

O soluție bioinspirată pentru generarea unui comportament autonom local, reflex, de evitare a obstacolelor

Dobrea Monica-Claudia, Dobrea Dan Marius

Abstract: Scopul principal al prezentei lucrări îl constituie găsirea unei soluții online, simple pentru problema de evitare a obstacolelor întâlnită în cazul roboților mobili. Soluția propusă de noi este una ce permite robotului dezvoltarea unui comportament autonom, local, de evitare a obstacolelor de fiecare dată când comanda motorie de nivel înalt care-l conduce (de ex., comanda *mergi înainte* sau *înapoi*) îl pune pe acesta într-o situație de coliziune iminentă. Soluția propusă de noi pentru un robot cu 36 de senzori infraroșii, distribuiți uniform în jurul acestuia, este una foarte simplă, bazată pe o rețea neuronală artificială minimală antrenată cu un algoritm derivat din algoritmul clasic de propagare înapoi a erorii (în engl. *backpropagation*, **BP**). Cu o încărcare computațională mică, algoritmul de învățare on-line implementat de noi s-a dovedit a fi unul foarte eficient atât în mediile aglomerate, statice cât și în cele dinamice. Rezultatele raportate aici au fost obținute în mediul de simulare MobotSim 1.0.03 – un simulator 2D, configurabil, destinat roboților mobili cu roți, cu acționare diferențială.

1. Introducere

În cazul roboților mobili o abilitate fundamentală a acestora constă în capacitatea autonomă locală a acestora de a evita obstacole. Până în prezent, numeroase metode [Tan, 2008], [Gutnisky, 2004] au fost propuse în vederea implementării acestui tip de comportament. În timp ce parte dintre aceste metode au fost considerate în cadrul unui simplu comportament de deplasare la întâmplare, alte metode au fost discutate și implementate în cadrul unor comportamente mai complexe cum ar fi cele în care roboților li se cere să ajungă într-o locație specificată, să-și planifice calea sau să urmeze fie pereții, fie o țintă mobilă, fie o linie marcată pe podea etc.

În cele ce urmează, problema evitării locale a obstacolelor este una tratată ca parte integrantă a unui proiect mai larg, al cărui scop final este acela de a construi un scaun cu roțile, inteligent, controlat mental de persoane cu deficiențe. Scaunul inteligent va asista utilizatorii în navigare; astfel, de fiecare dată și imediat după ce utilizatorii vor da o comandă (de exemplu, comanda *mergi înainte*), scaunul inteligent o va executa pe aceasta însă, va avea simultan grijă să evite, în mod adecvat și autonom, obstacolele aflate în cale.

Spre deosebire de multe alte abordări existente, abordarea noastră ia în considerare nu strategii de planificare a evitării obstacolelor ci însăși comportamentul reflex de apărare care, la oameni, este declanșat atunci când o modificare neașteptată și bruscă apare în mediul imediat înconjurător – mai exact, în așa-numitul spațiu peripersonal (SPP) [Graziano, 2006]. Rolul cheie al acestui mecanism de adaptare rapidă – mecanism implementat la nivelul cel mai de jos al execuției motorii – constă în plasarea, cu prioritate, a robotului în afara oricărui pericol de ciocnire cu obstacole dinamice și/sau statice. Conform unei arhitecturi reactive de tip de jos-în sus, decizia de la nivelul cel mai de jos al evitării unui obstacol este luată, mai departe, în

considerare de către nivelurile superioare de control motor ce sunt, la rândul lor, dedicate unor scopuri din ce în ce mai abstracte. În cazul nostru, următorul scop ce trebuie urmărit în cadrul proiectului global este acela ca robotul să aibă, pe lângă capacitatea de a evita (în mod autonom), și capacitatea de a ocoli (în același timp) obstacolul întâlnit astfel încât să fie păstrată, în continuare, direcția inițială de mers. În lucrarea prezentată în cele ce urmează doar comportamentul autonom reflex, de evitare a obstacolului (dezvoltat la nivelul motor cel mai de jos), integrat cu comanda superioară (în particular, cea care dictează o anumită direcție globală de mers) va face obiectul studiului nostru, urmând ca, pe viitor, să fie luată în studiu și problema ocolirii obstacolului (scop motor intermediar ce implică revenirea robotului pe direcția inițială de mers).

2. Câteva considerații tehnice și biologice

2.1 Considerații biologice

În vederea dezvoltării unui comportament reflex (bio-inspirat) de evitare a obstacolelor au fost avute în vedere o serie de rezultate biologice și psihologice raportate în literatură. Dintre acestea, elementele cheie folosite în implementarea noastră sunt cele prezentate mai jos.

