

Deep Reinforcement Learning. Studiu de caz: Deep Q-Network

Mihnea Horia VREJOIU

Institutul Național de Cercetare-Dezvoltare în Informatică – ICI București

mihnea.vrejoiu@ici.ro

Rezumat: Inteligența artificială (*Artificial Intelligence* – AI) a ajuns astăzi poate cel mai de actualitate subiect în numeroase domenii. Unul dintre obiectivele principale ale AI este crearea de agenți complet autonomi capabili să interacționeze cu lumea înconjurătoare și să învețe prin încercare și eroare comportamente optime în diferite contexte, perfectibile în timp. Între metodele AI de învățare automată (*machine learning*), învățarea cu întărire / prin recompensă (*reinforcement learning* – RL) prin interacțiuni repetate cu mediul înconjurător cu urmărirea unui scop are un rol deosebit de important, pe lângă învățarea supervizată și respectiv nesupervizată. Totuși, metodele clasice de RL au limitări importante în scalabilitatea la probleme cu dimensionalitate mai mare. În ultimii ani, tehnologiile de învățare supervizată și nesupervizată bazate pe *deep learning* cu utilizarea rețelelor neuronale adânci (*deep neural networks*), având proprietăți remarcabile de aproximare a funcțiilor complicate pe spații multidimensionale, precum și de învățare de reprezentări ierarhice caracteristice extrase automat direct din date, cu reducere dimensională semnificativă, au cunoscut o dezvoltare explozivă producând rezultate impresionante comparabile cu, sau chiar peste performanța umană în domenii diverse ca: recunoașterea obiectelor / imaginilor, recunoașterea vorbirii, traducerea automată etc. Combinarea metodelor RL cu cele *deep learning* a condus la ceea ce se numește astăzi *deep reinforcement learning* (DRL), oferind noi posibilități în realizarea de agenți autonomi performanți în spații multidimensionale. Lucrarea de față își propune o prezentare succintă a domeniului DRL și studierea și analizarea în detaliu a uneia dintre primele metode DRL de succes și anume Deep Q-Network, dezvoltată de Google DeepMind.

Cuvinte cheie: învățare cu/prin întărire/recompensă, agent, stare, acțiune, politică, diferență temporală, Q-learning, învățare profundă, rețea neuronală adâncă, rețea neuronală convoluțională.

Deep Reinforcement Learning. Case Study: Deep Q-Network

Abstract: Artificial Intelligence (AI) became today perhaps the most up-to-date topic in many areas. One of the main goals of AI is to create completely autonomous agents able to interact with the surrounding world and learn by trial and error optimal behaviors in different contexts, perfectible in time. Among the machine learning methods of AI, reinforcement learning (RL) by repetitive interactions with the environment while targeting a purpose plays a particularly important role, besides supervised and unsupervised learning. However, classical RL methods have important limitations in scalability to higher-dimensionality problems. In recent years, supervised and unsupervised learning technologies based on deep learning, using deep neural networks with remarkable properties of approximating complex functions on multi-dimensional spaces, as well as the learning of characteristic hierarchical representations automatically extracted directly from data, with significant dimensional reduction, have had an explosive development, producing astonishing results comparable with, or even surpassing human performance in areas such as object / image recognition, speech recognition, automatic translation etc. The combination of RL with deep learning methods has led to what is now called deep reinforcement learning (DRL), providing new possibilities for producing autonomous agents in multidimensional spaces. This paper is proposing a brief presentation of the DRL field, while also studying and analyzing in detail one of the first successful DRL methods, namely Deep Q-Network developed by Google DeepMind.

Keywords: reinforcement learning, agent, state, action, policy, temporal difference, Q-learning, deep learning, deep neural network, convolutional neural network.

1. Introducere

Atunci când se pune problema reprezentării și înțelegerii mecanismelor de învățare, atât la nivel empiric și intuitiv dar și în majoritatea teoriilor despre învățare și inteligență, una dintre primele metode avute în vedere apare ca fiind învățarea bazată pe interacțiunea cu mediul înconjurător prin percepții senzoriale și acțiuni. În acest mod pot fi acumulate informații despre diverse legături cauzale, despre rezultate și consecințe ale acțiunilor, precum și despre ce anume

trebuie făcut pentru atingerea unui obiectiv/scop, acest tip de interacțiune constituind pentru om una dintre principalele surse de cunoștințe și cunoaștere despre mediul înconjurător și despre sine.

Crearea de agenți complet autonomi cu capacități de a interacționa cu mediul înconjurător și de a-și antrena singuri, pe bază de încercare și eroare, comportamentul cel mai adecvat (perfectibil în timp) în diferite situații, constituie unul dintre obiectivele principale în domeniul inteligenței artificiale. Aria de aplicabilitate acoperă de la roboți care pot „simți” și interacționa cu lumea înconjurătoare, până la agenți software care interacționează cu limbajul natural și multimedia (Arulkumaran et al., 2017 [1]). Metoda de învățare prin interacțiuni repetate, cu urmărirea atingerii unui anumit scop, este denumită *reinforcement learning* – RL (învățare cu întărire / prin recompensă) și reprezintă un alt tip de învățare automată (*machine learning*), pe lângă învățarea supervizată și cea nesupervizată.

RL furnizează mediul matematic principal pentru învățarea autonomă bazată pe experiență. În RL, un agent interacționează pas cu pas prin acțiuni cu mediul înconjurător, acțiunea sa la fiecare pas conducând la modificarea stării acestui mediu (tranziția la o nouă stare) care îi întoarce ca răspuns (*feedback*) agentului câte o mărime scalară denumită recompensă (*reward*). Scopul agentului este maximizarea recompensei cumulate pe termen lung, după o secvență de pași. În acest sens, RL este asemănătoare controlului optimal, utilizând tehnici specifice programării dinamice (procese de decizie Markov) dar fără un model matematic exact al mediului și funcționând pe spații cu dimensiuni mari. Prin învățare în RL este stabilită o așa numită politică pentru o problemă dată, prin care se definește comportamentul agentului la un moment dat. Astfel, o politică π asociază/mapează stările mediului înconjurător percepute de agent cu acțiunea cea mai potrivită a fi întreprinsă de către acesta din urmă în fiecare din stările respective. Agentul va acționa conform politicii și/sau ocazional va încerca aleator și alte acțiuni care eventual ar putea îmbunătăți politica.

Deși cu ajutorul metodelor RL clasice au fost obținute unele rezultate foarte bune, acestea s-au dovedit totuși a fi destul de limitate, mai ales în privința scalabilității pentru probleme cu număr crescut de dimensiuni unde apar dificultăți ținând de complexitate. Astfel de inconveniente au putut fi depășite însă odată cu apariția succeselor spectaculoase obținute în ultimii ani în diferite alte domenii (cum sunt recunoașterea imaginilor, înțelegerea vorbirii, traducerea automată etc.) de *deep learning* cu utilizarea rețelelor neuronale adânci (*deep neural networks*) având proprietăți remarcabile de aproximare a funcțiilor, precum și de învățare de reprezentări ierarhice caracteristice extrase automat direct din date, cu reducere dimensională semnificativă. Astfel, utilizarea algoritmilor de *deep learning* în conjuncție cu metodele de RL a condus la ceea ce se numește astăzi generic *deep reinforcement learning* (DRL), oferind noi posibilități în realizarea de agenți autonomi performanți în spații multidimensionale de stări și acțiuni în perspectiva creării de sisteme capabile să învețe cum să se adapteze în lumea reală. Cu siguranță DRL va fi o componentă importantă în realizarea sistemelor generale de IA (Lake et al., 2016 [9]).

