

- noțiuni / concepte
- clasificare

# Inteligență Artificială

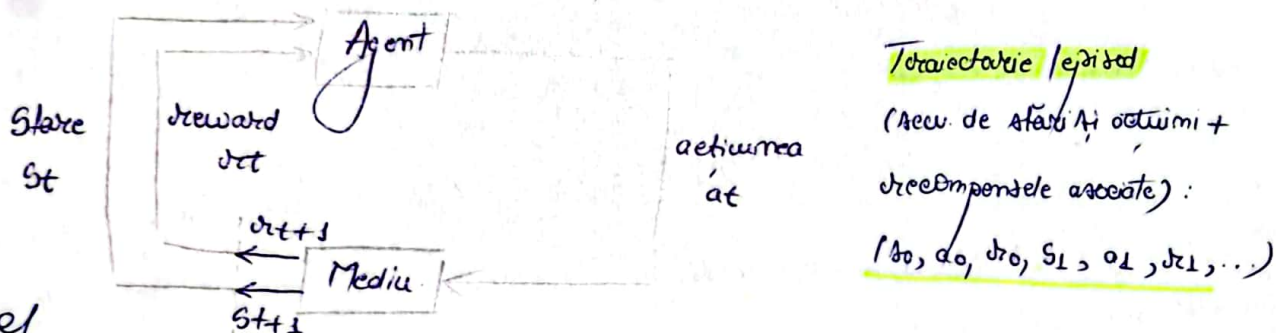
## Cursul 4

### Reinforcement Learning

Agentul trebuie să învețe un comportament față unei formă de instruire:

- la care miște datele
- face o serie de acțiuni
- la primii un feedback sub forma unor recompense (positive / negative)

În final, agentul trebuie să găsească politica pentru a trece dintr-o stare în alta.



Scopul este să maximizeze recompensele, deci vom să determinăm acțiune de stare care ne duce spre această stare finală pe baza datelor de antrenare (S, A, R)

Deep Reinforcement Learning - Aplicații

- AlphaGo (programul a învățat să joace Go (2600) → prima ca unghi o imagine cu un scar.
- 2016 → Go
- 2017 → AlphaZero a învățat să joace fără informații din jocuri ale oamenilor.
- 2017 → AlphaZero a învățat doar pe baza experienței proprii, în 8 ore ( timp relativ mic)

## Procese de decizie Markov

Mediu  $\rightarrow$  **determinist** / **stocastic** (cum ardeam probabilitatea de a ajunge într-o stare) (starea în care ajungem în ultima urnă de ieșire)

**stocastic** / nu întotdeauna ajungem în starea dorită, ci eu o anumită probabilitate

pe baza unui model de tranziție  $P(s' | s, a)$

Mediul stocastic îl vom folosi mai departe.

**Presupunerea Markov** constă în cât de mult ne uităm înapoi.

ex. pp Markov de ordinul  $T$  ne uităm doar la starea anterioară

$$P(s_t | s_{t-1} \dots s_0) = P(s_t | s_{t-1}) \quad (\text{starea curentă depinde de un nr. finit de stări anterioare})$$

**MDP** = problemă pt un mediu stocastic în care se aplică modelul de tranziție al lui Markov

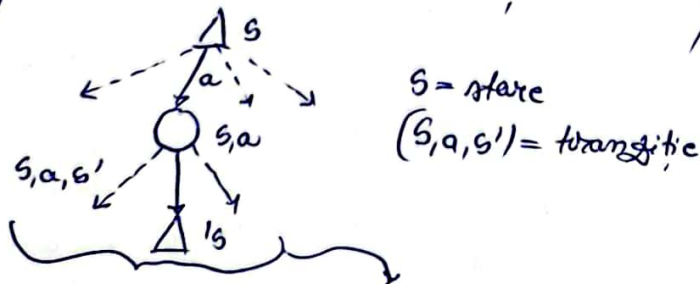
este format din:

- stări  $s \in S$
- modelul de tranziție  $P(s' | s, a)$  (probabilitatea de a ajunge într-o stare aplicând o acțiune pe starea curentă)
- funcția de recompensă  $R(s)$

Rezoluția MDP presupune găsirea unei soluții / politici. Acestea sunt alg de căutare nedeterministice.