La om, majoritatea celulelor sistemului nervos central (SNC) sunt formate înainte de naștere însă cea mai mare parte a conexiunilor dintre celulele nervoase se realizează după naștere, în timpul copilăriei timpurii. Modul în care se realizează aceste conexiuni este unul dictat, în mod esențial, de: (a) interacțiunea constantă a copilului cu mediul său înconjurător și, mai mult, (b) de cronologia experiențelor acumulate de copil (astfel, de exemplu, o experiență timpurie, de interacțiune a copilului cu mediul său, este de așteptat să aibă o influență mai mare în dezvoltarea sistemului nervos central al copilului, prin aceea că ea dă naștere unor conexiuni nervoase care vor dicta ulterior maniera în care copilul va realiza noi achiziții – fie acestea sub forma unor informații, abilități noi etc.).

Reflexele – principalele răspunsuri motorii involuntare de apărare: (1) pot fi învățate sau (2) pot fi comportamente învățate ce permit realizarea, într-un mod automat, a unor activități mai complexe. Acestea din urmă: (i) sunt fie învățate de fiecare individ, în mod independent, fie sunt predate de către terțe persoane, (ii) vin din experiență (fie într-o manieră de tip încercare și eroare, fie folosind memoria experiențelor trecute și observarea terților) și (iii) pot fi rafinate printr-o practică continuă, în scopul unei mai bune adaptări la condițiile de mediu aflate într-o permanentă schimbare.

La apărarea corpului și menținerea unei margini de siguranță a acestuia își aduc aportul două arii corticale importante, având răspunsuri neuronale cu latență mică. Aceste arii corticale sunt: a) aria ventrală intra-parietală (VIP) și b) aria frontală polisenzorială (PZ) [Graziano, 2006]. În ariile VIP și PZ (arii unde neuronii răspund la tipuri similare de stimuli senzoriali și a căror stimulare conduce la ieșiri de tip defensiv¹ similare [Graziano, 2002]): i) cei mai mulți neuroni sunt bimodali și trimodali, răspunzând simultan la stimuli vizuali, tactili și auditivi, ii) independent de mărimea stimulilor, circa jumătate din celulele VIP

¹ Mișcările evocate sunt mișcări de evitare, retragere sau de protejare a părții corpului pe care se află localizat câmpul tactil receptiv.

răspund cel mai bine la stimuli vizuali situați la mai puțin de 30 cm de corp iar multe dintre aceste celule răspund doar la stimuli aflați la foarte puțini centimetri de corp; în mod similar, în aria PZ circa 46% dintre neuroni prezintă un răspuns puternic și susținut doar pentru acei stimuli vizuali situați la mai puțin de 5 cm de corp în timp ce alți 40% dintre neuroni răspund la stimuli vizuali situați la mai puțin de circa 20 cm de corp; *iii*) răspunsurile celulelor multimodale pun în evidență un gradient al ratei de descărcare ca o funcție de distanța la stimulii localizați în spațiul peripersonal (de exemplu, răspunsurile descresc neliniar ca o funcție de distanța la stimuli); *iv*) unii neuroni VIP multisenzoriali primesc și intrare vestibulară și, din acest motiv, se consideră că aceștia ar fi implicați în detectarea direcției mișcării autogenerate a subiectului; *v*) neuronii PZ proiectează direct la nivel spinal, fiind astfel implicați în ieșirea motorie defensivă.

Comenzile motorii ce coboară de la nivel cortical sunt direct responsabile de inițierea și generarea mișcărilor voluntare. Reflexele complexe (așa cum este, de exemplu, reflexul de evitare a obstacolelor, care cel mai frecvent sub-servește alte mișcări voluntare cu scop precis) apar prin proiecția intrărilor senzoriale la nivelul circuitelor neuronale spinale. Proiecția intrărilor senzoriale se realizează prin intermediul inter-neuronilor spinali iar circuitele spinale la nivelul cărora se face proiecția sunt frecvent reprezentate de generatorii centrali de pattern² (GCP). În această rețea de conexiuni, rolul inter-neuronilor spinali este acela de a integra intrările descendente, venite de la creier, cu intrările senzoriale aferente și, în consecință, acela de a adapta reflexele și activitatea neuronilor motori spinali la diferitele cerințe funcționale apărute.

Teoria celor doi factori ai lui Mowrer privind evitarea postulează faptul că învățarea procesului de evitare implică *două etape*:

- (1) *O primă etapă*, în care subiectul care învață experimentează *condiționarea clasică/bazată pe sentimentul de frică generat*. Mai exact, un stimul de avertizare, așa cum este o distanță mică (sau un anumit prag) până la cel mai apropiat obstacol, este asociat cu o situație neplăcută, așa cum este coliziunea; în acest mod, un stimul inițial neutru devine un stimul condiționat, SC, capabil să producă un puternic răspuns condiționat, de frică.
- (2) *O a doua etapă*, în care subiectul experimentează *condiționarea operantă*. În această etapă subiectul adoptă o acțiune de răspuns la *stimulul condiționat aversiv* și, astfel, elimină – prin *întărire negativă* (în engl., negative reinforcement) – *evenimentul aversiv*. Această ultimă etapă corespunde însuși procesului de învățare a evitării iar în această etapă comportamentul de evitare nu este întărit prin evitarea situației neplăcute (ciocnirea de obstacol) ci, prin terminarea stimulului condiționat aversiv – adică, a stimulului care a evocat sentimentul de frică . Cu alte cuvinte, *stimulul aversiv este cel care întărește acele răspunsuri care îl elimină pe el însuși*.