Lucrarea de față își propune o prezentare succintă a domeniului DRL, cu studierea și analiza în detaliu a uneia dintre primele metode de DRL de succes și anume Deep Q-Network (DQN), dezvoltată de Mnih et al., 2015 [15] de la Google DeepMind, ca metodă de referință în domeniu care îmbină o metodă RL (*Q-learning*) cu *deep learning* utilizând o rețea neuronală convoluțională. În continuare, lucrarea este structurată după cum urmează. Secțiunea 2 trece sumar în revistă elementele definitorii principale pentru RL, cu enumerarea în final a câtorva metode de RL de referință, menționarea limitărilor acestora, precum și a îmbunătățirilor produse de noile metode de RL care valorifică avantajele aduse de revoluția *deep learning*, fiind menționate câteva aplicații cu succes ale acestora. În secțiunea 3 este analizată metoda DQN, cu detalierea arhitecturii, algoritmului și tehnicilor speciale utilizate, a experimentelor efectuate și a rezultatelor obținute. Lucrarea se încheie cu secțiunea 4 în care sunt enunțate câteva concluzii relevante.

2. Reinforcement Learning

Reinforcement learning – RL (învățarea cu întărire / prin recompensă) reprezintă un tip de metode de învățare automată (*machine learning*) din cadrul inteligenței artificiale, prin interacțiuni (continue) cu mediul cu urmărirea unui scop. În general, aceeași denumire – RL – referă totodată atât o problemă, o clasă de metode de soluționare potrivite pentru aceasta, precum și domeniul care

studiază aceste probleme și metodele de soluționare a lor. De aceea trebuie avută permanent în vedere distincția conceptuală între acestea, în particular între probleme și metode de soluționare, pentru a se evita potențiale confuzii. Practic, RL nu este definită prin descrierea/caracterizarea metodelor/algoritmilor de învățare, ci prin caracterizarea problemei de învățare, esențială fiind perceperea aspectelor cele mai importante ale problemei reale puse în fața agentului interacționând cu mediul său înconjurător pentru a atinge scopul urmărit. Orice metodă/algoritm care rezolvă o astfel de problemă se consideră a fi metodă/algoritm de RL (Sutton & Barto, 1998 [20]).

În RL, un agent autonom controlat de un algoritm de învățare automată descoperă prin încercări care acțiuni ale sale în raport cu mediul sunt cele mai recompensatoare în fiecare moment/situație. Agentul observă/percepe starea mediului său înconjurător s_t la momentul t , interacționează cu acest mediu efectuând acțiunea a_t în starea s_t , determinând tranziția mediului la momentul $t+1$ într-o nouă stare s_{t+1} care furnizează agentului o recompensă scalară r_{t+1} (*reward, reinforcement*) ca *feedback*. De cele mai multe ori, o acțiune nu are efect doar asupra recompensei imediate, ci și asupra situației/stării următoare și, prin aceasta, asupra recompenselor ulterioare (viitoare) corespunzătoare. Căutarea de tip încercare-eroare (*try-error search*) și recompensa întârziată (*delayed reward*) reprezintă cele două caracteristici definitorii în RL. Secvența de stări (și acțiuni) parcursă de un agent în interacțiunea sa cu mediul este denumită **traietorie**.

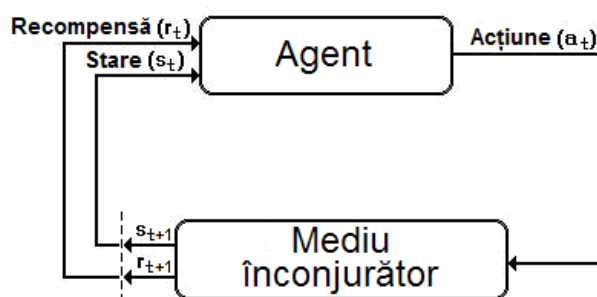


Figura 1. Interacțiunea agentului cu mediul (adaptare după: Sutton & Barto, 1998 [20])

Pe lângă **agentul** capabil să învețe din experiența proprie luând decizii de acțiune în funcție de aceasta și **mediul înconjurător** cu care agentul respectiv interacționează, un sistem RL mai este compus din următoarele elemente:

- a) o **politică** (*policy*),
- b) un semnal de **recompensă** (*reward signal*),
- c) o **funcție de evaluare** (*value function*) și
- d) opțional, un **model al mediului** (*environment model*).

a) O **politică** π asociază/mapează stările mediului înconjurător percepute de agent cu acțiunea cea mai potrivită a fi întreprinsă de către acesta din urmă în fiecare din stările respective, determinând astfel comportamentul agentului RL la un moment dat. În unele cazuri, politica poate fi reprezentată printr-o tabelă de asocieri (*lookup table*) sau printr-o funcție simplă, dar în alte cazuri poate implica calcule extensive cum sunt procesele de căutare. Politicile pot fi în general stohastice, precizând probabilități asociate fiecărei acțiuni posibile în fiecare stare.

b) La fiecare pas (în timp), urmare a unei acțiuni a agentului RL acesta primește de la mediul ca *feedback* o valoare numerică R denumită **recompensă** (*reward*). Un semnal de recompensă definește practic obiectivul unei probleme RL, indicând evenimentele „bune” și respectiv „rele” pentru agent ca și caracteristici definitorii imediate ale problemei abordate de către acesta. Pe baza semnalului de recompensă se poate realiza ajustarea/modificarea politicii: dacă o acțiune selectată de politică este urmată de o penalizare (sau recompensă slabă), atunci politica poate fi modificată să selecteze pe viitor o altă acțiune în respectiva situație. Semnalele de recompensă sunt adesea funcții stohastice (probabilistice) de starea mediului și de acțiunile întreprinse.

c) Scopul agentului este maximizarea recompensei cumulate pe termen lung. O **funcție de evaluare** determină valoarea (utilitatea) unei stări $V(s)$ ca măsură a recompensei totale pe care agentul se așteaptă să o acumuleze în viitor pornind de la starea respectivă. Astfel, valoarea indică atractivitatea pe termen lung a unei stări a mediului ținând cont și de stările probabile să urmeze și de recompensele disponibile în acestea, în timp ce recompensa determină atractivitatea intrinsecă imediată a stării. O stare s poate avea o recompensă imediată $R(s)$ scăzută, dar totuși o valoare $V(s)$ mare dacă este urmată de alte stări cu recompense în general mari, sau vice-versa. În evaluarea și luarea deciziilor se țin cont de valori, fiind alese acele acțiuni care conduc la stări cu valoarea cea mai mare (și nu neapărat cu recompensa cea mai mare), care indică obținerea cantității celei mai mari de recompensă pe termen lung. În timp ce recompensele sunt oferite direct de mediu, valorile trebuie estimate și re-estimate din secvența de observații pe care le face un agent de-a lungul existenței sale. Analog valorii unei stări $V(s)$ este definită și valoarea unei acțiuni pentru o anumită stare (*action-value*, *quality* sau *Q-value*), $Q(s, a)$. La limită, valoarea unei stări este maximul dat de valoarea acțiunii optime din starea respectivă: $V(s) = Q^*(s, a)$. Funcția valorilor acțiunii Q asociază/mapează perechi stare-acțiune cu combinația cea mai valoroasă între recompensa imediată și toate recompensele viitoare (cu ponderi diminuate la fiecare pas în timp) ce pot fi obținute de acțiuni ulterioare pe traiectorie. O metodă eficientă de estimare a valorilor stărilor sau, respectiv, acțiunilor constituie practic componenta cea mai importantă a majorității algoritmilor RL.

d) Un al patrulea element component al unor sisteme RL poate fi reprezentat de un **model al mediului** care simulează comportamentul acestuia, permițând inferențe despre cum se va comporta mediul. De exemplu, pentru o pereche stare-acțiune dată, modelul poate „prezice” tranziția la următoarea stare, precum și următoarea recompensă. Modelele sunt utilizate pentru planificare, prin care se înțelege orice modalitate de decizie de acțiune considerând situațiile viitoare posibile înainte ca acestea să fie efectiv experimentate. Metodele de rezolvare a problemelor RL care utilizează modele și planificare se numesc metode bazate pe model (*on-model*), spre deosebire de metodele simple fără model (*off-model*) care învață politici optime exclusiv din încercări și greșeli. Sunt utilizate și metode RL mixte în care se învață din încercări și erori și simultan se învață și un model al mediului care este folosit pentru planificare.