Alg. sunt de tip programare dinamică (dar pt dimensiuni reduse ale problemei)

ex. identificarea, pentru fiecare stare, care este acțiunea cea mai profitabilă.



Pentru această pb. se poate asocia un **graf** / **arbore de căutare**, unde fiecărei stări  $i$  se poate asocia un arbore de căutare.



Într-o problemă de căutare identificăm starea inițială și finală.

Într-un proces de decizie Markov vom să identificăm o politică optimă  $\pi$  (strategie):

$$\pi: S \rightarrow A, \pi(s) = \text{acțiunea recomandată în starea } s.$$

Pentru a vedea ce e optimă este o politică, vom calcula utilitatea.

Utilitatea = tot câștigul obținut (printr-o drumă ponderată a stărilor mai apropiate)

Câștigul / Recompensa = date de la o stare la alta.

Ouzgent limit:  $U_k([s_0, s_1, \dots, s_{N+k}]) = U_k([s_0, s_1, \dots, s_N])$ ,  $\forall k \geq 0$ .

Politică nu este staționară într-un ouzgent, adică se poate schimba în timp.

Ouzgent infinit  $\rightarrow$  poate avea oricâte stări.  
politică optimă este staționară

Vom avea 2 chipuri de recompense pt. un ouzgent infinit

recompensă aditivă:

$$U_k([s_0, s_1, s_2, \dots]) = R(s_0) + R(s_1) + \dots$$

$\gamma \in [0, 1]$  = discountul (factor de discount)

recompense actualizate (pe aste o vom folosi)

$$U_k([s_0, s_1, s_2, \dots]) = R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots$$

stările mai apropiate de starea inițială au discount mai mare (conține mai mult)

ex. 

$s_0$	$\leftarrow$	$\leftarrow$	$\leftarrow$	$ $	$1$
a	b	c	d	e	

 (recompensele imediate  $R(b) = R(c) = R(d) = 0$ ),  $\gamma = 1$

Vom calcula politicile:  $U_b, U_c, U_d$ .

$$U_b = R(b) + \gamma R(s_0) \text{ (la stânga)} = 1$$

$$\downarrow$$
  

$$= R(b) + \gamma R(c) + \gamma^2 R(d) + \gamma^3 R(e) \text{ (la dreapta)} = 1$$

$\rightarrow$  politică spune să mergem la stânga

$$U_c = R(c) + \gamma R(b) + \gamma^2 R(s_0) \text{ (la stânga)}$$

$$\downarrow$$
  

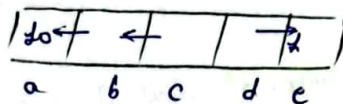
$$R(c) + \gamma R(d) + \gamma^2 R(e)$$

Dacă  $\gamma = 0,1$

$$U(b) = \begin{cases} R(b) + \gamma R(a) = 0,1 \cdot 10 = 1 \text{ (stânga)} \\ \gamma^2 R(c) + \gamma^2 R(d) + \gamma R(c) + R(b) = 9,01 \text{ (dreapta)} \end{cases}$$

$$U(c) = \begin{cases} \gamma^2 R(a) = 0,01 \cdot 10 = 0,1 \text{ stânga} \\ \gamma^2 R(c) = 0,01 \text{ dreapta} \end{cases}$$

$$U(d) = \begin{cases} \gamma^2 R(a) = 0,001 \cdot 10 = 0,01 \text{ (stânga)} \\ \gamma^2 R(c) = 0,1 \end{cases}$$



Dacă recompensele sunt mărginite și  $\gamma < 1$ :  $U^\pi([A_0, A_1, A_2, \dots]) =$

$$= \sum_{t=0}^{\infty} \gamma^t R(A_t) = \sum_{t=0}^{\infty} \gamma^t R_{\max} = \frac{R_{\max}}{1-\gamma}$$

Pentru a identifica o politică optimă, vom folosi noțiunea de **utilitate așteptată**:

$$U^\pi(s) = E \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t) \right] \text{ (dacă suntem într-o stare, avem mai multe Acv. de stări în care putem ajunge \(\Rightarrow\) considerăm media)}$$