Teoria ecologică a lui Gibson privind percepția vizuală [Shumway-Cook, 2007] reclamă faptul că pentru a genera acțiuni adecvate oamenii au nevoie de informația perceptuală (nu senzorială!) legată de factorii de mediu care sunt importanți pentru task-ul motor de executat. În cazul nostru ne-am folosit de o

² Aceștia sunt implicați în generarea secvențelor motorii stereotype.

percepție a adâncimii medii până la obstacolele din jur, alături de percepțiile legate de direcția mișcării și viteza de deplasare.

2.2 Considerații tehnice și practice

În cele prezentate în continuare, implementarea propusă pentru procesul de învățare a evitării obstacolelor va trebui privită doar ca un model foarte reducăționist al analogului său uman, fără a se pierde însă și din relevanța sa.

Dintre modelele bio-inspirate propuse în literatură ca soluții pentru învățarea evitării obstacolelor, *modelele de învățare cu întărire* (în engl. *reinforcement learning*, RL), precum și *modelele de tip comportament operant* (în engl. *operant behavior*, OB) par a fi cele mai promițătoare.

Așa cum am precizat anterior, în teoria celor doi factori ai lui Mowrer vorbim de *învățare prin evadare* (învățăm să sfârșim ceva aversiv) și *învățare prin evitare* (învățăm să prevenim ceva aversiv). Ca și în *învățarea prin evadare*, cele două tipuri de modele, RL și OB:

- (1) nu necesită o cunoaștere completă a mediului și/sau cunoașterea acțiunii ce trebuie luată în fiecare context particular de mediu (oricum, nu putem vorbi de o soluție unică pentru acțiunea ce poate fi abordată la un moment dat);
- (2) aceste două modele necesită, mai presus de toate, cât mai multe interacțiuni posibile ale robotului cu mediul său înconjurător; din acest punct de vedere, aceste modele se încadrează în clasa *algoritmilor de învățare de tip on-line*;
- (3) interacțiunea cu mediul (static și/sau dinamic) are la bază tehnica de învățare *bazată pe încercare și eroare*;
- (4) ambele modele furnizează robotului mobil autonom o *evaluare permanentă* și corespunzătoare a performanțelor sale, evaluare ce se face în termeni de *pedeapsă și recompensă*.

Cu toate acestea, în timp ce rezultatele obținute cu fiecare dintre aceste două modele par a fi foarte promițătoare, fiecărui model, în parte, fie îi lipsesc unele detalii importante, fie acesta se confruntă cu unele probleme practice, de implementare; iar toate acestea fac ca, în final, soluția obținută să fie una departe de a exprima un comportament asemănător celui întâlnit la subiecții umani. În particular, în cazul metodei RL și a celei mai larg utilizate implementări a ei (vorbim aici de *algoritmul Q-learning*³ [Watkins, 1992]) menționăm cel puțin următoarele dezavantaje:

- (i) În primul rând, această metodă implică manipularea unui tabel foarte mare. Acest tabel este unul utilizat în actualizarea valorilor-Q. Astfel, spre exemplu, pentru un robot cu numai 8 senzori, cu 5 acțiuni posibile de executat și cu un domeniu al valorilor de intrare de [0, 1022] pentru fiecare sensor în parte, avem nevoie de un tabel-Q cu nu mai puțin de $1.1995 \times 10^{24} \times 5$ intrări. Recent, o alternativă la tabelele de tip look-up (utilizate

³ În care se estimează pentru fiecare pereche posibilă [stare, acțiune] un semnal mediu de tip recompensă numerică.

inițial pentru a stoca valorile-Q) o reprezintă rețelele neuronale artificiale (RNA). Acestea din urmă sunt utilizate în principal datorită celor două capacități ale lor, respectiv: **a)** capacitatea de a oferi o reprezentare mai compactă a valorilor-Q și **b)** capacitatea de a interpola valorile-Q pentru perechile stare-acțiune care nu au fost vizitate niciodată. Totuși, în cazul utilizării diferitelor paradigme RNA noi probleme pot să apară [Tan, 2008], așa cum este, spre exemplu, cazul instabilității raportate pentru arhitecturile de tip perceptron multistrat (MLP), antrenate cu algoritmul BP. Mai precis, este foarte dificil de a garanta faptul că învățarea noilor pattern-uri nu erodează cunoștințele acumulate anterior.

