

Control prin învățare

Master ICAF, An 1 Sem 2

Lucian Bușoniu, Ștefan Pîrje



Partea III

Învățarea prin recompensă



Partea III în plan

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



Gama 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)

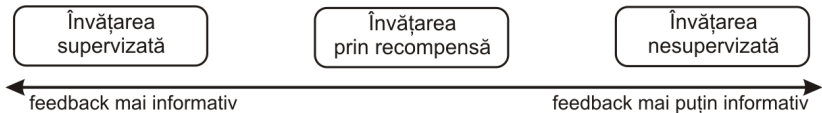


Conținut partea III

- 1 Monte Carlo, MC
- 2 Nevoia de explorare
- 3 Diferențe temporale, TD
- 4 Accelerarea metodelor TD
- 5 Încheiere



Învățarea prin recompensă în gama învățării automate



- Supervizată: pentru fiecare eșantion de antrenare, **ieșirea corectă** este cunoscută
- Nesupervizată: doar intrări, **fără ieșiri**; găsește structuri în date
- Prin recompensă: acțiunile corecte nu sunt disponibile, **doar recompense**

Important: învățarea prin recompensă găsește un **control optimal pentru un sistem dinamic!**



- 1 Monte Carlo, MC
- 2 Nevoia de explorare
- 3 Diferențe temporale, TD
- 4 Accelerarea metodelor TD
- 5 Încheiere

Reamintim: Iterația pe legea de control

Iterația pe legea de control cu funcții Q

inițializează legea de control h_0 arbitrar

repeat la fiecare iterație ℓ

1: **evaluarea legii de control:** găsește Q^{h_ℓ}

2: **îmbunătățirea legii de control:**

$$h_{\ell+1}(x) \leftarrow \arg \max_u Q^{h_\ell}(x, u)$$

until convergență la h^*



Evaluarea legii de control

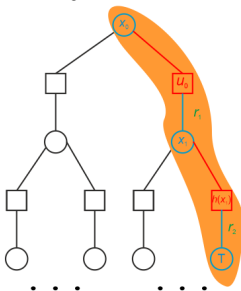
Pentru a găsi Q^h :

- Până acum: metode bazate pe model
- Învățare prin recompensă: **modelul nu este disponibil**
- Învăță Q^h din date sau prin **interacțiune online cu sistemul**



Evaluarea Monte-Carlo a legii de control

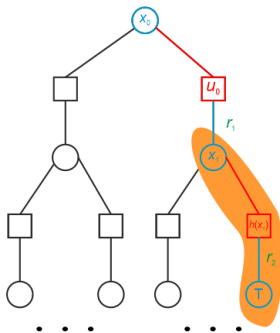
Reamintim: $Q^h(x_0, u_0) = \sum_{k=0}^{\infty} \gamma^k r_{k+1}$



- Traiectorie de la (x_0, u_0) la x_K (**terminal**) utilizând $u_1 = h(x_1)$, $u_2 = h(x_2)$ etc.
- ⇒ $Q^h(x_0, u_0) =$ returnul de-a lungul traiectoriei:

$$Q^h(x_0, u_0) = \sum_{j=0}^{K-1} \gamma^j r_{j+1}$$

Evaluarea Monte-Carlo a legii de control (cont.)



- În plus, la fiecare pas k :

$$Q^h(x_k, u_k) = \sum_{j=k}^{K-1} \gamma^{j-k} r_{j+1}$$

Iterația Monte Carlo pe legea de control

Iterația Monte Carlo pe legea de control

for fiecare iterație ℓ **do**

efectuează N traiectorii aplicând h_ℓ

resetează la 0 acumulator $A(x, u)$, counter $C(x, u)$

for fiecare pas k din fiecare traiectorie i **do**

$$A(x_k, u_k) \leftarrow A(x_k, u_k) + \sum_{j=k}^{K_i-1} \gamma^{j-k} r_{i,j+1} \text{ (return)}$$

$$C(x_k, u_k) \leftarrow C(x_k, u_k) + 1$$

end for

$$Q^{h_\ell}(x, u) \leftarrow A(x, u) / C(x, u)$$

$$h_{\ell+1}(x) \leftarrow \arg \max_u Q^{h_\ell}(x, u)$$

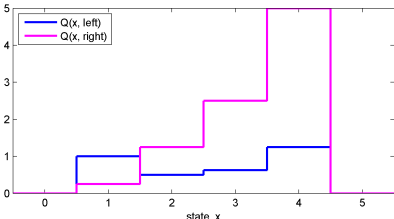
end for

De notat: **atingerea stării terminale** trebuie garantată!



