

Control prin învățare

Master ICAF, An 1 Sem 2

Lucian Bușoniu

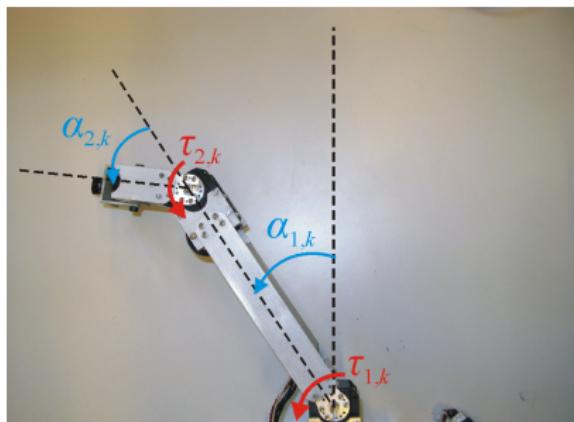


Partea V

Învățarea prin recompensă online cu
aproximare

Recap: Nevoia de aproximare

- În aplicații reale de control, x, u **continue!**



- Reprezentarea prin tabel **imposibilă**
- Aproximarea** funcțiilor de interes
 $Q(x, u)$, $V(x)$, $h(x)$ necesară

Partea 5 în plan

- Problema de învățare prin recompensă
- Soluția optimală
- Programarea dinamică (variabile discrete)
- Învățarea prin recompensă (variabile discrete)
- Tehnici de aproximare
- Programarea dinamică cu aproximare (var. continue)
- **Partea V: Învățarea prin recompensă cu aproximare (var. continue)**
- Planificarea online (var. continue și discrete)

Recap: Partea 4 – Algoritmi offline

pornind de la:

- model f, ρ
- sau date $(x_s, u_s, r_s, x'_s), s = 1, \dots, n_s$

- ➊ găsește soluție aproximată $\widehat{Q}(x, u)$, $\widehat{h}(x)$, etc.
- ➋ controlează sistemul folosind soluția găsită

Algoritmi exemplificări:

- iterația fuzzy Q
- iterația Q bazată pe date
- LSPI, iterația lege de control CMMP

Partea 5 în plan: Categorii 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)

Există mulți algoritmi, doar câțiva sunt selectați pentru discuție



Conținut partea 5

- 1 Învățarea Q și SARSA cu aproximare
- 2 Actor-critic
- 3 LSPI online
- 4 Accelerarea RL cu aproximare
- 5 Perspective

1 Învățarea Q și SARSA cu aproximare

- Învățarea Q approximată
- SARSA approximată

2 Actor-critic

3 LSPI online

4 Accelerarea RL cu aproximare

5 Perspective

Reamintim: Învățarea Q

Învățarea Q cu ε -greedy

for fiecare traiectorie **do**

 initializează x_0

repeat la fiecare pas k

$$u_k = \begin{cases} \arg \max_u Q(x_k, u) & \text{cu prob. } (1 - \varepsilon_k) \\ \text{aleatoare} & \text{cu prob. } \varepsilon_k \end{cases}$$

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

$$Q(x_k, u_k) \leftarrow Q(x_k, u_k) + \alpha_k \cdot$$

$$[r_{k+1} + \gamma \max_{u'} Q(x_{k+1}, u') - Q(x_k, u_k)]$$

until traiectoria terminată

end for

Diferența temporală: $[r_{k+1} + \gamma \max_{u'} Q(x_{k+1}, u') - Q(x_k, u_k)]$



Învățarea Q aproximată

- Învățarea Q **scade diferența temporală**

$$Q(x_k, u_k) \leftarrow Q(x_k, u_k) + \alpha_k [r_{k+1} + \gamma \max_{u'} Q(x_{k+1}, u') - Q(x_k, u_k)]$$

- $r_{k+1} + \gamma \max_{u'} Q(x_{k+1}, u')$ înlocuiește **idealul** $Q^*(x_k, u_k)$
 [Vezi și Bellman: $Q^*(x, u) = \rho(x, u) + \gamma \max_{u'} Q^*(x', u')$]

⇒ Ideal, scade eroarea $[Q^*(x_k, u_k) - Q(x_k, u_k)]$

Învățarea Q aproximată (continuare)

Aproximare: folosim $\hat{Q}(x, u; \theta)$, actualizăm **parametri**

- **Gradient pe eroarea** $[Q^*(x_k, u_k) - \hat{Q}(x_k, u_k; \theta)]$:

$$\begin{aligned}\theta_{k+1} &= \theta_k - \frac{1}{2} \alpha_k \frac{\partial}{\partial \theta} \left[Q^*(x_k, u_k) - \hat{Q}(x_k, u_k; \theta_k) \right]^2 \\ &= \theta_k + \alpha_k \frac{\partial}{\partial \theta} \hat{Q}(x_k, u_k; \theta_k) \cdot \left[Q^*(x_k, u_k) - \hat{Q}(x_k, u_k; \theta_k) \right]\end{aligned}$$