- (ii) Un al doilea dezavantaj important este dat de numărul foarte limitat al acțiunilor posibil de executat de către robot – număr ce este determinat, în principal, ca urmare a constrângerilor de calcul.
- (iii) Un alt dezavantaj major îl constituie numărul mare de parametri necunoscuți atât ai algoritmului Q-learning (de ex., rata de actualizare, parametrul de temperatură inițial, funcția de recompensă etc.), cât și ai implementării RNA; pentru mai multe dezavantaje vezi și [Gutnisky, 2004].

Spre deosebire de *metodele RL*, unde într-o primă fază de explorare (*fază de învățare*), în mod independent de starea curentă, robotul explorează mediul prin selectarea unor *acțiuni non-greedy*⁴ (folosind, în acest sens, distribuția de probabilitate Boltzmann), în *condiționarea operantă* [Gutnisky, 2004] stimulii primiți de robot sunt utilizați nemijlocit pentru a învăța ce acțiuni să realizeze mai mult (este vorba de acțiunile ce au primit recompensă) și ce acțiuni să realizeze mai puțin (respectiv, acțiunile care au fost penalizate). În consecință, în timp ce ambele tipuri de metode sunt considerate metode de tip on-line, *operarea în timp real* rămâne o caracteristică doar pentru metodele de tip OB. În plus, în [Gutnisky, 2004] găsim raportat faptul că nici una dintre metodele menționate aici nu conduce la performanțe maxime; o explicație posibilă pentru aceasta ar putea fi faptul că în faza de antrenare termenul de recompensă sau pedeapsă este unul utilizat în detrimentul semnalului de intrare senzorial care este unul mult mai adecvat dat fiind marele său potențial informativ.

3. *Procesul de învățare (bio-inspirată) a evitării reflexe a obstacolelor*

În cele ce urmează o nouă metodă bio-inspirată, de învățare a evitării reflexe a obstacolelor va fi prezentată, alături de implementarea ei practică pe o platformă robotică. Platforma robotică cu roți și cu acționare diferențială este una dotată cu 36 de senzori infraroșii, IR, distribuiți în mod echidistant în jurul ei (vezi Figura 1.a). Un controller bazat pe o rețea minimală de tip MLP și antrenată cu un nou algoritm bio-inspirat (**BBP**) – derivat din algoritmul BP – va constitui, în cele ce urmează, elementul cheie folosit în vederea întrunirii criteriului de învățare on-line și în timp real, întâlnit de altfel la subiecții umani.

⁴ *Strategia (metoda) greedy* presupune efectuarea unei alegeri. Dintre toți pașii următori posibili de ales, se alege acel pas care asigură un maximum de “câștig”, de unde și numele metodei: *greedy* = lacom.

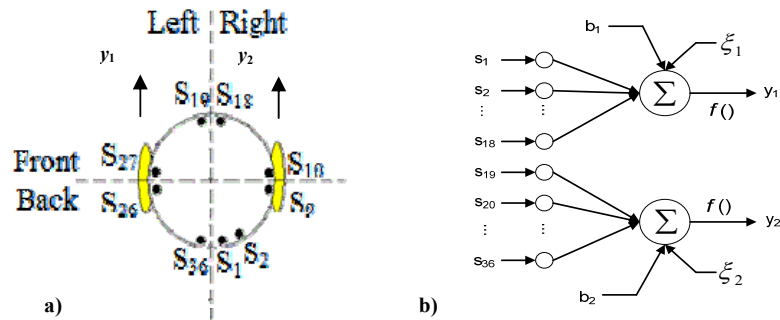


Figura 1: a) Dispunerea senzorilor IR, b) arhitectura rețelei de tip MLP.

Arhitectura minimală a rețelei de tip MLP, **Figura 1.b**, a constat în 36 de intrări (corespunzător celor 36 de valori, s_i , citite de la senzorii IR) și în 2 ieșiri (y_1 , y_2 , ce furnizează comenzile pentru cele două motoare ce acționează roțile diferențiale ale robotului). Raza scurtă de acțiune a senzorilor IR, de doar 30 cm, i-a făcut pe aceștia să fie foarte adecvați în cadrul mecanismelor de adaptare rapidă la mediu.

Într-o primă etapă (de preprocesare), valorile senzorilor au fost liniarizate iar, mai apoi, ele au fost normalizate în intervalul $[0, 0.9]$, cu 0 desemnând lipsa oricărui obstacol în vecinătatea imediată a robotului și 0.9 denotând coliziune. În cadrul topologiei rețelei MLP, funcțiile de transfer adoptate pentru cele două elemente de procesare de ieșire au fost de tip *tanh*; acest din urmă tip de funcții prezintă, suplimentar, avantajul că limitează valorile de ieșire în intervalul $[-1, +1]$.