Robot menajer: Monte Carlo, demo

Monte Carlo, trial 70 [piter 7 done, peval 10]



Nevoia de explorare

$$Q^h(x, u) \leftarrow A(x, u) / \mathbf{C(x, u)}$$

Cum asigurăm $C(x, u) > 0$ – **informație** despre fiecare (x, u) ?

- 1 **Stări inițiale** x_0 selectate reprezentativ
- 2 **Acțiuni:**
 u_0 reprezentative, câteodată diferite de $h(x_0)$
și în plus, posibil:
 u_k reprezentative, câteodată diferite de $h(x_k)$



Dilema explorare-exploatare

- **Explorarea** necesară:
acțiuni diferite de legea curentă de control
- **Exploatarea** cunoștințelor curente necesară:
legea de control trebuie aplicată

Dilema explorare-exploatare

– esențială în toți algoritmi de RL

(nu doar în MC)

Explorare-exploatare: strategia ϵ -greedy

- O soluție simplă: **ϵ -greedy**

$$u_k = \begin{cases} h(x_k) = \arg \max_u Q(x_k, u) & \text{cu probabilitatea } (1 - \epsilon_k) \\ \text{o acțiune aleatoare} & \text{cu probabilitatea } \epsilon_k \end{cases}$$

- Probabilitatea de explorare $\epsilon_k \in (0, 1)$ scade de obicei în timp



Strategia softmax

- Selecția acțiunii:

$$u_k = u \text{ cu proba. } \frac{e^{Q(x_k, u)/\tau_k}}{\sum_{u'} e^{Q(x_k, u')/\tau_k}}$$

unde $\tau_k > 0$ este **temperatura de explorare**

- Pentru $\tau \rightarrow 0$, recuperăm selecția greedy;
 $\tau \rightarrow \infty$ duce la distribuție aleatoare uniformă
- Comparativ cu ε -greedy, acțiunile mai bune au o probabilitate mai mare de a fi aplicate chiar și în explorare



Explorare bazată pe bandiți



Într-o stare unică, explorarea este modelată ca **bandit cu mai multe brațe**:

- Acțiunea j = braț cu distribuție a recompensei ρ_j , valoare așteptată μ_j
- Cel mai bun braț (acțiunea optimală) are valoarea așteptată μ^*
- La pasul k , tragem brațul (încercăm acțiunea) j_k , obținând $r_k \sim \rho_{j_k}$
- **Obiectiv:** După n încercări, regret minim: $\sum_{k=1}^n \mu^* - \mu_{j_k}$

Algoritmul UCB

Un algoritm des utilizat: după n pași, alegem brațul cu cea mai mare **margine superioară de încredere** (en. upper confidence bound):

$$b(j) = \hat{\mu}_j + \sqrt{c \frac{\log n}{n_j}}$$

unde:

- $\hat{\mu}_j$ = media recompenselor observate pentru brațul j până în prezent
- n_j = de câte ori a fost tras brațul j
- c constantă acordabilă, de exemplu $3/2$

Acestea sunt doar câteva metode simple, există multe altele, de exemplu explorare bayesiană, recompense intrinseci, inițializare optimistă etc.



- 1 Monte Carlo, MC
- 2 Nevoia de explorare
- 3 Diferențe temporale, TD**
 - Introducere
 - SARSA
 - Învățarea Q
- 4 Accelerarea metodelor TD
- 5 Încheiere



Diferențe temporale (TD)

În rest, actualizarea rămâne aceeași, dar explicităm returnul estimat:

$$\begin{aligned} Q(x_k, u_k) &\leftarrow Q(x_k, u_k) + \alpha_k [\hat{R}_k - Q(x_k, u_k)] \\ &= Q(x_k, u_k) + \alpha_k [r_{k+1} + \gamma Q(x_{k+1}, h(x_{k+1})) - Q(x_k, u_k)] \end{aligned}$$

- [...] este **diferența temporală** între două estimări ale $Q(x_k, u_k)$, utilizând informații la pași de timp consecutivi
- Actualizări bazate pe date, fără model (ca în MC):
 r_{k+1}, x_{k+1} observate de ex. în timpul interacțiunii online
- Actualizările estimează $Q(x_k, u_k)$ utilizând o altă estimare, $Q(x_{k+1}, h(x_{k+1}))$): **bootstrapping**
- Programarea dinamică face, de asemenea, bootstrapping, dar utilizând un model