Obiectivul agentului este acela de a învăța o politică (strategie de control) π care să maximizeze recompensa cumulată așteptată din pașii următori (aportul recompenselor individuale fiind diminuat la fiecare astfel de pas cu un factor subunitar γ). Fiind dată o stare s , o politică π întoarce o acțiune de executat a . O politică optimală π^* maximizează recompensa așteptată. În acest sens, RL este similar cu controlul optimal. Formalismul matematic al problemelor RL are la bază idei din teoria sistemelor dinamice și programarea dinamică (Bellman, 1952; 1957 [4; 5]), cum este cazul controlului optimal în procesele de decizie Markov (*Markov Decision Processes* - MDP) incomplet cunoscute. În esență, pentru a învăța, agentul trebuie:

- să aibă capabilitatea de a „simți” (într-o anumită măsură) starea mediului înconjurător;
- să poată efectua acțiuni prin care să afecteze/modifice această stare;
- să aibă unul sau mai multe scopuri/obiective legate de starea mediului respectiv.

Procesele de decizie Markov includ aceste trei aspecte în formele lor cele mai simple fără a le trivializa pe niciunul. În RL, în general, nu este disponibil un model al dinamicii tranzițiilor de stare, agentul trebuind să învețe despre consecințele acțiunilor sale în mediu prin încercare și eroare. Fiecare interacțiune cu mediul aduce informație pe care agentul o utilizează pentru a-și actualiza cunoștințele. Tabelul 1 prezintă sintetic o paralelă între MDP și RL.

Tabelul 1. Paralelă sintetică între MDP și RL

<p>MDP sunt definite prin:</p> <ul style="list-style-type: none"> • mulțimea de stări S, mulțimea de acțiuni A; • modelul de tranziții $T(s_i, a, s_{i+1})$ cunoscut (determinist sau stohastic); • funcția de recompensă $R(s)$ cunoscută; • se calculează o politică optimă π^* ($\pi^*(s)$ = acțiunea recomandată în starea s). 	<p>RL se bazează pe MDP, dar:</p> <ul style="list-style-type: none"> • modelul de tranziții este necunoscut; • funcția de recompensă este necunoscută; • se învață o politică optimă.
--	---

Metodele de învățare RL pot fi pasive, în care agentul execută o politică fixă pe care doar o evaluează (*on-policy*), sau active, fără politică fixă explicită (*off-policy*), în care agentul își actualizează politica pe măsură ce învață. Necesitatea compromisului între explorare și exploatare (*exploration vs. exploitation*) este specifică metodelor de învățare RL active. Pe de o parte, agentul trebuie să aleagă acele acțiuni dovedite din experimentări anteriorare ca eficiente în aducerea de recompensă (exploatare a experienței acumulate), dar pe de altă parte, pentru descoperirea acestora, este necesară experimentarea de acțiuni neîncercate anterior, în căutarea uneia mai „profitabile” (explorare pentru acumularea de experiență nouă). Aceasta este „dilema” agentului la un moment dat: trebuie să favorizeze explorarea stărilor și acțiunilor necunoscute la momentul respectiv, sau exploatarea stărilor și acțiunilor despre care știe deja că aduc recompense mari? În general, agentul trebuie să mixeze acțiuni de ambele tipuri și, în mod progresiv, să le favorizeze pe acelea care apar ca fiind cele mai bune (e.g. metoda ϵ -greedy prin care se alege cu probabilitate ϵ o acțiune aleatoare și cu probabilitate $1 - \epsilon$ o acțiune optimă cunoscută, ϵ scăzând de la valoarea 1,0 la o valoare reală mică subunitară pozitivă de prag pe parcursul unui număr mare de pași, după care nu mai scade sub această valoare). Totodată, într-o sarcină stohastică, fiecare acțiune trebuie încercată de mai multe ori, pentru a se obține o estimare de încredere privind recompensa așteptată în urma acesteia.

Dintre metodele RL cele mai utilizate amintim (Sutton & Barto, 1998 [20]): programarea dinamică adaptivă, metoda Monte Carlo, metoda diferenței temporale (*temporal difference* – TD), SARSA (*State-Action-Reward-State-Action*), *Q-learning* (Watkins, 1989 [34]; Watkins & Dayan, 1992 [35]), metode de tip actor-critic etc. Totuși, deși au condus la anumite succese, abordările clasice de RL prezintă limitări serioase în ceea ce privește scalabilitatea la un număr crescut de dimensiuni, fiind restricționate de probleme ținând de complexitate (de memorie, computațională, de exemple).

Progresele spectaculoase din ultimii ani în domeniul *deep learning* bazate pe proprietățile excepționale ale rețelelor neuronale adânci (*deep neural networks*) de aproximare a funcțiilor și de găsire automată de reprezentări compacte de dimensiuni reduse (caracteristici) ale datelor de mari dimensiuni, au oferit noi mijloace pentru depășirea acestor limitări. Utilizarea algoritmilor *deep learning* în RL definește astăzi ceea ce numim *deep reinforcement learning* (DRL), deschizând noi perspective în realizarea de agenți autonomi performanți în spații multidimensionale de stări și acțiuni. O primă aplicație revoluționară utilizând DRL a putut învăța să joace cu performanțe comparabile cu, sau chiar peste cele umane o gamă de jocuri video Atari 2600 direct din pixelii din imaginile acestora (Mnih et al., 2015 [15]), demonstrând că agenți RL pot fi antrenați direct din date brute multidimensionale (doar pe baza unui semnal de recompensă. Un al doilea succes important a fost reprezentat de dezvoltarea sistemului DRL hibrid AlphaGo, care a învins campionul mondial la Go (Silver et al., 2016 [19]), utilizând rețele neuronale antrenate prin învățare supervizată și RL, combinate cu un algoritm tradițional de căutare euristică.

Algoritmi DRL au fost de asemenea deja utilizați într-o gamă largă de aplicații de robotică, în care politicile de control al roboților pot fi învățate direct din imagini capturate din lumea reală (Levine et al., 2016 [10; 11]). Mai mult, DRL a permis crearea de agenți capabili să „învețe să învețe” – meta-învățare – (Duan et al., 2016 [8]; Wang et al., 2017 [32]), ajungând astfel să poată generaliza la medii vizuale complexe pe care nu le-au mai „văzut” anterior (Duan et al., 2016 [8]).

DRL oferă posibilități reale pentru crearea de sisteme capabile să învețe cum să se adapteze în lumea reală, permițând automatizarea a tot mai multe sarcini fizice concrete, de la gestionarea consumului de energie (Tesauro et al., 2008 [25]) la apucarea și depunerea obiectelor (Levine et al., 2016 [11]). Totodată, RL fiind o metodă generală pentru abordarea problemelor de optimizare prin încercare și eroare, DRL a fost deja utilizată în diverse sarcini legate de învățarea automată, de la proiectarea de modele avansate (*state-of-the-art*) de traducere automată (Zoph & V Le, 2017 [36]), la crearea de noi funcții de optimizare (Li & Malik, 2017 [12]). DRL va fi o componentă importantă în realizarea sistemelor generale de IA (Lake et al., 2016 [9]).