Comenzile globale de deplasare a robotului – codate, după cum urmează: *înainte* (1,1), *înapoi* (-1,-1), *la stânga* (0,0.5), *la dreapta* (0.5, 0) și *stop* (0,0) – au fost luate în considerare, la nivelul rețelei MLP, prin intermediul bias-urilor (b_1 , b_2) ale neuronilor de ieșire. Direcțiile de deplasare înainte și înapoi ale roților robotului au fost și ele codate prin valori pozitive și, respectiv, negative în timp ce virajele la dreapta și, respectiv, la stânga, au fost obținute prin comenzi diferențiale adecvate transmise la nivelul celor două motoare. Informația de viteză, y_j , a fost una furnizată ca o valoare proporțională în intervalul 0 (*stop*) și 1/-1 (*viteză maximă înainte/înapoi*). La calculul activării neuronilor de ieșire ai rețelei s-a adăugat, de asemenea, un termen reprezentând un mic zgomot, ξ , cu rol în a preveni robotul de a se bloca în anumite condiții particulare de mediu.

În lumina celor prezentate mai sus, o analogie grosieră a sistemului nostru robotic cu sistemul neuro-motor uman poate arăta astfel: stratul de intrare al rețelei neuronale este omologul ariei frontale polisenzoriale PZ, neuronii de ieșire sunt inter-neuronii care integrează informația senzorială primară (ce a trecut prin aria PZ) cu comanda corticală iar sistemul de acționare a roților (controller-ul implementat cu o rețea MLP) ar corespunde generatorilor centrali de pattern-uri.

În forma sa standard (1), algoritmul BP impune existența unui profesor care știe sau care poate calcula ieșirea dorită pentru fiecare intrare (condiție de mediu) dată. Totuși, acesta nu este și cazul nostru întrucât noi nu putem defini valorile dorite pentru comenzile y_1 , y_2 , acestea din urmă depinzând de factori multipli care, adesea, sunt și imprevizibili.

$$BP: \Delta w_{ij} = -\eta \frac{\delta E}{\delta w_{ij}} = -\eta(d_j - y_j)s_i(1 - y_j^2), \quad \Delta b_j = -\eta(d_j - y_j)(1 - y_j^2) \quad (1)$$

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

În relația (1) w_{ij} sunt ponderile rețelei neuronale artificiale, E reprezintă eroarea iar η este rata de învățare (valoarea ei a fost aleasă fixă, și anume 0.7).

Pentru a putea face, totuși, uz de algoritmul BP (algoritm ce descrie cel mai bine modul în care copiii învață exersând, într-o manieră de tip *încercare și eroare*) am ținut cont și exploatat, în mod corespunzător, trei dintre teoriile psihologice, și anume: **i) teoria ecologică a percepției vizuale, ii) teoria motivației interne [Oudeyer, 2007] și iii) teoria evitării a lui Mowrer** (teoria celor doi factori). În cele din urmă am ajuns la o nouă paradigmă, și anume, am substituit *paradigma centrată pe profesor* (algoritmul BP) cu o *nouă paradigmă, centrată pe elev*. Relația de calcul aferentă noului algoritm, BBP (algoritm pe care noi l-am numit, oarecum impropriu, algoritmul BP-bioinspirat) este cea prezentată în relația (2).

$$Forward : y_j = \tanh \left(\sum_i w_{ij} s_i + b_j + \xi_j \right), \quad i = \overline{1,18} \text{ pt. } j = 1 \text{ si } i = \overline{19,36} \text{ pt. } j = 2$$

$$BBP : \Delta w_{ij} = -\eta T_j s_i (1 - y_j^2), \quad \Delta b_j = 0$$

$$T_j = E_j \text{ Crit }_j$$

$$E_j = A_j - A_j^d = A_j, \quad A_j = \frac{\sum_i s_i}{n_j}, \quad i = \overline{1,18} \text{ pt. } j = 1, i = \overline{19,36} \text{ pt. } j = 2 \quad (2)$$

$$\text{Crit }_j = \text{sign}(\Delta E_j) \text{sign}(y_j) = \text{sign}(\Delta A_j) \text{sign}(y_j)$$

În noua relație de actualizare a valorilor ponderilor (2), variația ponderilor w_{ij} diferă față de variația ponderilor din relația (1) doar prin termenul compozit T_j pentru care termenul echivalent în relația (1) este $(d_j - y_j)$. În relația (1) acest termen din urmă este o măsură a cât de mult din valoarea erorii se datorează ieșirii j a RNA. Spre deosebire de algoritmul BP standard (1), în care acest termen al erorii are o formă analitică clară ce depinde de valoarea dorită, d_j , și de ieșirea actuală, y_j , a RNA, în algoritmul BBP modul de calcul al termenului său echivalent, T_j , este unul ce ține cont:

- pe de o parte, de *maniera* particulară în care „creierul elevului” *percepe stimulii din mediul imediat înconjurător* (în particular, vorbim aici de termenul E_j) și,
- pe de altă parte, de *modul* în care elevul *manipulează*, la nivel cortical, *aceste percepții* (reflectat, la nivel de relație, prin termenul Crit_j).

