

CAPITOLUL VII

ÎNVĂȚARE

7.1 Problematika învățării în Inteligența Artificială

Problemele de învățare reprezintă o direcție importantă de cercetare în Inteligența Artificială, cea a *mașinilor auto-instruibile (învățării automate)*. Învățarea automată reprezintă studiul modelelor de sisteme, care pe baza unui set de date (de instruire), își îmbunătățesc performanțele prin experiențe și prin învățarea unor cunoștințe specifice domeniului, prin experiențe.

Încercarea de modelare a raționamentului uman conduce la apariția noțiunii de *raționament inteligent*. *Raționamentul* este un proces de a deduce concluzii; raționamentul *intelligent* vrem să simuleze raționamentul efectuat de oameni.

Majoritatea sistemelor cu Inteligență Artificială sunt sisteme deductive, capabile să tragă concluzii (să facă inferențe) pe baza cunoașterii lor inițiale sau furnizate, nefiind însă capabile să achiziționeze și să genereze singure cunoștințe.

Capacitatea de a învăța fiind strâns legată de comportamentul inteligent, o direcție de cercetare în Inteligența Artificială este implementarea în mașini a capacității de învățare. Realizarea unor astfel de mașini auto-instruibile devine deosebit de importantă în realizarea sistemelor complexe bazate pe cunoaștere, cum ar fi de exemplu sistemele expert.

După cum am arătat până acum, *inteligenta* unui agent este încorporată în partea sa de program. La un moment dat, agentul va alege cel mai bun mod de acțiune, în sensul în care a fost programat să o facă.

În situații în care programul are informații (cunoștințe) incomplete despre mediul în care agentul acționează, *învățarea* este singura modalitate prin care agentul poate să achiziționeze cunoștințele de care are nevoie. Astfel, *învățarea* conferă *autonomia* agentului, în sensul în care a fost descrisă în capitolul I. Deasemenea, conferă o modalitate de a construi sisteme performante – furnizând o experiență dobândită prin învățare în domeniul aplicației.

Știința mașinilor cu învățare automată (mașini auto-instruibile) se ocupă cu dezvoltarea de algoritmi care:

- capătă experiență când operează într-un anumit mediu;
- recunosc propria stare internă (poate chiar și datele și programul) pentru a deveni mai eficienți în mediul pentru care au fost pregătiți;
- au abilitatea de a-și evalua propriile performanțe fără a fi necesară intervenția unor stimuli externi.

Studiul sistemelor inteligente cuprinde atât mașinile cu învățare automată, cât și tehnologiile de genul sistemelor expert și rețelelor neuronale, care necesită uriașe eforturi umane pentru dezvoltarea seturilor de reguli euristice, respectiv a datelor necesare pentru antrenare.

7.1.1 Modelul general al unui agent care învață

Un agent care învață poate fi împărțit în patru componente conceptuale. Distincția cea mai importantă care trebuie făcută este între *elementul de învățare*, care este responsabil de a face îmbunătățiri, și *elementul de performanță*, care este responsabil de a selecta acțiuni externe. Elementul de performanță este ceea ce până acum am considerat a fi întregul agent: primește percepțiile și alege acțiunile. Elementul de învățare vede cum se “descurcă” agentul și stabilește cum ar trebui modificat elementul de performanță pentru a acționa mai bine în viitor.

Modelul general al unui agent care învață este descris în Figura 7.1.

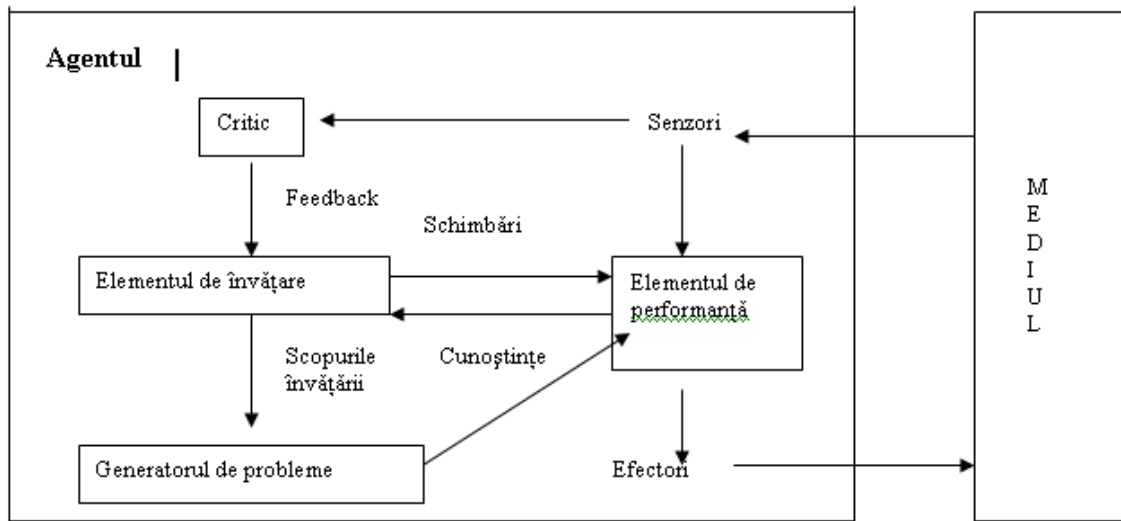


Figura 7.1. Modelul general al unui agent care învață

Într-un astfel de model, *criticul* îi spune elementului de învățare cât de bine se descurcă agentul. Criticul se folosește de un standard fix de performanță. Acest lucru este necesar deoarece percepțiile nu dau nici o indicație despre succesul agentului. De exemplu, un program de șah poate primi o percepție care indică faptul că i-a dat mat adversarului său, dar este nevoie de un standard al performanței pentru a ști dacă acesta este un lucru bun deoarece percepția nu i-o spune. Este important ca standardul de performanță să fie o măsură fixă aflată în exteriorul agentului, altfel agentul ar putea ajusta standardul de performanță pentru ca acesta să arate că este bun comportamentul său.

Generatorul de probleme este responsabil de a sugera acțiuni care ar putea conduce la noi informații. Dacă elementul de performanță ar face totul, fără a fi întrebat generatorul de probleme, acesta ar alege acțiunile cele mai bune având în vedere ceea ce cunoaște agentul la un moment dat. Dar dacă agentul ar explora puțin și ar alege acțiuni suboptimale pe termen scurt, ar putea descoperi acțiuni mult mai bune pe termen lung. Sarcina generatorului de probleme este de a sugera aceste acțiuni de explorare.

Să considerăm exemplul unui agent care este șofer de taxi. Elementul de performanță este alcătuit din colecția de cunoștințe și proceduri pe care le are agentul pentru a putea selecta acțiunile necesare conducerii mașinii (întoarcerea, accelerarea, frânarea, etc.). Agentul conduce mașina folosind acest element de performanță. Elementul de învățare formulează țelurile, de exemplu, să învețe cum se comportă mașina pe un drum umed. Criticul observă lumea și trimite informațiile la elementul de învățare. De exemplu, dacă șoferul trece rapid de pe o bandă pe alta fără a semnaliza, criticul observă iritarea celorlalți participanți la trafic și îi transmite aceste observații elementului de învățare. Acesta formulează o regulă care spune că este greșită schimbarea benzii fără semnalizare și elementul de performanță se va modifica prin adăugarea acestei noi reguli. Din când în când generatorul de probleme dă agentului o sugestie. De exemplu: dacă este aglomerată strada, generatorul de probleme i-ar putea sugera agentului să încerce să meargă pe o stradă paralelă și să vadă dacă nu cumva aceasta este mai rapidă decât calea aleasă inițial.

Elementul de învățare mai este responsabil și de îmbunătățirea eficienței elementului de performanță. De exemplu, dacă șoferului de taxi i se va cere să meargă pe o rută nouă, acesta își va consulta harta pentru a găsi drumul cel mai scurt. Când i se va solicita din nou acest traseu, șoferul nu va mai căuta pe hartă căci deja va cunoaște acest traseu.

Componentele elementului de performanță

Elementul de performanță al unui agent care învață are următoarele componente:

1. o mapare directă între condiții din stări și acțiuni;
2. o metodă de a deduce proprietăți relevante ale mediului din secvența perceptuală;
3. informații despre modul în care lumea (mediul) evoluează;
4. informații despre rezultatele posibilelor acțiuni pe care agentul le poate efectua;
5. informații de utilitate indicând utilitatea stărilor mediului;
6. informații acțiuni-valoare indicând utilitatea unor acțiuni particulare în stări particulare;
7. scopuri indicând clase de stări a căror realizare va maximiza utilitatea agentului.

Fiecare dintre cele șapte componente ale elementului de performanță pot fi descrise matematic sub forma unor **funcții**. Spre exemplu componenta 3 poate fi interpretată ca o funcție care asignează unei stări din mediu (starea curentă) o altă stare a mediului (starea următoare); componenta 7 poate fi interpretată ca o funcție care asignează unei stări a mediului o valoare booleană (adevărat sau fals) după cum starea este sau nu o stare finală (scop).

Ca urmare, orice sarcină de învățare poate fi văzută de fapt ca fiind învățarea reprezentării unei anumite funcții.

Feed-back-ul disponibil

Pentru anumite componente ale elementului de performanță, cum ar fi componenta pentru precizarea rezultatelor unei acțiuni (componenta 4), feed-back-ul disponibil îi spune, în general, agentului, care este rezultatul corect (ieșirea corectă).

Orice situație în care pot fi observate atât semnalele de intrare cât și semnalele de ieșire ale unei componente se numește **învățare supervizată** (*Supervised Learning*) (în general, ieșirile sunt furnizate de către un profesor). Aceasta constă în faptul că agentul crede că o anumită acțiune are un anumit rezultat iar mediul emite un stimul prin care îi spune agentului care este rezultatul corect al acțiunii sale.

Pe de altă parte, în învățarea componentei *condiție-acțiune* (componenta 1), agentul primește o evaluare a acțiunii sale, dar nu i se spune care este ieșirea (rezultatul corect). Acest tip de învățare se numește **învățare prin întărire** (*Reinforcement Learning*), iar evaluarea primită de agent se numește **recompensă** (*reinforcement*).

Învățarea în care nu se dă agentului nici o indicație referitoare la rezultatele corecte se numește **învățare nesupervizată** (*Unsupervised Learning*).

7.1.2 Învățarea unui domeniu

În învățarea unui domeniu specific, sistemul conține numeroase concepte predefinite, structuri de reprezentare a cunoștințelor, restricțiile domeniului și reguli euristice. Se așteaptă ca sistemul să deducă noi atribute și concepte în procesul de învățare (inducție constructivă). Un sistem bazat pe această abordare este dezvoltat pentru un domeniu particular și nu poate fi folosit direct în alt domeniu. Tipurile de învățare folosite sunt: *învățarea din exemple*, *învățarea prin analogie* și *învățarea prin observație și descoperire*.

S-au construit sisteme ce combină această abordare și calculul neuronal. Separând mecanismele de inferență generale de cunoștințele specifice domeniului s-au realizat sisteme ce pot fi aplicate la o largă varietate de domenii diferite.

7.1.3 Tipuri de învățare

În procesul de învățare, informația furnizată de un profesor sau de mediu este transformată într-o formă nouă în care ea este memorată pentru uzul ulterior. Natura acestei transformări determină tipul de strategie. Astfel putem distinge *învățarea mecanică*, *învățarea prin instruire*, *învățarea prin analogie* și *învățarea prin deducție*. Ultima se subdivide în *învățarea din exemple* și *învățarea din observație* (*învățarea nesupervizată*) și prin *descoperire*.

În învățarea mecanică informația furnizată de profesor nu suferă nici o transformare esențială, fiind memorată direct. Problema de bază este cum să fie indexată și memorată cunoașterea în vederea regăsirii ulterioare. În învățarea prin instruire, transformările de bază sunt *selecția* și *reformularea* (mai ales la nivel sintactic) a informației furnizate de profesori. În învățarea prin deducție, pe baza cunoașterii date, sunt formulate inferențe deductive care au proprietatea de a conserva adevărul, concluziile utile fiind memorate. Dacă procesele de transformare presupun *generalizarea* informației primite și *selectarea* celui mai plauzibil sau dezirabil rezultat, avem de-a face

cu un proces de *inferență inductivă*. În învățarea *prin analogie* se combină învățarea deductivă și cea inductivă. Descrierile din diferite puncte de vedere sunt combinate pentru a determina o substructură comună, care să servească drept bază pentru analogii. Găsirea substructurii comune presupune inferență inductivă, în timp ce realizarea proiecției analogice este o formă de deducție.

7.2 Învățarea în Inteligența Artificială Tradițională

7.2.1 Învățarea Supervizată

7.2.1.1 Învățarea inductivă

La învățarea supervizată elementul de învățare primește valorile corecte (sau aproape corecte) ale unei funcții pentru anumite semnale de intrare particulare. Elementul de învățare modifică reprezentarea funcției astfel încât să se potrivească cu informațiile pe care le primește. De exemplu, se dă o pereche de forma $(x, f(x))$ unde x este semnalul de intrare, iar $f(x)$ este valoarea funcției f aplicată lui x . Sarcina inducției este următoarea: fiind dată o colecție de valori ale funcției f , să se găsească funcția h care aproximează funcția f . Funcția h se numește *ipoteză*. Funcția f nu se cunoaște, dar sunt mai multe variante de alegere a funcției h .

Să presupunem că avem un agent reflex care învață de la un profesor. Prin funcția *Reflex_Learning_Element*, descrisă în Figura 7.2.a), agentul, prin elementul de învățare, își reactualizează variabila globală, *exemple*, care conține o listă de perechi (*acțiune*, *stare*). Percepția ar putea fi o poziție de pe tabla de șah și acțiunea ar putea fi cea mai bună mutare care îl poate conduce pe agent la câștigarea partidei. Când elementul de performanță (descrie în Figura 7.2.b) prin funcția *Reflex_Performance_Element*) va primi o percepție despre care a învățat, va alege acțiunea corespunzătoare. În celelalte cazuri va apela algoritmul de învățare *Inducție*, care va formula o ipoteză h pe care agentul o va folosi pentru a alege acțiunea următoare.

După cum se vede din Figura 7.2, elementul de învățare, memorează fiecare pereche (percepție, acțiune). Elementul de performanță repetă acțiunea care a mai făcut-o când a primit același stimul sau caută o acțiune. Setul *exemple* e o variabilă globală care apare atât în elementul de învățare cât și în cel de performanță.

funcția <i>Reflex_Learning_Element</i> (percepție, acțiune) intrări: <i>percepție, acțiune</i> <i>exemple</i> \leftarrow <i>exemple</i> \cup {(percepție, acțiune)} sf- <i>Reflex_Learning_Element</i>	a
funcția <i>Reflex_Performance_Element</i> (percepție) returnează acțiune dacă (percepție, a) în <i>exemple</i> atunci returnează a altfel h \leftarrow Inducție(<i>exemple</i>) sf-dacă returnează h(percepție) sf- <i>Reflex_Performance_Element</i>	b

Figura 7.2. Modelul unui agent reflex care învață de la un profesor

7.2.1.2. Învățarea prin rețele neuronale

Calculul neuronal implică două aspecte fundamentale: învățarea și reprezentarea cunoașterii. Abordarea neuronală a acestor aspecte este diferită de abordarea standard întâlnită în modelele bazate pe reprezentarea simbolică a cunoștințelor. Calculul neuronal este caracterizat prin posibilitatea învățării bazată pe informații parțiale sau pe date ce conțin erori. Acest lucru este posibil datorită caracterului robust (în raport cu informațiile ce reprezintă erori) al învățării neuronale.

Există un consens aproape unanim în a accepta faptul că tehnicile de învățare reprezintă calea prin care mașinile pot deveni capabile să realizeze sarcinile dificile ale inteligenței artificiale. În mod evident un program inteligent, care învață prin instruire și din experiență, poate fi capabil să realizeze sarcini mai dificile decât un program bazat pe o listă de posibilități care a fost stabilită și analizată de programatorul uman.