- Folosește **estimare** pentru $Q^*(x_k, u_k)$:

$$\theta_{k+1} = \theta_k + \alpha_k \frac{\partial}{\partial \theta} \widehat{Q}(x_k, u_k; \theta_k) -$$

$$\left[r_{k+1} + \gamma \max_{u'} \widehat{Q}(x_{k+1}, u'; \theta_k) - \widehat{Q}(x_k, u_k; \theta_k) \right]$$

(diferență temporală aproximată)

Învățarea Q aproximată: algoritm

Învățarea Q aproximată cu explorare ε -greedy

for fiecare trajectorie **do**

initializează x_0

repeat la fiecare pas k

$$u_k = \begin{cases} \arg \max_u \widehat{Q}(x_k, u; \theta_k) & \text{cu prob. } (1 - \varepsilon_k) \\ \text{aleatoare} & \text{cu prob. } \varepsilon_k \end{cases}$$

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

$$\theta_{k+1} = \theta_k + \alpha_k \frac{\partial}{\partial \theta} \hat{Q}(x_k, u_k; \theta_k).$$

$$\left[r_{k+1} + \gamma \max_{u'} \hat{Q}(x_{k+1}, u'; \theta_k) - \hat{Q}(x_k, u_k; \theta_k) \right]$$

until traiectoria terminată

end for

Desigur, **explorarea** necesară și în cazul aproximat



Reamintim: Maximizare

Soluția 1:

- Legea de control nu este reprezentată explicit
 - Acțiuni greedy calculate la cerere din \hat{Q}

Soluția 2:

- Legea de control aproximată explicită

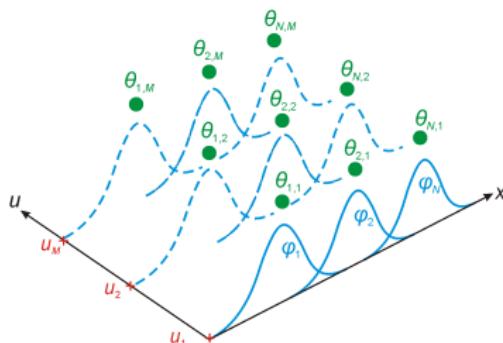
Maximizare în Învățarea Q aproximată

- Acțiuni greedy calculate la cerere din \hat{Q} :

$$\dots \max_u \hat{Q}(x, u; \theta) \dots$$

⇒ Soluția 1: Legea de control reprezentată implicit

- Aproximatorul funcției Q trebuie să garanteze
soluție eficientă pentru max
- Ex. acțiuni discrete & funcții de bază în x



Învățarea Q aprox.: demo mers robotic (E. Schuitema)

Aproximator: **tile coding**



1 Învățarea Q și SARSA cu aproximare

- Învățarea Q aproximată
- **SARSA** aproximată

2 Actor-critic

3 LSPI online

4 Accelerarea RL cu aproximare

5 Perspective

SARSA aproximată

Reamintim SARSA clasică:

$$Q(x_k, u_k) \leftarrow Q(x_k, u_k) + \alpha_k [r_{k+1} + \gamma Q(x_{k+1}, u_{k+1}) - Q(x_k, u_k)]$$

Aproximare: similar cu învățarea Q

- actualizăm parametri
- bazat pe **gradientul** funcției Q
- și diferența temporală aproximată

$$\theta_{k+1} = \theta_k + \alpha_k \frac{\partial}{\partial \theta} \hat{Q}(x_k, u_k; \theta_k) \cdot$$

$$\left[r_{k+1} + \gamma \hat{Q}(x_{k+1}, u_{k+1}; \theta_k) - \hat{Q}(x_k, u_k; \theta_k) \right]$$

SARSA

SARSA aproximată

for fiecare traекторie **do**

 initializează x_0

 alege u_0 (ex. ε -greedy din $Q(x_0, \cdot; \theta_0)$)

repeat la fiecare pas k

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

 alege u_{k+1} (ex. ε -greedy din $Q(x_{k+1}, \cdot; \theta_k)$)

$$\theta_{k+1} = \theta_k + \alpha_k \frac{\partial}{\partial \theta} \hat{Q}(x_k, u_k; \theta_k) \cdot$$

$$\left[r_{k+1} + \gamma \hat{Q}(x_{k+1}, u_{k+1}; \theta_k) - \hat{Q}(x_k, u_k; \theta_k) \right]$$

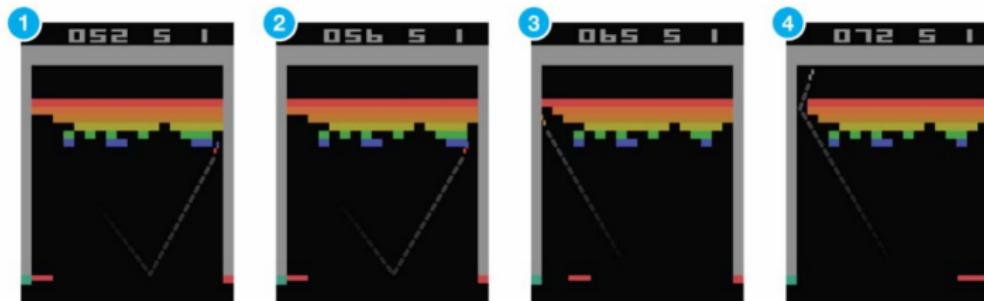
until traectoria terminată

end for

Învățarea Q cu “deep neural networks”