Acești ultimi termeni nou introduși vor fi explicați mai târziu în cadrul acestei lucrări.

În mod natural și firesc creierul uman folosește diferența dintre modul în care lumea este percepută și modul în care ea este de așteptat să fie percepută (corespunzător scopurilor noastre) ca o informație de eroare funcție de care încearcă, apoi, să corecteze actul motor (mișcarea).

Pentru a genera acțiuni adecvate oamenii utilizează informația perceptuală (nu pe cea senzorială) legată de factorii de mediu care prezintă importanță pentru task-ul motor ce trebuie făcut. În cazul nostru, *percepția stimulului aversiv*, A_j – obținută ca *percepția adâncimii medii*⁵ până la obstacolele aflate în partea contralaterală dreaptă/stângă a robotului – a fost cea utilizată în faza de învățare (2). Pentru fiecare neuron de ieșire al RNA s-a calculat o valoare diferită A_j corespunzătoare, scopul fiind acela de a facilita învățarea reflexului de evitare a obstacolului pe ambele părți (în acest mod se asigură evitarea obstacolului pe partea cea mai potrivită). În relația (2), parametrul $(n_j)_{j=1,2}$ reprezintă numărul senzorilor de valoare non-nulă dispuși pe partea contralaterală a robotului. Valoarea acestui parametru variază de la un moment de timp la altul, funcție de condițiile de mediu (respectiv, prezența sau absența obstacolelor în spațiul peripersonal al robotului).

În cadrul algoritmului de învățare BBP, *stimulul condiționat aversiv* dat de percepția A_j , crește cu cât adâncimea medie până la obstacolele aflate în spațiul peripersonal (SPP) drept/stâng scade și, respectiv, dispare (devine zero) atunci când valorile citite de la nivelul tuturor senzorilor robotului devin nule. Acest din urmă caz corespunde percepției dorite a mediului, $A_j^d = 0$; în această situație ideală robotul este situat la cel puțin 30 cm distanță de orice obstacol înconjurător. În mod corespunzător, *eroarea apreciată de elev* (robot), E_j , devine, în final (conform relației (2)), egală cu percepția A_j care nu reprezintă altceva decât însăși stimulul aversiv care trebuie eliminat. Această eroare calculată de către robot (și care ia valori în intervalul $[0, 0.9]$) este una conformă cu paradigma fundamentală care definește eroarea BP⁶.

În continuare, corespunzător aceleiași etape a condiționării operante, elevul (robotul) va întări doar acele acțiuni care elimină stimulii aversivi. În acest sens, robotul evaluează la fiecare pas consecințele acțiunilor sale ca răspuns la stimulii din mediu și, apoi, *generează intern o recompensă scalară*, $\text{sign}(A_j[n] - A_j[n-1])$; aceasta din urmă trebuie înțeleasă ca o măsură a *progresului înregistrat în procesul de învățare*. Această măsură calitativă, împreună cu percepția direcției de mișcare a roților, formează așa-numitul termen pe care noi l-am denumit *critic*, **Crit**. Termenul **Crit** este cel care controlează maniera în care rețeaua neuronală artificială își actualizează parametrii. Mai exact, actualizarea ponderilor rețelei MLP în sensul menținerii sensului de mers de la pasul anterior⁷ este încurajată ori de câte ori robotul se mișcă astfel încât să elimine/diminueze stimulii aversivi, $A_j[n] \leq A_j[n-1]$, și această actualizare a ponderilor se face în sensul schimbării sensului de mers de la pasul anterior în caz contrar.

La oameni am văzut că, într-o primă fază (numită și de *condiționare clasică*), comportamentul de evitare este unul bazat, în principal, pe *sentimentul de frică generat de consecințele evenimentului aversiv* (în cazul nostru, coliziunea). În momentul, însă, în care comportamentul de evitare începe să fie unul realizat, în mod repetat, cu succes, sentimentul de frică începe să dispară (subiectul uman începe să aibă controlul asupra situației) iar procesul de învățare a comportamentului de evitare încetează.

⁵ A se vedea semnificația dată valorilor citite de la senzori.

⁶ Atunci când sistemul adaptiv reușește să rezolve cu succes problema, eroarea calculată devine zero; în caz contrar, eroarea măsoară distanța dintre rezultatele dorite și ieșirile curente ale sistemului adaptiv.

⁷ Vorbim aici de sensul de rotire (înainte sau înapoi) a fiecărei roți în parte.