Există o analogie între neuronii umani și cei artificiali. Această analogie este sintetizată mai jos:

Neuron uman	Neuron Artificial
Neuron	Element de procesare
Dendrite	Conexiuni
Corp celular	Funcție de activare
Axon	Element de ieșire
Sinapse	Ponderi

Totuși, analogia nu este perfectă. Neuronii umani și activitatea neuronală sunt mult mai complecși decât poate sugera studierea neuronilor artificiali. Neuronii umani nu efectuează o simplă sumă ponderată, iar mecanismul de funcționare al dendritelor în sistemul biologic este mult mai elaborat.

Învățarea în rețele neuronale poate fi văzută în două moduri. Din punct de vedere computațional este vorba despre o metodă de reprezentare a funcțiilor folosind rețele de elemente computaționale aritmetice simple și despre metode de învățare a acestor reprezentări din exemple. Din punct de vedere biologic, învățarea în rețele neuronale și

rețele “de încredere” poate fi văzută ca un model matematic pentru reprezentarea operațiilor creierului. Elementele computaționale aritmetice simple corespund *neuronilor* (celulele care execută procesul de informare în creier), iar rețeaua corespunde unei colecții de neuroni interconectați. Din acest motiv rețelele se numesc *rețele neuronale*. Rețelele neuronale oferă una dintre acele șanse pentru înțelegerea multor fenomene psihologice, fenomene care apar din structura specifică și din operațiile creierului.

O *rețea neuronală* este compusă dintr-un număr de noduri sau *unități* conectate prin *legături*. Fiecare legătură are asociată o anumită *ponderare*. Învățarea are loc prin reactualizarea ponderilor. Unele elemente sunt conectate la mediul extern și pot fi proiectate ca elemente de intrare sau ieșire. Ponderile sunt modificate în încercarea de a determina comportamentul rețelei să fie cât mai aproape de semnalele de intrare date de mediu.

Funcția de activare neuronală

Neuronii unei rețele neuronale își transformă intrarea în ieșire folosind o funcție numită funcție de activare. Această ieșire constituie parte din intrarea tuturor neuronilor conectați la neuronul respectiv. Învățarea unei rețele neuronale depinde de valorile ponderilor conexiunilor interneuronale dar și de funcția de activare (numită și funcție de transfer a informației) specifică neuronilor.

În cazul neuronilor ascunși, funcția de activare are scopul de a introduce non-liniaritatea în rețeaua neuronală. Rețelele neuronale cu mai multe straturi dau rezultate foarte bune și datorită non-liniarității.

Câteva din cele mai importante funcții de activare neuronală sunt prezentate mai jos:

▪ Funcția identitate

Această funcție lasă intrarea nemodificată, ieșirea neuronului fiind astfel identică cu intrarea sa. Uneori însă, avem de-a face cu extinderi ale acestei funcții, intrarea fiind multiplicată cu o constantă pentru a determina valoarea de ieșire.

$$g(x)=x$$

▪ Funcția prag

Este cunoscută și sub denumirea de funcție binară de activare, deoarece poate lua doar două valori: 0 și 1 și este folosită de obicei în rețelele neuronale cu un singur strat.

$$g(x)=\begin{cases} 1 & \text{if}(x \geq \theta) \\ 0 & \text{if}(x < \theta) \end{cases}$$

▪ Funcția sigmoidală

Este indicată folosirea acestui tip de funcție de activare la rețelele care folosesc un algoritm de învățare cu propagare înapoi a erorii și în aplicațiile în care ieșirile sunt numere reale între 0 și 1.

$$g(x) = \frac{1}{1 + e^{-x}}$$

▪ Funcția sigmoidală bipolară

După cum sugerează și numele, această funcție este asemănătoare funcției sigmoidale, fiind o extensie a acesteia. Este aplicabilă problemelor pentru care ieșirea ia valori reale între -1 și 1.

$$g(x) = \frac{1 - e^{-x}}{1 + e^{-x}}$$

Se observă că funcțiile de activare liniare, dau posibilitatea neuronului să aibă mai mult de două ieșiri. Ieșirea neuronului este proporțională cu activarea sa. Pentru neuronii a căror funcție de transfer este sigmoidală, valorile de ieșire variază continuu, dar nu liniar.

Rețele neuronale cu transmitere înainte (feedforward)

Rețelele neuronale cu transmitere înainte sunt cele mai populare și cele mai folosite modele în aplicațiile practice. Mai sunt cunoscute și sub denumirea de perceptroni multi-strat. O rețea neuronală cu transmitere înainte este o rețea de elemente similare neuronilor, organizate în straturi, cu o funcție de activare diferențiabilă, de obicei funcția sigmoidală. Algoritmul cu propagare înapoi a erorilor ajustează ponderile pe baza ideii de minimizare a erorii pătratice. Funcția de activare diferențiabilă permite algoritmului cu propagare înapoi să ajusteze ponderile pentru toate straturile de neuroni. Având mai multe noduri în fiecare strat, problemele n-separabile ca și sau-exclusiv (sau XOR), care nu se puteau rezolva doar cu perceptronul, pot fi rezolvate cu rețele cu transmitere înainte. În figura de mai jos se poate observa o rețea neuronală cu transmitere înainte în care, de la intrare până la ieșire, fiecare nod este conectat cu toate nodurile din straturile adiacente:

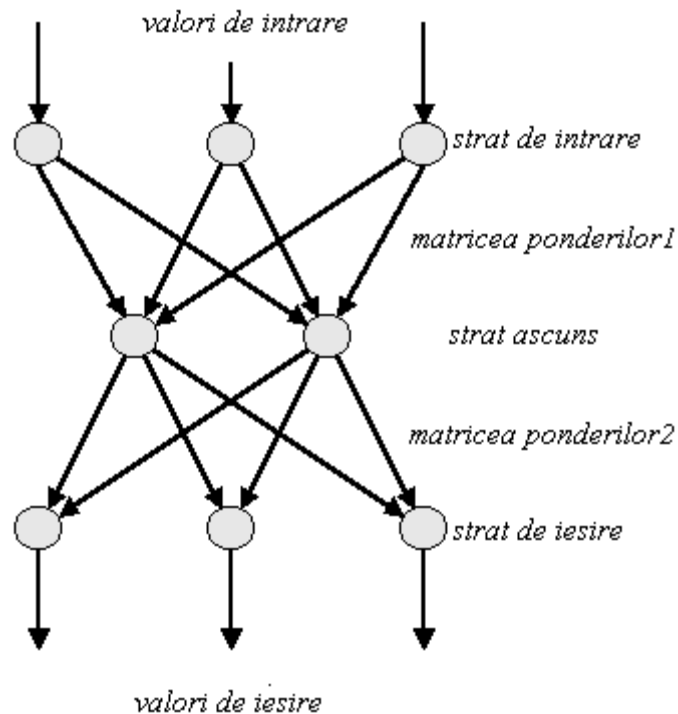


Figura 7.3. Rețea cu transmitere înainte

Fiecare neuron din stratul neuronal de intrare are o valoare de activare care reprezintă un stimul extern. Un neuron de intrare trimite valoarea sa de activare fiecărui neuron din stratul ascuns cu care este conectat. Fiecare neuron al stratului ascuns calculează propria sa valoare de activare depinzând de valorile de activare pe care le primește de la neuronii de intrare. Acest semnal este apoi trimis neuronilor din stratul de ieșire sau următorului strat ascuns de neuroni, în cazul în care există mai multe straturi ascunse. Acești neuroni din straturile ascunse calculează valorile de activare în același fel și le trimit mai departe vecinilor. În cele din urmă semnalul provenind de la neuronii de intrare se propagă prin rețea pentru a determina valorile de activare pentru toți neuronii de ieșire.

Activările neuronilor sunt determinate de ponderile conexiunilor dintre neuroni. Ponderile pot fi pozitive sau negative. O pondere negativă reprezintă inhibarea neuronului care primește informația, de către cel care o trimite. Valoarea de activare pentru fiecare neuron este calculată conform unei funcții de activare simple. Funcțiile de activare variază în detaliu, dar sunt toate conforme aceluiași plan de bază. Funcția însumează contribuțiile tuturor neuronilor din stratul anterior, contribuția unui neuron fiind ponderea conexiunii dintre neuronul din stratul anterior și cel care primește înmulțită cu valoarea de activare a neuronului care trimite. Această sumă este de obicei modificată mai departe, de exemplu prin scalarea la o valoare între 0 și 1, sau prin setarea valorii de activare la 0 dacă nu este atinsă o anumită valoare prag. Deoarece toți neuronii calculează aceeași funcție de activare, realizările intelectuale depind în primul rând de valorile ponderilor dintre neuroni.

Așadar în rețelele neuronale cu transmitere înainte, activarea este transmisă direct de la stratul de intrare la cel ascuns și apoi la cel de ieșire. În modelele mai realiste ale

creierului, însă, există mai multe straturi ascunse de neuroni și de conexiuni recurente care transmit semnalele înapoi de la nivelele mai înalte, la cele mai de jos. Această recurență este necesară pentru a putea explica anumite procese cognitive, cum ar fi memoria pe termen scurt. În rețelele neuronale cu transmitere înainte, prezentarea repetată a aceluiași intrări produce de fiecare dată aceeași ieșire, dar chiar și cele mai simple organisme învață să ignore prezentarea repetată a aceluiași stimul.

Algoritmul de Propagare Înapoi (Backpropagation)

Algoritmul backpropagation este cel mai folosit algoritm de antrenare pentru rețelele neuronale multistrat.

Rețelele neuronale multistrat sunt formate în mod normal din trei sau patru straturi, având un strat de intrare, unul de ieșire și de obicei un strat ascuns, deși în cazuri mai rare sunt necesare două straturi ascunse. Termenul de strat de intrare este impropriu, deoarece acestui strat nu-i este aplicată funcția de activare neuronală, ci valorile de intrare sunt trimise stratului ascuns fără a fi modificate în vreun fel. În calcularea valorii fiecărui neuron din stratul ascuns și din cel de ieșire trebuie mai întâi să calculăm produsul vectorial dintre intrările neuronului și ponderile asociate acestora, iar apoi să aplicăm rezultatului o funcție (de exemplu sigmoidală) pentru a calcula activarea neuronului. Ieșirea neuronului este tocmai această activare.

Rețeaua învață apoi prin modificarea tuturor ponderilor. Calculând derivatele parțiale ale erorii în raport cu fiecare pondere, vom observa direcția de îndreptare a erorii rețelei. Dacă adunăm apoi la veche pondere derivatele corespunzătoare cu semn schimbat, eroarea va scădea. Acest lucru are sens pentru că dacă derivata este pozitivă, înseamnă că eroarea crește odată cu ponderea și de aceea firesc este să adunăm la ponderea respectivă o valoare negativă. Vice versa dacă derivata este negativă. Deoarece calculul acestor derivate parțiale și aplicarea lor fiecărei ponderi se face începând de la startul de ieșire la cel ascuns și apoi de la stratul ascuns la cel de intrare (acest lucru este necesar deoarece pentru modificarea acestui set de ponderi avem nevoie de derivatele parțiale calculate pentru stratul anterior) algoritmul poartă numele de algoritmul cu propagare înapoi (sau backpropagation algorithm).

Antrenarea cu backpropagation

Așadar antrenarea folosind backpropagation presupune transmiterea exemplelor de antrenare ca vectori de intrare prin rețea, calcularea erorii la ieșire și ajustarea ponderilor rețelei pentru a minimiza eroarea.

Antrenarea poate fi oprită dacă eroarea scade sub un anumit prag (de exemplu e bun un prag de 0.001 pentru eroarea pătratică. Acesta variază de la o problemă la alta și în unele cazuri s-ar putea să nici nu se obțină o eroare pătratică mai mică sau egală cu 0.001). Antrenarea excesivă poate avea rezultate dăunătoare deoarece rețeaua poate deveni prea adaptată la exemplele din setul de antrenare și va fi incapabilă să clasifice corect exemple din afara acestuia.

Exemplele din afara setului de antrenare formează setul de validare. Acestea sunt de fapt exemplele pe care putem măsura performanța rețelei. Nu putem aprecia performanța doar bazându-ne pe succesul rețelei în învățarea unui anumit set de

antrenare. Trebuie facute teste pentru a confirma faptul că rețeaua este capabilă să clasifice și exemple din afara setului de antrenare.

Descrierea algoritmului

Primul pas constă în inițializarea ponderilor cu valori aleatoare, de obicei între -0.5 și 0.5, sau între -1 și 1, iar apoi transmiterea vectorului de intrare prin rețea și calcularea valorii corespunzătoare fiecărui neuron din rețea. Aceasta se face, așa cum am mai amintit, prin calcularea produsului vectorial dintre intrările neuronului și ponderile corespunzătoare acestora și apoi aplicarea funcției sigmoideale rezultatului. Fie x vectorul activărilor (ieșirilor) stratului anterior (sau vectorul intrărilor stratului curent). Atunci are loc:

$$o = \sigma(\vec{w}\vec{x})$$
$$\sigma(y) = \frac{1}{1+e^{-y}}$$

Vectorul w este vectorul ponderilor conexiunilor dintre stratul curent și stratul anterior.

Cel de-al doilea pas este calcularea erorii pătratice a rețelei. Aceasta se calculează prin însumarea pătratelor erorilor pentru fiecare neuron din stratul de ieșire. Vectorul țintă care apare în calculul erorii, este asociat exemplului de antrenare (vectorului de intrare).

$$E(\vec{w}) = \frac{1}{2} \sum_{k \in \text{outputs}} (t_k - o_k)^2$$

Valorile t denotă valorile din vectorul țintă, iar o reprezintă activările neuronilor din stratul de ieșire.

Cel de-al treilea pas constă în calcularea erorii pentru fiecare neuron de ieșire, indicată mai jos prin δ .

$$\delta_k \leftarrow o_k(1 - o_k)(t_k - o_k)$$

Această eroare este de fapt derivata parțială a erorii în raport cu fiecare pondere.

Pasul al patrulea este calcularea erorii fiecărui neuron din stratul ascuns.

$$\delta_h \leftarrow o_h(1 - o_h) \sum_{k \in \text{outputs}} w_{kh} \delta_k$$

Eroarea unui neuron din stratul ascuns depinde de erorile calculate pentru neuronii de ieșire.

Pasul al cincilea constă în calcularea valorilor cu care se vor ajusta ponderile. η reprezintă rata de învățare. O rată de învățare mică poate asigura o convergență mai stabilă, pe când o valoare mare a acesteia poate grăbi convergența.

$$\Delta w_{ji} = \eta \delta_j x_{ji}$$

Valoarea x reprezintă neuronul conectat cu neuronul curent prin conexiunea având ponderea w .

Ultimul pas reprezintă adunarea la fiecare pondere a valorilor calculate anterior.

$$w_{ji} \leftarrow w_{ji} + \Delta w_{ji}$$

După terminare, se revine la pasul întâi. Algoritmul se oprește când eroarea devine mai mică sau egală cu pragul stabilit, sau când se atinge un număr maxim de iterații.

Inconveniente ale algoritmului backpropagation

- Timpul de antrenare mare datorat numărului mare de iterații.

- Modificarea lentă a ponderilor.
- Posibilitatea blocării într-un minim local.
- Posibilitate ca rețeaua să nu converge din cauza numărului mare de perechi de vectori de antrenare.

Aplicații ale rețelelor neuronale

Rețelele neuronale sunt utilizate cu succes în acele domenii în care programarea clasică, algoritmică, sau chiar inteligența artificială tradițională cu metode simbolice, au eșuat, sau s-au dovedit ineficiente. Astfel, rețelele neuronale sunt aplicate problemelor al căror algoritm sau model nu este cunoscut, sau a căror rezolvare, folosind algoritmi cunoscuți, durează foarte mult.

Aplicațiile pentru care rețelele neuronale dau rezultatele cele mai bune, cuprind:

- Aproximări de funcții.
- Predicții.
- Clasificări.
- Asocieri de date.
- Optimizări.
- Filtrări de date.
- Luări de decizii.
- Recunoaștere de voce.

și multe alte probleme care intră în categoria recunoașterii unor modele sau tipare.