Diferențe temporale pentru evaluarea h

Diferențe temporale pentru evaluarea h

for fiecare traiectorie **do**

inițializează x_0 , alege acțiunea inițială u_0

repeat la fiecare pas k

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

alege acțiunea **următoare** $u_{k+1} \sim h(x_{k+1})$

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

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

until traiectoria terminată

end for

Observație: am înlocuit $h(x_k + 1)$ cu u_{k+1} , aleasă cu h

Explorare-exploatare

alege acțiunea următoare $u_{k+1} \sim h(x_{k+1})$

- Informații despre $(x, u) \neq (x, h(x))$ necesare
⇒ **explorare**
- h trebuie urmărită
⇒ **exploatare**
- Ex. ε -greedy:

$$u_{k+1} = \begin{cases} h(x_{k+1}) & \text{cu prob. } (1 - \varepsilon_{k+1}) \\ \text{aleatoare} & \text{cu prob. } \varepsilon_{k+1} \end{cases}$$

- 1 Monte Carlo, MC
- 2 Nevoia de explorare
- 3 Diferențe temporale, TD**
 - Introducere
 - **SARSA**
 - Învățarea Q
- 4 Accelerarea metodelor TD
- 5 Încheiere



Îmbunătățirea optimistă a legii de control

- Legea de control neschimbată pentru N traiectorii
- ⇒ Algoritmul învață încet
- Îmbunătățirea legii de control după fiecare traiectorie
= **optimist**
- Vom folosi de asemenea și explorare ϵ -greedy



Metoda Monte Carlo optimistă

Metoda Monte Carlo optimistă

inițializează la 0 acumulator $A(x, u)$, counter $C(x, u)$

for fiecare traiectorie **do**

efectuează traiectoria, ex. aplicând ε -greedy:

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

for fiecare pas k al fiecărei traiectorii **do**

$$A(x_k, u_k) \leftarrow A(x_k, u_k) + \sum_{j=k}^{K-1} \gamma^{j-k} r_{j+1}$$

$$C(x_k, u_k) \leftarrow C(x_k, u_k) + 1$$

end for

$$Q(x, u) \leftarrow A(x, u) / C(x, u)$$

end for

- h implicit, greedy în Q
- actualizarea $Q \Rightarrow$ implicit îmbunătățirea h



Optimism în TD

- Algoritm precedent: h fixată
- Îmbunătățirea h : cel mai simplu, după fiecare tranziție
⇒ interpretare: iterație pe legea de control
optimistă la nivel de tranziție
- h implicit, greedy în Q
(actualizarea $Q \Rightarrow$ implicit îmbunătățirea h)



SARSA

SARSA cu ϵ -greedy

for fiecare traiectorie **do**

inițializează x_0

$$u_0 = \begin{cases} \arg \max_u Q(x_0, u) & \text{cu prob. } (1 - \epsilon_0) \\ \text{aleatoare} & \text{cu prob. } \epsilon_0 \end{cases}$$

repeat la fiecare pas k

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

$$u_{k+1} = \begin{cases} \arg \max_u Q(x_{k+1}, u) & \text{cu prob. } (1 - \epsilon_{k+1}) \\ \text{aleatoare} & \text{cu prob. } \epsilon_{k+1} \end{cases}$$

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

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

until traiectoria terminată

end for



Originile numelui SARSA

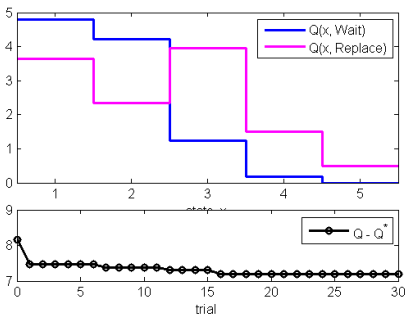
$(X_k, U_k, r_{k+1}, X_{k+1}, U_{k+1}) =$
(Stare, Acțiune, Recompensă, Stare, Acțiune) = SARSA



Înlocuirea unei mașini: SARSA, demo

Parametri: $\alpha = 0.1$, $\varepsilon = 0.3$ (constanți), 20 pași pe traiectorie
 $x_0 = 1$

