate" de la traiectoriile pe care le-au început. Este însă adevărat că o dată cu creșterea gradului de acoperire al vocabularelor, perturbarea care poate să fie făcută de un modul către altul (prin transmiterea unei propoziții nedorite în pasul respectiv) este mai ușor să apară.

6.4.4. O partiționare bazată pe întregul vocabular al sistemului

Urmărind manual liniile de raționament, am observat ca mai mult de jumătate din cele 19 WL-PHINT-uri neantrenate ca traiectorii la nivel conexionist (T2P-uri), nu puteau să constituie corecții ale liniilor de raționament incomplete (la nivel simbolic) din cauza faptului că vocabularul lor era diferit de cel al WL-PHINT-urilor învățate. Am arătat în subcapitolul anterior care au fost criteriile de alegere și partiționare a WL-PHINT-urilor propuse pentru antrenare. Experimentele au arătat că este bine ca numărul de module să fie mic, vocabularele să fie cât mai apropiate (cu un grad mare de acoperire), iar traiectoriile să fie cât mai puțin trunchiate (trunchierea poate duce în plus și la reducerea vocabularului). Dar cele 19 traiectorii au fost mereu aceleași. Am mai făcut însă un experiment, cu alt set de WL-PHINT-uri, alese din setul inițial de 38, astfel încât vocabularul obținut prin reuniunea acestor WL-PHINT-uri, să fie exact vocabularul complet al sistemului (120 de cuvinte). Acest lucru s-a putut obține, fără

a include în setul de antrenare toate cele 38 de WL-PHINT-uri; au fost necesare doar 25. Puteau fi și mai puține pentru a obține vocabularul, dar am ales traiectoriile în așa fel încât lungimea maximă a T2P-urilor obținute să fie de 6 (cu două excepții, de lungime 8), iar trunchierile să fie pe cât posibil evitate.

De această dată, am încercat să grupez WL-PHINT-urile doar în funcție de lungime, nu în funcție de vocabular (avantajul experimentului anterior era că în cazul acesta vocabularul era minimal). Chiar dimpotrivă, am ales pentru același modul WL-PHINT-uri cu vocabulare cât mai diferite, pentru a avea prin reuniune, un vocabular cât mai mare. Folosind această metodă, obțineam cel mai mare grad de acoperire între vocabulare. Ținând cont și de efortul computațional necesar pentru antrenare (nu era recomandat ca modulele să aibă un vocabular mai mare de 30 de cuvinte) și de echilibrarea acestui efort între module, configurația modulelor a fost următoarea:

nume	număr traietorii	lungime traietorii	cardinalitate vocabular	număr de unități (total + I/O)	necesar de memorie (MBytes)
MC1	2	8	24	340 + 180	1.7
MC2	3	6	27	300 + 190	1.3
MC3	3	5	27	270 + 190	1.2
MC4	3	5	22	250 + 150	0.9
MC5	3	4	25	230 + 190	0.8
MC6	3	6	29	330 + 220	1.8
MC7	3	4	24	270 + 180	1.1
MC8	3	4	27	280 + 200	1.3
MC9	3	4	28	250 + 200	1.1

Se poate observa un număr mic de traiectorii/modul, relativ la un vocabular/modul destul de bogat. Însușind numărul de cuvinte din vocabularele tuturor modulelor, obținem un total de 223 de cuvinte, care raportat la numărul total de 120, dă un factor de acoperire global de 1.86, mai mare decât la încercările anterioare (unde a fost de aproximativ 1.4). Raportarea se face de această dată față de numărul total al cuvintelor, față de cazul anterior, unde se făcea față de cele 88 de cuvinte care compuneau vocabularul reunit al celor 19 WL-PHINT-uri folosite ca bază de învățare. Pentru a vedea dacă acoperirea este și locală trebuie facute și calculele pentru cele C_9^2 perechi de module. Factorul de acoperire între două module este egal cu raportul dintre numărul total de cuvinte din cele două module și numărul de cuvinte din mulțimea reunită a vocabularelor. Dacă acoperirea este nulă, acest factor este egal cu 1; iar în caz contrar crește până la maximum 2 (când acoperirea este totală, modulele având vocabulare identice). Acoperirea ideală am considerat-o ca fiind egală cu 1.5 (media

celor două extreme). Făcând calculele pentru toate perechile de module (28 de perechi posibile), am obținut valori apropiate de 1.5, cu câteva excepții (1.31 pentru MC4 cu MC5 - cel mai mic grad de acoperire; 1.77 pentru MC3 și MC8 - cel mai mare grad de acoperire). Trebuie observat că factorul de acoperire global și cel local sunt diferite ca și exprimare matematică, cel global variind între 1.0 și un număr real egal cu numărul de module (atunci când toate modulele cunosc toate cuvintele), iar factorul local între 1.0 și 2.0.

Din cauza faptului că numărul de traiectorii/modul este mic, antrenarea rețelelor recurente a fost relativ ușoară, chiar dacă dimensiunea rețelelor a fost mare, datorată vocabularelor consistente (mult peste limita propusă inițial de 22 cuvinte). Din totalul de linii de raționament incomplete (38), 25 au fost antrenate, iar pe restul de 13 urma să se facă testul legat de factorul de generalizare. Rezultatele, folosind toate cele patru protocoale de comunicare, sunt date în tabelul următor:

protocol	generalizare (nr.traiectorii, procent - din total 13)	deteriorare la nivel simbolic (din total 55)	deteriorare la nivel conexiunist (din total 25 antrenate)	deteriorare globală (din total 80 de linii)
I	7 (54.8%)	0	0	0
II	8 (61.5%)	2 (3.6%)	7 (28%)	9 (11.2%)
IIIa	12 (92.3%)	15 (27.3%)	11 (44%)	26 (32.5%)
IIIb	9 (69.2%)	10 (10.5%)	9 (36%)	19 (23.7%)

Cu toate că factorul de deteriorare crește foarte mult pentru protocoalele II, IIIa și IIIb, rezultatele sunt foarte bune relativ la factorul de generalizare (pentru toate protocoalele, mai ales la IIIb). **Factorul de deteriorare este mare pentru că numărul de module este prea mare.** Câștigând multe linii de raționament, pierdem însă un număr și mai mare, singura excepție fiind protocolul cel mai simplu (I). Remarcabil este rezultatul obținut pentru protocolul IIIa, unde aproape toate (12 din 13) liniile de raționament incomplete sunt rezolvate de nivelul conexiunist. Dar în același timp, deteriorarea globală este cea mai mare (aproximativ o treime din liniile complete inițiale sunt pierdute).