Cea mai comună utilizare a rețelelor neuronale este pentru proiecția a ceea ce urmează să se întâmple. În multe aplicații predicțiile pot ajuta la stabilirea priorităților. De exemplu în camera de urgență a unui spital este foarte important de stabilit cine are nevoie mai rapid de ajutor medical pentru a permite o operație cu succes. Practic, fiecare organizație trebuie să își stabilească anumite priorități care să guverneze alocarea resurselor acesteia. Această nevoie de proiecție a viitorului a dus la crearea rețelelor pentru predicții.

O altă utilizare importantă și foarte frecventă a rețelelor neuronale este pentru clasificări. O rețea care poate face clasificări, poate fi folosită în industria medicală pentru a prelucra atât rezultate de laborator, cât și înregistrări ale simptomelor pacienților pentru a determina cea mai probabilă boală.

Recent, au fost descrise în presă câteva aplicații de recunoașteri de modele. Unul din acestea, era un sistem care poate detecta bombe în bagaje în aeroporturi, identificând modelele din niște senzori specializați. Un alt articol a raportat cum un medic a antrenat o rețea neuronală feedforward cu date colectate în camere de urgență de la oameni care simțeau că au un atac de cord și dădea ca rezultat probabilitatea ca atacurile de cord să fie reale și nu alarme false.

În continuare vom prezenta câteva exemple de aplicații ale rețelelor neuronale:

▪ Aplicarea rețelelor neuronale în recunoașteri de caractere

Ideea recunoașterii de caractere a devenit foarte importantă, deoarece dispozitivele pentru scrisul de mână au devenit tot mai populare. Utilizarea rețelelor feedforward în recunoașterea caracterelor scrise de mână este mai veche. Ca și în majoritatea cazurilor de învățare supervizată, este prezentat

rețelei un model de intrare, în acest caz matricea caracterului scris de mână, împreună cu ieșirea dorită, aici litera sau cifra corectă. Antrenarea se face cu modele de scris de mână ale utilizatorului.

- **Aplicarea rețelelor neuronale în comprimarea imaginilor**

Pentru că pot accepta un șir mare de date de intrare și le pot procesa rapid, rețelele neuronale sunt utile și în comprimarea imaginilor. O arhitectură a unei rețele neuronale pentru comprimare de imagini are dimensiunea stratului de intrare egală cu cea a stratului de ieșire și un strat intermediar cu o dimensiune mai mică. Rata de comprimare este tocmai raportul dintre dimensiunea stratului de intrare și dimensiunea stratului intermediar. Scopul acestor rețele pentru comprimarea datelor este de a recrea intrarea. Așadar, în antrenare, imaginea folosită ca intrare, va fi și ieșirea dorită. Imaginea va fi o matrice de dimensiunea 256x256, iar intrările rețelei vor fi matrici de dimensiuni 8x8 din imaginea inițială, alese aleator. Acestea sunt prezentate rețelei repetat, iar ponderile sunt ajustate continuu, până când rețeaua reproduce cât mai fidel imaginea.

Odată terminată antrenarea, reconstrucția imaginii este demonstrată în faza de testare. În acest caz, rețelei îi sunt prezentate tot părți de dimensiuni 8x8 din imagine, dar selectate secvențial, de la stânga la dreapta și de sus în jos. Pentru fiecare astfel de matrice de 8x8, este calculat și afișat pe ecran rezultatul pentru a se putea observa vizual performanța rețelei.

- **Rețelele neuronale în predicții financiare**

Rețelele neuronale s-au dovedit a fi instrumente puternice în predicțiile de pe piața de capital. Multe companii au obținut profituri foarte mari folosind rețele neuronale și le-au folosit cu multă ușurință pentru predicții ale prețurilor acțiunilor.

Ideea predicțiilor pe piața financiară nu este nouă. Oamenii de afaceri încearcă deseori să anticipeze prețul pieței, interpretând parametri externi cum ar fi indicatorii economici, opinia publică și climatul politic curent. Întrebarea este dacă rețelele neuronale pot descoperi anumite tendințe în datele de intrare, pe care oamenii s-ar putea să nu le observe și să folosească aceste tendințe cu succes în predicțiile lor. Au fost obținute rezultate bune folosind o rețea relativ simplă, cu doar 6 indicatori financiari ca intrări și cu trei straturi neuronale. Rețeaua a fost antrenată cu algoritmul backpropagation pe un volum mare de date anterioare. Rezultate și mai bune au fost obținute cu o rețea neuronală având 2 straturi ascunse și mai mult de 6 variabile de intrare.

- **Aplicarea rețelelor neuronale în medicină**

Rețelele neuronale sunt un domeniu de cercetare important și în medicină și se crede că vor avea o mare aplicabilitate în sistemele biomedicale în următorii ani. Pentru moment, cercetătorii se ocupă cu modelarea anumitor

părți ale corpului uman și recunoașterea unor boli. Rețelele neuronale sunt ideale în stabilirea de diagnostice, deoarece nu au nevoie de un algoritm precis pentru a identifica o anumită boală. Ele învață din exemple, așa că nu sunt necesare detalii legate de modul de recunoaștere a bolii. Este nevoie doar de un set de exemple reprezentativ pentru toate variațiile bolii. Acestea trebuie alese cu multă grijă pentru ca sistemul să fie cât mai eficient și de încredere.

Unul dintre domeniile care a intrat în atenție în ultima vreme, este diagnosticarea cardiopulmonară. Modul în care rețelele neuronale funcționează în acest domeniu sau în alte domenii ale diagnosticării medicale, este prin compararea mai multor modele diferite. Datele de intrare pot include pulsul, presiunea arterială, respirația, vârsta, sexul, activitatea fizică, etc. pentru diferiți pacienți. Rețeaua neuronală poate învăța prin studierea diferitelor condiții și modele, determinându-le să formeze un tablou conceptual complet, putând apoi să diagnosticheze pacientul bazându-se pe aceste modele.

Rețelele neuronale mai sunt frecvent folosite și în diagnosticarea cancerului la sân. Celulele canceroase sunt examinate în mod tradițional la microscop, de către om, iar acesta decide gradul cancerului prezent. Însă oamenii sunt inconsistenți în astfel de judecăți. Rețelele neuronale au obținut rezultate corecte între 52% și 89% din cazurile noi, neprezentate anterior rețelei.

▪ **Rețele neuronale folosite în împrumuturi**

Acordarea de împrumuturi este un domeniu în care rețelele neuronale pot ajuta oamenii, deoarece nu este un domeniu bazat pe criterii bine definite. Băncile vor să aibă cât mai mulți bani, iar o modalitate de a-i avea este de a diminua rata eșecurilor folosind rețele neuronale pentru a decide dacă să acorde sau nu un împrumut. Rețelele neuronale sunt cu atât mai utile în acest domeniu cu cât nici un proces nu va garanta o precizie de 100%. Chiar o precizie de 85-90% ar fi o îmbunătățire față de metodele folosite de om. În unele bănci rata de eșec în cazul împrumuturilor aprobate folosind rețele neuronale este mai mică decât în cazul celor acordate folosind cele mai bune metode tradiționale.

Procesul funcționează prin analizarea eșecurilor anterioare, iar deciziile curente se bazează pe experiența trecută. Singura problemă este justificarea deciziei luate față de cel care solicită împrumutul. Însă explicarea modului în care a învățat rețeaua și pe ce caracteristici se bazează deciziile luate, este dificilă.

▪ **Folosirea rețelelor neuronale în marketing**

Rețelele neuronale au fost folosite și în aplicații de marketing. Una din aceste aplicații este cea de împărțire a locurilor pentru o companie aeriană. Aplicația este un sistem format din mai multe tehnologii inteligente. În acest sistem este integrată o rețea neuronală feedforward care a fost antrenată

folosind algoritmul backpropagation. Mediul s-a schimbat rapid și constant și de aceea abordarea neuronală adaptivă a fost o soluție bună. Sistemul a fost folosit pentru a monitoriza și a recomanda rezervările pentru fiecare zbor. Astfel de informații au un impact direct asupra profitului unei companii aeriene și dau un avantaj tehnologic pentru utilizatorii sistemului.

7.2.1.3 Învățarea prin arbori de decizie

Învățarea prin arbori de decizie este una dintre cele mai simple și totuși una dintre cele mai de succes metode de învățare. Este o bună introducere pentru aria învățării inductive și este ușor de implementat.

Arborii de decizie primesc la intrare o situație, descrisă printr-un set de proprietăți, și au ca ieșire o decizie de forma "Da / Nu". De aceea arborii de decizie reprezintă funcții booleene. Fiecare nod intern din arbore corespunde unui test al valorii unei proprietăți și fiecare ramură a nodului este etichetată cu o valoare posibilă a testului. Fiecare nod frunză din arbore specifică o valoare booleană ce va fi returnată dacă se va merge pe un drum în arbore ce conduce la acel nod.

Ca un exemplu, să considerăm următoarea problemă: să se aștepte sau nu pentru o masă în restaurant. Predicatul scop este *VomAștepta* iar definiția acestui predicat este exprimată printr-un arbore de decizie. Pentru a putea defini această problemă de învățare mai întâi trebuie să vedem care dintre proprietăți sau *attribute* sunt valabile pentru a putea descrie exemplele din domeniu. Să presupunem că avem următoarea listă de attribute:

1. *Alternativă*: dacă este un restaurant acceptabil în apropiere.
2. *Vineri / Sâmbătă*: dacă este vineri sau sâmbătă.
3. *Bar*: dacă restaurantul are bar unde se poate aștepta.
4. *Apetitul*: dacă clienților le este foame.
5. *Clienți*: câte persoane sunt în restaurant (valorile sunt: *Nimeni*, *Câțiva*, *Plin*).
6. *Preț*: prețurile din restaurant. (valorile sunt: \$, \$\$, \$\$\$)
7. *Ploaie*: dacă plouă afară.
8. *Rezervare*: dacă s-a făcut rezervare sau nu.
9. *Tipul*: tipul restaurantului (valorile sunt Francez, Chinezesc, Italian).
10. *PerioadaDeAșteptare*: perioada de așteptare pentru o masă. (valorile sunt: 0-10 minute, 10-30 minute, 30-60 minute, mai mult de o oră).

Arborele poate fi exprimat ca o conjuncție de implicații individuale care corespund drumurilor prin arbore până la nodurile frunză care au valori "Da". De exemplu, drumul pentru un restaurant plin de clienți, cu un timp de așteptare de 10-30 de minute, când agentului nu îi este foame, este exprimat prin următoarea propoziție logică:

$$\forall \text{Clienți}(r, \text{Plin}) \wedge \text{PerioadaDeAșteptare}(r, 10-30) \wedge \text{Apetit}(r, N) \Rightarrow \text{VomAștepta}(r)$$

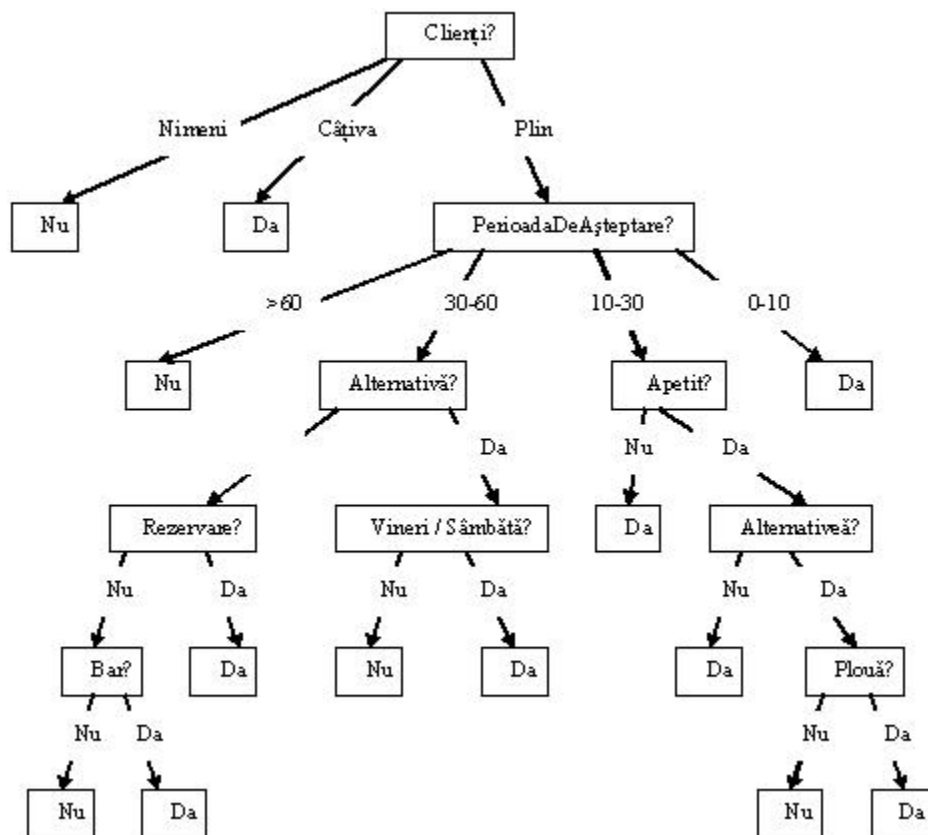


Figura 7.4 Un arbore de decizie pentru a decide dacă se așteaptă sau nu pentru o masă

Arborii de decizie sunt buni pentru anumite tipuri de funcții și nepotriviti pentru alte tipuri de funcții. Nu există nici o reprezentare eficientă pentru toate tipurile de funcții. În cazul în care considerăm un set de funcții booleene cu n atribute sunt tot atâtea tabele de adevăr pentru aceste funcții. Un tabel de adevăr are 2^n coloane deoarece fiecare intrare este descrisă de n atribute, ceea ce înseamnă că sunt 2^n funcții diferite cu n atribute.

De exemplu, dacă avem numai șase atribute booleene, sunt aproximativ $2 * 10^{19}$ funcții diferite. Va fi nevoie de niște algoritmi ingenioși pentru a găsi ipoteze bune într-un spațiu atât de mare.

Ideea de bază pentru algoritmul *Învățarea prin arbori de decizie* este ca prima oară să fie testat cel mai important atribut. Prin cel mai important atribut se înțelege acel atribut care realizează cea mai semnificativă diferență în clasificarea exemplelor. În acest mod, se speră găsirea unei clasificări corecte cu un număr mic de teste, aceasta însemnând că drumurile în arbore vor fi scurte și întregul arbore va fi mic.

Definiție: Se numește *clasificare* a unui exemplu valoarea predicatului scop. Dacă predicatul scop este adevărat pentru un exemplu, acesta se numește *exemplu pozitiv*, altfel se numește *exemplu negativ*.

În Figura 7.5 sunt descrise cinci exemple pentru problema restaurantului. Exemplele pozitive sunt cele pentru care valoarea predicatului *VomAștepta* este adevărat (X_1, X_3, X_4) iar cele negative sunt cele pentru care această valoare este fals (X_2, X_5, X_6). Setul complet de exemple se numește *setul de instruire*.

Exemple	Atribute										Scop
	Alternativă	Bar	Vineri	Apetit	Clienți	Preț	Plouă	Rez.	Tip	Timp	VomAștepta
X_1	da	nu	nu	da	câțiva	\$\$\$	nu	da	Francez	0-10	da
X_2	da	nu	nu	da	plin	\$	nu	nu	Chinezesc	30-60	nu
X_3	nu	da	nu	nu	câțiva	\$	nu	da	Italian	0-10	da
X_4	da	nu	da	da	lin	\$\$\$	nu	da	Chinezesc	10-30	da
X_5	da	nu	da	nu	plin	\$\$	nu	da	Francez	>60	nu
X_6	nu	nu	nu	nu	nimeni	\$	nu	nu	Chinezesc	0-10	nu

Figura 7.5 Exemple pentru problema restaurantului

Avem cele cinci exemple de instruire pe care le-am clasificat în exemple pozitive și exemple negative. Apoi trebuie să decidem care este atributul cel mai semnificativ. Figura 7.5 arată că atributul *Clienți* este un atribut important, deoarece dacă valoarea lui este *câțiva* sau *nimeni* atunci ne rămâne un set pentru care putem da un răspuns; pentru atributul *plin* mai avem nevoie de câteva teste. Se observă că atributul *Tip* este un atribut irelevant, deoarece nu reduce domeniul căutărilor. Se cercetează astfel toate atributele și se alege ca rădăcină a arborelui acel atribut care este cel mai important. Am ales atributul *Clienți*.