În cazul nostru, pentru a putea surprinde în cadrul soluției propuse de noi și acest comportament întâlnit la oameni, o modelare a sentimentului de frică menționat mai sus apare, practic, ca o condiție necesară. În lipsa, însă, a unei astfel de modelări (reamintim că algoritmul BBP implementează condiționarea operantă, nu și etapa premergătoare, cea a condiționării clasice), soluția adoptată – și rezultată din practică – a constat în impunerea ca după primii 300 de pași de antrenare (în care regula de actualizare BBP s-a aplicat la fiecare pas), actualizarea ponderilor rețelei MLP să aibă loc în continuare doar atunci când cel puțin o intrare a RNA avea o valoare mai mare de 0.8. Această din urmă situație s-ar putea traduce prin aceea că robotul ajunge să se confrunte cu o situație nemaiîntâlnită până atunci (o situație nouă⁸), pe care nu o poate rezolva și pentru care este necesară o nouă învățare.

În mediul de simulare, implementarea soluției s-a făcut ținând cont și de caracteristicile și limitările tehnice ale robotului fizic pe care se urmărește implementarea, într-o etapă următoare, a soluției simulate. Din aceste considerente, viteza maximă de deplasare a robotului a fost setată la 0.3 m/s iar timpul scurs între fiecare acțiune (mișcare) a robotului și următoarea citire a valorilor senzorilor a fost aleasă de 400 ms. Valorile inițiale ale ponderilor rețelei nu au fost generate în mod aleator ci, dimpotrivă, ele au fost setate la zero pentru a simula procesul de conectare neurală progresivă ce are loc în sistemul nervos al copiilor începând de la naștere și ținând pe tot parcursul procesului de dezvoltare a controlului neural.

4. Rezultate și discuții

Testarea algoritmului de învățare a evitării obstacolelor propus mai sus s-a făcut în mediul de simulare MobotSim 1.0.03, cu medii atât statice cât și dinamice (de exemplu, cu doi și, respectiv, cinci roboți mobili diferiți). Un exemplu de comportament de evitare a obstacolelor învățat cu noul algoritm este și cel prezentat în **Figura 2.a** (comportament obținut după primii 2527 de pași de antrenare) sau în **Figura 2.b** (comportament obținut după primii 70517 de pași de antrenare).

În urma analizei modului în care robotul a învățat comportamentul de evitare putem trage următoarele concluzii. La fel ca și la oameni:

- (1) Metoda propusă a permis *implementarea online și în timp real* a procesului de învățare. Aceasta trebuie privită și ca o consecință a faptului că algoritmul este unul foarte simplu, ce nu implică costuri computaționale foarte mari.
- (2) *Învățarea are loc relativ repede* (după doar una sau două coliziuni) și ea este foarte durabilă, fără să presupună alte coliziuni ulterioare.
- (3) Pentru procesul de învățare *foarte importante sunt primele interacțiuni ale robotului cu mediul său inconjurător*, aceste interacțiuni fiind cele care conturează, în linii mari, comportamentul de evitare de mai târziu.
- (4) O consecință majoră a soluției propuse pentru algoritmul de evitare a obstacolelor constă în *construirea unei zone de siguranță* în jurul robotului. Acest fapt vine să confirme ipoteza că “spațiul personal” al omului este rezultatul unui mecanism defensiv al cărui rol este acela de monitorizare a obiectelor aflate în jurul corpului, obiecte care pot aduce

⁸ Nici o situație dintre cele întâlnite de robot până atunci nu seamănă (în anumite limite) cu cea prezentă.

atingere integrității acestuia.

- (5) Robotul reacționează prompt indiferent de direcția din care vine stimulul mobil (în cazul nostru, stimulul îl reprezintă un alt robot aflat în mișcare).
- (6) Maniera în care este generată mișcarea⁹ împreună cu modul în care sunt calculate cele două ieșiri ale RNA¹⁰ (reprezentând comenzile pentru motoare) oferă suportul real pentru obținerea unei varietăți foarte mari de mișcări. Așa se explică marea flexibilitate a mișcării robotului obținută utilizând algoritmul nou propus.
- (7) Preluarea controlului în mod automat, fie de către comanda „corticală”, fie de către comportamentul autonom local, s-a făcut în mod corespunzător, ori de câte ori contextul de mediu a cerut-o.