6.5. Relevanța experimentelor

După efectuarea acestor experimente, am tras o serie de concluzii, care pot fi utile pentru cel care dezvoltă un astfel de sistem:

- **soluția ideală** este ca nivelul conexiunist să fie compus dintr-un singur modul, care să cunoască toate cuvintele din vocabularul nivelului simbolic și să fie antrenat cu cât

mai multe traiectorii identificate ca și WL-PHINT-uri. În acest mod se poate obține o perfectă sincronicitate cu nivelul simbolic iar factorul de deteriorare la nivel conexiunist poate fi redus la 0. Am presupus că în cazul în care la un pas de raționament avem o singură propoziție rezultată la nivel conexiunist și folosim unul din protocoalele de tip III, posibilitatea ca aceasta să perturbe procesul de raționament este foarte mică.

- pentru ca **nu există niciodată soluții ideale**, proiectantul este constrâns să dezvolte mai multe module conexiuniste. Relația combinatorială între vocabular și dimensiunea rețelelor, complexitatea polinomială de grad mare a algoritmilor de antrenare impun acest lucru. Numărul de module trebuie să fie însă cât mai mic, dus până la limita dimensională a modulelor și a timpilor lor de învățare.
- **vocabularele** trebuie să aibă un grad de acoperire locală (două câte două) cât mai aproape de 1.5. Acest lucru va facilita creșterea factorului de generalizare. Am observat că sistemul funcționează cel mai bine când vocabularul de la nivelul conexiunist este identic cu vocabularul de la nivel simbolic.
- **protocolul de comunicare** între cele două nivele trebuie să fie cât mai simplu (ca și protocolul I). Protocoalele care implică un grad mai mare de transmitere a informației și păstrarea propozițiilor pe blackboard-uri timp de mai mulți pași, conduc de regulă la degradarea unor linii de raționament care în cazul protocoalelor simple nu erau afectate. Chiar dacă factorul de generalizare crește spectaculos, trebuie întotdeauna verificat dacă numărul de linii de raționament "câștigate" nu este mai mic decât numărul de linii de raționament "pierdute".

Experimentele prezentate mai sus sunt pe departe exhaustive. Din cauza naturii excesive a antrenării rețelelor recurente este foarte dificil a se face în timp rezonabil, de către un singur experimentator, fără resurse deosebite (un accelerator matematic, ar fi indicat) un număr suficient de mare de experimente. Se poate observa că sunt mai multe dimensiuni ale problemei:

- număr de module simbolice
- număr de module conexiuniste
- vocabular la nivel simbolic
- vocabular la nivel conexiunist
- număr de linii de raționament
- număr de WL-PHINT-uri identificate
- număr de WL-PHINT-uri antrenate
- vocabularele asociate grupelor de WL-PHINT-uri
- gradul de acoperire global și cel local al vocabularelor modulelor conexiuniste.

Toate acestea ar trebui investigate sistematic, prin variația uneia și fixarea celorlalte pe diverse valori. Se poate observa însă că volumul de experimente este foarte mare. Am încercat să variez acele dimensiuni care mi s-au părut esențiale (vocabular, număr de

module conexioniste, grad de acoperire) făcând totuși un număr cât mai mare posibil de experimente.

Un aspect important care ar urma să fie investigat este metoda de a comunica între modulele conexioniste prin vectori de activare și nu prin structuri simbolice. Și din punct de vedere al plauzibilității cognitive și al analogiei cu creierul uman, acest mod de implementare ar fi fost mult mai potrivit. Metoda de a transla în fiecare pas de raționament, în timpul transmiterii informației de la modul la modul, informația din starea ei neuronală, distribuită, în structuri simbolice, este o metode foarte rigidă. Cele mai moderne teorii (vezi [McLelland&O'Reilly94], [Fahlman97]) susțin că cea mai bună modelare a creierului uman este prin utilizarea unor rețele recurente conectate prin trunchiuri de transmitere a informației cu dimensiune mult mai mică decât a rețelelor (bottleneck) unde informația poate fi la nevoie restabilită în formă simbolică. Realizarea acestor bottleneck prin șiruri de așteptare care conțin grafuri conceptuale, este cel puțin artificială. Ar fi necesar aplicarea conceptului de "comunicare laterală" introdus în schemele ierarhice multi-rețele de [Miiikulainen96]. Comunicarea între modulele conexioniste s-ar realiza prin legături tipic conexioniste, existând și posibilitatea introducerii unor matrici de ponderi, cu o dimensiune mică, care să indice gradul de legătură și influență laterală. Acesta ar fi evident în raport direct cu gradul de similaritate între vocabulare. Problema cea mai spinoasă este firește modul în care se vor calcula aceste ponderi. Blackboardul BS trebuie să rămână, ca un segment al punții de legătură între cele două nivele, dar nu ar mai avea un rol în comunicarea între modulele conexioniste.

Firește, asocierea unor cuvinte din limbajul natural ar fi cel mai mare câștig. Dar soluția pentru introducerea unei semantici "naturale" este lărgirea vocabularului. Dar am văzut că aceasta duce la o modularizare excesivă, ideal fiind de fapt utilizarea unei singure rețele, sau mai multe, dar legate prin comunicarea laterală (formând de fapt una singură, cu o topologie mai deosebită). Ca de obicei, rezolvarea oricărei probleme tehnice presupune un "trade-off", câștigul într-o parte însemnând pierdere în alta.

Concluzii

C1. Un sumar al rezultatelor

Teza de doctorat prezintă o metodologie de îmbunătățire a performanțelor calitative în sistemele de raționament compoziționale eterogene. Motivele pentru care am propus această arhitectură specifică sunt prezentate dealungul celor 6 capitole. Un scop mai general, aflat deasupra scopului general pragmatic, cu un caracter aplicativ de dezvoltare strictă a metodologiei, este investigarea unor modele de integrare simbolic-conexionistă (neurosymbolic integration).

Metodologia se poate aplica la orice sistem compozițional eterogen fără supervizare, la care problema întreruperii liniilor de raționament datorată "starvation"-ului se rezolvă prin intervenția manuală a proiectantului. Soluția oferită de mine este de a adăuga un subsistem, numit nivel conexionist, compus din unul sau mai multe module bazate pe rețele neurodinamice discrete. Aceste rețele pot fi antrenate cu traiectorii formate din puncte în spațiul euclidian, care sunt analoge printr-o transformare bijectivă, cu o linie de raționament a sistemului compozițional, formată dintr-o secvență ordonată de structuri simbolice.

Pașii pentru utilizarea metodologiei pot fi rezumați succint în felul următor:

- Sistemul compozițional pur simbolic este testat de către utilizator și proiectant, care furnizează "hint"-urile necesare deblocării raționamentului pentru anumite linii specifice unde apar întreruperi datorate fenomenului de "starvation". Cu cât sistemul este mai exhaustiv testat, cu atât posibilitatea de a lăsa puncte de blocare nerezolvate este mai mică, dar este imposibil de a le rezolva pe toate.

