

# Control prin învățare

Master ICAF, An 1 Sem 2

Lucian Bușoniu, Ștefan Pîrje



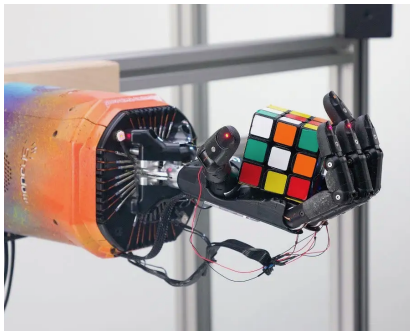
## Partea V

# Învățarea prin recompensă online cu aproximare



## Recap: Nevoia de aproximare

- În aplicații reale de control,  $x$ ,  $u$  **continue!** (sau discrete cu multe valori)



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

## Recap: Partea 4 – Algoritmi offline

pornind de la:

– model  $f, \rho$

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

- 1 găsește o soluție aproximată  $\hat{Q}(x, u), \hat{h}(x)$ , etc.
- 2 controlează sistemul folosind soluția găsită

Algoritmi exemplificați:

- iterația Q cu interpolare
- iterația Q bazată pe date



## Partea V 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 V în plan: Categoriile de algoritmi

După utilizarea unui model:

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

După nivelul de interacțiune:

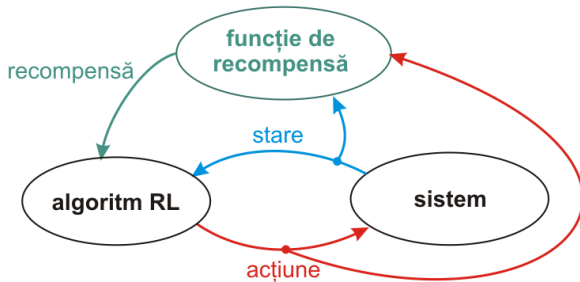
- **Offline:** algoritmul rulează în avans
- **Online:** algoritmul controlează direct sistemul

Exact vs. cu aproximare:

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



# Principiul RL



- Urmărim acum cu adevărat procedura RL de interacțiune online
- Există mulți algoritmi; discutăm câțiva

# Conținut partea 5

- 1 TD cu aproximare
- 2 RL prin gradientul legii de control (policy gradient)



- 1 TD cu aproximare
  - SARSA aproximată
  - Învățarea Q aproximată
  - Maximizare și discuție
- 2 RL prin gradientul legii de control (policy gradient)



- 1 TD cu aproximare
  - SARSA aproximată
  - Învățarea Q aproximată
  - Maximizare și discuție
- 2 RL prin gradientul legii de control (policy gradient)



# Reamintim SARSA clasică

## SARSA cu $\epsilon$ -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**



# Reamintim: Derivarea actualizării SARSA (on-policy)

- Actualizare:

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

- $\hat{R}_k$  este un estimat cu bootstrap (folosind ecuația Bellman) al return-ului Monte Carlo  $R_k$  din  $(x_k, u_k)$  pentru legea curentă de control  $h$
- $R_k$  este el însuși un eșantion al  $Q^h(x_k, u_k)$ , deci în final rulăm o versiune estimată a actualizării ideale:

$$Q(x_k, u_k) \leftarrow Q(x_k, u_k) + \alpha_k [Q^h(x_k, u_k) - Q(x_k, u_k)]$$



# SGD pentru cazul aproximat

- Extindem ideea de mai sus la aproximatorul parametric  $\hat{Q}(x, u; \theta)$
- Nu mai putem actualiza direct  $Q$ , în schimb actualizăm  $\theta$  folosind **coborârea pe gradient stohastică (SGD)** pe eroarea de aproximare pătratică:

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

- Înlocuim  $Q^h(x_k, u_k)$  cu eşantionul Monte Carlo  $R_k$ , apoi  $R_k$  cu estimatul cu bootstrap  $\hat{R}_k = r_{k+1} + \gamma \hat{Q}(x_{k+1}, u_{k+1}; \theta_k)$ :