3. Deep Q-Network (DQN)

Modul de reprezentare și stocare a funcțiilor de evaluare și/sau a politicilor constituie una din provocările majore la aplicarea RL în probleme din lumea reală. O reprezentare exhaustivă printr-o tabelă de asocieri (*lookup table*) este posibilă numai în cazuri simple, pentru care mulțimea stărilor este finită și suficient de restrânsă. În general, trebuie utilizată o aproximare printr-o funcție parametrică, liniară sau neliniară după caz, pe baza unor caracteristici (*features*) conținând informația necesară și care să fie disponibile pentru sistemul de învățare. Astfel, majoritatea aplicațiilor RL depind în mare măsură de seturi de caracteristici atent și migălos manufacturate anterior pe baza cunoștințelor și intuiției umane despre fiecare problemă specifică abordată.

Odată cu popularizarea algoritmului *backpropagation* ca metodă de învățare a reprezentărilor interne (a caracteristicilor) specifice problemei / sarcinii (Rumelhart et al., 1986 [17]), au început să fie utilizate pentru aproximarea funcțiilor în RL rețelele neuronale artificiale (RNA) cu mai multe straturi ascunse. Au apărut astfel aplicații cu performanțe comparabile cu, sau chiar peste cele umane ca de exemplu: TD-Gammon (Tesauro et. al., 1992; 1994; 1995; 2002 [21; 22; 23; 24]) în abordarea jocului de table (*backgammon*), sistem care a fost ulterior adaptat și la crearea strategiei de pariare „Daily Double” utilizate de IBM Watson în Jeopardy! (Tesauro et al., 2012; 2013 [26; 27]) alături de alte strategii avansate. Majoritatea acestor aplicații necesitau ca intrare a rețelei un set de caracteristici specifice, special manufacturate pentru problema dată.

Mnih et al. (2013; 2015) [14; 15] de la Google DeepMind au realizat o demonstrație impresionantă a modului în care o rețea neuronală artificială adâncă poate automatiza procesul de generare/extragere a caracteristicilor. Aceștia au dezvoltat un agent RL numit „Deep Q-Network” (DQN), combinând pentru prima dată *Q-learning* cu o rețea neuronală adâncă de tip convoluțional (*convolutional neural network*) specializată în procesarea imaginilor (Vrejoiu, 2019 [31]) și au arătat că acest agent poate atinge un nivel înalt de performanță pentru o colecție de probleme diferite fără a necesita câte un set de caracteristici prefabricate specifice fiecăreia dintre acestea.

Rețeaua convoluțională este utilizată pentru aproximarea funcției optimale de evaluare a valorii așteptate a acțiunilor (*optimal action-value function*):

$$Q^*(s, a) = \max_{\pi} E[R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots \mid s_t = s, a_t = a, \pi], \quad (1)$$

dată de valoarea așteptată maximă a sumei recompenselor R , diminuate cu factorul de reducere γ la fiecare pas ulterior în timp, care pot fi obținute cu o politică $\pi = P(a|s)$ după ce se face o observație s și se execută o acțiune a . Ponderile rețelei convoluționale reprezintă parametrii funcției de evaluare a valorii acțiunilor (*action-value function*) care aproximează funcția optimală, $Q(s, a; w) \approx Q^*(s, a)$.

Ideea de bază a mai multor algoritmi de RL utilizată și de DQN constă în îmbunătățirea estimării funcției valorii acțiunii Q – practic a aproximării parametrice a acesteia – prin actualizări iterative pe baza estimării curente (*bootstrapping*) utilizându-se ecuația Bellman (Bellman, 1952 [4]) pentru valorile din oricare două iterații succesive:

$$Q_{i+1}(s, a; w) = E_{s'} [r + \gamma \max_{a'} Q_i(s', a'; w) \mid s, a]. \quad (2)$$

DQN a fost antrenat să joace prin interacțiune cu un emulator un număr de 49 de jocuri video Atari 2600 diferite. Practic, DQN a învățat câte o politică diferită pentru fiecare dintre cele 49 de jocuri – ponderile rețelei convoluționale fiind resetate la valori aleatoare înainte de învățarea fiecărui nou joc, – dar a utilizat aceeași intrare (pixeli imagine), aceeași arhitectură a rețelei și

aceleași valori ale parametrilor (e.g. mărimea pasului, factorul de reducere, parametrii de explorare și alți parametri specifici implementării) pentru toate jocurile. Chiar dacă aceste jocuri au în comun faptul că sunt jucate prin urmărirea fluxurilor de imagini video, ele sunt foarte diferite în alte aspecte: acțiunile au efecte diferite, dinamica tranziției stărilor e diferită și necesită politici diferite pentru a se obține scoruri mari. Totuși, DQN a atins niveluri de joc comparabile sau superioare nivelului uman pentru o mare parte dintre aceste jocuri, rețeaua convoluțională învățând să transforme intrarea comună tuturor jocurilor (cadrele video) în caracteristici specifice, specializate în reprezentarea valorii acțiunilor necesare pentru a juca la nivelul înalt atins de DQN pentru majoritatea jocurilor.

În perioada 1977-1992, Atari Inc. a produs și comercializat cu mare succes diverse versiuni ale consolei de jocuri video pentru acasă Atari 2600, care a introdus sau popularizat numeroase jocuri video Arcade clasice (e.g.: Pong, Breakout, Space Invaders, Asteroids etc.), încă jucate cu entuziasm și astăzi. În paralel, aceste jocuri au reprezentat de asemenea și excelente medii de testare pentru dezvoltarea și evaluarea metodelor de RL (Diuk et al., 2008 [7]; Naddaf, 2010 [16]; Cobo et al., 2011 [6]; Bellemare et al., 2012 [2]). Pentru a facilita utilizarea jocurilor Atari 2600 în testarea și evaluarea algoritmilor de învățare și planificare, Bellemare et al., 2013 [3] au dezvoltat emulatorul și mediul de învățare disponibil public denumit Arcade Learning Environment (ALE).

Toate acestea au constituit premisele care au condus la alegerea colecției de jocuri Atari 2600 pentru demonstrația grupului Mnih et al., care a fost de asemenea inspirat și de performanța la nivel uman atinsă de TD-Gammon. Astfel, DQN se aseamănă până la un punct cu TD-Gammon în utilizarea unei rețele neuronale artificiale adânci ca metodă de aproximare de funcții neliniare cu forma cu semi-gradient a algoritmului TD al *Q-learning*, fără model și fără politică fixă explicită (*model-free* și *off-policy*) și gradientii calculați cu *backpropagation*.

Arhitectura de bază a DQN, reprezentată schematic în Figura 2, este una de tip rețea neuronală convoluțională adâncă, cu reducere dimensională (*sub-sampling*) în fiecare strat convoluțional și hărți de caracteristici (*feature maps*) compuse din neuroni conectați fiecare doar la câte un grup local de neuroni din stratul anterior, care alcătuiesc câmpul lor receptor (*receptive field*). Arhitectura exactă este prezentată în cele ce urmează, detaliile fiind sintetizate (și) în Tabelul 2. La intrare, rețeaua DQN a utilizat cadrele 210×160 pixeli cu 128 de culori (cu componente R, G, B) la 60 Hz ale celor 49 de jocuri Atari, dar nu ca atare în mod direct. Din considerente de reducere a memoriei și procesărilor necesare, aceste cadre au fost mai întâi preprocesate fiind transformate în matrici de 84×84 de valori ale luminanței ($Y = 0,2126 \cdot R + 0,7152 \cdot G + 0,0722 \cdot B$). Totodată, deoarece stările complete ale multora din jocurile Atari nu sunt integral observabile dintr-un singur cadru video, dar și pentru a putea fi surprinsă dinamica elementelor de interes și viteza acestora, au fost grupate „stive” (*stacks*) de câte 4 astfel de matrici în niveluri de gri provenite din cele mai recente 4 cadre video, acestea constituind practic intrarea rețelei, de dimensiune 84×84×4.