SARSA, trial 30 completed



- 1 Monte Carlo, MC
- 2 Nevoia de explorare
- 3 Diferențe temporale, TD**
 - Introducere
 - SARSA
 - Învățarea Q**
- 4 Accelerarea metodelor TD
- 5 Încheiere



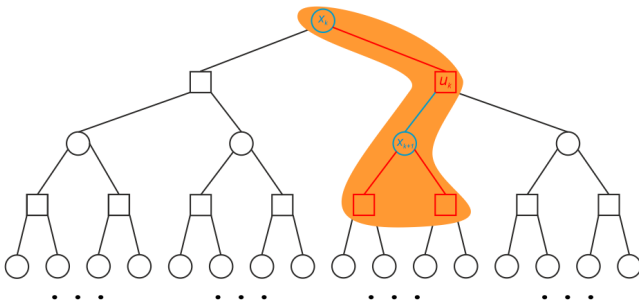
Estimarea Q^* prin bootstrapping

Ecuția de optimalitate Bellman:

$$Q^*(x, u) = E_{x'} \left\{ \tilde{\rho}(x, u, x') + \gamma \max_{u'} Q^*(x', u') \right\}$$

conduce la estimatul:

$$\hat{Q}_k = r_{k+1} + \gamma \max_{u'} Q(x_{k+1}, u')$$



Actualizare TD pentru Q^*

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



Învățarea Q

Învățarea Q cu ϵ -greedy

for fiecare traiectorie **do**

inițializează x_0

repeat la fiecare pas k

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

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

On-policy / off-policy

SARSA: **on-policy**

- Estimează permanent funcția Q a legii de control curente

Învățarea Q: **off-policy**

- Indiferent de legea de control curentă, estimează funcția Q optimală



Diferențe temporale: Discuție

Avantaje

- Simplu de înțeles, implementat
- Complexitate scăzută \Rightarrow execuție rapidă

SARSA vs. învățarea Q

- SARSA mai puțin complex decât învățarea Q (fără max în actualizarea funcției Q)

Secvențele α_k, ε_k **influențează semnificativ** performanța

Dezavantaj principal

- Necesită multe date



- 1 Monte Carlo, MC
- 2 Nevoia de explorare
- 3 Diferențe temporale, TD
- 4 Accelerarea metodelor TD**
 - Motivare
 - Reluarea experienței
 - Return pe n pași
- 5 Încheiere



Nevoia de a accelera metodele TD

Dezavantaj principal: învață încet – **necesită multe date**

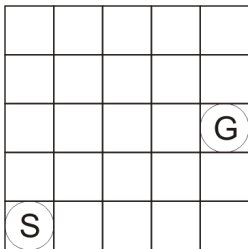
În practică, datele costă:

- timp
- profit (performanță scăzută datorită explorării)
- uzură a sistemului

Accelerarea RL = **eficientizarea folosirii datelor**

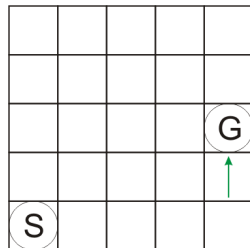
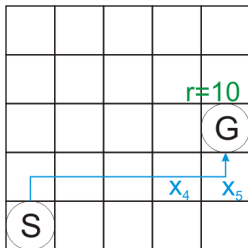


Exemplu: Navigare 2D



- Navigare într-o lume 2D discretă de la **Start** la **Goal**
- Singura recompensă = 10 la atingerea G (stare terminală)

Exemplu: TD



- Alegem SARSA, $\alpha = 1$; inițializăm $Q = 0$
- Actualizări de-a lungul traiectoriei din stânga:

...

$$Q(x_4, u_4) = 0 + \gamma \cdot Q(x_5, u_5) = 0$$

$$Q(x_5, u_5) = 10 + \gamma \cdot 0 = 10$$

- O nouă tranziție de la x_4 la x_5 necesară pentru a propaga informația la x_4 !