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



# Semigradien

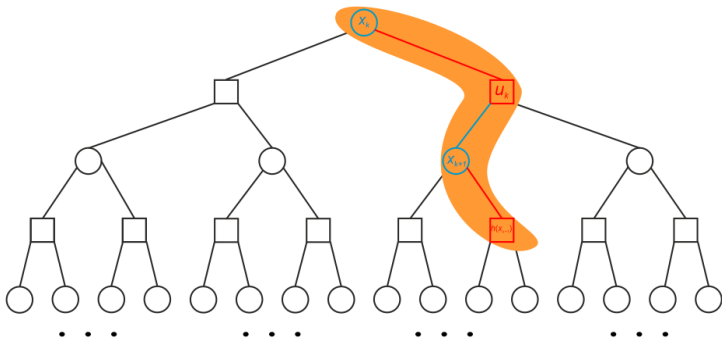
$$\theta_{k+1} = \theta_k + \alpha_k \nabla_{\theta} \widehat{Q}(x_k, u_k; \theta_k) \left[ r_{k+1} + \gamma \widehat{Q}(x_{k+1}, u_{k+1}; \theta_k) - \widehat{Q}(x_k, u_k; \theta_k) \right]$$

- Actualizarea finală nu este un gradient complet, deoarece  $\widehat{R}_k$  depinde de  $\theta_k$  prin  $\widehat{Q}(x_{k+1}, u_{k+1}; \theta_k)$ , dar doar al doilea termen al erorii este diferențiat!
- ⇒ Astfel de metode se numesc de **semigradien**.
- Termenul din paranteza pătrată este o **diferență temporală aproximată**.



# Ilustrație

Ilustrația grafică este similară cazului clasic:



dar acum prezența aproximării necesită utilizarea gradientului

## Obiectiv și comparație cu metodele fitted

$$\theta_{k+1} \leftarrow \theta_k + \alpha_k \nabla_{\theta} \widehat{Q}(x_k, u_k; \theta_k) \left[ Q^h(x_k, u_k) - \widehat{Q}(x_k, u_k; \theta_k) \right]$$

minimizează următorul obiectiv sub distribuția generată de eșantioanele  $(x_k, u_k)$ :

$$\mathcal{L}_{\text{eval}} = \mathbb{E} \left\{ \left[ Q^h(x, u) - \widehat{Q}(x, u; \theta) \right]^2 \right\}$$

- Comparăm cu obiectivul “fitted”, unde distribuția este implicit tot cea a eșantioanelor:  $\sum_{s=1}^{n_s} \left[ \widehat{R}_s - \widehat{Q}(x_s, u_s; \theta) \right]^2$
- Dacă  $\sum_{k=0}^{\infty} \alpha_k^2$  este finită și  $\sum_{k=0}^{\infty} \alpha_k \rightarrow \infty$ , SGD converge către un minim local al  $\mathcal{L}_{\text{eval}}$ .
- Acest lucru rămâne valabil când  $Q^h(x_k, u_k)$  este înlocuit cu eșantionul Monte Carlo  $R_k$ , dar nu și pentru estimatul cu bootstrap  $\widehat{R}_k$ , din cauza actualizărilor semigradient!



# SARSA aproximată semigradient

## SARSA aproximată

**for** fiecare traiectorie **do**

inițializează  $x_0$

alege  $u_0$  (ex.  $\varepsilon$ -greedy din  $\widehat{Q}(x_0, \cdot; \theta_0)$ )

**repeat** la fiecare pas  $k$

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

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

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

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

**until** traiectoria terminată

**end for**

**Explorarea** este desigur necesară și în cazul aproximat.



## Actualizarea semigradient pentru învățarea Q

Reamintim actualizare din învățarea Q clasică:

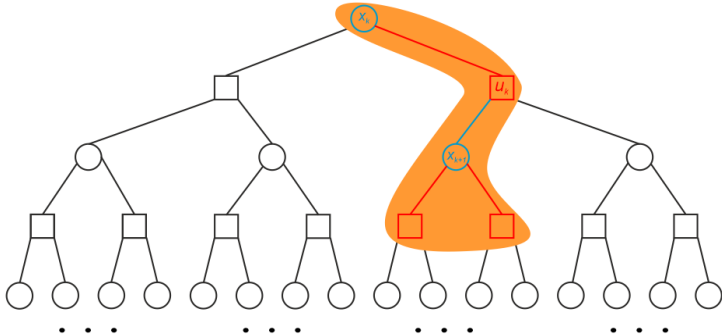
$$\begin{aligned} Q(x_k, u_k) &\leftarrow Q(x_k, u_k) + \alpha_k [r_{k+1} + \gamma \max_{u'} Q(x_{k+1}, u') - Q(x_k, u_k)] \\ &\approx Q(x_k, u_k) + \alpha_k [Q^*(x_k, u_k) - Q(x_k, u_k)] \end{aligned}$$

