

- clasificări

- noțiuni/concepte

- important pt motare.

Inteligință Artificială

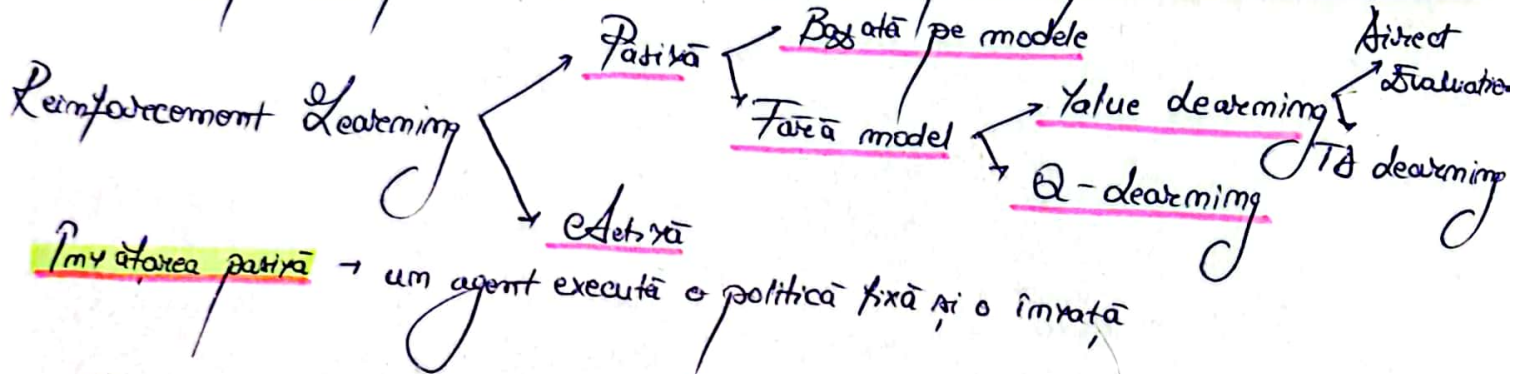
Cursul 9

Procese de decizie Markov

- avem un set de stări S și un set de acțiuni A
- modelul de tranziție $P(S'|S, a)$ este cunoscut
- fct. de recompensă $R(s)$ este cunoscută
- calcularea politicii optime

Învățare bazată pe reinforcement

- se bazează pe procese de decizie Markov.
- modelul de tranziție nu este cunoscut
- fct. $R(s)$ nu este cunoscută
- învățare politică optimă



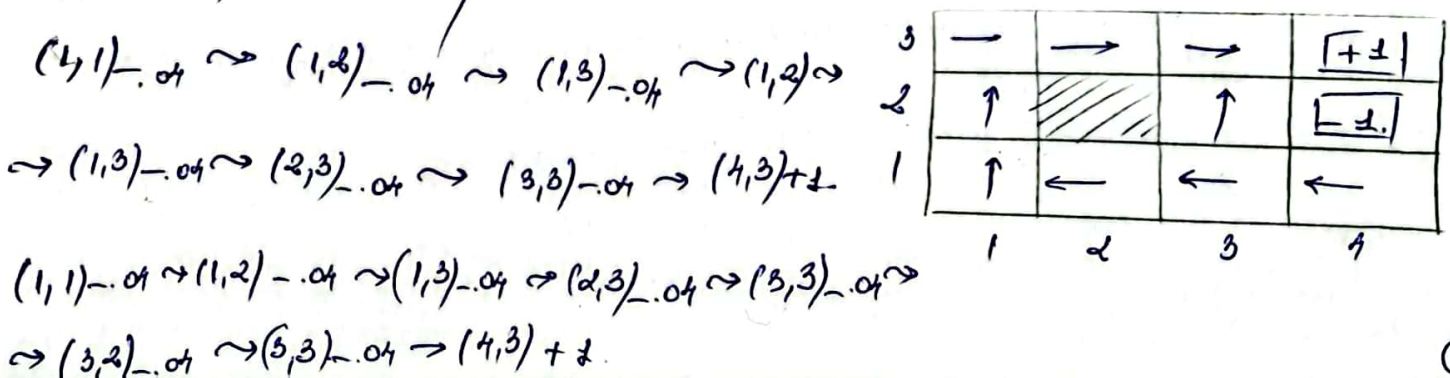
Învățarea pasivă → un agent execută o politică fixă și o învață

Scopul este să învațăm cât de bună este politica π .

Pașul 1) Învățăm utilitatea $U^\pi(s)$ pt fiecare stare.

Pașul 2) Se aplică politica și se învață din experiență

Pașul 3) Evaluăm politica



$(1,1) - .04 \leadsto (2,1) - .04 \leadsto (3,1) - .04 \leadsto (3,2) - .04 \leadsto (4,3) - 1$.

Învățarea bazată pe model \rightarrow învățăm modelul probabilist (MDP)
aplicăm tehnicile de la procesele de decizie Markov.
(Rezolvă MDP)

Cea mai simplă metodă de a învăța este să învățăm modelul probabilist și recompensele.

MDP / Programarea Dinamică Stochastică \rightarrow pt. rezolvarea procesului de decizie Markov.

Vom estima $P(s'|s, \pi(s))$ și $R(s)$ din interacțiuni.

$$U^\pi(s) = R(s) + \gamma \sum_{s'} P(s'|s, \pi(s)) U^\pi(s') \quad (\text{ecuație liniară})$$

Se va folosi un tabel de probabilități pt. a calcula procesul de învățare al modelului.

Acești jocuri de dimensiuni mari, nu este recomandată folosirea MDP (jocul de table).