**Politică optimă ( $\pi^*$ )** = maximizează valoarea utilității așteptate:  $\pi^* = \underset{\text{utilitatea adevărată}}{\operatorname{argmax}_{\pi}} U^\pi(s)$

$U^{\pi^*}(s)$  **utilitatea adevărată** este suma discounturilor pe care un agent o are dacă execută politica optimă

$$\pi^*(s) = \underset{a \in A(s)}{\operatorname{argmax}} \sum_{s'} P(s'/s, a) U(s')$$

**Ecuația Bellman**: utilitatea unei stări este recompensa imediată pt. acea stare,  $R(s)$ , plus utilitatea maximă a stării următoare.

$$U(A) = R(A) + \gamma \max_a \sum_{s'} P(s'/s, a) \cdot U(s')$$

Utilitatea stării (1, 1).

Mediul este stochastic  $\Rightarrow$  probabilitatea nu este mereu 1.

$$U(1, 1) = -0,64 + \gamma \max [0,8 U(1, 2) + 0,1 U(1, 2) + 0,1 U(2, 1) +$$

+ 0,1 U(1, 1), (stângă)

$$0,9 U(1, 1) + 0,1 U(1, 2), (stângă)$$

$$0,9 U(1, 1) + 0,1 U(2, 1), (jos)$$

$$0,8 U(2, 1) + 0,1 U(1, 2) + 0,1 U(1, 1)] \text{ dreapta}$$

Acțiunea cea mai bună este sus.

Metode de rezolvare ale proceselor de decizie Markov (prin programare dinamică)

**Metoda 1.** iterarea valorilor până se stabilizează.

Pașul 1) Inițializează utilitățile cu val. arbitrare

Pașul 2) Actualizează utilitatea fiecărei stări din utilitățile vecinilor:

$$U_{i+1}(s) = R(s) + \gamma \max_a \sum_{s'} P(s'/s, a) \cdot U_i(s')$$

Pașul 3) Se oprește pentru fiecare A simultan, până când se stabilizează.

ex (de pe yut):  $U_{i+1}(s) = \max_a \sum_{s'} P(s'/s, a) [R(s') + \gamma U_i(s')]$

↓  
de la la  
cramen.

$$U_1(\text{cool}) = \frac{P(\text{Cool} | \text{Cool}, a = \text{Slow})}{0,5} \left[ \frac{1}{2} R(\text{Cool}) + \frac{0}{2} U_1(\text{Cool}) \right] + \frac{0}{2} = -1, a = \text{Slow}$$

$$\frac{P(\text{Warm} | \text{Cool}, a = \text{Fast})}{0,5} \left[ \frac{2}{2} R(\text{Cool}) + \frac{0}{2} U_1(\text{Warm}) \right] + \frac{0}{2}$$

$$\frac{P(\text{Cool} | \text{Cool}, a = \text{Fast})}{0,5} \left[ \frac{2}{2} R(\text{Cool}) + \frac{0}{2} U_1(\text{Cool}) \right] = -3, a = \text{Fast}$$

$$U_1(\text{Warm}) = \frac{0,5}{1} [1] + 0,5 \cdot -1, a = \text{Slow}$$

$$1 \cdot (-10) = -10, a = \text{Fast}$$



## Metoda II Iterarea politicilor (până politica converge)

Pașul 1) evaluarea politicii: acțiunea e fixată de politică.

$$U_i(s) = R(s) + \gamma \sum_{s'} P(s'|s, \pi_i(s)) U_i(s')$$

pt. actualizarea val. utilitatilor.  
sistem de ecuații liniare cu m necunoscute.  
→ se poate calcula ușor.

Vom aplica alg. Value Iteration:  $U_{i+1}(s) = R(s) + \gamma \sum_{s'} P(s'|s, \pi_i(s)) U_i(s')$

Pașul 2) în bunăstărea politicii:  $q_i^*(s) = \max_a \sum_{s'} P(s'|s, a) \cdot U(s')$

Apoi acțiunea optimă este diferită de politica curentă, atunci se actualizează politica.

Extensii ale MDP - variante POMDP (Partial Observability) → vom avea observații parțiale / cu zgomot, adică nu știm în ce stări vom putea ajunge, ci se bazează pe observațiile făcute de agent, estimate cu probab  $P(o|s)$