În cazul aproximat, folosim gradientul:

$$\begin{aligned} \theta_{k+1} &= \theta_k - \frac{1}{2} \alpha_k \nabla_{\theta} \left[ Q^*(x_k, u_k) - \hat{Q}(x_k, u_k; \theta_k) \right]^2 \\ &= \theta_k + \alpha_k \nabla_{\theta} \hat{Q}(x_k, u_k; \theta_k) \left[ Q^*(x_k, u_k) - \hat{Q}(x_k, u_k; \theta_k) \right] \\ &\approx \theta_k + \alpha_k \nabla_{\theta} \hat{Q}(x_k, u_k; \theta_k) \cdot \\ &\quad \left[ r_{k+1} + \gamma \max_{u'} \hat{Q}(x_{k+1}, u'; \theta_k) - \hat{Q}(x_k, u_k; \theta_k) \right] \end{aligned}$$



# Ilustrație



# Învățarea Q aproximată semigradient: algoritm

## Învățarea Q aproximată cu explorare $\varepsilon$ -greedy

**for** fiecare traiectorie **do**

inițializează  $x_0$

**repeat** la fiecare pas  $k$

alege  $u_k$  (ex.  $\varepsilon$ -greedy din  $\widehat{Q}(x_k, \cdot; \theta_k)$ )

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

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

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

**until** traiectoria terminată

**end for**



- 1 TD cu aproximare
  - SARSA aproximată
  - Învățarea Q aproximată
  - Maximizare și discuție
- 2 RL prin gradientul legii de control (policy gradient)



# 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 explicit



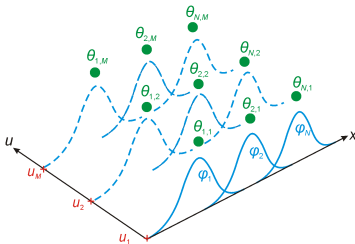
# Maximizare în învățarea Q aproximată

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

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

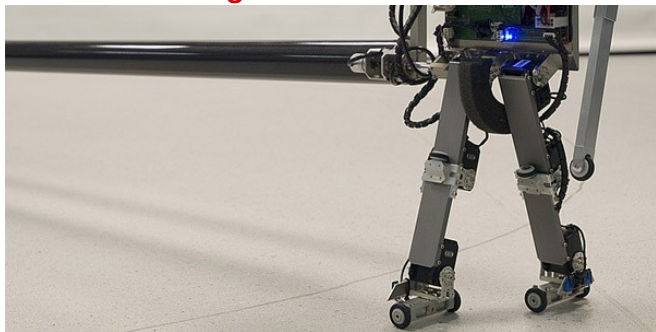
⇒ Soluția 1: Legea de control reprezentată implicit

- 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**



## Discuție despre metodele TD aproximare

- Complexitate redusă
- Convergența este garantată pentru **versiuni modificate**
- Ratele de explorare și de învățare trebuie să fie **reglate cu grijă** pentru toate metodele
- La fel ca în cazul clasic, metodele TD aproximare învață lent, deci trebuie **accelerate**
- Reluarea experienței și return-urile pe  $n$  pași sunt aproape direct aplicabile (ultimele pentru SARSA, dar pot fi extinse la Q-learning off-policy)



- 1 TD cu aproximare
- 2 RL prin gradientul legii de control (policy gradient)



# Reprezentarea legii de control

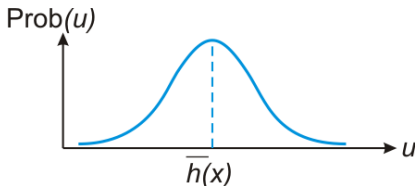
- Tipul 2: Legea de control **aproximată explicit**
- Reamintim avantajele: mai ușor de tratat acțiuni continue și de integrat cunoștințe a priori
- De exemplu, lege de control reprezentată folosind atribute (funcții de bază):

$$\bar{h}(\mathbf{x}; \mu) = \sum_{i=1}^n \phi_i(\mathbf{x}) \mu_i$$



# Lege de control cu explorare

- RL online  $\Rightarrow$  policy gradient trebuie să exploreze



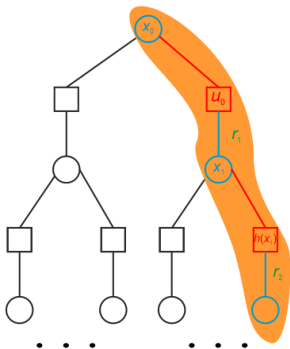
- Explorarea gaussiană** aplică  $u$  în  $x$  cu probabilitatea:

$$P(u|x) = \mathcal{N}(\bar{h}(x; \mu), \sigma) =: \hat{h}(x, u; \vartheta)$$

unde  $\vartheta$  conține pe lângă  $\mu$  și covarianțele în matricea  $\sigma$