După ce a fost ales primul atribut fiecare ieșire este la rândul ei un arbore de decizie, dar cu mai puține exemple și mai puține atribute. Sunt patru cazuri:

- Dacă mai sunt exemple negative și pozitive trebuie ales cel mai bun atribut pentru a le împărți. Se alege *Apetit*.
- Dacă toate exemplele rămase sunt fie pozitive, fie negative, am terminat: se poate răspunde prin *da* sau *nu*. Exemple de felul acesta sunt cazurile *nimeni* și *câțiva*.
- Dacă nu au mai rămas exemple în setul de instruire, înseamnă că nu a mai fost observat nici un exemplu și se va returna o valoare implicită calculată din majoritatea clasificărilor nodurilor părinți.
- Dacă nu au mai rămas atribute, dar nu se poate lua o decizie atunci înseamnă că avem exemple cu aceeași descriere dar clasificări diferite. Acest lucru se întâmplă când unele date sunt incorecte sau când atributele nu oferă destule informații pentru a putea descrie complet situația.

S-ar putea trage concluzia că algoritmul de învățare nu se comportă conform așteptărilor în învățarea funcției corecte. Această concluzie este greșită deoarece algoritmul de învățare ia în considerare exemplele, nu cunoaște funcția corectă.

Algoritmul de învățare în arbori de decizie este descris Figura 7.6.

```

funcția Învățare_Arb(exemple, atribute, implicit) returnează un arbore de decizie
intrări: exemple           //setul de exemple
          atribute          //setul de atribute
          implicit         //valoarea implicită pentru predicatul scop
dacă nu mai există exemple atunci Învățare_Arb ← implicit altfel
dacă toate exemplele au aceeași clasificare atunci Învățare_Arb ← clasificarea altfel
  dacă nu mai există atribute
    atunci
      Învățare_Arb ← Majoritatea_Valorilor(exemple)
    altfel
      cel_mai_bun ← Alege_Atribut(atribute, exemple)
      arbore ← un nou arbore de decizie cu rădăcina cel_mai_bun
      pentru fiecare valoare  $v_i$  din cel_mai_bun execută
        exemplei ← {elementele din exemple cu cel_mai_bun =  $v_i$ }
        subarbore ← Învățare_Arb(exemplei, atribute – cel_mai_bun,
                                   Majoritatea_Valorilor(exemplei))
        @adaugă o nouă ramură la arbore cu eticheta  $v_i$  și subarborele subarbore
      sf-pentru
    sf-dacă
  sf-dacă
sf-dacă
  Învățare_Arb ← arbore
sf-Învățare_Arb

```

Figura 7.6. Algoritmul de învățare în arbori de decizie

7.2.2 Învățarea Prin Întărire

În această secțiune vom aborda problematica generală a învățării prin întărire, fără a intra în amănunte legate de algoritmi specifici.

Problema învățării prin întărire este problema generală de îmbunătățire a comportamentului unui agent artificial pe baza unui feed-back al performanței sale.

Vom studia în acest capitol modul în care agenții pot învăța în condiții mai puțin favorabile, când ei nu primesc exemple și pornesc fără un model al mediului înconjurător și fără a avea o funcție de performanță. *Învățarea prin întărire* (*reinforcement learning*) este modalitatea de a îmbunătăți comportarea unui agent pe baza unui anumit feed-back primit în legătura cu performanța sa.

Învățarea prin întărire este o abordare a învățării care combină două discipline în scopul rezolvării unor probleme pe care nici una nu o poate rezolva separat.

- *Programarea dinamică* este un domeniu al matematicii utilizat îndeosebi în probleme de optimizare și control.
- *Învățarea supervizată* este o metodă generală de instruire a unor rețele (care aproximează funcții – cum ar fi rețelele neuronale) pentru a reprezenta funcții.

În esență, învățarea supervizată necesită perechi (intrare-ieșire) pentru funcția ce urmează a fi învățată (un set de întrebări cu răspunsurile corespunzătoare).

Într-un sistem de învățare prin întărire calculatorului i se dă pur și simplu sarcina pe care trebuie să o realizeze. Sistemul va învăța cum să-și realizeze sarcina (să-și atingă scopul) prin încercări și eventuale feedback-uri (răspunsuri) pe care le primește de la mediu.

Un exemplu este un agent care învață să joace șah fără a fi supervizat, adică fără a avea niște exemple de situații de joc împreună cu cele mai bune mutări alese în situațiile respective. În situația în care nu există un supervizor care să-i furnizeze exemplele, ce va face agentul? Încercarea aleatoare a mutărilor, ar fi total inefficientă, atâta timp cât numărul de mutări la un moment dat este mare, procesul devenind exponențial. Fără un feedback despre ce e bine și ce e rău, agentul nu ar putea să decidă la un moment dat cea mai bună mutare care ar trebui aleasă. Există totuși un fel de feedback pe care-l obține agentul la sfârșitul jocului – chiar în absența unui supervizor –, când percepe dacă a câștigat sau a pierdut. Acest feedback poartă numele de *recompensă* (*reward*) sau *întărire*, iar agentul îl primește fie la sfârșit, într-o stare finală, fie pe parcurs când are informații exacte despre ce a făcut bine sau nu.

Rolul *învățării prin întărire* este de a folosi recompense pentru învățarea cu succes a unui agent funcțional.

Învățarea prin întărire poate fi văzută ca un microcosmos pentru toate problemele de Inteligență Artificială. Un agent primește niște percepții de la mediul în care se află, le marchează ca fiind utilități pozitive sau negative și apoi se decide ce acțiune să execute (Figura 7.7).

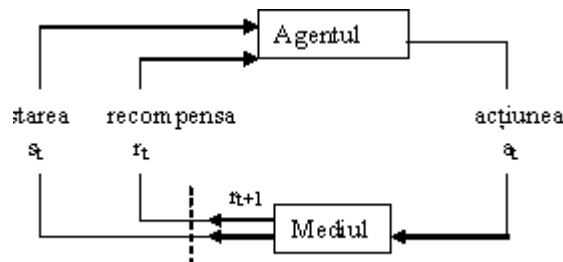


Figura 7.7. Interacțiunea dintre agent și mediu în învățarea prin întărire

La fiecare pas de timp t , agentul primește informații despre starea în care se află, $s_t \in S$, unde S este mulțimea stărilor posibile, și pe baza acestei informații selectează acțiunea $a_t \in A(s_t)$ unde $A(s_t)$ este mulțimea de acțiuni disponibile din starea s_t . La următorul pas, în parte ca o consecință a acțiunii efectuate, agentul primește o recompensă numerică $r_{t+1} \in R$, și se află într-o nouă stare s_{t+1} . La fiecare pas de timp agentul efectuează o mapare de la stări la probabilitățile de a alege o acțiune posibilă. Această mapare se numește *politica* agentului notată cu π_t , unde $\pi_t(s, a)$ este probabilitatea ca $a = a_t$ și $s = s_t$. Metodele de învățare prin întărire specifică cum ar trebui agentul să-și modifice politica ținând cont de experiența acumulată.

Învățarea prin întărire ne arată cum poate învăța un agent și cum poate deveni expert într-un mediu necunoscut, când i se dau numai recompense ocazionale. Scopul Învățării prin Întărire este de a folosi recompensa primită pentru ca agentul să învețe să facă acțiuni corecte. Recompensa e un număr și valoarea acesteia variază de la pas la

pas. Scopul agentului este de a maximiza totalul recompenselor primite. Agentul trece printr-o succesiune de stări și primește o anumită recompensă. Aceasta poate fi recompensa de câștig (+1 de exemplu), de eșec (-1) sau nulă când nu s-a ajuns la o stare terminală (0).

Agentul învață să-și maximizeze recompensele primite. Dacă dorim ca agentul să facă un anumit lucru trebuie să stabilim recompensele astfel încât agentul să ajungă acolo unde vrem maximizând aceste recompense. Dar recompensa care îi va fi atribuită agentului pentru o acțiune trebuie să-i spună acestuia ce dorim ca el să obțină nu cum să obțină. Învățarea prin Întărire e sinonimă cu învățarea prin interacțiune. În timpul învățării sistemele adaptive încearcă aplicarea de acțiuni asupra mediului în care se află, iar apoi sunt recompensate în funcție de acțiunea pe care o execută. Algoritmii de învățare prin întărire rețin selectiv ieșirile care maximizează recompensa primită. În multe medii complexe învățarea prin întărire este singura metodă de a antrena programe să îndeplinească sarcini cu un grad mare de dificultate. De exemplu, la jocuri este extrem de greu pentru un antrenor uman să dea evaluări exacte și consistente ale unui număr mare de poziții, care ar fi necesare pentru a antrena o funcție de evaluare direct din exemple. În loc de acest lucru, programului i se poate spune când a câștigat sau a pierdut, și acesta poate folosi aceste informații pentru a-și construi o funcție de evaluare care estimează aproape exact probabilitatea de reușită sau eșec dintr-o anumită poziție. Un agent primește percepții de la mediul în care se află, le clasifică în utilități pozitive și negative, iar apoi se decide ce acțiune să aleagă.

O problemă de învățare prin întărire are trei părți fundamentale:

- *mediul* – reprezentat sub formă de stări.
- *funcția de întărire* – funcția exactă a recompenselor pe care agentul caută să le maximizeze. Agentul va învăța să execute acele acțiuni care să maximizeze suma recompenselor primite pe drumul de la o stare inițială spre o stare finală. Există trei clase utilizate pentru construirea unor funcții de întărire corespunzătoare scopului urmărit.
 - **Probleme de evitare** sau de **recompensă amânată** – recompensele sunt 0, cu excepția stărilor terminale. Semnul recompensei într-o stare terminală va indica dacă acea stare e finală (succes) sau trebuie evitată (eșec) – probleme de teoria jocurilor.
 - **Timp minim de obținere a soluției** – scopul este ca agentul să aleagă acele acțiuni care să genereze cel mai scurt drum către o stare finală. În acest caz funcția de întărire va fi -1 pentru toate tranzițiile dintre stări, cu excepția tranziției spre o stare finală, caz în care recompensa va fi 0. Maximizarea recompenselor se va reduce în acest caz la minimizarea numărului de tranziții (a lungimii drumului spre o stare finală).
 - **Jocuri** – în cazul jocurilor cu doi sau mai mulți adversari, se va folosi o funcție de întărire alternantă. Sistemul va învăța să genereze comportări optime ale jucătorilor, în sensul maximizării sau minimizării funcției de întărire. Problemele de joc se reduc de fapt la căutări în arbori (de joc) în care alternează două nivele: cel maximizant corespunzător primului jucător, respectiv cel minimizant corespunzător celui de-al doilea jucător.
- *funcția de utilitate (valoare)* – **Strategia de control** va desemna ce acțiuni vor trebui efectuate în fiecare stare. **Valoarea** (utilitatea) unei stări se va defini ca fiind suma

recompenselor primite pe drumul de la starea curentă către o stare finală. Scopul va fi maximizarea sumei recompenselor obținute pe drumul de la o stare inițială arbitrară către o stare finală, în stările accesibile pe acest drum. *Funcția de utilitate* va fi deci o asignare de valori stărilor și va putea fi aproximată folosind orice tip de aproximator de funcții (perceptron, etc.).

Marea problemă ce se pune în general în cercetările de învățare prin întărire este: Care este șablonul unui algoritm care să găsească valorile optime ale funcției de utilitate? Dintr-un anumit punct de vedere *învățarea prin întărire* este o nouă exprimare a problemelor de Inteligență Artificială. Un agent într-un mediu primește semnale, le atașează unora utilități pozitive sau negative și apoi va trebui să decidă ce acțiune să aleagă. Sarcina de *învățare* poate varia în funcție de anumite situații:

- Mediul poate fi accesibil sau inaccesibil. Într-un mediu accesibil, stările pot fi identificate cu percepții, pe când într-un mediu inaccesibil agentul va trebui să încerce să rețină un model al mediului;
- Agentul poate începe cu o anumită cunoaștere despre mediu și despre efectele acțiunilor sale; sau va trebui să învețe acest model la fel ca și informațiile de utilitate;
- Recompensele pot fi recepționate doar în stările finale sau în orice altă stare;
- Recompensele pot fi componente ale funcției de utilitate (puncte pentru un agent care joacă ping-pong sau dolari pentru un agent de pariuri) pe care agentul încearcă să o maximizeze, sau pot fi sugestii despre funcția de utilitate (“mutare bună” sau “mutare rea”);
- Agentul poate învăța **pasiv** sau **activ**. Un agent pasiv se rezumă doar la a observa mediul și a învăța utilitatea de a fi în diverse stări; un agent activ va trebui să și acționeze conform cu informația ce a învățat-o și poate să-și utilizeze generatorul de probleme pentru a sugera explorarea unor părți necunoscute ale mediului.

Pe de altă parte, deoarece agentul va primi recompense referitoare la utilități, există două configurații de bază:

- Agentul învață o funcție de utilitate (asociată stărilor) și o folosește pentru a selecta acțiunea care să maximizeze performanța preconizată a rezultatelor.
- Agentul învață **valorile acțiunilor (action-value function)**, furnizând performanța preconizată în urma alegerii unei anumite acțiuni într-o anumită stare. Acest mod de învățare se numește **Q-învățare**.

Un agent care învață funcții de utilitate va trebui să aibă un model al mediului, pentru a putea aplica deciziile luate, va trebui să cunoască stările în care se ajunge în urma aplicării acțiunilor. În schimb, un agent care învață valorile acțiunilor, nu are nevoie de un model al mediului; atâta timp cât își cunoaște mutările permise, poate compara direct valorile acestora fără a trebui să cunoască stările rezultate în urma aplicării mutărilor. Din acest punct de vedere agenții din a doua categorie sunt mai ușor de realizat decât cei din prima categorie. Pe de altă parte însă, pentru că nu cunosc unde conduc acțiunile lor, agenții care învață valori de acțiune, nu pot reveni asupra acțiunilor lor, ceea ce este o serioasă restricție asupra abilității lor de a învăța.

Un exemplu clasic de învățare prin întărire este **problema banditului cu n brațe**, problemă ce va fi descrisă în continuare.

Să considerăm următoarea problemă de învățare: suntem puși adesea în situația de a alege între n variante sau acțiuni. După fiecare acțiune primim o recompensă numerică aleasă dintr-o distribuție de probabilitate staționară care depinde de acțiunea selectată. Obiectivul este de a maximiza recompensa totală obținută după o anumită perioadă de timp (de exemplu după 100 de acțiuni selectate). Fiecare selectare a acțiunii se numește *joc*.

Fiecare acțiune are o recompensă acordată dacă acțiunea este selectată. Aceasta se numește *valoarea* acțiunii respective. Dacă s-ar cunoaște valoarea fiecărei acțiuni problema banditului cu n brațe ar fi foarte ușor de rezolvat. Dacă fiecare acțiune are asociată o anumită recompensă, atunci există o acțiune a cărei valoare a recompensei este cea mai mare. Selectarea acțiunii a cărei recompensă estimată este cea mai mare se numește *selecția greedy* (și se realizează prin metoda *greedy*). Dacă se selectează acțiunea care maximizează recompensa atunci se *exploatează* cunoștințele obținute despre valorile acțiunilor. Dacă nu se selectează acest tip de acțiune spunem că se *exploarează*. Exploatarea are avantajele ei, dar explorarea ar putea produce un rezultat mai bun pe termen lung. Recompensa este mai mică pe termen scurt, în timpul explorării, dar pe termen lung, după ce au fost descoperite acțiunile mai bune acestea pot fi exploatate.

Deoarece nu este posibil să exploatăm și să explorăm cu un singur selector de acțiuni, apare un conflict între explorare și exploatare.

În cazuri practice, decizia de a explora sau exploata depinde de valorile precise ale estimărilor, de incertitudini și de numărul de jocuri rămase de jucat. Sunt mai multe metode de combinare a exploatării cu explorarea pentru formulări particulare ale problemei banditului cu n brațe și ale problemelor înrudite cu aceasta.

În problemele învățării prin întărire, agenților li se face o descriere a stării curente și trebuie să aleagă următoarea acțiune dintr-un set de acțiuni posibile pentru a maximiza feedback-ul primit după fiecare acțiune.