- Se grupează toate liniile de raționament identificate în primul pas, în funcție de lungimea lor și în funcție de gradul de acoperire al vocabularelor. Se caută ca acesta din urmă să fie cât mai mare.
- Se împart liniile astfel încât să rezulte un număr cât mai mic de grupe, dar din punct de vedere al antrenării, niciuna din grupe să nu ducă prin lungime sau vocabular la module conexiuniste care să conțină rețele cu un necesar computațional exagerat. Dacă este necesar, traiectoriile prea lungi pot fi alocate pe două module consecutive (modul de prefață + modul principal).
- Se trasează pentru fiecare linie de raționament, o traiectorie discretă în spațiul euclidian, astfel încât acesta să fie netedă și centrată, iar punctele sale să fie în relație corectă cu punctele altor traiectorii similare.
- Se antrenează rețelele modulelor conexiuniste, în următoarea ordine: prima dată FFO cu algoritmul Cascor, se stabilește apoi mărimea stratului ascuns pentru FFI (aceiași ca la FFO) și se antrenează rețeaua neurodinamică.
- Se atașează nivelul conexiunist sistemului compozițional simbolic și se testează în ordine, cele patru protocoale de comunicare (în ordinea dată în teză: I, II, IIIa, IIIb) și se evaluează experimental factorul de generalizare și cel de deteriorare. Se alege în funcție de acestea, cel mai favorabil protocol.

O cale necesară de extindere a metodologiei este automatizarea procesului de identificare a traiectoriilor cu puncte de blocare și o partajare asistată a acestora pe modulele, astfel încât să existe o optimizare a procesului de antrenare. Acest aspect devine și mai important dacă se pune problema ca nivelele să poată fi extinse prin adăugarea unor noi module. Adăugarea de noi cunoștințe în sistem prin atașarea unor noi module la nivel simbolic, va duce foarte probabil la apariția unor noi linii de raționament cu puncte de blocare. Ar fi de dorit, ca atașarea unor noi module la nivel conexiunist care să rezolve aceste puncte de blocare să fie cât mai automatizată, existând și posibilitatea ca unele linii noi care trebuiesc antrenate să fie atașate unor module existente, care vor fi reantrenate. Această posibilitate apare firește doar în cazul în care modulul respectiv mai are capacități de antrenare nefolosite, iar vocabularul atașat este suficient pentru linia respectivă.

Ca și o concluzie mai generală pot spune că cercetarea pentru stabilirea metodologiilor potrivite pentru implementarea sistemelor cognitive este extrem de laborioasă, fapt datorat în primul rând interacțiunii software-mediului de lucru și de interdependența unor aspecte legate de metodologia în sine. Aici nu s-a dovedit deloc fructuoasă izolarea componentelor problemei și cercetarea concentrată pe subprobleme. Deși am folosit rezultate notabile ale cercetării și cercetarea concentrată pe subprobleme și din domeniul conexiunist, cred că nici unul din cele două domenii, în mod separat, nu

vor putea realiza "Modelul", implementabil în software, pentru un sistem cognitiv.

Metodologia propusă de mine pentru implementarea sistemelor cognitive este un exemplu de abordare inginerescă, constructivistă și empirică (la modul non-trivial). Dintr-un anumit punct de vedere, metodologia propune o *ipoteză științifică*. Mai exact: Se pot simula procese cognitive robuste prin cooperarea a două nivele de raționament complet diferite ca natură? (simbolic și conexionist). Dată fiind complexitatea covârșitoare a subiectului, este extrem de dificil, atât experimental cât și teoretic, de a demonstra dacă ipoteza este falsă sau nu.

Rezultatul cel mai important obținut este descoperirea că linii de raționament cu puncte de blocare neantrenate, dar potențial posibile, apar la funcționare ca și confirmare a puterii de generalizare a grupului de module conexioniste utilizate. Interpretarea acestui rezultat este dificilă și chiar riscantă, dimensiunea sistemului folosit la experimente nefiind suficient de mare pentru a se face supoziții legate de eventualele performanțe ale unui "real-life system", sau de plauzibilitatea psihologică, atât de căutată în cercetarea din domeniul integrării neuro-simbolice. Utilitatea acestui rezultat va fi găsită mai mult ca sigur în continuarea dezvoltării metodologiilor de acest tip, ca și fenomen relevant și analizat pe larg.

Consider că alte două rezultate sunt importante:

- demonstrarea experimentală a generalizării traiectoriilor simbolice (subcapitolul 4.4.3.), împreună cu schema de evitare a bifurcațiilor (subcapitolul 4.5.3.)
- detectarea artificialității distribuției nivelului conexionist în module separate.

Deasemenea, este important de știut că se pot executa raționamente robuste cu o rețea neurodinamică. Întrebarea este: de ce nu s-a încercat până acum și de ce nu se aplică metoda pe larg? (deși o demonstrație teoretică a acestei posibilități - pentru rețele liniare - există din 1986, vezi [Smolensky86]). Răspunsul: complexitatea computațională a antrenării unei rețele recurente (sistem dinamic neliniar - nu există algoritmi de antrenare pentru rețele recurente linare) este foarte mare, în special pentru un sistem viabil pentru piață. Dar la fel de important este de remarcat că algoritmi folosiți nu sunt intractabili, o estimare personală, pe baza experienței dobândite, fiind de 10^{20} , iterații necesare pentru construirea unor sisteme comerciale bazate pe această metodologie. În acest caz, întrebarea devine: când vom avea la dispoziție, cu un preț acceptabil, o putere de calcul suficientă (eventual dedicată pentru taskul de antrenare). Probabil că într-un viitor destul de apropiat. Tot ca o estimare rezultată în urma antrenării rețelelor, ordinul de mărime actualmente fezabil pentru aplicarea metodologiei este de 10^{16} , în condițiile unei dotări modeste (stații de lucru la nivelul anului 1995). Tot din considerente computaționale, metoda distribuției nivelului conexionist este ad-hoc și a fost necesară pentru a se realiza experimentele

practice, imposibile în cazul unei rețele de mari dimensiuni.

C2. Contribuții originale

Cea mai elocventă trecere în revistă a contribuțiilor originale este o lectură a celor 9 articole publicate în perioada 1995-1997, care sunt trecute în bibliografia tezei. Alte 5 lucrări prezente în bibliografie, din perioada imediat premergătoare (1993-1994) prezintă deasemenea idei originale, dar nu sunt legate de conținutul tezei decât prin unele metodologii folosite, având un aspect mult mai aplicativ. Trei din lucrări [Szirbik&a195], [Szirbik95b] [Szirbik97] sunt publicate în străinătate și au fost recenzate conform unei proceduri stricte, impuse de editurile respective. Recenziile indică o notă bună pentru factorul de originalitate, fapt semnalat de toți referenții acestor lucrări.