- Astfel, reprezentăm o lege de control stohastică, incluzând direct explorarea în parametri

# Traectorie



- Traectoria  $\tau := (x_0, u_0, \dots, x_k, u_k, \dots)$  generată cu  $\hat{h}$ ; recompensele rezultate  $r_1, \dots, r_{k-1}, \dots$
- Considerăm un MDP determinist pentru simplitate. Return-ul de-a lungul traiectoriei:

$$R(\tau) = \sum_{k=0}^{\infty} \gamma^k r_{k+1} = \sum_{k=0}^{\infty} \gamma^k \rho(x_k, u_k)$$

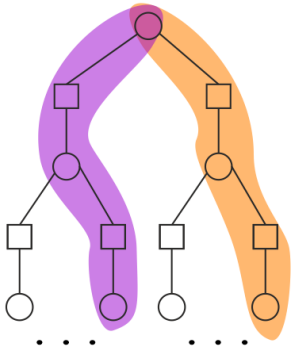
- Probabilitatea traiectoriei  $\tau$  dați fiind parametrii legii de control  $\vartheta$ :

$$P_{\vartheta}(\tau) = \prod_{k=0}^{\infty} \hat{h}(x_k, u_k; \vartheta)$$

unde  $x_{k+1} = f(x_k, u_k)$



# Obiectiv



Considerăm  $x_0$  fixat, pentru simplitate

**Obiectiv**  
**Maximizarea return-ului așteptat** din  $x_0$  al  
 legii de control  $\hat{h}(\cdot, \cdot; \vartheta)$ :

$$J_\vartheta := E_\vartheta \{R(\tau)\} = \int R(\tau)P_\vartheta(\tau)d\tau$$

# Ideea principală

**Ascensiune pe gradient** pentru maximizarea

$J_{\vartheta}$ :

$$\vartheta \leftarrow \vartheta + \alpha \nabla_{\vartheta} J_{\vartheta}$$



# Derivarea gradientului

Dorim să ușurăm implementarea actualizării folosind date:

$$\begin{aligned} \nabla_{\vartheta} J_{\vartheta} &= \int R(\tau) \nabla_{\vartheta} P_{\vartheta}(\tau) d\tau \\ &= \int R(\tau) P_{\vartheta}(\tau) \nabla_{\vartheta} \log P_{\vartheta}(\tau) d\tau \\ &= \mathbb{E}_{\vartheta} \left\{ R(\tau) \nabla_{\vartheta} \log \left[ \prod_{k=0}^{\infty} \hat{h}(x_k, u_k; \vartheta) \right] \right\} \\ &= \mathbb{E}_{\vartheta} \left\{ R(\tau) \sum_{k=0}^{\infty} \nabla_{\vartheta} \log \hat{h}(x_k, u_k; \vartheta) \right\} \end{aligned}$$

Unde:

- am folosit “likelihood ratio trick”

$$\nabla_{\vartheta} P_{\vartheta}(\tau) = P_{\vartheta}(\tau) \nabla_{\vartheta} \log P_{\vartheta}(\tau)$$

- am înlocuit integrala cu o așteptare și am substituit  $P_{\vartheta}(\tau)$
- am înlocuit logaritmul produsului cu suma logaritmilor



# Implementarea gradientului

- Există multe metode de estimare a gradientului, bazate de ex. pe Monte-Carlo
- De exemplu, REINFORCE folosește legea de control curentă pentru a executa  $n_\tau$  traiectorii-eșantion, fiecare cu lungimea finită  $K$ , și estimează:

$$\widehat{\nabla}_{\vartheta} J_{\vartheta} = \frac{1}{n_{\tau}} \sum_{s=1}^{n_{\tau}} \left[ \sum_{k=0}^{K-1} \gamma^k r_{s,k} \right] \left[ \sum_{k=0}^{K-1} \nabla_{\vartheta} \log \hat{h}(x_{s,k}, u_{s,k}; \vartheta) \right]$$