În afară de agent și mediu, pot fi identificate patru principale elemente ale sistemului de învățare prin întărire: o politică, o funcție recompensă, o funcție de evaluare și un model al mediului.

- **Politica (π)** definește comportamentul agentului la un moment dat, este o funcție de la stările percepute ale mediului la acțiunile ce trebuie efectuate din aceste stări. Ea corespunde la ceea ce în psihologie se numește set de reguli stimul-răspuns. În unele cazuri politica poate fi o funcție simplă, pe când în altele poate implica calcule extrem de complicate (de exemplu într-un proces de căutare). Politica este nucleul unui agent de învățare prin întărire în sensul că este suficientă pentru a determina comportamentul acestuia. O politică π este o mapare de la starea $s \in S$ și acțiunea $a \in A(s)$ la probabilitatea $\pi(s,a)$ de a efectua acțiunea a din starea s .
- **Funcția recompensă** definește scopul, atribuie fiecărei stări percepute (sau pereche stare-acțiune) un singur număr, o recompensă, semnificând dezirabilitatea stării (perechii) respective. Obiectivul fundamental al unui agent de învățare prin întărire este maximizarea recompensei primite pe termen lung. Funcția recompensă semnalează agentului care sunt acțiunile bune și care sunt cele greșite. Recompensele sunt caracteristicile fundamentale

ale problemei cu care se confruntă agentul și de aceea funcția de recompensă trebuie să fie fixă. Ea poate fi însă folosită pentru a modifica politica. De exemplu, dacă o acțiune selectată de politică este urmată de o recompensă mică, atunci politica se poate schimba astfel încât să fie aleasă o altă acțiune în momentul în care apare din nou aceeași situație.

Funcția recompensă este reprezentată de funcția $Q: S \times A \rightarrow R$ și estimează recompensa pentru fiecare pereche stare-acțiune. Deci, valoarea efectuării acțiunii a din starea s sub politica π , notată $Q^\pi(s, a)$, este recompensa așteptată când se pleacă din starea s și se urmează politica π :

$$Q^\pi(s, a) = E_\pi\{R_t \mid s_t = s, a_t = a\} = E_\pi\{\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \mid s_t = s, a_t = a\}$$

Q se mai numește funcție acțiune-valoare pentru politica π .

- În timp ce funcția recompensă semnalează ce e bun într-un moment imediat, **funcția de evaluare** semnalează ce este bun pe termen lung. Valoarea unei stări este recompensa totală pe care agentul se poate aștepta să o primească de-a lungul timpului, plecând din acea stare. O stare poate produce o recompensă imediată mică, dar să aibă totuși o valoare mare deoarece este de obicei urmată de stări ce produc recompense mari (sau reciproc). Între recompense și valori există o strânsă legătură în sensul că fără recompense nu ar exista valori și unicul motiv pentru care valorile sunt estimate este obținerea de recompense mai mari. Cu toate acestea, valorile sunt cele importante atunci când se iau decizii. Se aleg acțiunile care duc în stări cu valoare mare, nu cu recompensă mare, deoarece aceste acțiuni vor duce, pe termen lung, la o recompensă totală maximă. Din păcate, este mult mai greu a determina valorile decât recompensele. Recompensele sunt date direct de mediu, pe când valorile trebuie estimate și reestimate din secvența de observații pe care agentul le face de-a lungul întregii sale vieți. De fapt, cea mai importantă componentă a majorității algoritmilor de învățare prin întărire este o metodă de estimare eficientă a valorilor.

Valoarea unei stări s sub o politică π , notată $V^\pi(s)$ este recompensa așteptată plecând din s , având în vedere faptul că agentul urmează politica π , iar t este un pas oarecare de timp:

$$V^\pi(s) = E_\pi\{R_t \mid s_t = s\} = E_\pi\{\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \mid s_t = s\}$$

- Cel de-al patrulea element al unui sistem de învățare prin întărire este un **model al mediului**, o reprezentare internă, de obicei simplificată, a mediului utilizată de agent pentru a prezice ce s-ar putea întâmpla ca rezultat al acțiunilor viitoare. Un agent mai complicat ar putea utiliza acest model pentru a-și planifica viitorul curs al acțiunilor.

Mediul celui care învață poate fi modelat de un timp discret, de o stare finită și de un proces de decizie Markov care poate fi reprezentat prin cvadruplul (S, A, P, r) unde S este un set de stări, A este un set de acțiuni, iar $P: S \times S \times A \rightarrow [0, 1]$ indică probabilitatea de a trece din starea s_1 în starea s_2 întreprinzând acțiunea a și $r: S \times A \rightarrow R$ este o funcție cu rezultat scalar, funcția recompensă. Fiecare agent păstrează o politică π care translatează starea curentă în acțiunea dorită a fi realizată în acea stare.

Q-Learning

Au fost propuse mai multe strategii pentru învățarea prin întărire care pot fi folosite de agenți în dezvoltarea unei politici de maximizare a rezultatelor pozitive acumulate de-a lungul timpului. Pentru a evalua sistemul de clasificare paradigmatic pentru învățarea prin întărire în sistemele multiagent, acesta este comparat cu algoritmul Q-Learning, care este conceput pentru a găsi o politică π care maximizează $V^\pi(s)$ pentru toate stările s din S . Valorile Q sunt definite ca $Q^\pi(s, a) = V^{a:\pi}$, unde cu $a:\pi$ se notează secvența de alegere a acțiunii a în starea curentă, urmată de alegerea acțiunilor bazate pe politica π . Acțiunea a în starea s este aleasă în așa fel încât să se maximizeze recompensa:

$$V^{\pi^*}(s) = \max_{a \in A} Q^{\pi^*}(s, a), \quad \forall s \in S.$$

Politica π^* se numește politică optimă și este mai mare sau egală decât toate celelalte politici. Deși pot fi mai multe astfel de politici toate au aceeași funcție de evaluare, numită funcție valoare optimă, notată V^* , definită prin:

$$V^*(s) = \max_{\pi} V^\pi(s), \quad \forall s \in S.$$

Politicile optime au și aceeași funcție acțiune-valoare, notată Q^* , definită prin:

$$Q^*(s, a) = \max_{\pi} Q^\pi(s, a), \quad \forall s \in S, a \in A(s).$$

Dacă o acțiune a într-o stare s produce o întărire a lui R și o tranziție la starea s' , atunci valoarea Q corespunzătoare se modifică astfel:

$$Q(s, a) \leftarrow (1 - \beta)Q(s, a) + \beta(R + \gamma * \max_{a' \in A} Q(s', a')),$$

unde β este o constantă mică, numită rată de învățare.

7.2.3 Învățarea nesupervizată

Învățarea care are loc fără ca agentul să primească vreo o sugestie despre semnalele de ieșire corecte se numește **învățare nesupervizată**.

În ciuda numeroaselor sale aplicații, învățarea supervizată a fost criticată ca fiind neplauzibilă din punct de vedere biologic. Aceste critici pornesc de la observația că este greu de conceput existența în creier a unui mecanism de instruire care ar compara ieșirile dorite cu cele reale, propagând înapoi în toată rețeaua de neuroni corecțiile efectuate. Dacă acesta ar fi mecanismul de instruire se pune problema provenienței formelor de ieșire dorite.

Pe de altă parte s-a pus în evidență o autoorganizare a creierului încă din stadiile timpurii ale dezvoltării sale. Această autoorganizare nu poate fi explicată satisfăcător ca rezultatul unui mecanism de învățare supervizată. Învățarea nesupervizată este un model de învățare mult mai plauzibil pentru sistemele biologice. În general, un model de învățare nesupervizată este unul în care ajustarea ponderilor nu se bazează pe compararea cu răspunsuri ideale predeterminate. Mulțimea de instruire constă doar în vectorii de intrare.

În general algoritmi de clasificare nesupervizată modifică ponderile rețelei pentru a produce vectori de ieșire care sunt consistenți. Prin aceasta înțelegem faptul că aplicarea a doi vectori de intrare care sunt suficient de asemănători va genera aceeași formă de ieșire (sau două forme foarte apropiate). Procesul de învățare grupează vectorii similari în clase. Prezintănd rețelei un vector dintr-o clasă aceasta va genera un vector de ieșire specific. Răspunsul pe care îl va produce o clasă de vectori de intrare nu poate fi însă determinat înainte de încheierea procesului de învățare. Prin urmare ieșirile unei astfel de rețele trebuie, în general, să fie interpretate, adică să primească o formă inteligibilă după procesul de învățare.

7.2.3.1 Clusterizare

Clusterizarea, cunoscută și sub numele de *clasificare nesupervizată* este procesul împărțirii unui grup (set) de date în grupuri de obiecte similare în raport cu o mulțime de atribute relevante ale obiectelor analizate. Aceste grupuri vor fi numite *clusteri*.

Clusterizarea, din perspectiva modelării datelor, este în strânsă legătură cu matematica, statistica și analiza numerică. Din perspectiva învățării automate, clusterii corespund unor *șabloane ascunse*, iar procesul de căutare (identificare) a clusterilor este învățare nesupervizată.

Dintr-o perspectivă practică procesul de clusterizare joacă un rol deosebit de important în aplicații de *Data Mining* (procesul de descoperire de șabloane în date) cum ar fi: explorări de date științifice, regăsirea informației (*Information Retrieval*), identificarea de șabloane în texte (*Text Mining*), aplicații cu baze de date spațiale, analiză a informațiilor Web, marketing, diagnostic medical și multe alte aplicații.

În continuare vom prezenta pe scurt problematica generală a clusterizării.

Se consideră un set de date $O = \{O_1, O_2, \dots, O_n\}$ care constă într-o mulțime de puncte (obiecte, instanțe, modele, tuple, tranzacții, etc). Fiecare obiect O_i este caracterizat în raport cu o mulțime $A = \{A_1, A_2, \dots, A_m\}$ de atribute, $O_i = \{O_{i1}, O_{i2}, \dots, O_{im}\}$, unde O_{il} este valoarea atributului A_l pentru obiectul O_i .

Scopul clusterizării este asignarea obiectelor unei mulțimi finite de k submulțimi, numite *clusteri*, astfel încât:

$$O = C_1 \cup C_2 \cup \dots \cup C_k, \quad C_i \cap C_j = \emptyset, \forall i \neq j$$

Scopul partiționării obiectelor este ca obiectele din cadrul unui cluster să fie mai similare între ele decât cu obiectele din alt cluster. Similaritatea și disimilaritatea între obiecte este calculată folosind *metrice* sau *semi-metrici* aplicate valorilor atributelor care caracterizează obiectele. Spre exemplu, distanța între două obiecte poate fi exprimată ca fiind *distanța euclidiană* între acestea (una din metricile des folosite în clusterizare):

$$d(O_i, O_j) = d_E(O_i, O_j) = \sqrt{\sum_{l=1}^m (O_{il} - O_{jl})^2}$$

Similaritatea între două obiecte va fi calculată ca inversa distanței între cele două obiecte:

$$\text{sim}(O_i, O_j) = \frac{1}{d(O_i, O_j)}$$

Fără a intra în amănunte legate de algoritmi de clusterizare, amintim principalele metode de clusterizare existente:

- metode ierarhice;
- metode de partiționare (algoritmul *k-means* este unul dintre algoritmi de clusterizare din această grupă);
- clusterizare bazată pe constrângeri;
- algoritmi de clusterizare folosiți în învățare automată:
 - rețele neuronale
 - metode evolutive.

7.3 Învățarea în Inteligența Artificială Distribuită

7.3.1 Introducere

Învățarea și inteligența sunt strâns legate una de alta. Se spune de obicei că un sistem capabil să învețe merită să fie numit inteligent, și reciproc, de la un sistem considerat inteligent se așteaptă, printre altele, să fie capabil să învețe. Învățarea presupune întotdeauna auto-îmbunătățirea comportamentului viitor pe baza experiențelor precedente. Mai precis, din punct de vedere al inteligenței artificiale (AI), învățarea poate fi definită astfel:

Achiziționarea de noi cunoștințe și aptitudini cognitive și încorporarea acestora în activitățile viitoare ale sistemului cu condiția ca achiziționarea și încorporarea să fie făcute de către sistem și să ducă la o îmbunătățire a performanțelor acestuia.

Învățarea în sistemele multiagent reprezintă un subdomeniu important atât al inteligenței artificiale distribuite (IAD), cât și al mașinilor auto-instruibile (machine learning - ML). Mașinile auto-instruibile reprezintă o direcție distinctă de cercetare în IA și se ocupă de aspectele computaționale ale învățării atât în sistemele naturale cât și în sistemele tehnice.

Intersecția dintre IAD și ML o constituie o tânăra, dar importantă arie de cercetare și aplicații. Comunitățile IAD și ML au ignorat în mare măsură această arie o perioadă îndelungată (există excepții de ambele părți care doar confirmă regula). Pe de o parte, munca în IAD era în principal preocupată de sistemele multiagent a căror organizare structurală și comportamental funcțională era determinată în detaliu și de aceea era mai mult sau mai puțin fixată. Pe de altă parte, munca în ML trata la început învățarea ca un proces centralizat și izolat care are loc în sistemele inteligente stand-alone. Această ignorare reciprocă a IAD și ML a dispărut, și astăzi aria învățării în sistemele multiagent se bucură de o largă atenție într-o constantă creștere. Există două motive principale care justifică această atenție, ambele ilustrând importanța unificării IAD și ML:

- există o mare nevoie să echipăm sistemele multiagent cu abilități de învățare;

- o privire de ansamblu asupra ML care cuprinde nu doar învățarea monoagent, ci și cea multiagent poate duce la o înțelegere îmbunătățită a principiilor generale, fundamentale ale învățării atât în sistemele naturale cât și în cele computaționale.

Primul motiv se bazează pe faptul că sistemele multiagent trebuie să reacționeze în medii complexe, mari, deschise, dinamice și imprevizibile. Pentru astfel de medii este deosebit de dificil și uneori chiar imposibil să se specifice sistemelor corect și complet o prioritate la momentul creerii lor și înainte de a fi folosite. Pentru asta este nevoie, de exemplu, să se cunoască prioritatea pentru condițiile ambientale care se vor ivi, agenții care vor fi disponibili în acel moment și cum vor reacționa și interacționa agenții disponibili ca răspuns la aceste condiții. Singurul mod posibil de a face față acestei dificultăți este înzestrarea agenților individuali cu abilitatea de a-și îmbunătăți propriile performanțe și pe cele ale sistemului în general. Al doilea motiv reflectă faptul că învățarea în sistemele multiagent nu este doar o sporire a învățării în sistemele stand-alone și nici doar suma activităților de învățare izolate a câtorva agenți. Învățarea în sistemele multiagent cuprinde învățarea în sisteme stand-alone pentru că un agent poate învăța individual și complet independent de alți agenți. Mai mult, învățarea în sistemele multiagent extinde învățarea în sisteme stand-alone. Asta pentru că activitățile de învățare ale unui agent individual pot fi influențate considerabil (întârziate, accelerate, redirectate sau făcute posibile) de alți agenți și pentru că mai mulți agenți pot lucra într-un mod distributiv și interactiv ca un întreg coerent. Această viziune extinsă a învățării este diferită calitativ de viziunea tradițională asupra ML și are capacitatea de a provoca impulsuri în cadrul cercetării care să conducă la noi tehnici și algoritmi în domeniul mașinilor auto-instruibile.

Învățarea în sistemele multiagent este un fenomen cu mai multe laturi și de aceea nu ne surprinde faptul că mulți termeni ce se găsesc în literatură se referă la acest tip de învățare accentuând diferite laturi. Exemple de astfel de termeni sunt: învățarea reciprocă, învățarea prin cooperare, învățarea prin colaborare, învățarea în echipă, învățarea socială, învățarea împărțită, învățarea pluralistă, învățarea organizatorică. Scopul acestei secțiuni este de a explicita mai multe astfel de laturi, oferind o caracterizare generală a învățării în sistemele multiagent. Aceasta se face prin descrierea, din punct de vedere al sistemelor multiagent, a categoriilor principale de învățare și a trăsăturilor de bază prin care abordările în învățare pot diferi.