Rețeaua DQN este compusă din 3 straturi ascunse convoluționale succesive, urmate de un al patrulea strat ascuns complet conectat, urmat de stratul de ieșire la rândul său complet conectat la acesta din urmă. Cele 3 straturi convoluționale succesive (Conv1 – Conv3) utilizează: 32 filtre 8×8 cu pas (*stride*) 4, 64 filtre 4×4 cu pas 2 și respectiv 64 filtre 3×3 cu pas 1, producând: 32 de hărți de caracteristici de dimensiuni 20×20, 64 de hărți de caracteristici de dimensiuni 9×9 și respectiv 64 de hărți de caracteristici de dimensiuni 7×7. Funcția de activare a fiecărui neuron din fiecare hartă de caracteristici este ReLU (*Rectified Linear Unit*: $\max(0, x)$). Cei 3.136 (7×7×64) neuroni din cel de-al treilea strat convoluțional sunt fiecare conectați la fiecare dintre cei 512 neuroni ReLU din stratul complet conectat (*fully connected*) FC1, care la rândul lor sunt fiecare conectați la toate cele 18 unități ale stratului de ieșire FC2, câte una pentru fiecare acțiune posibilă într-un joc Atari. Nivelurile de activare ale ieșirilor DQN sunt date de valorile optime ale acțiunilor, $Q^*(s, a)$, estimate/prezise pentru starea reprezentată de intrarea respectivă a rețelei.

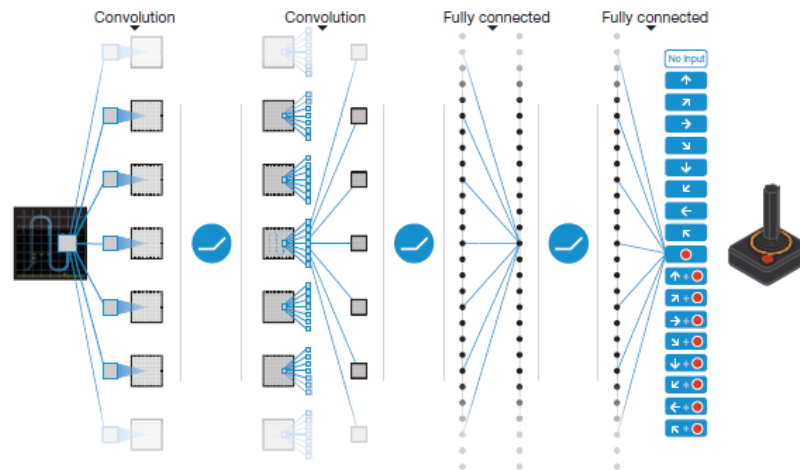


Figura 2. Ilustrarea schematică de principiu a rețelei neuronale convoluționale DQN (sursa: Mnih et al., 2015 [15])

Tabelul 2. Detalierea sintetică a arhitecturii DQN

Strat	Intrare	Dim. filtre	Pas	Nr. filtre	Activare	Ieșire
Conv1	$84 \times 84 \times 4$	8×8	4	32	ReLU	$20 \times 20 \times 32$
Conv2	$20 \times 20 \times 32$	4×4	2	64	ReLU	$9 \times 9 \times 64$
Conv3	$9 \times 9 \times 64$	3×3	1	64	ReLU	$7 \times 7 \times 64$
FC1	$7 \times 7 \times 64$			512	ReLU	512
FC2	512			18	Liniară	18

Principalul avantaj al acestui tip de arhitectură este reprezentat de faptul că permite estimarea valorilor pentru toate acțiunile posibile într-un singur pas de parcurgere înainte (*forward*) a rețelei. Asocierea unităților de ieșire cu acțiunile diferite de la un joc la altul și, deoarece numărul acțiunilor posibile a variat între 4 și 18, nu toate unitățile de ieșire au avut un rol funcțional în toate jocurile. Rețeaua poate fi gândită ca 18 rețele separate, câte una pentru estimarea valorii optime pentru fiecare acțiune posibilă. În realitate, aceste rețele partajază straturile inițiale, dar unitățile de ieșire învață să utilizeze caracteristicile extrase de acestea în moduri diferite.

La învățare, semnalul de recompensă utilizat în DQN a fost standardizat pentru toate jocurile printr-un singur parametru asociat fiecărui pas, care indică modul în care se modifică scorul jocului de la un pas la următorul în timp prin valorile:

- +1 pentru cazurile în care scorul crește,
- -1 pentru acelea în care scorul scade,
- 0 dacă scorul rămâne nemodificat,

ceea ce a permis buna funcționare pentru toate jocurile în ciuda plajelor variate de evoluție a scorurilor acestora. Această modalitate de trunchiere standardizată a recompenselor limitează amplitudinea derivatelor erorii (la calculul gradientilor în *backpropagation*) și permite utilizarea aceleiași rate de învățare α pentru mai multe jocuri.

DQN utilizează un algoritm fără model (*model-free*), folosind direct cadrele video (imaginile) provenite de la emulator fără estimare explicită a dinamicii recompenselor și tranzițiilor $P(r, s'|s, a)$, precum și fără politică fixă explicită (*off-policy*), învățând o politică (*greedy*) $a = \operatorname{argmax}_a(Q(s, a'; w))$ urmărind o distribuție de comportamente care să asigure explorarea adecvată a spațiului stărilor. Practic, o politică ϵ -greedy asigură selecția acestei distribuții la fiecare pas, cu probabilitate $1 - \epsilon$ de alegere aleatoare uniformă a unei acțiuni optime deja experimentate anterior,

alternând cu selectarea cu probabilitate ε a unei acțiuni noi. Parametrul $\varepsilon \in [0, 1]$ descrește liniar de la 1,0 la 0,1 pentru primul milion de cadre, rămânând apoi nemodificat la valoarea 0,1 pentru tot restul sesiunii de învățare. Astfel, în primele etape ale învățării este favorizată explorarea situațiilor care inițial sunt toate noi, pentru ca treptat să ajungă să fie favorizată exploatarea experiențelor deja învățate anterior, permițându-se totuși ocazional și explorarea câte unei acțiuni potențial noi care poate eventual duce la o îmbunătățire a experienței. Valorile celorlalți diferiți hiperparametri și parametri de optimizare, precum mărimea pasului de învățare (*learning step size*), factorul de reducere (*discount rate*) și alții specifici implementării, au fost selectate empiric prin căutări informale pentru un număr mic (cinci) de jocuri selectate pentru a se vedea care valori au funcționat cel mai bine. Aceste valori au fost apoi păstrate fixe pentru toate celelalte jocuri.

După ce DQN selectează la un moment dat t o acțiune A_t , aceasta este executată de către emulatorul de jocuri care întoarce o recompensă R_{t+1} și următorul cadru video. Acest nou cadru este preprocesat și actualizează stiva ultimelor 4 cadre cele mai recente S_t care devine astfel următoarea intrare S_{t+1} pentru rețea.

Pentru ajustarea ponderilor rețelei, DQN utilizează următoarea formă cu semi-gradient de diferență temporală (*temporal difference* – TD) din *Q-learning*:

$$w_{t+1} = w_t + \alpha [R_{t+1} + \gamma \max_a Q(S_{t+1}, a; w_t) - Q(S_t, A_t; w_t)] \nabla Q(S_t, A_t; w_t) \quad (3)$$

unde w_t este vectorul ponderilor rețelei, A_t este acțiunea selectată la momentul t , S_t și S_{t+1} reprezintă stivele de imagini preprocesate aplicate la intrarea rețelei la momentele t și respectiv $t+1$, α este rata de învățare ($\alpha = 0,00025$), R_{t+1} este valoarea recompensei la momentul $t+1$, γ este factorul de reducere a recompensei de la fiecare pas următor ($\gamma = 0,99$), a sunt acțiunile posibile la momentul $t+1$, iar Q este valoarea funcției de evaluare a valorii acțiunii (*action-value function*) pentru o pereche stare-acțiune, aproximată parametric prin ponderile curente ale rețelei, w_t .