- Funcția Q reprezentată via o rețea neuronală $Q(x_{k+1}, \cdot; \theta_k)$
- Rețea neuronală “deep” cu multe nivele, cu structuri și funcții de activare specifice
- Rețeaua antrenată pentru a minimiza diferența temporală, similar cu algoritmul standard
- Antrenare pe mini-batch-uri de tranzitii, similar cu iterația Q bazată pe date \Rightarrow algoritmul combină RL online și offline

(DeepMind, [Human-level control through deep reinforcement learning](#), Nature 2015)



1 Învățarea Q și SARSA cu aproximare

2 Actor-critic

- Algoritm
- Exemplu

3 LSPI online

4 Accelerarea RL cu aproximare

5 Perspective

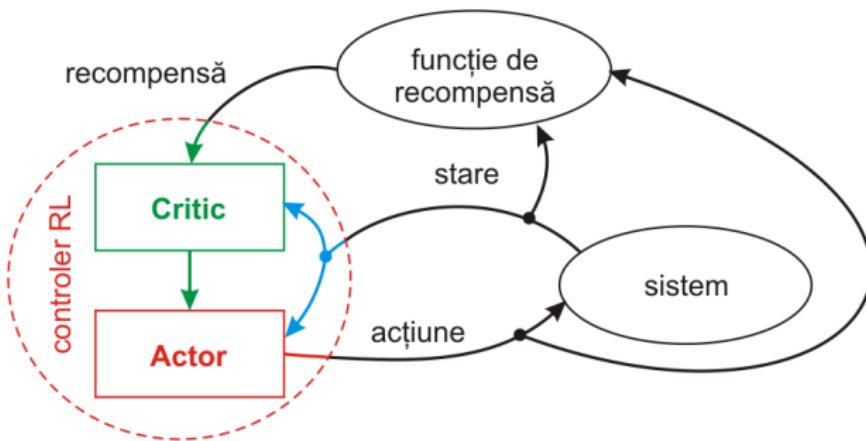
Lege de control explicită

- Soluția 2: Legea de control **aproximată explicită**: $\hat{h}(x; \vartheta)$

Avantaje:

- **Acțiuni continue** mai ușor de folosit
- Reprezentarea poate include mai ușor **cunoștințe expert**

Schema actor-critic



- **Actor**: legea de control $\hat{h}(x; \vartheta)$
- **Critic**: funcția V , $\hat{V}(x; \theta)$

Actualizarea criticului

- Gradient pe diferență temporală:

$$\begin{aligned}\theta &\leftarrow \theta + \alpha_k^{\text{critic}} \frac{\partial}{\partial \theta} \hat{V}(x_k; \theta) [r_{k+1} + \hat{V}(x_{k+1}; \theta) - \hat{V}(x_k; \theta)] \\ &= \theta + \alpha_k^{\text{critic}} \frac{\partial}{\partial \theta} \hat{V}(x_k; \theta) \Delta_k\end{aligned}$$

- Provine din ecuația Bellman pentru V^h :

$$V^h(x) = \rho(x, h(x)) + \gamma V^h(f(x, h(x)))$$

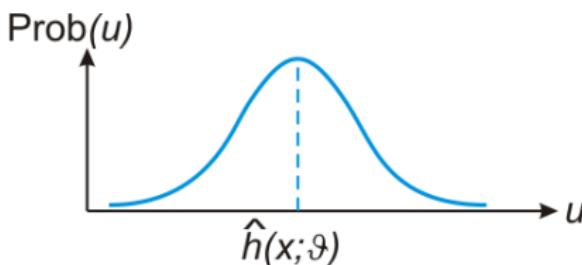
⇒ De fapt, evaluarea legii de control

Explorare-exploatare

- RL online \Rightarrow actor-critic trebuie să exploreze
- Ex. **explorare Gaussiană**

$$u_k = \hat{h}(x_k; \vartheta) + u_{\text{explor}}$$

unde termenul explorator u_{explor} Gaussian cu medie 0



Actualizarea actorului

$$\vartheta \leftarrow \vartheta + \alpha_k^{\text{actor}} \frac{\partial}{\partial \vartheta} \hat{h}(x_k; \vartheta) [u_k - \hat{h}(x_k; \vartheta)] \Delta_k$$

Intuiție:

- Dacă $\Delta_k > 0$, adică $r_{k+1} + \hat{V}(x_{k+1}; \theta) > \hat{V}(x_k; \theta)$, performanța mai bună decât cea așteptată
⇒ **apropie** de acțiunea u_k exploratoare
- Dacă $\Delta_k < 0$, performanța mai proastă
⇒ **îndepărtează** de u_k

