

Control prin învățare

Master ICAF, An 1 Sem 2

Lucian Bușoniu, Ștefan Pîrje



Partea VIII

Deep Reinforcement Learning



Partea VI în plan

- Problema de învățare prin recompensă
- Soluția optimală
- Programarea dinamică exactă
- Învățarea prin recompensă exactă
- Tehnici de aproximare
- Programarea dinamică cu aproximare (var. continue)
- Învățarea prin recompensă cu aproximare (var. continue)
- Rețele neurale profunde
- **Învățarea prin recompensă cu rețele neurale profunde**



Partea VI în plan: Categoriile de algoritmi

După utilizarea unui model:

- **Bazat pe model:** f, ρ cunoscute
- **Fără model:** doar date (învățarea prin recompensă)

După nivelul de interacțiune:

- **Offline:** algoritmul rulează în avans
- **Online:** algoritmul controlează direct sistemul

Exact vs. cu aproximare:

- **Exact:** x, u număr mic de valori discrete
- **Cu aproximare:** x, u continue (sau multe valori discrete)



- 1 Aproximarea funcției cu rețele neuronale
- 2 Iterația Q neurală bazată pe date
- 3 Deep Q-Networks
- 4 Perspective



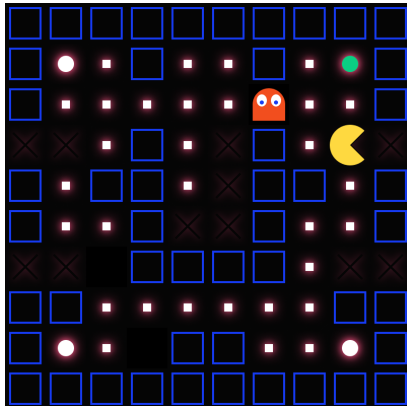
De ce rețele neurale în învățarea prin recompensă?



Motivație

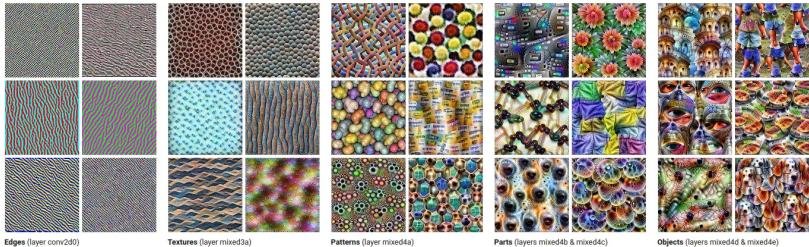
Chiar având acces la starea internă a jocului, ar trebui să extragem diverse atribute:

- distanțele până la fantome
- distanța până la cel mai apropiat power-up
- distanța până la pereți
- dacă Pac-Man este blocat
- distanța până la cele mai apropiate puncte
- puncte deja consumate



Aceste atribute sunt relevante doar pentru MDP-ul curent!

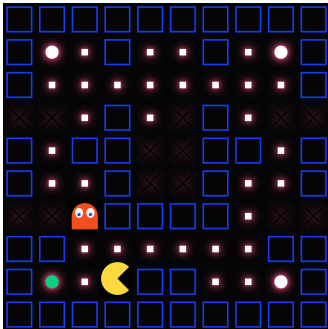
Rețele convoluționale: atribute automate!



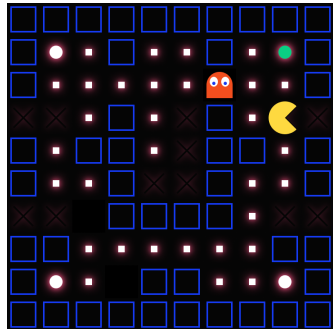
Reprezentări ierarhice învățate de o rețea convoluțională dintr-un set de date cu imagini naturale

Generalizare

Dorim Q antrenată pe $x \dots$



\dots să fie similară pentru alt x



- 1 Aproximarea funcției cu rețele neuronale
- 2 Iterația Q neurală bazată pe date**
- 3 Deep Q-Networks
- 4 Perspective



Reamintim: Iterația Q bazată pe date

Iterația Q bazată pe date

obține sau colectează setul de date

$$\mathcal{D} = \{(x_s, u_s, r_s, x'_s), s = 1, \dots, n_s\}$$

repeat la fiecare iterație ℓ

for $s = 1, \dots, n_s$ **do**

calculează ținta cu bootstrap pentru $\hat{Q}(x_s, u_s; \theta)$:

$$\hat{R}_s \leftarrow r_s + \gamma \max_{u'} \hat{Q}(x'_s, u'; \theta_\ell)$$

end for

$$\theta_{\ell+1} \leftarrow \arg \min_{\theta} \sum_{s=1}^{n_s} [\hat{R}_s - \hat{Q}(x_s, u_s; \theta)]^2 =: \mathcal{L}(\theta)$$