Accelerarea TD: 2 idei

- 1 Stochează și **reia experiența**
- 2 Return pe n pași

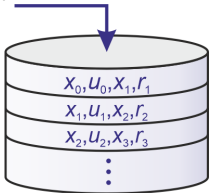


- 1 Monte Carlo, MC
- 2 Nevoia de explorare
- 3 Diferențe temporale, TD
- 4 Accelerarea metodelor TD**
 - Motivare
 - Reluarea experienței**
 - Return pe n pași
- 5 Încheiere



Reluarea experienței

- Stochează fiecare tranziție $(x_k, u_k, x_{k+1}, r_{k+1})$ (pentru SARSA, și u_{k+1}) într-o bază de date



- La fiecare pas, **reia** n tranziții din baza de date (pe lângă actualizările normale)

Învățarea Q cu reluarea experienței

Învățarea Q cu reluarea experienței

for fiecare traiectorie **do**

inițializează x_0

repeat la fiecare pas k

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

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

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

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

ReiaExperiența

until traiectoria terminată

end for



Procedura ReiaExperiența

ReiaExperiența

loop de n ori

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

$$Q(x, u) \leftarrow Q(x, u) + \alpha \cdot$$

$$[r + \gamma \max_{u'} Q(x', u') - Q(x, u)]$$

end loop

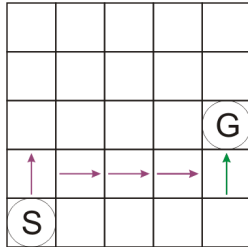
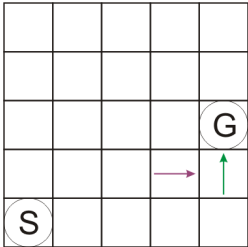
Direcția de reluare

Ordinea de reluare a tranzițiilor:

- 1 Înainte
- 2 Înapoi
- 3 Arbitrar



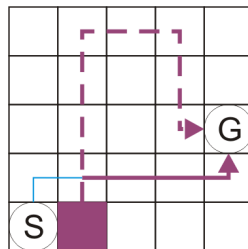
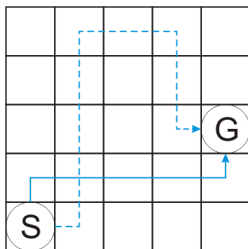
Exemplu: Influența direcției de reluare



- Verde: **actualizările normale**, mov: **reluarea experienței**
- Stânga: reluare înainte; dreapta: reluare înapoi
- **Reluarea înapoi** în general preferabilă



Exemplu: Agregarea informației

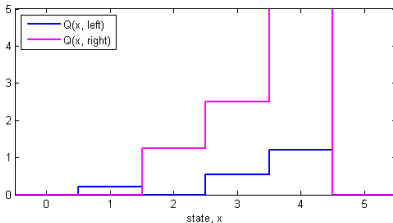


- Reluarea experienței **agregă informație** din mai multe traiectorii
- Căsuța indicată profită de informații de-a lungul ambelor traiectorii

Robot menajer: Q cu reluarea experienței, demo

Parametri: $\alpha = 0.2$, $\epsilon = 0.3$, $n = 5$, direcția înapoi
 $x_0 = 2$ or 3 (aleator)

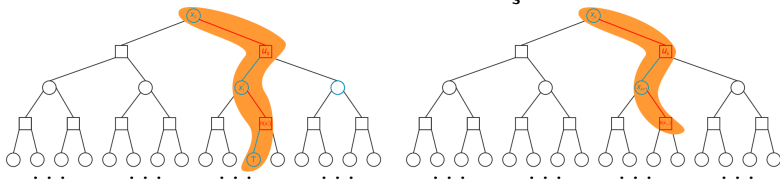
ER-Q-learning, trial 13, step 2 [replaying trial 8, step 2]



- 1 Monte Carlo, MC
- 2 Nevoia de explorare
- 3 Diferențe temporale, TD
- 4 Accelerarea metodelor TD**
 - Motivare
 - Reluarea experienței
 - Return pe n pași**
- 5 Încheiere



Reamintire: estimările returnului în MC și TD

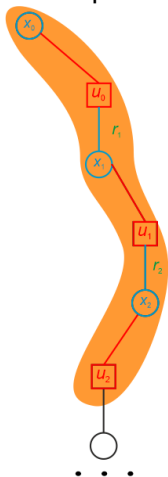


$$R_k = \sum_{j=k}^{K-1} \gamma^{j-k} r_{j+1}$$

$$\hat{R}_k = r_{k+1} + \gamma Q(x_{k+1}, h(x_{k+1}))$$

Există o variantă intermediară?