Ideea în sine, de a folosi rețele neurodinamice cu antrenare supervizată pentru raționamente simbolice robuste, este o idee totalmente originală și nu am găsit-o în literatura de specialitate până la nivelul anului 1997. Din discuțiile personale avute cu Prof. Scott Fahlman, de la Carnegie-Mellon University (Pittsburgh - Pennsylvania), unul din inițiatorii domeniului integrării neuro-simbolice în anii 80, am realizat că drumul pe care am mers a fost unic, fără precedent în istoria (foarte scurtă de altfel) a domeniului. Din alte discuții, cu Prof. Lee Giles de la Princeton University și Assoc. Prof. Riisto Miikulainen de la University of Texas (Austin), am aflat că au fost o serie de tentative nepublicate de a folosi rețele recurente cu reprezentarea locală a simbolurilor, dar s-a reușit doar simularea unor automate finite (cu un cost prohibitiv însă).

Pe capitole, mai pot fi evidențiate unele aspecte de originalitate:

- în capitolul de introducere, se face o taxonomizare a sistemelor și modelelor folosite în ultimii ani în domeniul integrării simbolic-conexioniste. Deși informația primară este culeasă din volumele unor conferințe indicate în bibliografie, modul în care am realizat clasificarea sistemelor hibride, este personal și diferit de cel propus de Melanie Hilario în 1996 [Hilario96].
- în capitolul 1, prezentarea unei rețele Towell ca o rețea neurodinamică fără recurență, cu o posibilă traiectorie trasabilă în spațiul stărilor dată de parcurgerea pașilor de raționament succesiv de straturile rețelei.
- în capitolul 2: demonstrarea formulei de antrenare pentru rețele neurodinamice continue, este făcută prin calcul matriceal și nu integral (deși în acest mod demonstrația este mai directă). Avantajul este apariția în clar a calculului matricei inverse, care ilustrează complexitatea task-ului de antrenare. La rețeaua discretă, o contribuție personală este

adăugarea pe stratul de intrare a unor unități speciale de "fan-out", doar pentru intrare, aceasta de obicei fiind furnizată direct unor unități din stratul de calcul (fără a mai fi ponderată). Avantajul a fost o antrenare mai rapidă în prima fază (de la o eroare de 25% la 10%), iar explicația este dată de reprezentarea mai facilă a unor dependențe liniare între intrare și ieșire prin aceste ponderi suplimentare.

- în capitolul 3: asocierea grafurilor conceptuale cu reprezentări conexioniste distribuite; metoda de translatare tensorială.

- începând cu capitolul 4, abordarea mea este diferită de toate abordările anterioare sau curente. Practic, toate ideile ilustrate sunt originale. În plus, metoda de paralelizare a algoritmului de antrenare, înregistrată ca și raport de cercetare la Vrije Universiteit din Amsterdam [Szirbik95a], este originală, contrazicând unele afirmații că accelerarea antrenării neuronale pe sisteme distribuite cu cuplare redusă (gen ZOO) este neperformantă.

C3. Posibile direcții de continuare a cercetării

Există două stadii pentru continuarea aceste cercetări. Primul este imediat și se referă pe de o parte la identificarea unor metode de automatizare a procesului de identificare a traiectoriilor de antrenare necesare pentru nivelul conexionist al unui sistem compozițional hibrid și pe de altă parte, la automatizarea procesului de creare și antrenare a unor module în funcție natura traiectoriilor detectate.

Dacă primul stadiu se referă mai mult la un proces de implementare al unor activități relativ bine determinate și are un caracter mai redus de cercetare, al doilea se referă la găsirea unor modele îmbunătățite pentru nivelul conexionist. Am afirmat în capitolul 6, că este ideal să folosim o singură rețea neurodinamică și că împărțirea pe module este artificială și dictată de considerente computaționale. O idee posibilă este transformarea setului de module bazate pe rețele recurente într-un singur tot, prin unirea rețelelor prin mănunchiuri de comunicare conexioniste (legături ponderate), reduse ca dimensiune, dar fără o semnificație simbolică. Ideea unei astfel de rețele (fără ca aceasta să aibă o aplicație imediată în procesarea simbolică) este întâlnită deja în literatură [Fahlman97], dar încă nu există algoritmi de antrenare propuși care să aibă performanțele de timp dorite. Deasemenea, pentru o astfel de topologie se poate încerca și utilizarea unor algoritmi de antrenare nesupervizați și aceasta este momentan calea pe care îmi voi continua cercetarea.

Bibliografie

- Ajjanagadde V., 1994, "Unclear Distinctions Lead to Unnecessary Shortcomings: Examining the role vs filler, rule vs fact and type vs predicate distinctions from a connectionist representation and reasoning perspective", Proc. of the Twelfth National Conference on Artificial Intelligence (AAAI-94).
- Allen J.F., 1987. Natural Language Understanding, capitol citat: p324-333, Benjamin-Cummings. Menlo Park, CA.
- Alexandre F., 1996, "Neural Computing and Artificial Intelligence" în Connectionist-Symbolic Integration: From Unified to Hybrid Approaches (editori R. Sun și F. Alexandre), p1-26, CSCSI Print, Montreal.
- Alonso-Betanzos A., Guijaro-Berdiñas B., Moret-Bonillo V., Lopez-Gonzales S., 1995, "The NST-EXPERT project: The Need to Evolve", Artificial Intelligence in Medicine, vol.7, no.4, p297-314, Elsevier New-York.
- Bal H.E., 1989, "The Shared Data Model as a Programming Language Paradigm for Parallel and Distributed Applications", PhD Thesis, Faculteit der Wiskunde en Informatika, Vrije Universiteit, Amsterdam. Olanda.
- Becraft W.R., Lee P.L., Newell R.B., 1991, "Integration of Neural Networks and Expert Systems" în Proc. of IJCAI'91, p832-837, Morgan-Kaufmann, San Mateo, CA.
- Beiu V., 1996, "Digital integrated circuit implementations - Chapter E1.4" în Handbook of Neural Computation (editori, E.Fiesler și R.Beale), Oxford University Press, NY.
- Bookman L.A., 1987, "A Microfeature Based Scheme for Modelling Semantics" în Proc. of IJCAI'87, Morgan Kaufmann, San Mateo, CA.
- Chang C.C., Kreisler H.J., 1973, Model theory, North Holland.

Clark A., 1989, "Connectionism and the multiplicity of mind", *Artificial Intelligence Review* 3/89, p49-65.

Cohen P.R., 1983, *Techniques for reasoning using context, self-reflection and partitioning*, teză de doctorat, Stanford University, CA, re-editată de Pitman Publishers, Londra, în 1985.

Craven M. W., Shavlik J. W., 1994, "Using Sampling and Queries to Extract Rules from Trained Neural Networks", în *Machine Learning: Proceedings of the Eleventh International Conference* (editori W. W. Cohen și H. Hirsch), Morgan Kaufman, San Mateo, CA.

Crevier D., 1993, *The Tumultuous History of the Search for Artificial Intelligence*, Basic Books.

Cristea L.L., Szirbik N., Holban S., 1997, "Some Complexity Issues Concerning a Nonlinear Dynamical System", în *Proc. of Intl. Conference on Analysis and Numerical Computation of Solutions of Nonlinear Systems* (Timișoara), p343-354, Editura Universității de Vest.