until terminare

\mathcal{L} funcție de pierdere



Iterația Q neurală bazată pe date

Utilizează NN pentru $\hat{Q}(x, u; \theta)$ și SGD pentru optimizare

Iterația Q neurală bazată pe date

inițializează $\hat{Q}(x, u; \theta)$ a.î. $Q(x, u; \theta) \approx 0, \forall x, u$

obține sau colectează setul de date

$\mathcal{D} = \{(x_s, u_s, r_s, x'_s), s = 1, \dots, n_s\}$

repeat la fiecare iterație ℓ

for $s = 1, \dots, n_s$ **do**

calculează ținta cu bootstrap pentru $\hat{Q}(x_s, u_s; \theta)$:

$\hat{R}_s \leftarrow r_s + \gamma \max_{u'} \hat{Q}(x'_s, u'; \theta_\ell)$

end for

rezolvă $\theta_{\ell+1} = \arg \min_{\theta} \mathcal{L}(\theta)$ **cu SGD pe batch-uri**
de b tranziții din \mathcal{D}

until terminare

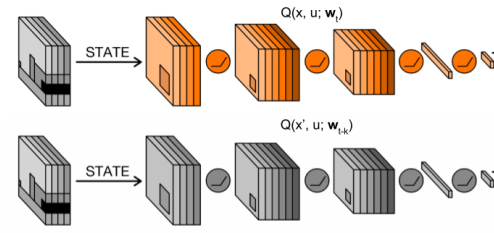
Problemă: actualizarea lui θ modifică țintele într-un mod corelat, ducând la acumularea erorilor.



- 1 Aproximarea funcției cu rețele neuronale
- 2 Iterația Q neurală bazată pe date
- 3 Deep Q-Networks**
- 4 Perspective



Deep Q-Networks. Fix 1: Rețea țintă



DQN introduce o rețea suplimentară pentru calcularea țintei de regresie, care “urmărește” rețeaua online

Această **rețea țintă** nu este antrenată separat – parametrii săi sunt copiați din rețeaua online la fiecare C pași

Deep Q-Networks. Fix 2: Reluarea experienței

DQN utilizează un buffer ciclic de reluare a experienței, care:

- permite antrenarea aproximatoarelor cu gradient stochastic descendent pe mini-batch-uri (ca în NFQI)
- permite îmbunătățiri rapide, optimiste ale legii de control fără a aștepta convergența NN



Alte ajustări importante

- recompensele sunt limitate la $[-1, 1]$
- un optimizator bun, precum Adam
- funcția de pierdere este Huber/Smooth L1 (L2 pătratic aproape de 0, L1 liniar pentru valori mari)
- în unele implementări se limitează și norma gradientului



Algoritm

Deep Q-Networks

init rețeaua Q $\hat{Q}(x, u; \theta)$ cu parametri aleatori θ

init rețeaua Q țintă $\bar{Q}(x, u; \bar{\theta})$ cu aceiași parametri $\bar{\theta} = \theta$

init buffer-ul \mathcal{D} cu capacitatea n_s

for fiecare traiectorie **do**

init x_0

repeat la fiecare pas k

alege u_k , de ex., ε -greedy din $Q(x_k, \cdot; \theta)$

aplică u_k , măsoară x_{k+1} , primește r_{k+1}

adaugă (x_k, u_k, r_k, x_{k+1}) în \mathcal{D} ; posibil eliminând eș. vechi

eșantionează un batch de b tranziții (x_s, u_s, r_s, x'_s) din \mathcal{D}

calculează țintele $\hat{R}_s = r_s + \gamma \max_{u'} \bar{Q}(x'_s, u'; \bar{\theta})$

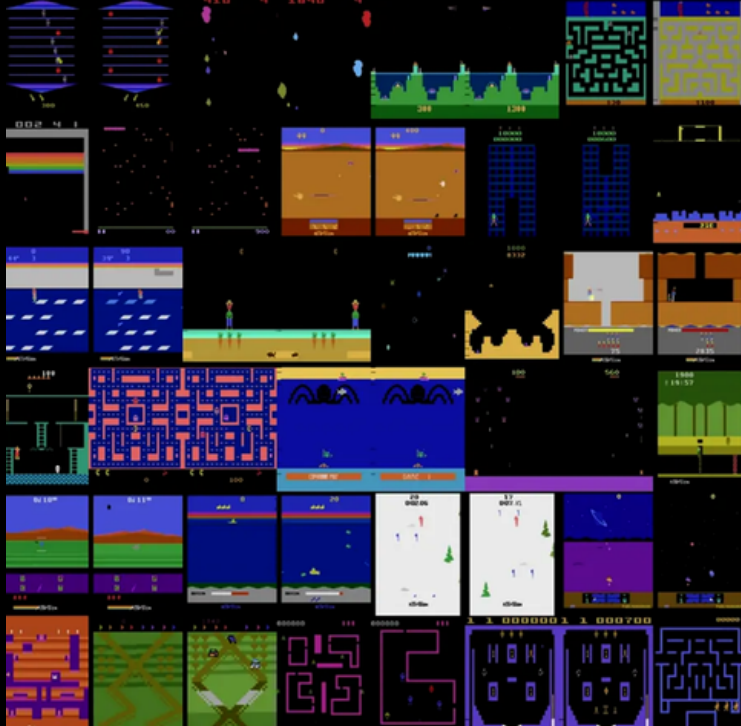
optimizează $\sum_{s=1}^b [\hat{R}_s - \hat{Q}(x_s, u_s; \theta)]^2$ față de θ