Algoritmul actor-critic

Actor-critic

```
for fiecare traiectorie do
    initializează  $x_0$ 
    repeat la fiecare pas  $k$ 
         $u_k \leftarrow \hat{h}(x_k; \vartheta) + \text{explorare}$ 
        aplică  $u_k$ , măsoară  $x_{k+1}$ , primește  $r_{k+1}$ 
         $\Delta_k \leftarrow r_{k+1} + \hat{V}(x_{k+1}; \theta) - \hat{V}(x_k; \theta)$ 
         $\theta \leftarrow \theta + \alpha_k^{\text{critic}} \frac{\partial}{\partial \theta} \hat{V}(x_k; \theta) \Delta_k$ 
         $\vartheta \leftarrow \vartheta + \alpha_k^{\text{actor}} \frac{\partial}{\partial \vartheta} \hat{h}(x_k; \vartheta) [u_k - \hat{h}(x_k; \vartheta)] \Delta_k$ 
    until traiectoria terminată
end for
```

- De notat: rate de învățare diferite actor & critic

Actor-critic – iterație legea de control optimistă

- Reamintim: Actualizare critic
= 1 pas de **evaluare** a legii de control
 - Actualizare actor
= **îmbunătățire** incrementală a legii de control
- ⇒ Actor-critic ≈ **iterație pe legea de control**
- Actualizările alternează la fiecare tranziție
- ⇒ Actor-critic ≈ iter. legea de control **optimistă**

1 Învățarea Q și SARSA cu aproximare

2 Actor-critic

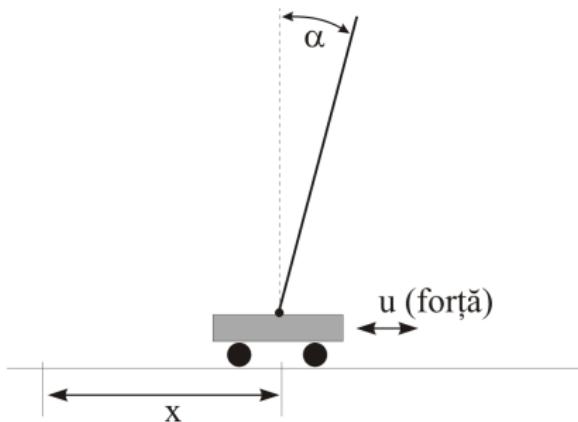
- Algoritm
- Exemplu

3 LSPI online

4 Accelerarea RL cu aproximare

5 Perspective

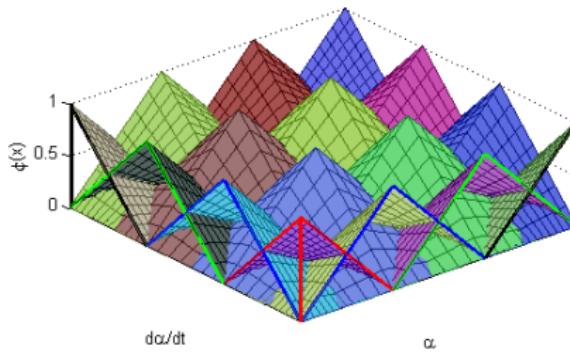
Exemplu: Pendulul inversat cu cărucior



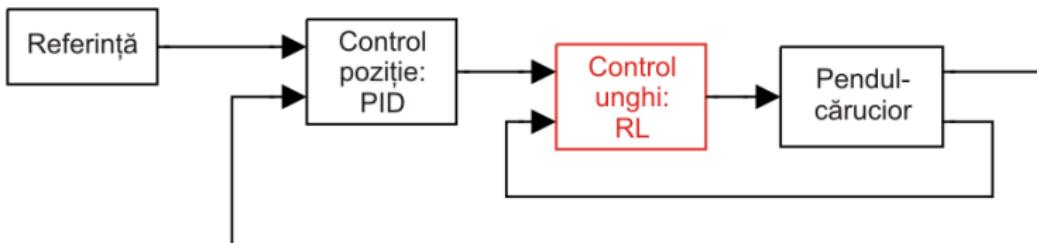
- Forța transmisă prin accelerarea căruciorului
- **Obiectiv:** căruciorul la poziție referință, menținând pendulul vertical
- Nu este nevoie de swingup

Pendul cu cărucior: Aproximator

Atât actor cât și critic: interpolare (aproximare “fuzzy”)