7.3.2 Învățarea Cooperativă în SMA

În ultimii ani a fost manifestat un interes deosebit pentru abordările descentralizate în rezolvarea problemelor complexe din lumea reală.

Învățarea în sistemele multiagent este aplicarea învățării automate în probleme care implică mai mulți agenți. Sunt două aspecte de bază ale învățării în sisteme multiagent, și anume:

- datorită faptului că învățarea în SMA abordează domenii de probleme ce implică mai mulți agenți, spațiul de căutare implicat poate fi deosebit de mare;
- învățarea în SMA *poate* implica *mai multe entități care învață*, fiecare dintre acestea învățând și adaptându-se în contextul celorlalți.

Anumite abordări din literatura de specialitate se referă la o singură entitate care învață pentru a îmbunătăți performanța întregii echipe de agenți. Această categorie de învățare se numește **învățare în echipă** deoarece procesul de învățare ajustează comportamentul echipei privită ca un întreg.

Alte abordări existente în literatură se referă la procese individuale de învățare pentru agenții individuali, urmărind în același timp și calitatea echipei de agenți. Această abordare se numește **învățare concurentă**.

Ca urmare, se disting două categorii principale de învățare în sistemele multiagent:

- învățarea centralizată (izolată);
- învățarea descentralizată (interactivă).

Pentru a clarifica ce tipuri de învățare sunt cuprinse în cele două categorii, introducem noțiunea de proces de învățare:

Procesul de învățare se referă la toate activitățile (ex.: planificare, deducție, pași de decizie) care sunt executate cu intenția de a atinge un anumit scop.

Învățarea este centralizată dacă procesul de învățare este executat sub toate aspectele lui de un singur agent și nu are nevoie de interacțiunea cu alți agenți; cel care învață se comportă ca și cum ar fi singur. Învățarea este descentralizată dacă mai mulți agenți sunt implicați în același proces de învățare. Asta înseamnă că activitățile ce constituie procesul de învățare sunt executate de agenți diferiți. În contrast cu învățarea centralizată, învățarea descentralizată se bazează sau chiar are nevoie de prezența mai multor agenți capabili să îndeplinească anumite activități.

În sistemele multiagent mai mulți indivizi care învață și încearcă să atingă diferite sau chiar aceleași scopuri pot fi activi în același timp. Similar, pot fi mai multe grupuri de agenți implicați în diferite procese de învățare descentralizată. Mai mult, scopurile urmărite de acești agenți pot fi diferite sau aceleași. Este important să vedem dacă un singur agent poate fi implicat în câteva procese de învățare centralizate și/sau descentralizate în același timp.

7.3.2.1 Învățarea în echipă

În *învățarea în echipă* este implicată o singură entitate care învață: dar această entitate descoperă o mulțime de comportamente mai degrabă pentru o echipă de agenți, decât pentru un singur agent.

Învățarea în echipă este o abordare simplă a învățării în SMA, deoarece se pot folosi tehnici standard de învățare (pentru un singur agent): există o singură entitate care efectuează procesul de învățare.

Un alt avantaj al a unei singure entități care învață este faptul că agenții tind să acționeze mai degrabă în scopul maximizării recompensei echipei, decât a recompenselor lor individuale. Acest lucru face ca agenții să se comporte mai degrabă egoist decât lacom.

Învățarea în echipă are, însă, și anumite dezavantaje: O problemă majoră este spațiul stărilor care crește foarte mult în timpul procesului de învățare. Spre exemplu, dacă agentul **A** poate fi într-una din 100 de stări și agentul **B** poate fi într-un a din alte 100 de stări, echipa formată din cei doi agenți poate fi în 10000 de stări posibile. Această explozie a spațiului de stări poate fi copleșitoare pentru metode de învățare care

explorează spațiul utilităților stărilor (cum ar fi învățarea prin întărire), dar poate să nu afecteze drastic tehnici care exponează spațiul comportamentelor (cum ar fi programarea evolutivă).

Un al doilea dezavantaj major al învățării în echipă este centralizarea algoritmului de învățare: toate resursele trebuie să fie disponibile într-un singur loc în care toate caculele sunt efectuate. Acest lucru poate fi un dezavantaj în domenii în care datele sunt distribuite.

Învățarea în echipă poate fi împărțită în două categorii: învățare *omogenă* și *eterogenă*. Entitățile omogene care învață dezvoltă un singur comportament care este folosit de fiecare membru din echipă. Entitățile eterogene pot dezvolta comportamente unice pentru fiecare agent. Există, deasemenea, o abordare hibridă între cele două.

7.3.2.2 Învățarea concurentă

Cea mai cunoscută alternativă la învățarea în echipă în cadrul sistemelor cooperative multiagent este învățarea concurentă, *în care fiecare agent* din echipă învață independent cum să își îmbunătățească performanțele și performanța echipei ca un întreg.

Cel mai mare avantaj al învățării concurente este faptul că împarte spațiul de căutare reunit al echipei în spații individuale de căutare, mai mici. Dacă problema poate fi descompusă astfel încât comportamentele individuale ale agenților sunt relativ disjuncte, atunci acest lucru poate duce la o reducere substanțială a spațiului de căutare și a complexității computaționale.

7.3.3 Trăsături diferențiale

Cele două categorii descrise anterior sunt generale și acoperă o varietate de forme de învățare care pot apărea în sistemele multiagent. Pentru a structura această varietate vom descrie câteva astfel de trăsături. Ultimele două trăsături, binecunoscute în ML, sunt la fel de potrivite pentru caracterizarea abordărilor centralizate și descentralizate. Celelalte sunt utile caracterizării învățării descentralizate.

(1) Gradul de descentralizare

Descentralizarea unui proces de învățare îi afectează:

- distribuția;
- paralelismul.

Extrema este atunci când un singur agent îndeplinește toate activitățile de învățare secvențial. Cealaltă extremă este când activitățile de învățare sunt distribuite și paralelizate la toți agenții din sistemul multiagent.

(2) Trăsături specifice interacțiunii

Există mai multe trăsături care pot fi folosite la clasificarea interacțiunilor necesare pentru a realiza un proces de învățare descentralizat. Iată câteva exemple:

- nivelul de interacțiune (variază de la observarea trecerii unui semnal la schimbul sofisticat de informații prin dialog și negociere);
- persistența interacțiunii (variază de la termen scurt la termen lung);
- frecvența interacțiunii (de la joasă la înaltă);
- tiparul interacțiunii (de la complet nestructurat la strict ierarhic);

- variabilitatea interacțiunii (de la fixă la schimbătoare).

Există situații în care învățarea are nevoie de o “interacțiune minimă” (ex.: observarea unui alt agent pentru o perioadă scurtă), în timp ce alte situații necesită “interacțiune maximă” (ex.: negocierea repetată pentru o perioadă îndelungată)

(3) Trăsături specifice implicării

Exemple de trăsături care pot fi folosite pentru a caracteriza implicarea unui agent în procesul de învățare:

- relevanța implicării;
- rolul jucat în timpul implicării.

În funcție de relevanță pot fi distinse două extreme: implicarea unui agent nu este o condiție pentru atingerea scopului deoarece activitățile sale pot fi executate și de un alt agent; și contrariul, scopul nu poate fi atins fără implicarea acestui agent. În funcție de rolul pe care un agent îl joacă în învățare, acesta se poate comporta ca un „generalist” dacă îndeplinește toate activitățile (în cazul învățării centralizate) sau ca un „specialist” dacă se specializează pe o anumită activitate (în cazul învățării descentralizate).

(4) Trăsături specifice scopului

Două exemple de trăsături ce caracterizează învățarea în sistemele multi agent, în funcție de scop sunt:

- tipul de îmbunătățire care se încearcă a fi obținut prin învățare;
- compatibilitatea scopurilor urmărite de agenți.

Prima trăsătură conduce la deosebirea importantă dintre învățarea care urmărește îmbunătățirea față de un singur agent (ex.: aptitudinile motorii sau abilitățile de deducție) și învățarea care urmărește îmbunătățirea față de mai mulți agenți lucrând în grup (de exemplu abilitățile de comunicare și negociere și gradul de coordonare și coerență). A doua trăsătură conduce la deosebirea importantă între scopurile conflictuale și complementare.

(5) Metoda de învățare

Se deosebesc următoarele metode sau strategii folosite de un agent:

- învățarea “pe de rost” (ex.: implantarea directă de cunoștințe și aptitudini fără a necesita deducție sau transformare din partea celui care învață);
- învățarea din instrucțiuni și ascultarea sfaturilor (ex.: operaționalizarea – transformarea într-o reprezentare internă și integrarea de cunoștințe și aptitudini anterioare – a noii informații ca o instrucțiune sau sfat care nu este direct executabil de către cel care învață);
- învățarea din exemple și antrenament (ex.: extragerea și rafinarea cunoștințelor și aptitudinilor ca un concept general sau un tipar standardizat de mișcare de la exemple pozitive și negative sau din experiență practică);
- învățarea prin analogie (ex.: păstrarea soluției de la transformarea cunoștințelor și aptitudinilor printr-o problemă rezolvată la o problemă nerezolvată, dar similară);
- învățarea prin descoperire (ex.: strângerea de cunoștințe și aptitudini noi făcând observații, experimentând, generând și testând ipoteze sau teorii pe baza rezultatelor experimentale).

(6) Feedback-ul învățării

Feedback-ul învățării indică nivelul de performanță obținut. Această trăsătură duce la următoarele deosebiri:

- învățarea supervizată (feedback-ul specifică activitatea dorită de cel care învață și obiectivul învățării este să potrivească această acțiune cât mai bine);
- învățarea prin întărire (feedback-ul specifică doar utilitatea activității celui care învață și obiectivul este să maximizeze această utilitate);
- învățarea nesupervizată (nu se oferă nici un feedback și obiectivul este să se găsească activități folositoare și dorite pe baza proceselor încercare – eroare și auto-organizare).

În toate aceste cazuri feedback-ul se presupune a fi oferit de mediul sistemului sau chiar de agenți. Asta înseamnă că mediul sau agentul care oferă feedback se comportă ca un “învățător” în cazul învățării supervizate, ca un “critic” în cazul învățării prin întărire și “pasiv” în cazul învățării nesupervizate.

Aceste trăsături caracterizează învățarea în sistemele multiagent din diferite puncte de vedere și pe diferite nivele. În particular, ele au un impact semnificativ asupra trebuințelor și abilităților agenților implicați în învățare. Sunt posibile diferite combinații cu diferite valori pentru aceste trăsături.

7.3.4 Problema Credibilității (The Credit Assignment Problem)

Problema principală cu care se confruntă fiecare sistem care învață este problema credibilității (**Credit Assignment Problem - CAP**). Acesta este problema asignării exacte a feedback-ului de crezare sau învinuire, pentru o schimbare globală a performanței (creștere sau descreștere), a fiecărei activități a sistemului care a contribuit la această schimbare.

Această problemă a fost, în mod tradițional, considerată în contextul sistemelor de sine stătătoare, dar există, deasemenea, și în contextul sistemelor multiagent. Considerând abordarea din Inteligența Artificială tradițională, conform căreia activitățile unui sistem inteligent sunt date de acțiuni externe efectuate de acesta și de inferențele interne și deciziile care implică aceste acțiuni, problema credibilității pentru sistemele multiagent poate fi descompusă în două probleme:

- problema *CAP inter-agenți* - asignarea creditului sau învinuirii pentru o schimbare globală a performanței datorate acțiunilor agenților;
- problema *CAP intra-agenți* - asignarea creditului sau învinuirii pentru o acțiune particulară a unui agent relativ la inferențele și deciziile sale interne.

Problema credibilității *inter-agenți* este o problemă dificilă în sistemele multiagent, deoarece o schimbare a performanței globale poate fi cauzată de acțiuni externe sau diferiți agenți distribuiți spațial și/sau logic.

Descrierea anterioară a problemei *CAP* este conceptuală, scopul fiind de a evidenția deosebirea clară între subproblemele inter și intra-agent. În practică, deosebirea nu este totdeauna evidentă. De altfel, abordările disponibile pentru învățare în sistemele multiagent nu fac diferența, în general, în mod explicit, între cele două subprobleme, sau se orientează doar pe una dintre subprobleme, simplificând-o mult pe cea de-a doua. Oricum, este util să fim conștienți de ambele subprobleme atunci când abordăm o problemă de învățare într-un sistem multiagent.

7.3.5 Învățarea și coordonarea activității

Această secțiune explică problema în care mai mulți agenți pot învăța să-și coordoneze activitățile (ex: pentru a-și împărți optim resursele sau pentru a maximiza profitul). Coordonarea activității se preocupă de dezvoltarea și adaptarea fluxului de date și de tiparele de control care îmbunătățesc interacțiunea dintre agenți. Anumiți agenți sunt nevoiți să interacționeze cu alți agenți care au diferite scopuri, abilități, structuri și durată de viață. Pentru a utiliza eficient oportunitățile prezentate și pentru a evita capcanele, agenții trebuie să învețe despre alți agenți și să-și adapteze comportamentul în funcție de structura grupului și de dinamică.

7.3.5.1 Învățarea izolată, concurentă și consolidată

Tehnicile de învățare prin întărire pot fi folosite de agenți pentru a dezvolta politici de selectare a acțiunii optimizând astfel feedback-ul mediului prin formarea unei corespondențe între percepție și acțiune. Un avantaj al acestor tehnici este că pot fi folosite în domenii în care agenții au foarte puțină experiență sau chiar deloc, și au informații puține despre aptitudinile și țelurile altor agenți. Lipsa acestor informații îngreunează problema coordonării. Aproape toate mecanismele de coordonare se bazează pe cunoștințele în domeniu și colaborarea dintre agenți. Ideea este că învățarea prin întărire se poate folosi pe post de tehnică nouă de coordonare în domenii în care tehnicile cunoscute sunt ineficiente.

Se pune întrebarea: să aleagă agenții să nu comunice în timp ce învață să se coordoneze? Cu toate că activitatea de comunicare este indispensabilă în activitățile de grup, nu garantează un comportament coordonat, consumă timp, și poate distra de la alte activități dacă nu este controlată cu atenție. De asemenea, agenții care depind de comunicare vor fi grav afectați dacă calitatea acesteia este compromisă (canale de comunicație întrerupte, informații eronate etc.). Câteodată, comunicarea poate fi riscantă sau chiar fatală (ca și în cazul luptelor când inamicul interceptează mesajele). Chiar și atunci când comunicarea se desfășoară normal, este bine să se facă doar dacă este absolut necesar. Acest design produce sisteme care nu inundă canalele de comunicare cu informație negarantată. Ca rezultat, agenții nu trebuie să se rătăcească printr-un labirint de informație inutilă pentru a localiza informația necesară.

În această formă de învățare izolată și concurentă, fiecare agent învață să-și îmbunătățească recompensa din mediu. Alți agenți din mediu nu sunt modelați explicit. Astfel, o altă întrebare interesantă ar fi dacă merită ca un astfel de agent să folosească același mecanism de învățare în ambele medii: cooperativ și non-cooperativ.

O prezumție de bază a majorității tehnicilor de învățare prin întărire este că dinamica mediului nu este afectată de alți agenți. Această prezumție este eronată în domenii cu mai mulți agenți care învață. Ne punem problema dacă tehnicile de învățare prin întărire vor fi adecvate pentru învățarea prin coordonare concurentă, izolată. Următoarele dimensiuni caracterizează învățarea prin întărire, concurentă și izolată (CIRL):

Cuplarea agenților. În unele domenii acțiunile unui agent afectează puternic planurile altor agenți (sisteme cuplate strâns), în timp ce în alte domenii acțiunile unui agent afectează mai puțin planurile altor agenți (sisteme cuplate slab).

Relațiile agenților. Agenții dintr-un sistem multiagent pot avea diferite relații:

- Acționează într-un grup pentru a rezolva o problemă comună (agenți cooperativi).
- Nu au o dispoziție unul față de celălalt, dar interacționează deoarece folosesc aceleași resurse (agenți independenți).
- Au interese opuse (agenți adversari).

Sincronizarea feedback-ului. În unele domenii agenții pot cunoaște efectele acțiunii lor imediat, în timp ce în altele pot primi feedback-ul acțiunii lor după o întârziere.