Crucianu M., Memmi D., 1992, "Implicit Structure Extraction with Connectionist Networks", *Proc. of NEuro '92* (Nanterre, Franța), p491-502.

DeFelipe J., 1997, "Microcircuits in the Brain" în *Biological and Artificial Computation: From Neuroscience to Technology* (editori, J.Mira, R.Moreno-Diaz, J.Cabestany), p1-14, Springer, Berlin.

Dietterch T.G., Flann N.S., 1995, "Explanation-based Learning and Reinforcement Learning: An Unified View" în *Proc. of 12th International Conference on Machine Learning*, p176-184, Tahoe City, CA, Morgan Kaufmann.

Dragomirescu M., Szirbik N., Dragomirescu C., Babii S., 1993a, "An Expert System for Infectious Diseases Research", *Proc. of CSCS-93*, vol.2, p321-325.

Dragomirescu M., Szirbik N., Dragomirescu C., Kerekes K., Andriuță A., Sălăgean L., Băbălău A., 1993b, "Sistem expert pentru modelarea procesului infecțios", *Conferința Națională de Informatică medicală* (Timișoara, nov.93), p10-18.

Dragomirescu M., Szirbik N., Dragomirescu C., Sălăgean L., 1994, "Modelarea procesului infecțios cu ajutorul rețelelor neuronale" în *Interdisciplinartatea medicinei interne* (editor I.Romoșan), p188-198, Editura Helicon-Banat, Timișoara.

Doya K., 1993, "Bifurcation of Recurrent Neural Networks in Gradient Descent Learning", TR-10-93, Department of Biology, University of California at San Diego, La Jolla, CA.

Dyer M.G., 1991, "Symbolic Neuro-Engineering for Natural Language Processing" în *Advances in Connectionist and Neural Computation Theory* (editori, J. Barnden și J. Pollack), Ablex Publishing, New-York, NY.

Edelman G., 1992, *On the Matter of the Mind*, Basic Books.

Elman J.L., 1990, "Finding Structure in Time", *Cognitive Science* 14 (2/1990).

-
- Fahlman S.E., 1988, "An Empirical Study of Learning Speed in Backpropagation Networks", TR-CMU-CS-88-162, Carnegie-Mellon University, Pittsburgh, PA.
- Fahlman S.E., 1989, "Faster-Learning Variations on Back-Propagation: QuickPropagation", TR-CMU-CS-89-45, Carnegie-Mellon University, Pittsburgh, PA.
- Fahlman S.E., Lebiere C., 1991, "The Cascade-Correlation Architecture", TR-CMU-CS-91-100, Carnegie-Mellon University, Pittsburgh, PA.
- Fahlman S.E., 1997, "Training Boltzmann Neural Networks Using Symbolic Features", acceptată la NIPS'97, vezi la <http://www.cs.cmu.edu/~sef>.
- Farhat N.H., 1997, "Dynamic Neural Networks for Cognition and Control", Proc. of NEURAP'97 (Marsilia, Franța), p1-22.
- Feldman J., Ballard D., 1982, "Connectionist Models and Their Properties", Cognitive Science, 6/82, p205-254.
- Fodor J.A., Pylyshyn Z.W., 1988, "Connectionism and Cognitive Architectures: A critical Analysis", Cognition 28, p3-72.
- Fu L.M., Fu L.C., 1990, "Mapping Rule-Based Systems into Neural Architectures", Knowledge-Based Systems, 3/90, p48-56.
- Gallant S.I., 1988, "Connectionist Expert Systems", Communications of the ACM, 2/88 (no.31).
- Giacometti A., 1992, Hybrid Expert Systems, PhD Thesis, Ecole Nationale Supérieure des Télécommunications, Paris.
- Giles C.L., Miller C.B., Chen D., Sun G.Z., Chen H.H., Lee Y.C., 1992, "Extracting and Learning an Unknown Grammar with Recurrent Neural Networks", în Advances in Neural Information Processing Systems 4, (editori J. E. Moody, S. J. Hanson și R. P. Lippmann), Morgan Kaufman, San Mateo, CA.
- Giles C.L., Chen D., Sun G.Z., Chen H.H., Lee Y.C., Goudreau M.W., 1995, "Constructive Learning of Neural Networks", IEEE Trans. on Neural Networks, vol.6, no.4, p829-836.
- Ginsberg A., 1990, "Theory Reduction, Theory Revision, and Retranslation" în Proc. of National Conference on Artificial Intelligence 1990, p777-782, U.S.A.
- Giunchiglia F., Weyhrauch R.W., 1988, "Meta-level Architectures and Multi-context Axiomatization" în Reasoning, Reflection and Non-monotonicity (editori, P.Maes și D.Nardi), p271-295, Elsevier, New-York.
- Goodman R.K., Miller J.W., Smith P., 1989, "An Information Theoretic Approach to Rule-Based Connectionist Expert Systems" în Advances in Neural Information Processing 1 (editor D.Touretzky), p256-263, Morgan Kaufmann, San Mateo, CA.
-

-
- Gutknecht M., Pfeifer R., Stolze M., 1991, "Cooperative Hybrid Systems" în Proc. of IJCAI'91, p824-829.
- Handelman D.A., Lane S.H., Gelfand J.J., 1989, "Integrating Knowledge-based System and Neural Network Techniques for Robotic Skill Aquisition" în Proj. of IJCAI'89, p193-198, Morgan Kaufmann, San Mateo, CA.
- Haykin S., 1994, *Neural Networks: A Comprehensive Foundation*, Macmillan College Publishing Comp, New-York, NY.
- Hayes S., Ciesielski V.B., Kelly W., 1992, "A Comparision of an Expert System with a NN for Respiratory System Monitory", TR92/1, Royal Melbourne Institute of Technology.
- Hecht-Nielsen R., 1990, "Neurocomputing", Addison-Wesley, New-York.
- Hendler J.A., 1989, "Problem Solving and Reasoning: A Connectionist Perspective" în *Connectionism in Perspective* (editori R.Pfeifer, Z.Schreter și F.Fogelman-Soulié), p229-243.
- Hilario M., 1996, "An Overview of Neurosymbolic Integration Process" în *Connectionist-Symbolic Integration: From Unified to Hybrid Approaches* (editori R. Sun și F. Alexandre), CSCSI, Montreal.
- Honavar V., 1994, "Symbolic AI and Numeric ANNs: Toward a Resolution of Dichotomy", TR94-14(CS), Iowa State University.
- Hölldobler S., Kurfess F., 1991, "CHCL - a Connectionist inference System" în *Parallelization in Inference Systems* (editori, B. Fronhofer și G. Wrightson), Springer-Verlag, Berlin.
- lenne. P., 1997, "Digital Connectionist Hardware: Current Problems and Future Challenges", în Proc. of IWANN'97 (Lanzarote, Spania), p688-713, Springer, Berlin.
- Jabri M., Picard S., Leong P., Chi Z., Flower B., Xie Y., 1992, "ANN Based Classification for Heart Defibrilators" în *Advances in Neural Information Processing 4*, p637-644, Morgan-Kaufmann, San Mateo, CA.
- Kay J., Phillips W. A., 1994, "Activation Functions, Computational Goals and Learning Rules for Local Processors with Contextual Guidance", Centre for Cognitive an Computational Neuroscience, University of Stirling, Technical Report CCCN-15.
- Khosla R., Dillon T., 1996, "Task Structure and Computational Level: Architectural Issues in Symbolic-Connectionist Integration" în *Connectionist-Symbolic Integration: From Unified to Hybrid Approaches* (editori R. Sun și F. Alexandre), CSCSI, Montreal.
- Kibler D., Langley P., 1988, "Machine Learning as an Experimental Science" în Proc. of Third European Working Session on Learning, p1-12, Edinburgh, UK.
- Kokinov B., 1996, "Micro-level Hybridization in DUAL" în *Connectionist-Symbolic Integration: From Unified to Hybrid Approaches* (editori R. Sun și F. Alexandre), CSCSI, Montreal.
-