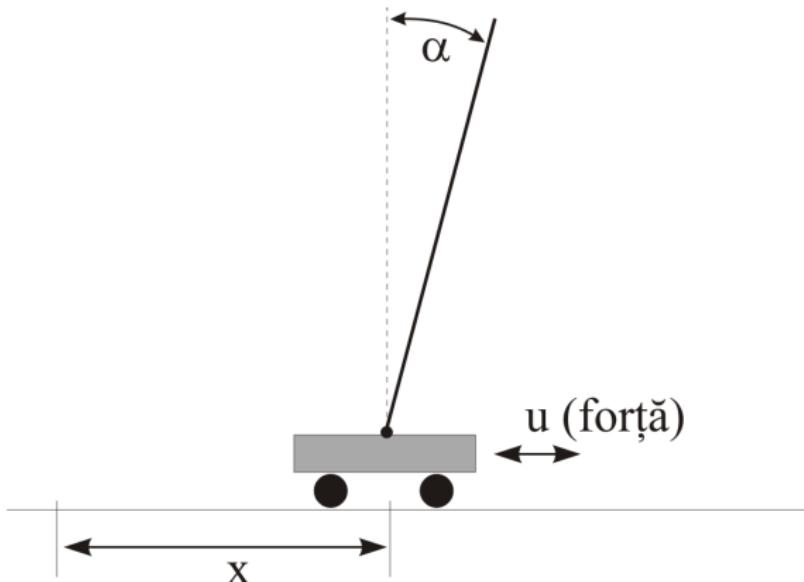


Schema de control



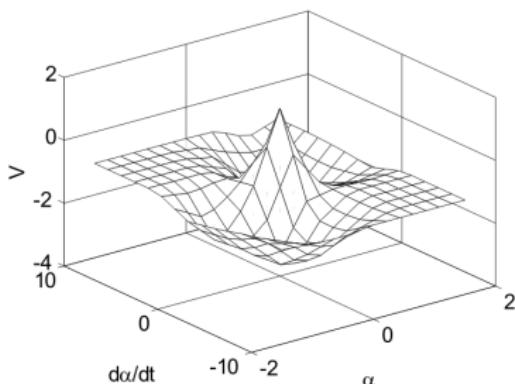
- Bucla externă, control pozitie: PID classic
- Bucla internă, control unghi: **actor-critic**

Pendulul cu cărucior: demo

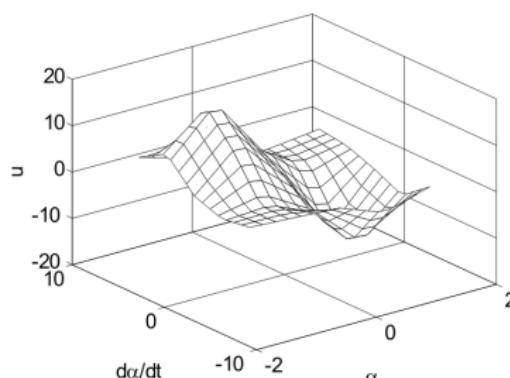


Rezultate

Suprafața criticului

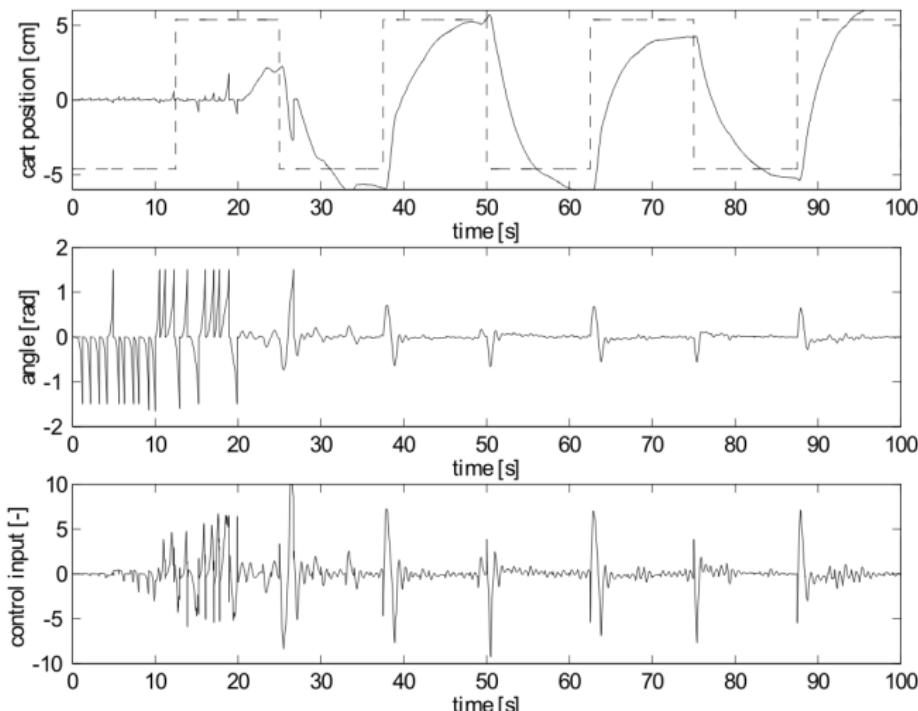


Suprafața actorului



Rezultate (continuare)

Traекторie de-a lungul învățării



- 1 Învățarea Q și SARSA cu aproximare
- 2 Actor-critic
- 3 LSPI online
- 4 Accelerarea RL cu aproximare
- 5 Perspective