Combinații comportamentale optime. Câte combinații comportamentale ale agenților participanți vor rezolva problema în mod optim? Numărul variază de la unu la infinit pentru diferite domenii.

Pentru a evalua aceste întrebări a fost folosită tehnica Q-Learning în trei domenii:

Împingerea cuburilor. Doi agenți individuali învață să împingă o cutie dintr-un loc inițial într-un loc țintă de-a lungul unei traiectorii date. Ambele situații, cooperativă (au aceeași țintă) și competitivă (au ținte diferite) au fost studiate. Feedback-ul se bazează pe devierea locului cutiei de la traiectoria stabilită. Caracteristicile domeniului sunt: învățarea concurentă de către doi agenți cu feedback de mediu apropiat; sisteme cuplate strâns; comportamente optime multiple.

Împărțirea resurselor. Fiindu-le date sarcini individuale, doi agenți trebuie să învețe să împartă o resursă pentru o anumită perioadă. Caracteristicile domeniului sunt: feedback-ul mediului întârziat, sisteme cuplate strâns; comportament optim unic.

Navigarea roboților Doi roboți învață să navigheze intersectându-și traiectoriile fără să se lovească. Caracteristicile domeniului sunt: feedback-ul mediului imediat, cuplare variabilă, comportamente optime multiple.

Concluzia a fost că CIRL conferă o nouă paradigmă sistemelor multiagent prin care atât prietenii cât și dușmanii pot obține o cunoaștere a coordonării foarte utilă. Nici cunoștințele vechi despre domeniu și nici un model explicit despre aptitudinile altor agenți nu este necesar. Limitarea constă în incapacitatea CIRL de a dezvolta o coordonare eficientă când acțiunile agenților sunt strâns legate, feedback-ul este întârziat și nu sunt decât una sau puține combinații optime. O posibilă soluție parțială la această problemă ar fi folosirea unei forme de învățare alternativă. În acest caz, fiecare agent poate să învețe o perioadă de timp, apoi să execute politica lui fără modificări, apoi să învețe din nou etc. Doi agenți pot să-și sincronizeze comportamentul astfel încât unul să învețe în timp ce altul acționează și invers. Chiar dacă sincronizarea perfectă nu este posibilă, învățarea alternativă este mai eficientă decât cea concurentă.

Alte observații interesante:

- În situații cooperative, agenții pot învăța politici complementare de a rezolva o problemă. Asta duce la specializare și nu la dezvoltarea comportamentului identic.
- Agenții pot transfera învățarea în situații similare (ex.: o dată ce învață să se coordoneze la o problemă, învață să se coordoneze mai repede la alta similară).

7.3.5.2 Învățând despre și de la alți agenți

În această secțiune sunt considerate scenarii în care agenții învață să-și îmbunătățească performanțele. Uneori această îmbunătățire a performanței ori creștere a recompensei mediului vine datorită altor agenți din acel mediu. Aici se pune accentul pe agenți care învață despre alți agenți pentru a profita de oportunități și pe întrebarea despre cum poate fi influențată învățarea unui agent de către alți agenți. Acest subcapitol se concentrează asupra prezumției comportamentului altor agenți (incluzând preferințele lor, strategii, intenții etc.), asupra îmbunătățirii comportamentului unui agent interacționând cu alți agenți, asupra dezvoltării unui punct de vedere comun asupra lumii.

Roluri organizaționale în învățare

Agenții din grupuri trebuie să învețe să-și distribuie diverse roluri pentru a se completa reciproc. Adaptând structura grupului și activitățile individuale ale membrilor la o manieră dependentă de situație, se permite unui grup să sporească performanțele sistemului și să întâmpine dificultățile neprevăzute. Nagendra Prasad, Lesser și Lander prezintă un formalism care combină negocierea bazată pe memorie și învățarea prin întărire pentru a permite membrilor grupului alegerea rolurilor organizaționale. Ei adresează problema agenților care învață să adopte roluri specifice într-un domeniu al rezolvării problemei în mod cooperativ. Fiecare agent se presupune capabil de a juca unul din rolurile situației. Țelul este ca un agent să fie capabil să aleagă cel mai potrivit rol.

Învățarea de a profita de condițiile de piață

Agenții informaționali care vând și cumpără informații într-o piață electronică trebuie să se poată adapta condițiilor mediului. Vidal și Durfee cercetează avantajele agenților care învață după modelele altor agenți. Ei cercetează folosirea agenților pentru a vinde și a cumpăra informații pe piețe electronice, precum bibliotecile digitale. Se presupune că astfel de medii sunt deschise și noi agenți (fie cumpărători, fie vânzători de informații) pot intra sau părăsi piața după voie. O abordare practică a acestui sistem ar fi să considerăm că fiecare agent este o entitate “egoistă” cu scopul de a maximiza utilitatea locală. Un mecanism economic este folosit pentru a controla transferul informațiilor între agenții care furnizează informațiile și agenții care au nevoie de ele. Calitatea informației disponibilă la diferiți vânzători poate varia și aprecierea și cumpărarea sunt lăsate în seama vânzătorilor și cumpărătorilor.

Se presupune că informația poate fi reprodusă arbitrar la costuri neglijabile și agenții au acces uniform la toți agenții din piață. Într-un asemenea scenariu, un vânzător trebuie să ofere servicii cu valoare adăugată pentru a-și diferenția produsele de alți vânzători. Într-o astfel de piață, un cumpărător anunță bunul de care are nevoie. Vânzătorii etalează prețuri pentru oferirea acelor bunuri. Cumpărătorul alege între aceste oferte și oferă prețul convenit. Vânzătorul oferă produsul cumpărătorului. Cumpărătorul poate evalua calitatea bunului numai după ce îl primește (nu poate înainte de a-l cumpăra). Profitul vânzătorului s ce vinde un produs g la un preț p este $p - c_s^g$, unde c_s^g este costul producerii acelui bun. Dacă acest bun este de calitatea q , valoarea lui pentru

un cumpărător b este $V_b^g(p, q)$. Într-o tranzacție, scopul cumpărătorului și al vânzătorului este de a maximiza valoarea și profitul.

Într-o astfel de piață sunt cercetate trei tipuri de agenți:

- Agenți nivel-0: Acești agenți nu modelează comportamentul altor agenți. Ei fixează prețurile de vânzare și cumpărare bazându-se pe experiență.
- Agenți nivel-1: Aceștia analizează comportamentul trecut al altor agenți și încearcă să prezică prețurile de vânzare-cumpărare ale acestora. Alți agenți sunt modelați ca agenți nivel-0 sau agenți care nu modelează alți agenți. Asta înseamnă ca dacă un agent nivel-1 A modelează un agent nivel-0 B , A nu consideră faptul că și B îl modelează pe A . Observăm faptul că agenții nivel-1 au informații despre agenți individuali din mediu, în timp ce agenții nivel-0 folosesc doar experiența trecută.
- Agenții nivel-2: Aceștia modelează alți agenți ca și agenți nivel-1. Acest lucru înseamnă că acești agenți văd alți agenți ca agenți care îi modelează pe alții ca și agenți nivel-0 sau agenți care nu au modelele altora.

Învățarea de a juca mai bine împotriva unui oponent

În domenii ca jocurile de masă, strategia clasică oferă o abordare conservatoare. Dacă strategia folosită de oponent se poate aproxima, atunci exploatarea slăbiciunilor strategiei poate fi folosită pentru a obține rezultate mai bune împotriva aceluși oponent. Într-o abordare similară în dezvoltarea jucătorilor care să exploateze punctele slabe ale adversarului, Sean și Arora au folosit principiul Utilității Maxime Așteptate (UMA) pentru exploatarea învățării modelelor adversarilor. În abordarea lor, sunt folosite probabilități condiționale pentru diferite mutări ale adversarului, corespunzătoare tuturor mutărilor din starea respectivă, pentru a calcula utilități așteptate pentru fiecare mutare posibilă. Este jucată mutarea cu utilitatea maximă. Astfel se dezvoltă un model probabilistic al strategiei adversarului, observând mutări ale adversarului în diferite intervale de discrepantă, măsurate de funcția de evaluare a jucătorului.

Să presupunem că jucătorul și adversarul aleg din mulțimile următoare de mutări, $\alpha = \{\alpha_1, \alpha_2, \dots\}$, respectiv $\beta = \{\beta_1, \beta_2, \dots\}$, iar utilitatea primită de A pentru perechea de mișcări (α_i, β_i) este $u(\alpha_i, \beta_i)$. Principiul UMA poate fi folosit pentru a alege mutări astfel: $\arg \max_{\alpha_i \in \alpha} \sum_{\beta_j \in \beta} p(\beta_j | \alpha_i) u(\alpha_i, \beta_j)$, unde $p(\beta_j | \alpha_i)$ este probabilitatea condiționată ca

adversarul să aleagă mutarea β_j știind că agentul alege mutarea α_i . Dacă strategia adversarului poate fi modelată perfect de către mecanismul de învățare, jucătorul UMA va putea exploata slăbiciunile adversarului.

Toate domeniile discutate implică învățarea izolată în sens distribuit. Unul sau mai mulți agenți pot învăța concurent în acel mediu. Agenții interacționează des și informația obținută este folosită pentru a crea modele pentru alți agenți. Deoarece fiecare agent învață separat, toți agenții trebuie să execute toate activitățile de învățare. Majoritatea mecanismelor de învățare sunt variante ale învățării prin întărire.

7.3.5.3 Învățarea și comunicarea

În acest subcapitol, accentul cade pe legătura dintre învățare și comunicare. Această legătură privește cerințele asupra abilității agenților de a interschimba eficient informația folositoare. Se identifică două mari direcții de cercetare:

- Învățarea de a comunica: învățarea este privită ca o metodă pentru a reduce sarcina comunicării între agenți.
- Comunicarea privită ca învățare: comunicarea este văzută ca o metodă de a schimba informații care permit agenților să continue sau să-și perfecționeze activitățile de învățare.

Prima direcție de cercetare pornește de la faptul că activitatea de comunicare este de obicei încheată și scumpă și de aceea ar trebui evitată sau redusă ori de câte ori este posibil. Ultima direcție se bazează pe premiza că învățarea (cum ar fi planificarea și luarea deciziilor) este limitată din start de informația disponibilă și pe care o poate procesa un agent. Ambele direcții doresc îmbunătățirea comunicării și a învățării în sistemele multiagent, și sunt legate de următoarele probleme:

- Ce să comunice (ce informații interesează pe ceilalți).
- Când să comunice (cât efort trebuie să investească un agent în rezolvarea problemei înainte de a cere ajutor).
- Cu cine să comunice (pe care agent îl interesează această informație, cărui agent să-i ceară ajutor).
- Cum să comunice (la ce nivel să comunice agenții, ce limbaj și protocol să fie folosit, schimbul de informații să se facă direct punct cu punct și difuzat, sau cu un mecanism “tablă”).

Aceste probleme trebuie adresate de proiectantul sistemului sau derivate de sistemul însuși. Următoarele două paragrafe ilustrează cele două direcții, descriind abordări reprezentative în “învățarea de a comunica” și “comunicarea privită ca învățare”.

Mai există un aspect care merită dezvoltat când vorbim de învățare și comunicare în sistemele multiagent. O condiție necesară pentru un schimb util de informație este existența unei ontologii comune. Desigur, comunicarea nu este posibilă dacă agenții atribuie semnificații diferite acelorași simboluri, fără să realizeze diferențele (sau fără să le poată detecta și manevra). Dezvoltarea unei înțelegeri comune a simbolurilor poate fi considerată o sarcină esențială a învățării în sistemele multiagent .

7.3.5.3.1 Reducerea comunicării prin învățare

În mediile complexe, în special, este nevoie de mecanisme pentru reducerea comunicațiilor transmise. Smith propune câteva mecanisme de acest gen, cum ar fi concentrarea pe adresare și contractarea directă care au ca scop înlocuirea comunicării punct cu punct cu difuzarea. Un dezavantaj al acestor mecanisme este că trebuie ca designerul sistemului să știe dinainte căile de comunicare directă, și deci tiparele comunicării rezultate pot fi inflexibile în medii non – statice. În cele ce urmează este descris un mecanism mult mai flexibil numit “învățarea adresată” (într-o formă simplificată).

Ideea de baza în învățarea adresată este de a reduce eforturile comunicării pentru anunțarea sarcinilor permițând agenților să acumuleze și sintetizeze cunoștințe despre

capacitatea celorlalți agenți de a rezolva sarcini. Cu ajutorul cunoștințelor asimilate, sarcinile pot fi atribuite direct fără a fi necesară difuzarea anunțurilor către toți agenții. Negocierea “bazată pe caz” este aplicată ca un mecanism bazat pe experiență pentru culegerea cunoștințelor și sintetizarea lor. Aceasta se bazează pe observația că oamenii rezolvă adesea o problemă pe baza unor soluții care au funcționat în trecut pentru probleme similare. Ea pune accentul pe construirea cazurilor, adică perechi problemă – soluție (*case based reasoning*). Când se ivește o problemă nouă se verifică dacă este complet necunoscută sau similară cu una cunoscută (refolosirea cazului). Dacă este necunoscută, o soluție trebuie concepută. Dacă sunt asemănări cu o problemă cunoscută, soluția acestei probleme poate fi folosită ca punct de plecare în rezolvarea celei noi (adaptarea cazului). Toate problemele întâlnite, împreună cu soluțiile lor, sunt memorate ca și cazuri într-o bază de cazuri. Acest mecanism poate fi aplicat pentru reducerea comunicării într-o rețea contractuală după cum urmează: Fiecare agent își păstrează propria bază de cazuri. Un caz se consideră că conține:

- detalii despre o sarcină;
- informații despre care agent a rezolvat această problemă și cât de bună sau rea a fost soluția.

Detaliile despre sarcina T_i sunt de forma:

$$T_i = \{A_{i1}V_{i1}, \dots, A_{im_i}V_{im_i}\},$$

unde A_{ij} este un atribut al lui T_i și V_{ij} este valoarea acestui atribut. Este nevoie, pentru a aplica negocierea bazată pe caz, de o măsură a similarității între sarcini. În cazul învățării adresate, această măsură este redusă la similaritatea dintre attribute și valorile lor. Mai precis, pentru fiecare două attribute A_{ir} și A_{js} , distanța dintre ele este definită astfel:

$$DIST(A_{ir}, A_{js}) = SIMILAR_ATT(A_{ir}, A_{js}) * SIMILAR_VAL(V_{ir}, V_{js}),$$

unde $SIMILAR_ATT$ și $SIMILAR_VAL$ reprezintă similaritatea dintre attribute și valori. Cum se definesc aceste două măsuri depinde de domeniul de aplicație și de cunoștințele despre attributele sarcinilor și valorile lor. În cea mai simplă formă, se definesc astfel:

$$SIMILAR_ATT(x, y) = SIMILAR_VAL(x, y) = \begin{cases} 1, & x=y \\ 0, & \text{altfel} \end{cases}$$

ceea ce înseamnă că similaritatea este egală cu identitatea. Cu ajutorul distanței $DIST$ între attribute, similaritatea dintre două sarcini T_i și T_j se va calcula astfel:

$$SIMILAR(T_i, T_j) = \sum_r \sum_s DIST(A_{ir}, A_{js})$$

Pentru fiecare sarcină T_i poate fi definită o mulțime de sarcini similare $S(T_i)$ specificând cerințele de similaritate între sarcini. Iată un astfel de exemplu:

$$S(T_i) = \{T_j : SIMILAR(T_i, T_j) \geq 0.85\},$$

unde sarcinile T_j sunt incluse în baza cazurilor agentului care caută cazuri similare. Să considerăm cazul în care un agent N trebuie să decidă atribuirea unor sarcini T_i altui agent. În loc să difuzeze anunțul T_i , N încearcă să aleagă unul sau mai mulți agenți pe care îi consideră potriviți pentru rezolvarea T_i calculând pentru fiecare agent M această potrivire:

$$SUIT(M, T_i) = \frac{1}{|S(T_i)|} \sum_{T_j \in S(T_i)} PERFORM(M, T_j)$$

unde $PERFORM(M, T_j)$ este o măsură bazată pe experiență, ce indică cât de bun sau cât de rău a fost executată T_j de către M în trecut. Astfel, agentul N trimite anunțul T_i celui mai potrivit agent sau agenți în loc să trimită tuturor.