- Koppel M., Feldman R., Segre A. H., 1994, "Bias-Driven Revision of Logical Domain Theories", *Journal of Artificial Intelligence Research* 1/94, p159-208.
- Hertz J., Krogh A., Palmer R.G., 1991, "Introduction to the Theory of Neural Computation", Addison-Wesley, Reading, MA.
- Krovvidy S., Wee W.G., "An Intelligent Hybrid System for Wastewater Treatment" în *Hybrid Architectures of Intelligent Systems*, cap.17, p358-377, CRC Press, Boca Raton, CA.
- Kunicky D.C., Hruska S.I., Lacher R.C., 1992, "Hybrid Systems: The Equivalence of Rule-Based Expert Systems and Artificial Neural Networks Inference", *International Journal of Expert Systems*, 4/92, p281-297.
- Labbi A., 1993 "Neural Networks for Decision Making in Dynamic Environments" în *Proc. of ICANN'93*. Amsterdam, Springer-Verlag, Berlin.
- Lendaris G.G., 1992, "A Neural Network Approach for Implementing Conceptual Graphs" în *Conceptual Structures: Current Research and Practice*, (editor T.E.Nagle) p165-177, Ellis Horwood Inc., Chichester, UK.
- Leonard J.J., Durrant-Whyte H.F., 1992, *Mobile Robot Navigation Using Neural Nets*, Kluwer Academic Publishers, Boston, MA.
- Leung K.S., Lam W., 1988, "Fuzzy Concepts in Expert Systems", *IEEE Expert*, 3/88, p43-56.
- Luger S., Stubblefield J., 1989, "Artificial Intelligence and the Design of Expert Systems", Benjamin/Cummings, New-York.
- MacKay D., 1992, "A Practical Bayesian Framework for Backpropagation Networks", *Neural Computation* 4/92, p448-472.
- Maclin R., Shavlik J. W., 1994, "Incorporating Advice into Agents that Learn from Reinforcements", *Proceedings of the Twelfth International Conference on Artificial Intelligence (AAAI-94)*.
- McClelland J.L., O'Reilly R.C., 1994, "Brain Conjunctive Encoding, Storage and Recall: Avoiding a Tradeoff", *TR-PDP-CNS-94-4*, Carnegie-Mellon University, Pittsburgh, PA.
- McCulloch W.S., 1988, *Embodiments of Mind*, MIT Press, Cambridge, MA.
- Memmi D., 1990, "Connectionism and Artificial Intelligence as Cognitive Models", *AI&Society* 4, p115-136.
- Miikulainen R., 1991, "Script Recognition with Hierarchical Feature Maps", *TR-UCLA-AI-91-10*, University of California at Los Angeles, CA.
- Miikulainen R., 1994, "Integrated Connectionist Models: Building AI Systems on Subsymbolic Foundations" în *AI and Neural Networks* (editori V.Honavar și L.Uhr), p483-508, Academic Press, San Diego, CA.
-

-
- Miikulainen R., 1996, "Subsymbolic Case-Role Analysis of Sentences with Embedded Clauses, TR-CS-45-96, University of Texas, Austin.
- Minsky M.L., Papert S.A., 1969, *Perceptrons*, MIT Press, Cambridge, MA.
- Minsky M.L., 1986, *The Society of Mind*, MIT Press, Cambridge, MA.
- Minsky M., 1991, "Local versus Analogical or Symbolic versus Connectionist or Neat versus Scruffy", *AI Magazine*, 3/91, p7-33.
- Nagayama I., Akamatsu N., 1994, "Learning of Chaos by Neural Networks and Bifurcation Phenomena", *Transactions of the Institute of Electronics, Information and Communication Engineers*, vol J77A, no.11, p1593-6.
- Newell A., Simon H.A., 1972, *Human Problem Solving*, Prentice-Hall, Engelwood Cliffs, NJ.
- Newell A., 1990, *Unified Theories of Cognition*, Harvard University Press, Cambridge, MA.
- Omlin C. W., Giles C. L., 1992, "Training Second-Order Recurrent Networks Using Hints", *Machine Learning: Proceedings of the Ninth International Conference*, edited by D. Sleeman and P. Edwards, Morgan Kaufman, San Mateo, CA.
- Opitz D.W., Shavlik J.W., 1997, "Connectionist Theory Refinement: Genetically Searching the Space of Network Topologies", *Journal of Artificial Intelligence Research* 6/97, p177-209.
- Orsier B., Labbi A., 1996, "NESSY3L: A NeuroSymbolic System with 3 Levels" în *Connectionist-Symbolic Integration: From Unified to Hybrid Approaches* (editori R. Sun și F. Alexandre), CSCSI, Montreal.
- Pazzani M., Brunk C., 1991, "Detecting and Correcting Errors in Rule-Based Expert Systems: An Integration of Empirical and Explanation-Based Learning", *Knowledge Acquisition* 2/91 (3), p157-173.
- Pearl J., 1988, *Probabilistic Reasoning in Intelligent Systems*, Morgan Kaufmann, San Mateo, CA.
- Pearlmutter B.A., 1990, "Dynamic Recurrent Neural Networks", TR-CMU-CS-90-191, Carnegie-Mellon University, Pittsburgh, PA.
- Pfeifer, W.T., 1996, "Symbols and Patterns: Beyond the Classical Information Processing", în *Encyclopedia of Computers* (editor, Marcel Dekker), vol.17, p253-275, Elsevier, New-York.
- Phan K.M., 1995, "A Neural Multi-Agent Approach for Modelling, Distributed Processing and Learning" în *Intelligent Hybrid Systems* (editori, S.Goonatilak și S.Kebal), John Wiley Comp., Chichester, UK.
- Pineda F.J., 1987, "Generalization of Back-Propagation to Recurrent Neural Networks", *Physical Review Letters* 19/87, p2229-32.
-