Reamintim: LSPI

Iterația pe legea de control CMMP (LSPI)

date fiind (x_s, u_s, r_s, x'_s) , $s = 1, \dots, n_s$

repeat la fiecare iterare

Evaluarea legii de control:

$A \leftarrow 0$, $B \leftarrow 0$, $b \leftarrow 0$

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

actualizează A , B , b folosind (x_s, u_s, r_s, x'_s)

end for

rezolvă $A\theta = \gamma B\theta + b$ găsind θ

Îmbunătățirea legii de control: $h(x) \leftarrow \arg \max_u \hat{Q}(x, u; \theta)$

until terminare

LSPI online

- Nu există date în avans, colectează **online**, interactiv
- Îmbunătățește legea de control **optimist**
- De notat: A, B, b refolosite chiar dacă h se schimbă!

$A \leftarrow 0, B \leftarrow 0, b \leftarrow 0$; initializează $h(x)$

for fiecare traекторie **do**

repeat fiecare pas k

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

 actualizează A, B, b folosind $(x_k, u_k, r_{k+1}, x_{k+1})$

if au trecut K tranziții **then**

 rezolvă $A\theta = \gamma B\theta + b$ găsind θ

 Îmbunătățește $h(x) \leftarrow \arg \max_u \hat{Q}(x, u; \theta)$

end if

until traекторia terminată

end for



LSPI: Explorare-exploatare

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

- Explorare necesară, ex. ε -greedy:

$$u_k = \begin{cases} h(x_k) & \text{cu prob. } 1 - \varepsilon_k \\ \text{o acțiune aleatoare} & \text{cu prob. } \varepsilon_k \end{cases}$$

Exemplu: Pendul inversat



- $x = [\text{unghi } \alpha, \text{ viteza } \dot{\alpha}]^\top$
- $u = \text{voltaj}$
- $\rho(x, u) = -x^\top \begin{bmatrix} 5 & 0 \\ 0 & 0.1 \end{bmatrix} x - u^\top 1 u$
- Factor de discount $\gamma = 0.98$

- **Obiectiv:** stabilizează orientat în sus
- Putere insuficientă \Rightarrow balansează înainte & înapoi

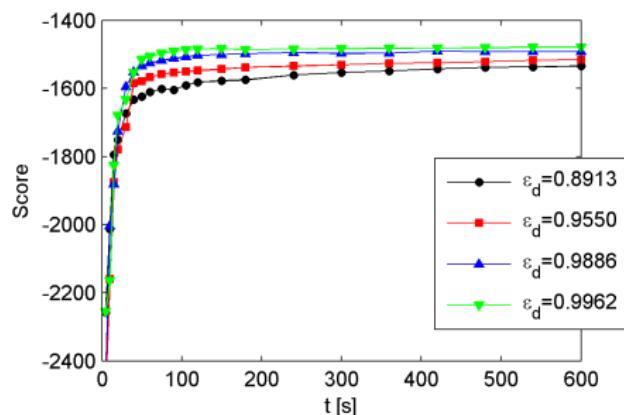
▶ Replay

Pendul inversat: LSPI online, demo

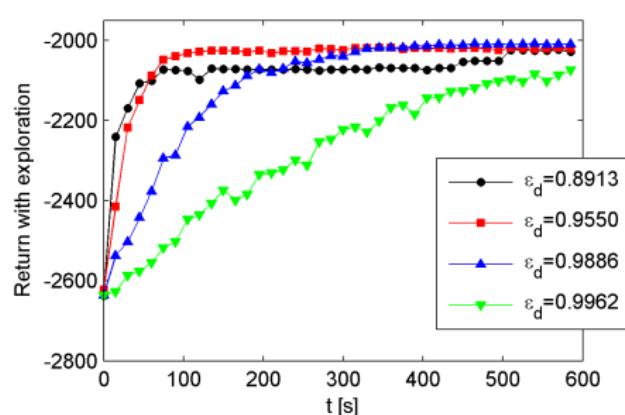


Câtă explorare?

Scorul legii de control învățate
(explorarea oprită)



Performanța **în timpul** învățării
(cu explorare)



⇒ Trebuie găsit un **echilibru**!