la fiecare C pași: $\bar{\theta} = \theta$

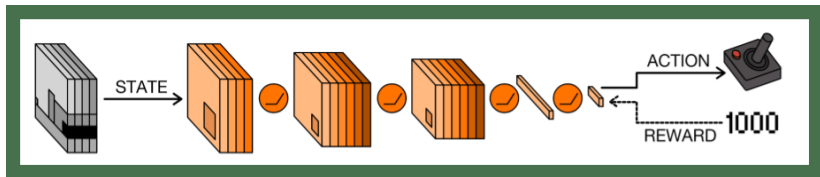
until traiectoria terminată

end for



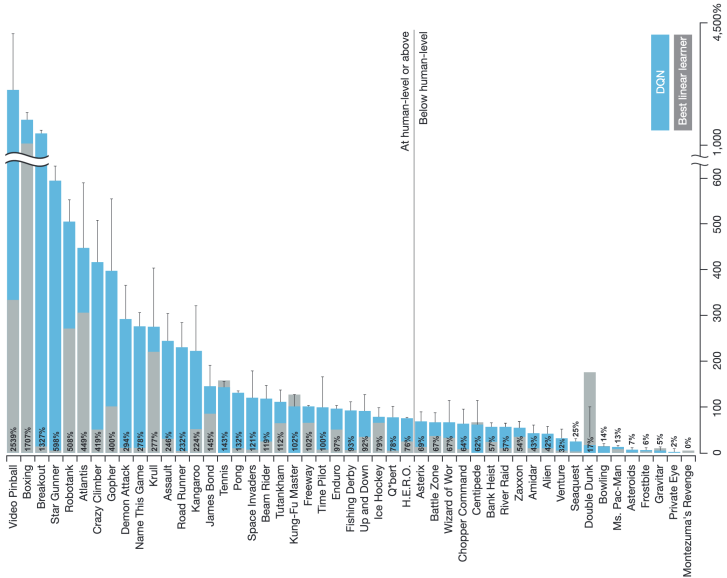


Arhitectura DQN pentru Atari



- primește ultimele 4 ecrane de joc în tonuri de gri
- trei straturi convoluționale ascunse
- un strat linear ascuns
- un strat linear de ieșire (uneori cu un bias comun pentru toate acțiunile)
- câte o ieșire pentru valoarea Q a fiecărei acțiune
- activări RELU

Performanța DQN relativ la un jucător uman



Dincolo de Deep Q-Networks

Multe direcții pentru îmbunătățirea DQN:

- obiective: Double-DQN, Dueling DQN, Munchausen-DQN, TD cu return pe n pași
- distribuții în loc de estimări punctuale: Categorical DQN, IQM
- eșantionare: Prioritized Experience Replay
- explorare: Random Network Distillation, Go-Explore, Bootstrapped DQN



- 1 Aproximarea funcției cu rețele neuronale
- 2 Iterația Q neurală bazată pe date
- 3 Deep Q-Networks
- 4 Perspective**



Conexiuni: Control optimal

- Controlează un sistem minimizând costul J
- Bazat pe model
- Posibil în timp continuu, orizont finit

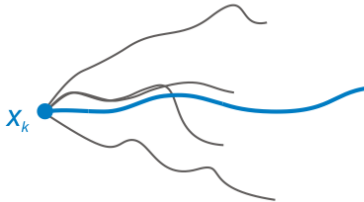
RL, DP **sunt control optimal!**

- Controlează un sistem maximizând returnul $V^h(x)$
- În timp discret, orizont infinit
- RL **fără model**, bazat pe date/interacțiune



Conexiuni: Control predictiv

- MPC: model-predictive control
- Variantă de control optimal, bazată pe model
- Principiu de bază: **receding horizon** (orizont alunecător)



Probleme deschise

RL & DP **în curs de dezvoltare**

Probleme deschise:

- Strategii de explorare
- Garanții de siguranță și stabilitate
- Stări care nu pot fi măsurate
- Sisteme multiagent
- Transfer, învățarea multi-task, învățarea continuă