- Pineda F.J., 1988, "Dynamics and Architecture for Neural Computation", *Journal of Complexity* 4.
- Pinker S., Prince A., 1988, "On Language and Connectionism: Analysis of Parallel Distributed Processing Model of Language Acquisition", *Cognition* 28, p73-193.
- Piramuthu S., Shaw J.M., 1994, "On Using a Decision Tree as Feature Selector for Feedforward Neural Networks" în *Proc. of International Symposium on Integrating Knowledge and Neural Heuristics* (Pensacola, FL), p67-74.
- Quillian R., Collins J., 1969, "A Graf Based Defined Lexicon as a Model for Plausible Representation in Human Mind", *Psychological Review* 80, p413-51.
- Rao S.S., Ramamurti V., 1993, "A Hybrid Technique to Enhance the Performance in Recurrent neural Networks", *Proc. of ICNN'93, San Francisco*, cat.no.93CH3274-8.
- Reeke G.N., Edelman G.M., 1988, "Real Brains and Artificial Intelligence" în *The Artificial Intelligence Debate: False Starts, Real Foundations* (editor S. Graubard), p144-173, MIT Press, Cambridge, MA.
- Riesbeck C.K., 1986, "Realistic Natural Language Processing" în *Issues in Cognitive Modelling* (editori, A.Aitkenhead și J.Slack), p193-206, Lawrence Erlbaum Assoc., London, UK.
- Roberts D.D., 1992, "The Existential Graphs" în *Semantic Networks in Artificial Intelligence*, p639-663, (editor F.Lehman), Pergamon Press, Oxford, UK.
- Robinson G., 1965, "The Resolution Principle". *AI Review* 4/65.
- Röbel A., 1994, "The Dynamic Pattern Selection Algorithm: Effective Training and Controlled Generalization of Recurrent Backpropagation Neural Networks", *Technical Report*, Institut für Augewandte Informatik, Technische Universität Berlin.
- Rumelhart D.E., Smolensky P., McClelland J.L., Hinton G.E., 1986, "Schemata and Sequential Thought Process in PDP Models" în *Parallel Distributed Processing (vol.2): Explorations in the Microstructure of Cognition* (editori, D.Rumelhart și J.McClelland), p7-57, Bradford Books, Cambridge, MA.
- Schank R.C., 1981, *Inside Computer Understanding*, Lawrence Erlbaum Assoc., Hillsdale, NJ.
- Scherer A., Schlegeler G., 1995, "A Multi-Agent Approach for the Integration of Neural Networks and Expert Systems" în *Intelligent Hybrid Systems* (editori, S.Goonatilak și S.Kebal), John Wiley Comp., Chichester, UK.
- Scott G., Shavlik J., Ray W., 1992, "Refining PID Controllers Using Neural Networks", *Neural Computation* 5/92, p746-757.
- Sharkey N. E., Reilly R., 1990, "Connectionist Natural Language Processing", *Parallel Processing Specialist Group Newsletter*, No. 13.
-

-
- Sharkey N., 1992, "The Gost in the Hybrid: a Study of Uniquely Connectionist Representations", *AISB Quarterly* 79, p10-16.
- Sharkey N., 1995, "Searching the Weight Space for Solutions" în *Current Trends in Connectionism* (editori, L.Niklasson și M.Boden), p103-120, Lawrence Erlbaum Publishers, Hove, UK.
- Shastri L., 1988, "A Connectionist Approach to Knowledge Representation and Limited Inference", *Cognitive Science*, 12/88, p331-392.
- Shavlik J.W., 1991, "Hybrid symbolic+connectionist reasoning automata", TR-118-91, University of Wisconsin, Madison.
- Shavlik J.W., 1997, "How to Close the Gap Between Neural and Symbolic Representations", *Machine Learning*, vol. 25 (Jan/Feb/Mar), p295-301.
- Shortliffe E.H., 1976, *Computer-Based Medical Consultations: MYCIN*, American Elsevier, New-York, NY.
- Sowa J.F., 1984, *Conceptual Structures: Information Processing in Mind and Machine*, Addison-Wesley, Reading, MA.
- Sowa, J.F., 1995, *Knowledge Representation: Logical, Philosophical and Computational Foundations*. PWS Publishing Comp. Boston, MA.
- Smolensky P., 1986, "Neural and Conceptual Interpretation of PDP Models" în *Parallel Distributed Processing (vol.2): Explorations in the Microstructure of Cognition* (editori, D.Rumelhart și J.McClelland), p390-431, Bradford Books, Cambridge, MA.
- Smolesky P., 1988, "On the Proper Treatment of Connectionism", *Behavioral and Brain Sciences*, 11/88, p1-74.
- Smolensky P., 1990, "Tensor Product Variable Binding and the Representation of Symbolic Structures in Connectionist Systems", *Artificial Intelligence* 46/1990, p159-216.
- Smolensky P., Legendre G., Miyata Y., 1992, "Principles for an Integrated Connectionist/Symbolic Theory of Higher Cognition", Computer Science Department, University of Colorado-Boulder, Technical Report CU-CS-600-92.
- Somlo G., 1996, "Improving Learning Speed in Recurrent Neural Networks", Research Progress Report-520 (Computer Science), Colorado State University, Fort Collins, CO.
- Sun R., 1991, "Integrating Rule and Connectionism for Robust Reasoning. A Connectionist Architecture with Dual Representation, PhD Thesis (TR-CS-91-160), Brandeis University, Waltham, MA.
- Sun R., 1992, "A Connectionist Model for Commonsense Reasoning Incorporating Rules and Similarities", *Knowledge Acquisition* 4, p293-321, Academic Press, Cambridge, UK.
-