(posibil adăugând un baseline pentru a reduce varianța)

- Comparat cu formula exactă:

$$\nabla_{\vartheta} J_{\vartheta} = \mathbb{E}_{\vartheta} \left\{ R(\tau) \sum_{k=0}^{\infty} \nabla_{\vartheta} \log \hat{h}(x_k, u_k; \vartheta) \right\}$$

- Gradientul  $\nabla_{\vartheta} \log \hat{h}$  preferabil direct calculabil (prin formulă)



# Policy gradient: algoritm

## Policy gradient cu REINFORCE

inițializează  $\vartheta$

**repeat**

execută  $n_\tau$  traiectorii-eșantion cu lungimea  $K$

estimează:

$$\widehat{\nabla}_{\vartheta} J_{\vartheta} = \frac{1}{n_{\tau}} \sum_{s=1}^{n_{\tau}} \left[ \sum_{k=0}^{K-1} \gamma^k r_{s,k} \right] \left[ \sum_{k=0}^{K-1} \nabla_{\vartheta} \log \hat{h}(x_{s,k}, u_{s,k}; \vartheta) \right]$$

actualizează parametrii:  $\vartheta \leftarrow \vartheta + \alpha \nabla_{\vartheta} J_{\vartheta}$

**until** finalizarea experimentului



## Scaun rulant asistat electric (cu Feng et al.)



- Sursă de putere hibridă: utilizator uman și baterie
- Scop: urmărirea referinței de viteză, optimizând asistența pentru:
  - (i) a atinge nivelului dorit de oboseală a utilizatorului
  - (ii) a minimiza consumului de baterie
- Provocare: utilizatorul are **dinamică necunoscută**

## PAW: Configurația experimentului

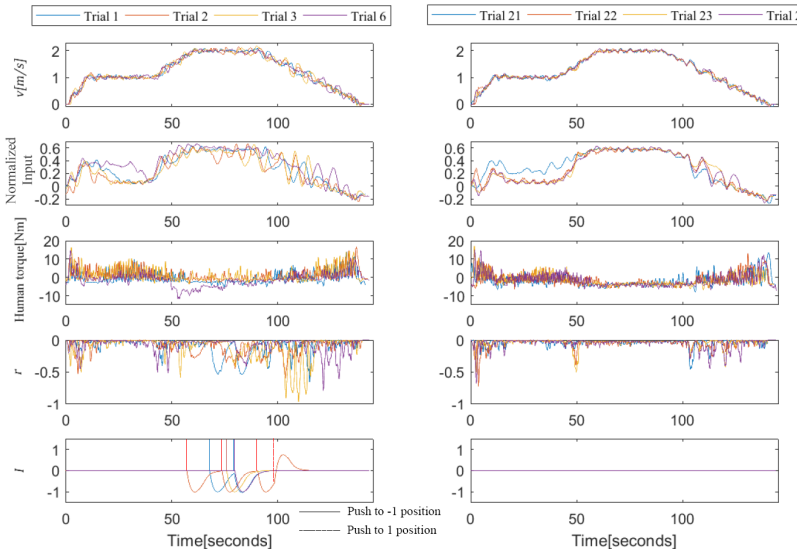
- Utilizatorul setează viteza, trage/împinge joystick-ul când e prea obosit/vrea mai mult efort
- Recompensa penalizează eroarea de viteză, semnalul joystick-ului  $I$ , și magnitudinea asistenței (pentru economisirea energiei)

$$r = -w_1(v - v_{\text{ref}})^2 - w_2 I^2 - w_3 U^2$$

- Controller PI cu reglarea constantelor folosind policy gradient (specific, algoritmul POWER)



# PAW: Resultate



## Terminologie engleză

coborâre pe gradient stohastică	=	<u>stochastic gradient descent, SGD</u>
semigradient	=	<u>semigradient</u>
diferență temporală aproximată	=	<u>approximate temporal difference</u>
gradientul pe legea de control	=	<u>policy gradient</u>



# Exerciții

- 1 Găsiți actualizările de tip semigradient ale parametrilor  $\theta$  pentru a aproxima  $V^h$  și  $V^*$ .
- 2 Generalizați actualizarea semigradient din SARSA cu return pe  $n$  pași.
- 3 Încercați să găsiți o formulă completă de coborâre pe gradient pentru SARSA, care ia în considerare dependența estimatului cu bootstrap de vectorul de parametri.