Comparație între algoritmii prezenți

Convergență

- Învățarea Q, SARSA, actor-critic:
convergență garantată pentru **variante modificate**
- LSPI online: nu există garanții

Complexitate

- Per iteratie, învățarea Q, SARSA, actor-critic < LSPI online

Reglaj parametri

- Explorarea **crucială** pentru toate metodele
- Ratele de învățare α delicate
(actor-critic are **două** rate de învățare)
- LSPI online: K – mai ușor de acordat

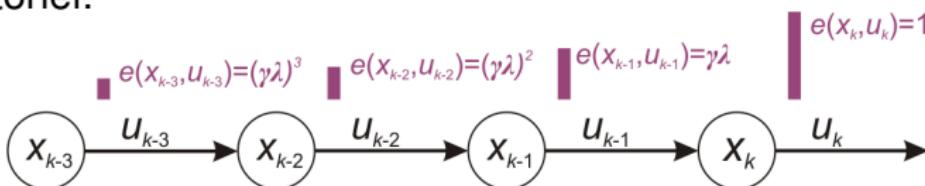
- 1 Învățarea Q și SARSA cu aproximare
- 2 Actor-critic
- 3 LSPI online
- 4 Accelerarea RL cu aproximare
 - Urme de eligibilitate
 - Reluarea experienței
- 5 Perspective

Motivare

- Dezavantaj metode TD approximate:
ca și în cazul discret, învăță încet
 - ⇒ Timp, uzură crescute, profituri scăzute
- **Accelerarea învățării** este necesară

Urme de eligibilitate

- Reamintim cazul discret – urmă $e(x, u)$ de-a lungul traiectoriei:



$$Q(x, u) \leftarrow Q(x, u) + \alpha_k \cdot e(x, u) \cdot$$

$$[r_{k+1} + \gamma \max_{u'} Q(x_{k+1}, u') - Q(x_k, u_k)] \forall x, u$$

- Când x, u continue, $e(x, u)$ nu se poate reprezenta direct

Urme de eligibilitate în cazul aproximat

- Idee: în actualizarea cu gradient, de ex. învățarea Q:

$$\theta_{k+1} = \theta_k + \alpha_k \frac{\partial}{\partial \theta} \hat{Q}(x_k, u_k; \theta_k) \cdot [r_{k+1} + \gamma \max_{u'} \hat{Q}(x_{k+1}, u'; \theta_k) - \hat{Q}(x_k, u_k; \theta_k)]$$

- ...trăiem gradientul $\frac{\partial}{\partial \theta_i} \hat{Q}(x_k, u_k; \theta_k)$ ca pe o **contribuție** a parametrului i la actualizarea curentă
- Luăm în considerare **contribuția cumulativă** (scăzând cu $\gamma \lambda$) până la pasul curent:

$$\begin{aligned} \theta_{k+1} &= \theta_k + \alpha_k e_{k+1} \cdot [r_{k+1} + \gamma \max_{u'} \hat{Q}(x_{k+1}, u'; \theta_k) - \hat{Q}(x_k, u_k; \theta_k)] \\ e_{k+1} &= \sum_{\ell=0}^k (\gamma \lambda)^{k-\ell} \frac{\partial}{\partial \theta} \hat{Q}(x_\ell, u_\ell; \theta_\ell) \end{aligned}$$

Urme de eligibilitate în Învățarea Q aproximată

Implementare iterativă în Învățarea Q:

Învățarea Q(λ) aproximată

for fiecare traiectorie **do**

 initializează x_0 , $e_0 = [0, \dots, 0]^\top$

repeat la fiecare pas k

$$u_k = \begin{cases} \arg \max_u \hat{Q}(x_k, u; \theta_k) & \text{cu prob. } (1 - \varepsilon_k) \\ \text{aleatoare} & \text{cu prob. } \varepsilon_k \end{cases}$$

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

 actualizează $e_{k+1} = (\gamma \lambda) e_k + \frac{\partial}{\partial \theta} \hat{Q}(x_k, u_k; \theta_k)$

$$\theta_{k+1} = \theta_k + \alpha_k e_{k+1} \cdot$$

$$\left[r_{k+1} + \gamma \max_{u'} \hat{Q}(x_{k+1}, u'; \theta_k) - \hat{Q}(x_k, u_k; \theta_k) \right]$$

until traiectoria terminată

end for

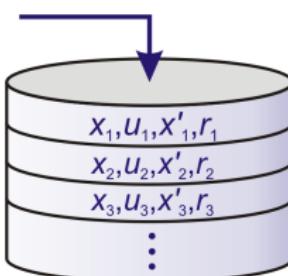
Urme de eligibilitate în alți algoritmi

- Urmele de eligibilitate se pot adapta simplu la SARSA și Învățarea criticului din actor-critic
- Ideea se extinde și la LSPI, dar mai complicat

- 1 Învățarea Q și SARSA cu aproximare
- 2 Actor-critic
- 3 LSPI online
- 4 Accelerarea RL cu aproximare
 - Urme de eligibilitate
 - Reluarea experienței
- 5 Perspective

Reluarea experienței

- Stochează tranzitiiile $(x_k, u_k, x_{k+1}, r_{k+1})$ într-o bază de date



- La fiecare pas, reia n tranzitii din baza de date pe lângă actualizările normale