- Sun G. Z., Giles C. L., Chen H. H., Lee Y. C., 1993, "The Neural Network Pushdown Automaton: Model, Stack and Learning Simulations", University of Maryland, Technical Report UMIACS-TR-93-77.
- Szirbik N., Babii S., Dragomirescu C., 1993, "A Method for Inserting PROLOG type Rules in a Neural Network", Proc. of CSCS-93 (București), vol.2, p202-206.
- Szirbik N., 1993, "Modus Ponens and Neural Networks" în Buletinul Științific și Tehnic al UTT. Tom 38(52), p145-149.
- Szirbik N., Somlo G., Buliga D., 1995, "Using the Conceptual Graph Model as Intermediate Representation for Knowledge Translation in Hybrid Systems" în Current Trends in Connectionism (editori, L.Niklasson și M.Boden), p141-152, Lawrence Erlbaum Publishers, Hove, UK.
- Szirbik N., Babii., 1995, "Intermediate Representation in Parallel-Distributed Processing", Proc. of CSCS-95 (București), vol.3, p243-252.
- Szirbik N., 1995a, "Recurrent Neural Networks Learning on the ZOO Multiprocessor: Speed-ups and Other Results", Research Progress Report, Faculteit der Viskunde en Informatika, Vrije Universiteit, Amsterdam, Olanda.
- Szirbik N., 1995b, "Towards a Model of Multi-Agent Connectionist Hybrid System" în Proc. of ANNES'95 (Dunedin, NZ), p273-6, IEEE Computer Soc.Press, Los Alamitos, CA.
- Szirbik N., Somlo G., 1995, "Parallelizing the Williams&Zipser Algorithm using PVM", Lucrare de diplomă (cu rezumat extins în engleză), Catedra de calculatoare, Universitatea "Politehnica" din Timișoara.
- Szirbik N., Babii S., 1996, "Modelling Connectionist Reasoning Patterns by Interacting Hybrid Reasoning Agents" în Buletinul Științific al UPT, Tom 41(55), p92-100.
- Szirbik N., 1996, "NeuroSymbolic Integration: A Way Towards Building Intelligent Systems" în Proc. of CONTI'96 (Timișoara), vol.1, p135-44.
- Szirbik N., 1997, "A Two-Level Heterogeneous Hybrid Model" în Proc. of IWANN'97 (Lanzarote, Spania), p644-650, Springer, Berlin.
- Tannenbaum A.S., Rennese R., 1989, "Experiences with the Amoeba Operating System", Report IR-194, Faculteit der Viskunde en Informatika, Vrije Universiteit, Amsterdam, Olanda.
- Thomason R.H., Touretzky D.S., 1990, "Inheritance Theory and Networks with Roles", TR-CMU-CS-90-139, Carnegie-Mellon University, Pittsburgh, PA.
- Todd B.S., 1991, "Bayesian Inference Theory", TR-PRG-95, Oxford University Computing Laboratory, Oxford, UK.
- Touretzky D.S., Geva S., 1987, "A Distributed Connectionist Representation for Concept Structures" în Proc. of Ninth Annual Conference of the Cognitive Science Society (Seattle,

WA).

Touretzky D.S., Hinton G.E., 1988, "A Distributed Connectionist Production System", *Cognitive Science*, 12/88, p423-466.

Touretzky D.S., 1989, "BoltzCONS: a Connectionist Symbolic Machine", TR-CMU-CS-89-182. Carnegie-Mellon University, Pittsburgh, PA.

Touretzky D.S., 1990, "Dynamic Symbol Structures in a Connectionist Network", *Artificial Intelligence*, 46 (1-2).

Towell G. G., Shavlik J. W., 1991, "The Extraction of Refined Rules from Knowledge-Based Neural Networks", *Machine Learning* 8/91.

Towell G. G., 1991, "Symbolic Knowledge and Neural Networks: Insertion, Refinement and Extraction", Ph.D. Thesis at the Computer Science Department of the University of Wisconsin, Madison.

Treur J., 1988, "Design of Modular and Interactive Knowledge-based Systems", Report P8826, Faculteit der Wiskunde en Informatika, Vrije Universiteit, Amsterdam, Olanda.

Treur J., 1990, "Modelling Non-classical Reasoning Patterns by Interacting Reasoning Modules", Report IR-236, Faculteit der Wiskunde en Informatika, Vrije Universiteit, Amsterdam, Olanda.

Treur J., 1992, "Towards Dynamic and Temporal Semantics of Meta-Level Architectures in DESIRE", Report IR-321, Faculteit der Wiskunde en Informatika, Vrije Universiteit, Amsterdam, Olanda.

Tsung F.S., Cottrell G.W., 1993, "Hopf Bifurcations and Hopf Hopping in Recurrent Neural Networks" în *Proc. of The IEEE Conference on Neural Networks*, vol.1, p39-45, IEEE Computer Soc.Press, New-York, NY.

Van Gelder T., 1995, "Modelling Connectionist and Otherwise" în *Current Trends in Connectionism* (editori, L.Niklasson și M.Boden), p217-36, Lawrence Erlbaum Publishers, Hove, UK.

Vilain M., McAllester D., 1989, "An Overview of NIKL: Tackling the Huge Number of Assertions" în *Annual Report of Research in Knowledge Representation*, (editori, B.Beranek și L.Newman), p45-79, BBN Press, Cambridge, MA.

Wang Y. Y., Waibel A., 1992, "Dialogue Processing with Neural Networks", Carnegie-Mellon University, Pittsburgh, Technical Report CMU-CS-92-165.

Watrous R., Towell G., Glassman M., 1993, "Synthesize, optimize, analyze, repeat (SOAR): Application of Neural Network Tools to ECG patient monitoring" în *Proc. of Symposium on Nonlinear Theory and its Applications*, p565-570, Honolulu, Hawaii.

Wermter S., Weber W., 1997, "SCREEN: Learning a Flat Syntactic and Semantic Spoken

Language Analysis Using Artificial Neural Networks", Journal of AI Research 6, 1/97, p35-85.

Wiggins S., 1990, Introduction to Applied Nonlinear Dynamical Systems, Bifurcations and Chaos, Springer-Verlag, New-York, NY.

Wilkins D.C., 1988, "Knowledge Base Refinement Using Apprenticeship Learning Techniques" in Proc. of National Conference on Artificial Intelligence 1988, p646-653.

Williams R. J., Zipser D., 1989, "A Learning Algorithm for Continually Running Fully Recurrent Neural Networks", Neural Computation, vol 1/89, p270-280.

Williams R.J., Zipser D., 1990, "Gradient Based Learning Algorithms for Recurrent Connectionist Networks", TR-NU-CCS-90-9, Northeastern University, Boston, MA.

Williams R.J., Peng J., 1990, "An Efficient Gradient Based Algorithm for On-line Training of Recurrent Network Trajectories", Neural Computation 2/90, p490-501.

Winograd T., 1980, "What Does It Mean to Understand Language ?", Artificial Intelligence 6/80, p231-263.

Notă autobiografică a autorului

Data nașterii: 31 mai, 1963

Locul nașterii: Timișoara

Clasele I-VIII: Colegiul "C.D. Loga", Timișoara, 1969-77

Liceul: Liceul de informatică "Grigore Moisil", Timișoara, 1977-1981

Stagiul militar: Rgt.56 Parașutiști, Caracal, 1981-1982

Facultatea: Electrotehnică, Specialitatea Calculatoare, Institutul Politehnic Traian Vuia Timișoara, 1982-1987.

Inginer stagiar la: I.A.M.S.A.T. București și Centrul de calcul al Universității Tehnice din Timișoara, 1987-1990

asistent universitar la: Catedra de calculatoare a Universității Tehnice din Timișoara, 1991-1994

înscriș la doctorat din: 1992

șef de lucrări la: Departamentul de Calculatoare al Universității "Politehnica" Timișoara, 1994-prezent

titular al cursurilor: Sisteme Expert (anul V), Modelul Computațional Conexiunist (Rețele Neuronale - curs pentru studenții de la Studii aprofundate)

~