7.3.5.3.2 Îmbunătățirea Învățării prin Comunicare

Deoarece un agent nu poate fi considerat omnipotent, în cele mai multe probleme nu poate fi considerat nici omniscient fără a încălca presupunerile realiste. Lipsa de informații pe care o poate avea un agent se poate datora:

- mediului în care este integrat (localizarea obstacolelor) și problemei pe care trebuie să o rezolve (specificarea stării obiectivului ce trebuie atins);
- altor agenți (aptitudinile, strategiile și cunoștințele lor);
- dependenței dintre diferite activități și efectele activității proprii față de ale altor agenți asupra mediului și asupra potențialelor activități viitoare (o acțiune a întreprinsă de agentul A poate împiedica un agent B să întreprindă acțiunea b și poate permite agentului C să facă acțiunea c).

Agenții care au un acces limitat la informații importante, pot eșua în învățarea unei sarcini date. Acest risc poate fi redus permițând agenților să interschimbe informații, deci să comunice între ei. Pot fi delimitate două forme de îmbunătățire a învățării prin comunicare:

- învățarea bazată pe comunicarea de nivel jos, adică, interacțiuni simple întrebare – răspuns în scopul schimbului de informații lipsă (cunoștințe și ipoteze);
- învățarea bazată pe comunicarea de nivel înalt, adică, interacțiune mult mai complexă, cum ar fi negocierea și explicarea în scopul combinării și sintetizării informației.

În timp ce prima formă de învățare comunicativă are ca rezultat împărțirea informației, a doua formă are ca rezultat împărțirea înțelegerii. În ambele forme, comunicarea este folosită ca mijloc de îmbunătățire a învățării. Lăsând deoparte acest punct de vedere, trebuie reținut faptul că doar acest tip de comunicare poate fi văzut ca învățare, pentru că este o realizare specifică a sistemului multiagent de acumulare a cunoștințelor. Dacă învățarea ar trebui să fie îmbogățită prin comunicare rămâne o întrebare dificilă. În lumina evaluării criteriilor standard pentru algoritmi de învățare – viteză, calitate, complexitate – această întrebare poate fi descompusă în trei întrebări:

- Cât de repede sunt dobândite rezultatele învățării cu și fără comunicare?
- Rezultatele învățării cu/fără comunicare au o calitate suficientă?
- Cât de complex este procesul de învățare cu sau fără comunicare?

Considerațiile de mai sus ne arată că activitatea de comunicare oferă numeroase posibilități pentru îmbunătățirea învățării, dar ea nu reprezintă leacul universal în rezolvarea problemelor de învățare în sisteme multiagent. Combinarea lor trebuie deci făcută cu grijă. În unele cazuri, se observă că, tocmai comunicarea poate introduce informații false sau incomplete în baza de informații a unui agent (datorită erorilor de transmisie) care îngreunează rezolvarea unei sarcini de învățare.

7.3.6 Învățarea despre alți agenți în sistemele multiagent dinamice

Un *sistem multiagent dinamic* este un sistem care se schimbă în timp, datorită interacțiunilor dintre mai mulți agenți, fiind un *sistem distribuit cu inteligență artificială*.

Un exemplu de sistem multiagent dinamic ar fi jocul de fotbal robotizat sau sistemele de informare în rețele.

Un sistem multiagent dinamic este o particularizare a unui sistem multiagent care poate fi folosit în rezolvarea problemelor complexe, al căror mediu este dinamic (se modifică în timp) și în care un rol important îl joacă interacțiunea între mai mulți agenți. Din acest punct de vedere, sistemele multiagent dinamice constituie un exemplu de sisteme rezolutive folosite în Inteligența Artificială Distribuită (IAD).

Pentru a face parte dintr-un sistem multi-agent dinamic un agent trebuie să fie capabil :

- să modeleze comportamentul altor agenți; un model corect este absolut necesar pentru succesul sau supraviețuirea agentului; de exemplu într-un joc de fotbal, abilitatea agentului de a prevedea mișcarea celorlalți agenți, va determina dacă agentul poate sau nu ajunge la minge;
- să se adapteze pe la schimbările mediului, adică agentul să-și poată modifica comportamentul și funcțiile (procedurile) de planificare a comportamentului în momentul în care primește noi informații.

Într-un astfel de sistem un aspect important este acela de a *învăța* despre alți agenți prin :

- folosirea unei metode de *învățare nesupervizate* - din moment ce agenții sunt autonomi și nu există nici un supervisor, *învățarea supervizată* nu este utilă într-un sistem multi-agent dinamic;
- agentul trebuie să-și formeze un *model* al celorlalți agenți și să folosească metode de învățare pentru a-și corecta acest model.

7.3.6.1 Aspecte formale în reprezentarea unui sistem multi-agent dinamic

Un *sistem multiagent dinamic* este format dintr-un număr $n \geq 2$ agenți care evoluează în timp (într-un interval de timp agenții își propun un scop).

Într-un astfel de sistem agenții pot acționa:

- *cooperativ* sau
- *non-cooperativ*.

Vom considera în continuare cazul non-cooperativ, în care fiecare agent își urmărește propriile obiective .

Sistemul pornește la momentul inițial **0**, deplasându-se în mod discret (individual) spre momentul **T**, sau până când atinge o stare care satisface o condiție de terminare. Agenții își evaluează rezultatele în baza stării finale.

Din punct de vedere formal, un sistem multi-agent dinamic este reprezentat sub forma $\langle S, A, h, U \rangle$, unde :

- $S = \prod S^i$ -spațiul reunit al stărilor; S^i -spațiul de stări perceput de agentul i ;
- $A = \prod A^i$ -spațiul reunit al acțiunilor; A^i -spațiul acțiunilor agentului i ;

- h este funcția de tranziție între stările mediului, $h: S \times A \rightarrow S$ funcție ce asignează unei perechi (stare, acțiune) o nouă stare;
- $U = (U^1, U^2, \dots, U^n)$ este vectorul funcțiilor de utilitate ale agenților.

Presupunem că spațiul stărilor și al acțiunilor sunt domenii continue.

La fiecare moment t , sistemul percepe starea $s_t \in S$, care e constituită din stări locale pentru fiecare agent individual $s_t = (s_t^1, s_t^2, \dots, s_t^n)$. Agenții își vor alege acțiunea $a_t^i \in A^i$, producând acțiunea reunită $a_t = (a_t^1, a_t^2, \dots, a_t^n)$.

Stările evoluează conform relației $s_{t+1} = h(s_t, a_t)$, unde $s_{t+1}^i = h^i(s_t^i, a_t^i)$. Altfel spus, starea agentului i la momentul de timp $t+1$ depinde de starea curentă a agentului și de acțiunea reunită curentă. Vom presupune că funcția de tranziție h este deterministă și staționară.

După cum s-a arătat, funcția de utilitate a agentului i este U^i , o funcție a stării finale a agentului s_T^i . Se va introduce noțiunea de **variație a utilității agentului i la momentul t** ca: $r_t^i = U^i(s_t^i) - U^i(s_{t-1}^i)$, fiind numită și **recompensa obținută de agentul i la momentul t** .

Se observă că $\sum_{t=1}^T r_t^i = \sum_{t=1}^T [U^i(s_t^i) - U^i(s_{t-1}^i)] = U^i(s_T^i) - U^i(s_0^i)$, iar $U^i(s_0^i)$ este o constantă independentă de acțiunile agentului i , deci maximizarea sumei recompenselor primite de agentul i este echivalentă cu maximizarea utilității corespunzătoare stării finale.

Din cele prezentate mai sus rezultă că scopul unui agent de învățare într-un sistem multiagent dinamic care evoluează de la momentul inițial **0** la un moment final **T**, este maximizarea sumei recompenselor pe care le primește de la starea inițială la starea finală.

Chiar dacă maximizăm recompensele pe o singură perioadă de timp,

$$r_{t+1}^i = U^i(h^i(s_t^i, a_t^i)) - U^i(s_t^i),$$

este nevoie să cunoaștem acțiunile celorlalți agenți, a_t^{-i} . Scopul în problema învățării este să prezicem aceste acțiuni – explicit sau implicit – în așa fel încât agentul să poată alege efectiv acțiunea lui.

7.3.6.2 Modelarea altor agenți

Pentru a simplifica problema învățării, presupunem că agentul ia deciziile considerând numai perioada curentă de timp. Pentru a maximiza recompensa curentă la momentul $t(1)$, agentul rezolvă :

$$\arg \max_{a_t^i \in A^i} U^i(h^i(s_t^i, a_t^i, a_t^{-i})).$$

O abordare generală este aceea de a face o estimare, a_t^{-i} , a acțiunilor celorlalți agenți și de a rezolva problema ca și cum estimarea ar fi corectă. Decizia se reduce astfel la modul în care se face estimarea.

Nivele de modelare

Cea mai directă cale de a estima pe a_t^{-i} se bazează pe o analiză în serii a observațiilor anterioare, $\{a_\tau^{-i}\}, \tau < t$. Agenții care folosesc această abordare nu raționează explicit asupra faptelor ce îl determină pe a_t^{-i} . Astfel de agenți se numesc agenți nivel-0. Un exemplu de estimare nivel-0 este considerarea lui a_t^{-i} ca o funcție liniară de observații importante,

$$a_t^{-j} = \sum_{k=1}^t w_k a_{t-k}^j.$$

Un agent care încearcă să modeleze politica altor agenți se numește agent nivel-1. Extensia naturală este să considerăm agenți care modelează alți agenți ca agenți nivel-1. Acești agenți se numesc agenți nivel-2. Ei modelează politica agentului j ca
Se pot defini agenți nivel-3, ..., agenți nivel-N repetând extinderea acestui model.

Agenți non-estimatori

Un agent nivel-0 poate învăța o politică, fără a considera explicit valoarea a_t^{-i} , caz în care politica sa, $f^i(s_t^i)$, este o funcție ce depinde doar de starea locală. Un exemplu de astfel de agent este agentul competitiv. A fi competitiv într-un context de piață înseamnă a presupune că efectul celorlalți asupra mediului este neglijabil, și deci specularea acțiunilor celorlalți nu aduce nici un avantaj. Acest tip de agent tinde să se comporte reactiv, cu toate că această caracterizare nu e la fel de precis definită ca și competiția în contextul de piață.

Ca și exemplu, fie un agent nivel-1 care modelează agenți nivel-0 non-estimatori. Pentru a estima politica agentului j se aplică o regresie liniară logică asupra datelor istorice disponibile, $\{(s_\tau^j, a_\tau^j) | \tau < t\}$. Dându-se starea curentă s_t^j , se iau cei mai apropiați k vecini (definiți prin distanța euclidiană), $\{s_{\tau_1}^j, \dots, s_{\tau_k}^j\}$ și se aplică regresia liniară pe punctele $\{(s_{\tau_1}^j, a_{\tau_1}^j), \dots, (s_{\tau_k}^j, a_{\tau_k}^j)\}$. Rezultă o estimare a parametrilor α și β în modelul

$$a^j = \alpha + \beta s^j.$$

7.3.7 Domenii de aplicare ale Învățării în Sistemele Multiagent

În ciuda faptului că este un domeniu destul de nou, domeniul Sistemelor Multiagent conține un număr foarte mare de domenii de probleme. În cele ce urmează vom prezenta câteva dintre problemele rezolvate folosind sisteme multiagent, majoritatea abordărilor de rezolvare folosind tehnici de învățare pentru a asigura autonomia agenților.

7.3.7.1 Agenți robotici

În ultima vreme, cercetările în domeniul roboticii multiagent cooperative au cunoscut o mare amploare. Operațiile care nu pot fi efectuate cu roboți reali, pot fi, acum, simulate; totuși, comunitatea roboticii încurajează încă validare rezultatelor pe roboți reali.

Urmărirea pradă-prădător

Este una dintre cele mai cunoscute medii în cercetările de învățare în SMA, fiind simplu de implementat. Jocurile de urmărire constă într-un număr de agenți (prădători) care vânează în mod cooperativ o pradă. Agenții individuali prădători sunt, în general, mai puțin rapizi decât prada, și, de cele mai multe ori, agenții simt prada doar dacă sunt în apropierea ei. Ca urmare, agenții au nevoie să coopereze în mod activ pentru a reuși să captureze prada.

Una din abordările de învățare folosite în literatură pentru a rezolva problema prăzii și prădătorilor este *învățarea bazată pe cazuri (case-based learning)*.

Explorare

Domeniul constă într-o hartă mare cu agenți (roboți) și elemente care trebuie găsite, spre exemplu comori. Sarcina este transportarea comorilor în locuri speciale. Eficiența unei abordări constă în cât de repede sunt găsite toate comorile și de numărul de comori descoperite într-un interval de timp fixat.

Fotbal

Jocul de fotbal, atât în simulare cât și cu roboți reali, este unul dintre cele mai utilizate domenii de cercetare în sistemele multiagent. Domeniul constă într-un teren de fotbal cu două porți, o minge și două echipe cu un anumit număr de agenți (de la 5 la 11, depinzând de dimensiunea roboților și de teren). Performanța unei echipe este în general bazată pe numărul de goluri marcate, dar pot fi luate în calcul și alte criterii de performanță. Cea mai cunoscută astfel de competiție, RoboCup, are diferite ligi împărțite pe continente și pe tipuri de roboți sau a mediilor de simulare (www.robocup.org).

7.3.7.2 Aplicații în lumea reală

În această secțiune vom prezenta câteva probleme derivate din domenii ale lumii reale, folosite în investigații anterioare în domeniul Sistemelor Multiagent. Cele mai multe astfel de domenii de probleme sunt probleme de logistică, planificare, satisfacere a constrângerilor, probleme care necesită decizii distribuite, luate în timp real. Deoarece aplicațiile sunt în general distribuite și foarte complexe, tehnicile de învățare au fost destul de rar aplicate acestora, ca urmare, le vom prezenta în ceea ce urmează ca fiind probleme provocatoare pentru învățare în SMA.

Monitorizarea distribuită a vehiculelor

Sarcina în acest domeniu este optimizarea traficului într-o rețea de intersecții. Fiecare intersecție este echipată cu un semnalizator (semafor) controlat de un agent. Agenții vor trebui să se coordoneze pentru a putea gestiona fluctuațiile în trafic.

Controlul traficului aerian

În scopuri de securitate, spațiul aerian este împărțit în regiuni tri-dimensionale folosite de controlori de trafic aerian pentru a ghida avioanele spre destinațiile lor. Fiecare astfel de regiune are o limită superioară (numită *capacitatea sectorului*) în număr de avioane de care se poate ocupa la un moment dat. Scopul final este de a ghida avioanele dintr-un sector într-altul de-a lungul unor rute de lungime minimă, asigurând în același timp verificarea constrângerilor existente; soluția trebuie să fie rapidă pentru a gestiona date în timp real.

Gestionarea rețelelor și rutare în rețele

Acest domeniu constă în într-un număr mare de rețele distribuite. Agenții sunt lansați în control cooperativ și distribuit și gestionează rețeaua, tratează avariile și își echilibrează sarcinile.

Gestionarea distribuirii electricității

Aici problema este de a menține o configurație de putere optimală în rețea pentru a servi toți clienții și a minimiza pierderile; în același timp se gestionează posibilele distrugereri în rețea, cereri diverse din partea clienților, deteriorarea materialelor, etc.

Îngrijire medicală distribuită

Acest domeniu aplică Inteligența Artificială pentru a asista personalul clinic în diagnosticare, decizii legate de terapii și teste, determină prescripții și efectuează alte sarcini medicale specifice. Problema este foarte potrivită pentru sisteme multiagent datorită descentralizării datelor și resurselor, caracterului stohastic și dinamic al datelor.

Producția industrială

Această problemă clasică de planificare și programare implică gestionarea procesului de producție a elementelor complexe prin secvențe de pași, fiind diferite constrângeri și costuri asociate fiecărui pas. Sarcina constă în construirea unui plan (o *secvență de producție*) care specifică ordinea operațiilor pentru diferite elemente astfel încât costurile de producție să fie minimizate, satisfăcând în același timp instrucțiunile clienților.