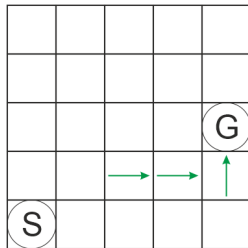
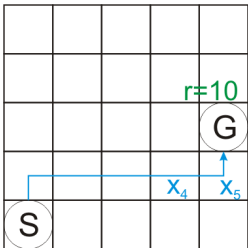
Compromis: return pe n pași



SARSA (on-policy):

$$\hat{R}_k = r_{k+1} + \gamma r_{k+2} + \dots + \gamma^{n-1} r_{k+n} + \gamma^n Q(x_{k+n}, u_{k+n})$$

Exemplu: Efectul returnului pe n pași



Pentru $n = 3$:

$$Q(x_5, u_5) = 10 + 0 \quad (\text{terminal})$$

$$Q(x_4, u_4) = 0 + \gamma 10 + 0 \quad (\text{terminal})$$

$$Q(x_3, u_3) = 0 + \gamma 0 + \gamma^2 10 + 0 \quad (\text{terminal})$$

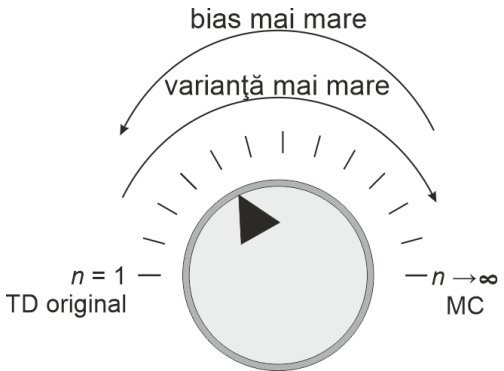
$$Q(x_2, u_2) = 0 + \gamma 0 + \gamma^2 0 + \gamma^3 0 \quad (\text{bootstrap})$$

...



TD versus MC

- $n = 1$ recuperează TD, $n \rightarrow \infty$ recuperează MC
- Valorile intermediare combină TD și MC, conducând la un compromis reglabil între bias și varianță



- 1 Monte Carlo, MC
- 2 Nevoia de explorare
- 3 Diferențe temporale, TD
- 4 Accelerarea metodelor TD
- 5 Încheiere**



Recapitulare: Metode în Partea III

Metode Monte Carlo, MC:

- Iterația Monte Carlo pe legea de control
- MC cu actualizări incrementale

Dilema explorare-exploatare:

- ϵ -greedy, folosit uzual
- Există multe alte soluții, de ex. UCB

Diferențe temporale, TD:

- TD pentru evaluarea legii de control
- Îmbunătățirea optimistă a legii de control
- SARSA
- Învățarea Q

Accelerarea TD:

- Reluarea experienței
- Return pe n pași



Terminologie engleză

metode Monte Carlo	=	<u>Monte Carlo methods</u>
dilema explorare-exploatare	=	<u>exploration-exploitation dilemma</u>
diferențe temporale	=	<u>temporal differences, TD</u>
învățarea Q	=	<u>Q-learning</u>
SARSA	=	<u>SARSA</u>
rata de învățare	=	<u>learning rate</u>
reluarea experienței	=	<u>experience replay</u>
return pe n pași	=	<u>n-step return</u>



Exerciții

- 1 Descreșterea exponentială a ratei de învățare, $\alpha_k = \alpha^k$, cu $\alpha \in (0, 1)$ o constantă, satisface condițiile de aproximare stohastică?
- 2 Este garantată convergența învățării Q când $\epsilon_k = \epsilon$, o constantă în $(0, 1)$? Dar pentru SARSA? Ce se întâmplă dacă folosim o descreștere exponentială $\epsilon_k = \epsilon^k$?
- 3 Ar avea sens un algoritm Monte-Carlo care îmbunătățește legea de control după fiecare tranziție (ca în TD)?
- 4 Ar propaga învățarea Q (fără return pe n pași, deoarece este mai complicat în cazul off-policy) informația mai rapid decât SARSA în exemplul traiectoriei pe grid?



Exerciții (cont.)

- 5 Presupunând că avem acces la un model **doar pentru îmbunătățirea legii de control**, furnizați alternative bazate pe funcția V pentru toți algoritmi din această parte a cursului. Urmați aceeași ordine ca pentru funcțiile Q :
- Estimări Monte-Carlo, bazate pe mediere și incrementale
 - Estimări și actualizări cu bootstrap
 - Evaluarea legii de control, SARSA și învățarea Q

Nu uitați să desenați arbori și să încercați nodurile relevante, acest lucru vă va ajuta să vizualizați metodele.