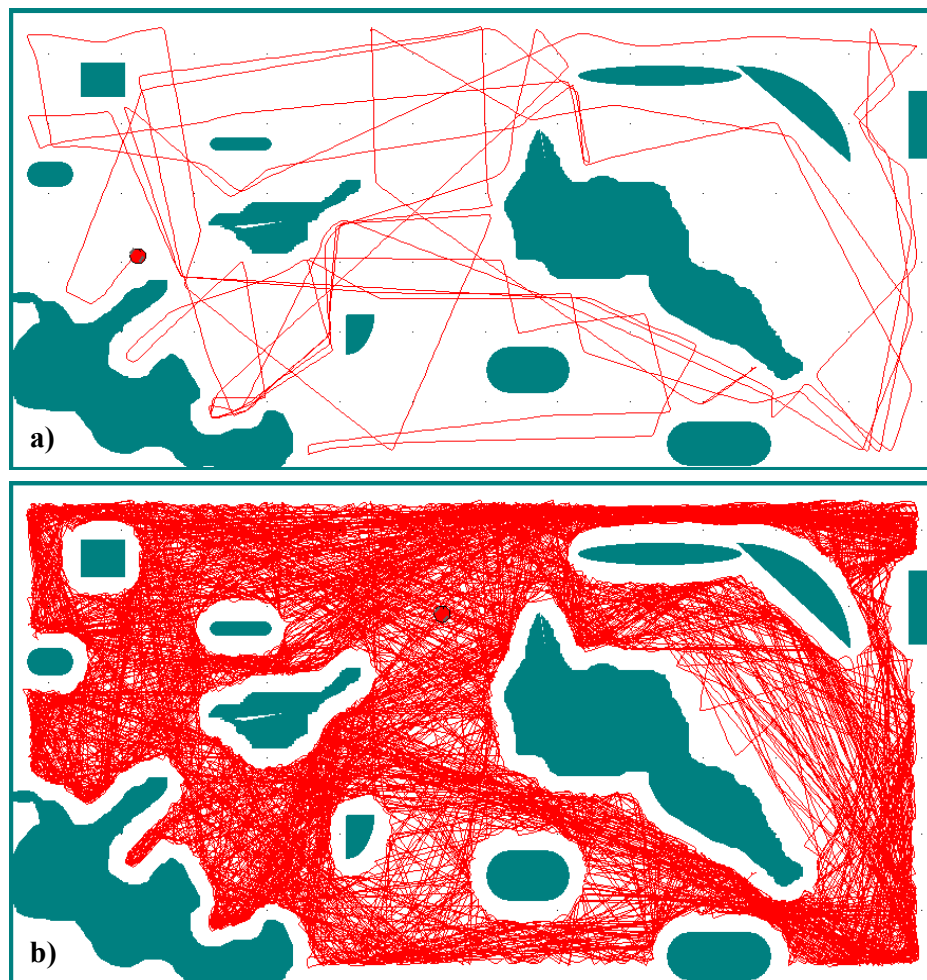


Figura 2: Comportamentul de evitare a obstacolelor învățat de robot: a) după primii 2527 pași de antrenare și b) după primii 70517 pași de antrenare ai algoritmului BBP.

⁹ Și anume, prin intermediul celor două roți acționate diferențial.

¹⁰ Avem în vedere aici valorile obținute – valori care baleiază, practic, intervalul [-1,+1].

5. Concluzii

În această lucrare s-a prezentat, pe scurt, un nou concept pentru controlul robotic de tip *comportament reactiv bio-inspirat*. Plauzabilitatea biologică a modelului, alături de simplitatea ridicată a tehnicii și a paradigmei de învățare, au făcut sistemul robotic astfel obținut să fie unul foarte eficient: *i*) acesta învață repede și de o manieră consistentă; *ii*) comportamentul autonom este obținut online, în timp real și în medii dinamice și nestructurate; *iii*) comanda „corticală” și comportamentul local autonom reușesc cu succes să preia, în mod automat și exclusiv (unul îl exclude pe celălalt) controlul sistemului, ori de câte ori contextul de mediu o cere.

Referințe

- 1 [Tan, 2008] Tan A. H., Lu N. și D. Xiao: *Integrating temporal difference methods and self-organizing neural networks for reinforcement learning with delayed evaluative feedback*, în IEEE Transactions on Neural Networks, 19: 230-244, 2008
- 2 [Gutnisky, 2004] Gutnisky D. A. și Zanutto B. S.: *Learning obstacle avoidance with an operant behavior model*, în Artificial Life, 10: 65-81, 2004
- 3 [Graziano, 2006] Graziano M. S. și Cooke D. F.: *Parieto-frontal interactions, personal space, and defensive behavior*, în Neuropsychologia, 44: 845–859, 2006
- 4 [Graziano, 2002] Graziano M. S., Taylor C. S. și Moore T.: *Complex movements evoked by microstimulation of precentral cortex*, în Neuron, 34: 841–851, 2002
- 5 [Shumway-Cook, 2007] Shumway-Cook A. și Woollacott M. H.: *Motor control: translating research into practice*, 3rd ed., Lippincott Williams & Wilkins, U.S.A., pp. 16, 2007
- 6 [Watkins, 1992] Watkins C. și Dayan P.: *Q-Learning*, în Machine Learning, 8: 279-292, 1992
- 7 [Oudeyer, 2007] Oudeyer P.-Y. și Kaplan F.: *What is intrinsic motivation? A typology of computational approaches*, în Frontiers in Neurorobotics, 1(6): 1-14, 2007