În DQN, gradientul este calculat cu *backpropagation*, considerându-se câte o rețea separată pentru fiecare acțiune, pentru ajustarea la momentul t aplicându-se *backpropagation* numai pe rețeaua corespunzătoare acțiunii A_t . A fost utilizată funcția *loss* Huber (pătratică pentru valori mici ale argumentului și liniară în rest) și metoda *stochastic gradient descent* (SGD) pe mini-loturi (*mini-batches*), cu ajustarea ponderilor numai după acumularea informației de gradient pe câte un lot restrâns de imagini (în acest caz, 32), față de procedura *gradient descent* standard cu ajustarea ponderilor la fiecare iterație (acțiune). Pentru accelerarea învățării a fost utilizat algoritmul de optimizare RMSProp (Tieleman & Hinton, 2012 [28]) bazat pe creșterea gradientului (*gradient-ascent*), cu determinarea dinamică a pasului de ajustare pentru fiecare pondere pe baza mediei mobile (*moving average*) a magnitudinilor gradientilor recentți pentru ponderea respectivă.

În *Q-learning*, atât intrările cât și țintele se modifică continuu de-a lungul procesului și fac ca învățarea să fie instabilă. Valorile țintă Q depind în fiecare iterație de valoarea Q la momentul respectiv (țintă nestaționară). Totodată, în cazul învățării secvențiale din stări consecutive, apropiate pe o anumită traiectorie, corelațiile inerente între acestea vor determina și creșterea valorii Q țintă odată cu creșterea valorii Q curente, ceea ce pune probleme de convergență. Pentru a depăși aceste inconveniente, metoda utilizată efectiv de grupul de la DeepMind pentru funcționarea DQN utilizează trei tehnici importante care modifică din punct de vedere funcțional procedura standard *Q-learning*.

Prima tehnică utilizată este denumită „reluarea experienței” (*experience replay*) și a fost introdusă inițial de Lin (1992) [13], și constă în păstrarea/stocarea inițială a experienței agentului la fiecare moment de timp într-o „memorie de reluare” (*replay memory*) și utilizarea fiecărei astfel de experiențe memorate pentru actualizarea ponderilor mai târziu, la momente ulterioare. Astfel, în DQN, după executarea de către emulatorul de jocuri a unei acțiuni A_t din starea reprezentată de stiva de imagini S_t urmată de întoarcerea recompensei R_{t+1} asociate stivei de imagini S_{t+1} , tuplul $(S_t, A_t, R_{t+1}, S_{t+1})$ este adăugat în această memorie de reluare – practic un *buffer* ciclic – în care se acumulează experiențele mai multor sesiuni/episoade pentru același joc (concret, cele mai recente 1 milion de cadre video). La învățare, la fiecare pas se fac mai multe actualizări *Q-learning*, fiecare pe baza câte unui mini-lot (*mini-batch*) de experiențe eșantionate aleator uniform din această memorie de reluare. Astfel, în loc ca S_{t+1} să devină noua S_t și să fie utilizată pentru următoarea

actualizare ca în *Q-learning*, o experiență neconectată din mini-lotul curent extras din memoria de reluare va fi cea utilizată la următoarea actualizare. Aceasta este posibil deoarece *Q-learning* este un algoritm fără politică fixă explicită (*off-policy*) și nu necesită să fie aplicat de-a lungul unor traiectorii conectate (secvențe de stări corelate). Această tehnică de reluare a experienței oferă câteva avantaje față de forma uzuală a *Q-learning*. În acest mod, DQN învață mai eficient din propriile experiențe prin utilizarea fiecărei experiențe stocate în memoria de reluare pentru (mai) multe actualizări. De asemenea, este redusă varianța deoarece actualizările succesive nu mai sunt corelate între ele ca în *Q-learning* standard. Totodată, este eliminată dependența experiențelor succesive de aceleași ponderi curente care se actualizează, ceea ce elimină o sursă de instabilitate.

A doua tehnică utilizată de DQN și care modifică *Q-learning* standard a vizat direct îmbunătățirea stabilității metodei. Valoarea țintă utilizată la o actualizare în *Q-learning* depinde de estimarea curentă a funcției de evaluare a valorii acțiunii (*action-value function*, Q) care este reprezentată în DQN printr-o aproximare parametrică, depinzând chiar de parametrii care se actualizează ai rețelei. Astfel, pentru actualizarea dată de relația (3) valoarea țintă este: $\gamma \max_a Q(S_{t+1}, a; w_t)$, dependența acesteia de ponderile w_t (care nu apare în cazul învățării supervizate simple în care țintele nu depind de parametrii care se actualizează) putând conduce la oscilații și/sau divergență (Tsitsiklis & Roy, 1997 [29]). Pentru a se evita astfel de situații, în DQN este utilizată o tehnică prin care valorile ponderilor rețelei valorii acțiunii la care s-a ajuns după un număr C ($C = 10.000$) de actualizări efectuate asupra lor sunt preluate/„înghețate” într-o (altă) rețea duplicat și menținute fixe în aceasta pentru următoarele C actualizări ș.a.m.d. Pentru fiecare actualizare, este utilizată ca valoare țintă *Q-learning* ieșirea corespunzătoare a acestei rețele duplicat pentru intrarea respectivă aplicată. Notând cu Q^- ieșirea rețelei duplicat/„țintă” și cu w^- ponderile fixe ale acesteia, regula de actualizare (3) devine:

$$w_{t+1} = w_t + \alpha [R_{t+1} + \gamma \max_a Q^-(S_{t+1}, a; w^-) - Q(S_t, A_t; w_t)] \nabla Q(S_t, A_t; w_t), \quad (3')$$

Reluarea experienței este principala „responsabilă” pentru îmbunătățirea performanței DQN. Rețeaua țintă aduce și ea o îmbunătățire semnificativă dar nu atât de critică, însă devine mai importantă când capacitatea rețelei este mică.

În fine, o a treia modificare a mai fost adusă *Q-learning* standard, având de asemenea drept scop îmbunătățirea stabilității și a constat în limitarea termenului de eroare TD din (3'), $\delta = R_{t+1} + \gamma \max_a Q^-(S_{t+1}, a; w^-) - Q(S_t, A_t; w_t)$, astfel încât să se mențină în intervalul $[-1, 1]$:

$$\delta' = \max(-1, \min(1, \delta)) \in [-1, 1]. \quad (4)$$

Această limitare corespunde utilizării unei funcții valoare absolută (modul) pentru erorile situate în afara intervalului $(-1, 1)$ deoarece funcția *loss* valoare absolută $|x|$ are derivata -1 pentru toate valorile x negative și derivata 1 pentru toate valorile x pozitive.