SARSA aproximată cu reluarea experienței

SARSA aproximată cu reluarea experienței

for fiecare traiectorie **do**

 inițializează x_0

 alege u_0

repeat la fiecare pas k

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

 alege u_{k+1}

$$\theta_{k+1} = \theta_k + \alpha_k \frac{\partial}{\partial \theta} \hat{Q}(x_k, u_k; \theta_k) \cdot$$

$$\left[r_{k+1} + \gamma \hat{Q}(x_{k+1}, u_{k+1}; \theta_k) - \hat{Q}(x_k, u_k; \theta_k) \right]$$

 adaugă $(x_k, u_k, x_{k+1}, r_{k+1})$ la baza de date

 ReiaExperiența

until traiectoria terminată

end for



Procedura ReiaExperiență

ReiaExperiență

loop de N ori

preia o tranziție (x, u, x', r) din baza de date

$$\theta_{k+1} = \theta_k + \alpha_k \frac{\partial}{\partial \theta} \hat{Q}(x_k, u_k; \theta_k) \cdot$$

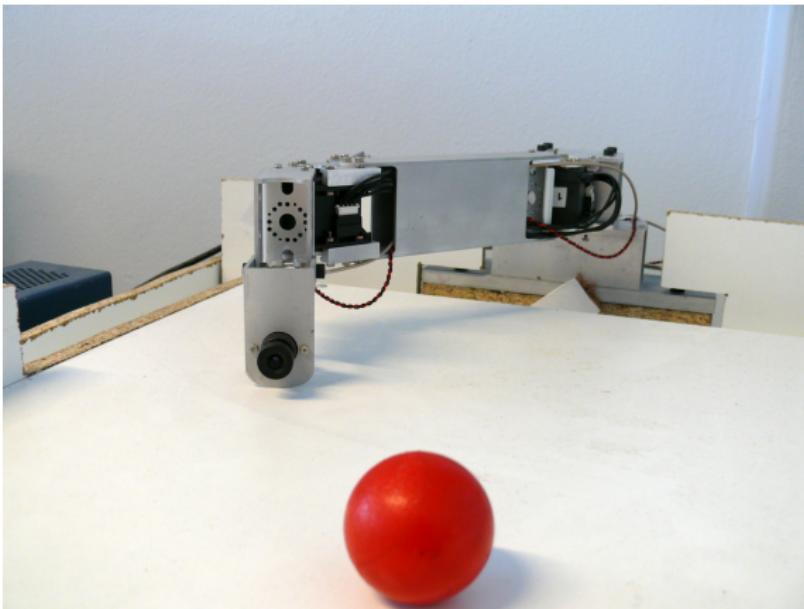
$$\left[r_{k+1} + \gamma \hat{Q}(x_{k+1}, u_{k+1}; \theta_k) - \hat{Q}(x_k, u_k; \theta_k) \right]$$

end loop

Pendul: RL cu reluarea exp., demo (Sander Adam)



Robot portar: RL cu reluarea exp., demo (S. Adam)



- 1 Învățarea Q și SARSA cu aproximare
- 2 Actor-critic
- 3 LSPI online
- 4 Accelerarea RL cu aproximare
- 5 Perspective

Conexiuni: Control optimal

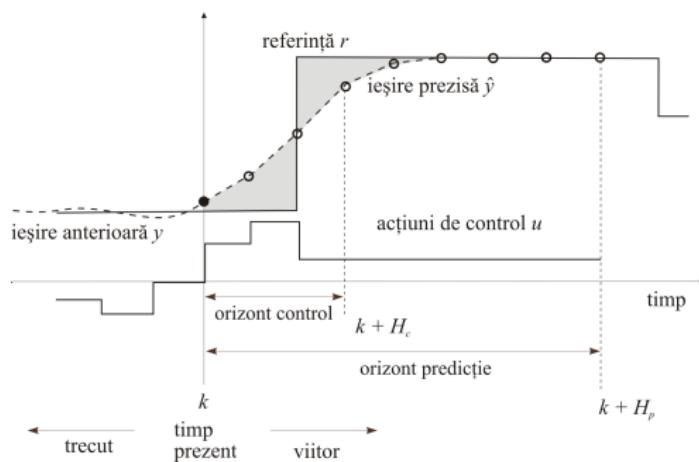
- Controlează un sistem minimizând costul J
- Bazat pe model
- Posibil în timp continuu, orizont finit

RL, DP **sunt control optimal!**

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

Conexiuni: Control predictiv

- Bazat pe model, clasic liniar
- Principiu de bază: **receding horizon**



Controlul predictiv este și el o variantă de control optimal, bazată pe model

Procese de fabricație

Focus curs: Controlul sistemelor generice

Procese de fabricație

- Legate de logistică și *operations research*
- Variabile **discrete** sunt importante posibil combinate cu variabile continue: hibrid
Aproximarea rămâne esențială
- Tranzitii de stare cu **durată variabilă**: procese de decizie semi-Markov

DP & RL din perspectiva operations research:

Warren B. Powell, *Approximate Dynamic Programming: Solving the curses of dimensionality*, ed. 2, Wiley, 2011.



Probleme deschise

RL & DP **în curs de dezvoltare**

Probleme deschise:

- Garanții de siguranță și stabilitate
- Stări și acțiuni cu dimensionalitate mare
- Stări care nu pot fi măsurate
- Strategii de explorare
- Sisteme multiagent