Învățarea fără model \rightarrow nu învățăm modelul probabilist, e direct politică.

Pașul 1) Estimarea directă a utilității (poate fi considerată o met. de învățare supervizată).

Utilitatea = suma tuturor recompenselor obținute de către agent (reward-to-go)

Dezavantajul: nu se știe cât de ec. Bellman, deci nu se știe cât de utilitatele stăruie

Successare \Rightarrow convergența este mai mare (se învață mai greu)

Împărțirea diferențelor Temporale

→ se folosește ec. Bellman pt. calcul utilitatilor. și actualizează toate stările direct afectate.

Ecuația diferențelor temporale - calculează diferența dintre utilitatea stării următoare și utilitatea stării curente.

$$U^{\pi}(s) \leftarrow U^{\pi}(s) + \alpha (R(s) + \gamma U^{\pi}(s') - U^{\pi}(s))$$

Actualizarea utilitatilor se face online, adică după fiecare tranziție, nu după fiecare episod.

Valoarea medie $U^{\pi}(s)$ va converge către valoarea corectă.

Această se realizează suficient înlocuiri, tranzițiile vor apărea de multe ori.

Această este o funcție care are de pe măsură ce nr. de vizitări a unei stări crește, atunci $U^{\pi}(s)$ converge către valoarea corectă: $\alpha(m) = \frac{1}{m}$ sau $\alpha(m) = \frac{1}{1+m} \in (0, 1]$

Examen: aplicarea împărțirii diferențelor temporale (slide 18) pt. calcularea utilitatilor

aplicăm formula $U^{\pi}(s) = U^{\pi}(s) + \alpha (R(s) + \gamma U^{\pi}(s') - U^{\pi}(s)) =$
 $= (1-\alpha) U^{\pi}(s) + \alpha (R(s) + \gamma U^{\pi}(s'))$

$B \xrightarrow{\text{east}} C, -2$

$$U^{\pi}(B) = \underbrace{U^{\pi}(B)}_0 + \underbrace{\alpha}_{0,5} (-2 + \gamma \underbrace{U^{\pi}(C)}_0 - \underbrace{U^{\pi}(B)}_0) = -1$$

$C \xrightarrow{\text{east}} D, -2$

$$U^{\pi}(C) = \underbrace{U^{\pi}(C)}_0 + \underbrace{\alpha}_{0,5} (-2 + \gamma \underbrace{U^{\pi}(D)}_8 - \underbrace{U^{\pi}(C)}_0) = \frac{8-2}{2} = 3$$

Stările

Tranzițiile observate

	A	
B	C	D
	E	

	0	
0	0	8
	0	

	0	
-1	0	8
	0	

	0	
-1	3	8
	0	

$$\gamma = 1, \alpha = \frac{1}{2}$$

B, east, C, -2

C, east, D, -2

Adaptive Dynamic Programming / PAA

- este bazat pe modele.
- actualizează doar un ducecat, nu pe tot.

Temporal Difference Learning

- nu are nevoie de model.
- converge mai greu.
- este o aproximare a PAA

Învățarea activă → agentul actualizează politica pe măsură ce învață.

Agentul poate face explorare sau exploatare. La început, acesta trebuie să exploreze deoarece este un mediu nou, iar apoi să exploateze ce a învățat. Între cele 2 trebuie să existe un echilibru deoarece altfel se va ajunge într-un optim local.

ϵ -greedy:

Fie $\epsilon \in [0, 1]$. Acțiunea următoare selectată va fi:

- o acțiune aleatoare, cu probabilitate ϵ
- acțiunea optimă, cu probabilitate $1 - \epsilon$

ϵ nu va scădea niciodată sub un anumit prag.

Alg. Q-Learning → cel mai folosit în RL, care ține cont de stări și de acțiuni.

$$Q(s, a) = R(s, a) + \gamma \sum_{s'} P(s' | s, a) \max_{a'} Q(s', a')$$

A fiecare epocă din (s, a, s', r) actualizăm valoarea Q .

Ec. de actualizare TD Q-Learning: $Q(s, a) = Q(s, a) + \alpha (R(s) + \max_{a'} Q(s', a') - Q(s, a))$, $\alpha \in (0, 1)$

α = coeficientul de înratare.

A $\xrightarrow{a_0}$ B, 2

s	a	s'	r
A	a_0	B	2
C	stop	A	0
B	stop	A	-2
B	a_0	C	-6
C	a_0	A	2
A	a_0	A	-2

$$Q(A, a_0) = Q(A, a_0) + \alpha \cdot (2 + \max_{a'} Q(B, a') - Q(A, a_0)) = 1$$

$$Q(C, a_0) = Q(C, a_0) + \alpha \cdot (0 + \max_{a'} Q(A, a') - Q(C, a_0)) = 0.5$$

SARSA - o variantă a alg. Q-Learning, dar nu mai avem acț. maxim.

$$Q(s, a) = Q(s, a) + \alpha (R(s) + Q(s', a) - Q(s, a))$$

Deep Reinforcement Learning \rightarrow folosește o rețea neurală care aproximează valoarea Q .

Ecuația de actualizare: $Q(s, a) = Q(s, a) + \alpha [r + \max_{a'} Q(s', a') - Q(s, a)]$

Funcția de pierdere: calculăm **MSE** (Mean Squared Error) a valorii prezise pentru Q .

Target value: $\text{target}(s') = r + \max_{a'} Q(s', a')$ (este estimat \rightarrow convergență slabă)

$$\text{Minimize } \text{loss}(s, a, s') = (Q(s, a) - \text{target}(s'))^2$$

O altă metodă este să se folosească metoda gradientului.

Pt lab: de implementat Q-Learning.