Algoritmul *deep Q-learning* cu reluarea experienței aplicat în DQN pe M episoade (sesiuni) ale unui joc poate fi descris după cum urmează (Φ este stiva ultimelor 4 cadre preprocesate 84×84):

```

Inițializare memorie de reluare D de capacitate N;
Inițializare funcție aproximare valori acțiuni Q cu ponderi aleatoare w;
Inițializare funcție valori țintă acțiuni Q^- cu ponderi w^- = w;
Pentru episod ep = 1 ÷ M execută:
{
  Inițializare secvență S1 = {x1} și stivă cadre preprocesate Φ1 = Φ(S1);
  Pentru pas t = 1 ÷ T execută:
  {
    Cu probabilitate ε: selectare acțiune aleatoare At,
    altfel: selectare At = argmaxa Q(Φ(St), a; w);
    Execuție acțiune At în emulator, cu întoarcere recompensă Rt și
    imagine xt+1;
  }
}

```

```

Setare  $S_{t+1} = S_t, A_t, x_{t+1}$  și preprocesare  $\Phi_{t+1} = \Phi(S_{t+1});$ 

Memorare tranziție  $(\Phi_t, A_t, R_t, \Phi_{t+1})$  în D;

Eșantionare mini-lot de 32 tranziții  $(\Phi_j, a_j, r_j, \Phi_{j+1})$  din D;

Pentru fiecare tranziție din mini-lot:
    dacă episodul se termină la pasul  $j + 1$ , atunci: setare  $y_j = r_j$ ,
    altfel: setare  $y_j = r_j + \gamma \max_{a'} Q^-(\Phi_{j+1}, a'; w^-);$ 

Execuție un pas gradient descent pe suma pe mini-lotul curent a
erorilor pătratice  $(y_j - Q(\Phi_j, a_j; w))^2$  relativ la parametrul
rețelei  $w$ , cu actualizare corespunzătoare a acestora;

La fiecare C pași: resetare  $Q^- = Q$  (cu ponderi  $w^- = w$ );
}
}

```

Pentru dezvoltarea și implementarea inițială a DQN a fost utilizat limbajul Lua pe o mașină Linux cu GPU și NVIDIA® CUDA® (v.5.5), LuaJIT și Torch 7.0, mnggraph, Xitari (o ramificație a Arcade Learning Environment, ALE – Bellemare et al., 2013 [3]) și AleWrap (o interfață Lua la Xitari) instalate. Codul sursă original este disponibil online pentru utilizare în scopuri necomerciale la adresa <https://sites.google.com/a/deepmind.com/dqn>.

Mnih et al., 2015 [15] au evaluat și analizat contribuția diferitelor tehnici caracteristice utilizate în proiectarea DQN amintite mai sus și modul cum afectează acestea performanța, prin efectuarea mai multor sesiuni/episoade de învățare pentru o selecție de 5 dintre jocurile Atari cu testarea tuturor celor 4 combinații posibile cu reluarea experienței și rețeaua duplicat țintă incluse sau neincluse. Chiar dacă rezultatele au diferit de la joc la joc, s-a constatat că atât fiecare dintre aceste caracteristici individual a condus la o îmbunătățire a performanței, cât mai ales utilizarea ansamblului tuturor acestora a îmbunătățit substanțial performanța. Totodată, a fost analizat rolul jucat de rețeaua neuronală convoluțională adâncă în abilitatea de învățare a DQN prin compararea cu o versiune utilizând o rețea cu doar un singur strat liniar, ambele primind la intrare aceleași cadre video preprocesate și grupate în stive de câte 4. Îmbunătățirea adusă de versiunea convoluțională adâncă față de cea liniară a fost remarcabilă pentru toate cele 5 jocuri de test. Per global, scorurile obținute de DQN au fost comparate cu cele obținute de sistemul de învățare cel mai performant până la momentul respectiv (Bellemare et al. 2013 [3]), de un specialist expert în testare uman, precum și de un agent selectând acțiunile aleator. DQN a învățat fiecare joc interacționând cu emulatorul pentru 50 de milioane de cadre, ceea ce corespunde la circa 38 de zile de experiență cu jocul. La începutul învățării fiecărui joc, ponderile rețelei DQN au fost resetate la valori aleatoare. Pentru evaluare, au fost mediate scorurile a câte 30 de sesiuni/episoade pentru fiecare joc, fiecare sesiune/episod începând dintr-o stare inițială aleatoare a jocului și durând până la maxim 5 minute. DQN a obținut performanțe mai bune ca cel mai bun sistem RL anterior pentru 40 de jocuri din cele 49 și mai bune ca expertul uman în 22 de jocuri, iar comparabil cu acesta (cu scoruri peste 75% din scorurile acestuia) în încă alte 7 jocuri din cele 49. În afara acestor performanțe remarcabile, ceea ce este cel mai important este faptul că același sistem de învățare a atins aceste niveluri de performanță pe o gamă variată de jocuri fără a se baza pe cunoștințe și modificări specifice pentru fiecare dintre acestea.

Ulterior, au fost aduse (tot la DeepMind) mai multe îmbunătățiri DQN, ca de exemplu: **Double DQN** (van Hasselt et al., 2016 [30]), utilizând practic rețeaua *online* pentru alegerea *greedy* a acțiunii și rețeaua țintă pentru estimarea valorii Q ; **Prioritized experience replay** (Schaul et al., 2016 [18]) acordând o pondere mai mare ca frecvență de apariție în eșantionările din memoria de reluare (*replay memory*) acelor tranziții cu diferențe mai mari față de țintă (cu eroare TD mai mare), care pot aduce un progres mai rapid al învățării; **Dueling DQN** (Wang et al., 2016 [33]) calculând valorile funcției $Q(s, a)$ prin compunerea funcției valorilor stărilor $V(s)$ cu funcția avantajului acțiunii $A(s, a)$; **Noisy Nets for exploration**, adăugând zgomot parametric stratului liniar final pentru facilitarea explorării, ca alternativă la selectarea ϵ -*greedy* a acțiunilor.

4. Concluzii

Crearea de agenți artificiali autonomi capabili să exceleze într-o colecție de sarcini dificile diverse a reprezentat un deziderat dificil de atins al inteligenței artificiale. Multă vreme, necesitatea manufacturării de reprezentări specifice pentru fiecare tip de problemă în parte, precum și problemele de scalabilitate în cazul spațiilor cu dimensionalități mari au constituit piedici serioase în abordările de învățare automată (*machine learning*) vizând atingerea acestui scop. Proprietățile și performanțele remarcabile ale *deep learning* utilizând rețele neuronale adânci (*deep neural network*) combinate cu metode de învățare cu întărire / prin recompensă (*reinforcement learning*) din interacțiuni repetate cu mediul au permis recent depășirea acestor neajunsuri. Deep Q-Network (DQN) este unul dintre primele exemple de referință în *deep reinforcement learning*.

Prin DQN s-a demonstrat că un același agent poate învăța caracteristici specifice problemei astfel încât să capete abilități peste, sau comparabile cu cele umane pentru o plajă relativ variată de sarcini. Chiar dacă jocurile Atari necesită aptitudini/abilități diverse, DQN a învățat să le joace cu performanțe remarcabile exclusiv prin observarea imaginilor video ale acestora. În acest context, utilizarea rețelei neuronale convoluționale a reprezentat practic alegerea naturală potrivită. Chiar dacă DQN nu este un agent care excelează simultan în toate sarcinile (învățarea fiind făcută separat pentru fiecare dintre jocuri), a dovedit cu succes faptul că prin utilizarea tehnologiei *deep learning* se poate reduce și posibil elimina necesitatea proiectării și ajustării pentru fiecare problemă specifică. Totuși, trebuie făcută precizarea că DQN nu reprezintă o soluție completă pentru problema învățării independente de sarcină (*task-independent learning*). Performanța DQN a fost chiar considerabil sub nivelurile de abilitate umană pe câteva dintre jocurile Atari 2600, unele dintre acestea necesitând planificare pe un orizont lung de timp (*deep planning*), ceea ce excede capacităților DQN. Pe de altă parte, trebuie menționat și faptul că abilitatea de învățare controlată prin exersare extensivă repetată, așa cum DQN a învățat să joace jocurile Atari, este doar una dintre multiplele modalități de învățare pe care oamenii le realizează în mod obișnuit.

Chiar și cu limitările amintite mai sus, realizarea DQN a produs un important pas înainte în domeniul învățării automate (*machine learning*), demonstrând noile posibilități deschise prin mixarea metodelor de învățare cu/prin întărire/recompensă (*reinforcement learning*) cu tehnologiile moderne extrem de performante de *deep learning* și capacitățile de reprezentare a caracteristicilor din datele de intrare și de aproximare a funcțiilor neliniare ale rețelelor neuronale adânci.

Mențiuni

Prezenta lucrare are la bază parte din activitățile și rezultatele etapei I a proiectului PN 1937-0601 derulat la ICI București (2019) în cadrul Programului național nucleu „SMARTIC” finanțat de Ministerul Cercetării și Inovării.

BIBLIOGRAFIE

1. Arulkumaran, K., Deisenroth, M. P., Brundage, M., Bharath, A. A. (2017). *A Brief Survey of Deep Reinforcement Learning*. IEEE Signal Processing Magazine, special issue on Deep Learning for image understanding (ArXiv extended version – *arXiv:1708.05866v2* [cs.LG]).
2. Bellemare, M. G., Veness, J. & Bowling, M. (2012). *Investigating contingency awareness using Atari 2600 games*. In Proceedings of the Twenty-Sixth AAAI Conference on Artificial Intelligence (AAAI-12), pp. 864–871. AAAI Press, Menlo Park, CA.
3. Bellemare, M. G., Naddaf, Y., Veness, J., Bowling, M. (2013). *The Arcade learning environment: An evaluation platform for general agents*. Journal of Artificial Intelligence Research, 47:253–279.
4. Bellman, R. (1952). *On the Theory of Dynamic Programming*. PNAS, 38(8):716–719.
5. Bellman, R. (1957). *Dynamic Programming*. Princeton Univ. Press, New York.

6. Cobo, L. C., Zang, P., Isbell, C. L., Thomaz, A. L. (2011). *Automatic state abstraction from demonstration*. In Proceedings of the Twenty-Second International Joint Conference on Artificial Intelligence (IJCAI-11), pp. 1243–1248. AAAI Press.
7. Diuk, C., Cohen, A., Littman, M. L. (2008). *An object-oriented representation for efficient reinforcement learning*. In Proceedings of the 25th International Conference on Machine Learning, pp. 240–247. ACM, New York.
8. Duan, Y., Schulman, J., Chen, X., Bartlett, P. L., Sutskever, I., Abbeel, P. (2016). *RL2: Fast Reinforcement Learning via Slow Reinforcement Learning*. NIPS Workshop on Deep Reinforcement Learning.
9. Lake, B.M., Ullman, T. D., Tenenbaum, J. B., Gershman, S. J. (2016). *Building Machines That Learn and Think Like People*. The Behavioral and Brain Sciences, page 1.
10. Levine, S., Finn, C., Darrell, T., Abbeel, P. (2016). *End-to-End Training of Deep Visuomotor Policies*. JMLR, 17(39):1–40.
11. Levine, S., Pastor, P., Krizhevsky, A., Quillen, D. (2016). *Learning Hand-Eye Coordination for Robotic Grasping with Deep Learning and Large-Scale Data Collection*. ISER.
12. Li, K., Malik, J. (2017). *Learning to Optimize*. International Conference on Learning Representations – ICLR (preprint in *arXiv:1606.01885*, 2016).
13. Lin, L.-J. (1992). *Self-improving reactive agents based on reinforcement learning, planning and teaching*. Machine Learning, 8(3-4):293–321.
14. Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., Riedmiller, M. (2013). *Playing Atari with deep reinforcement learning*. *arXiv:1312.5602*.
15. Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G. et al. (2015). *Human-Level Control through Deep Reinforcement Learning*. Nature, 518(7540):529–533.
16. Naddaf, Y. (2010). *Game-Independent AI Agents for Playing Atari 2600 Console Games*. Ph.D. Thesis, University of Alberta, Edmonton.
17. Rumelhart, D. E., Hinton, G. E.; Williams R. J. (1986). *Learning representations by back-propagating errors*. Nature 323, 533–536.
18. Schaul, T., Quan, J., Antonoglou, I., Silver, D. (2016). *Prioritized experience replay*. ICLR (ArXiv version – *arXiv: 1511.05952v4* [cs.LG]).
19. Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., van den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M. et al. (2016). *Mastering the Game of Go with Deep Neural Networks and Tree Search*. Nature, 529(7587):484–489
20. Sutton, R. S., Barto, A. G. (1998). *Reinforcement Learning: An Introduction*. MIT Press (second edition, 2018).
21. Tesauro, G. (1992). *Practical issues in temporal difference learning*. Machine Learning, 8(3-4):257–277.
22. Tesauro, G. (1994). *TD-Gammon, a self-teaching backgammon program, achieves master-level play*. Neural Computation, 6(2):215–219.
23. Tesauro, G. (1995). *Temporal difference learning and TD-Gammon*. Communications of the ACM, 38(3):58–68.
24. Tesauro, G. (2002). *Programming backgammon using self-teaching neural nets*. Artificial Intelligence, 134(1-2):181–199.
25. Tesauro, G., Das, R., Chan, H., Kephart, J., Levine, D., Rawson, F., Lefurgy, C. (2008). *Managing Power Consumption and Performance of Computing Systems using Reinforcement Learning*. NIPS.

26. Tesauo, G., Gondek, D. C., Lechner, J., Fan, J., Prager, J. M. (2012). *Simulation, learning, and optimization techniques in Watson's game strategies*. IBM Journal of Research and Development, 56(3-4):16–1–16–11.
27. Tesauo, G., Gondek, D. C., Lenchner, J., Fan, J., Prager, J. M. (2013). *Analysis of Watson's strategies for playing Jeopardy!* Journal of Artificial Intelligence Research, 47:205–251.
28. Tieleman, T., Hinton, G. (2012). *Lecture 6.5–RMSPop*. COURSEERA: Neural networks for machine learning 4.2:26–31.
29. Tsitsiklis, J., Roy, B. V. (1997). *An analysis of temporal-difference learning with function approximation*. IEEE Trans. Automat. Contr. 42, 674–690.
30. van Hasselt, H., Guez, A., Silver, D. (2016). *Deep Reinforcement Learning with Double Q-learning*. Proc. of AAAI (ArXiv version – arXiv: arXiv:1509.0646).
31. Vrejoiu, M. H. (2019). *Rețele neuronale convoluționale, Big Data și Deep Learning în analiza automată de imagini*. Revista Română de Informatică și Automatică, 29(1):91-114.
32. Wang, J. X., Kurth-Nelson, Y., Tirumala, D., Soyer, H., Leibo, J. Z., Munos, R., Blundell, C., Kumaran, D., Botvinick, M. (2017). *Learning to Reinforcement Learn*. CogSci.
33. Wang, Z., Schaul, T., Hessel, M., Van Hasselt, H., Lanctot, M., de Freitas, N. (2016). *Dueling Network Architectures for Deep Reinforcement Learning*. arXiv:1511.06581v3 [cs.LG].
34. Watkins, C. J. C. H. (1989). *Learning from Delayed Rewards*. Ph.D. Thesis, University of Cambridge.
35. Watkins, C. J. C. H., Dayan, P. (1992). *Q-learning*. Machine Learning, 8(3-4):279–292.
36. Zoph, B., V Le, Q. (2017). *Neural Architecture Search with Reinforcement Learning*. ICLR.



Mihnea Horia VREJOIU este cercetător științific gradul III în Departamentul “Sisteme inteligente distribuite intensive ca date” din ICI București. Domeniile și subiectele sale de expertiză și interes cuprind: vedere artificială (prelucrare și analiză de imagini, recunoașterea formelor, recunoașterea optică de caractere – OCR, recunoașterea numerelor de înmatriculare – LPR) și învățare automată (clasificatoare, memorii asociative).

Mihnea Horia VREJOIU is scientific researcher 3rd degree in the “Distributed and Data Intensive Intelligent Systems” Department at ICI Bucharest. His main areas and topics of expertise and interest cover: Artificial Vision (Image Processing and Analysis, Pattern Recognition, Optical Character Recognition – OCR, License Plate Recognition – LPR) and Machine Learning (classifiers, associative memories).