

# Gyakorló feladatok a 2. pótzh. előtt

## Megerősítéses tanulás



Pataki Béla,

BME I.E. 414, 463-26-79

[pataki@mit.bme.hu](mailto:pataki@mit.bme.hu),

<http://www.mit.bme.hu/general/staff/pataki>

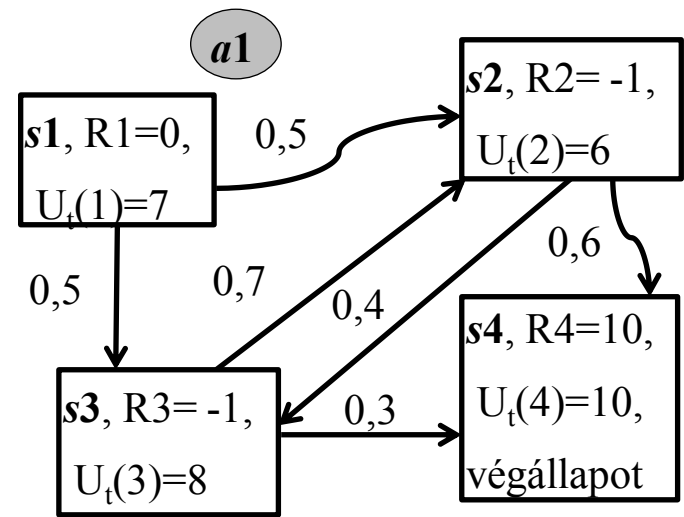
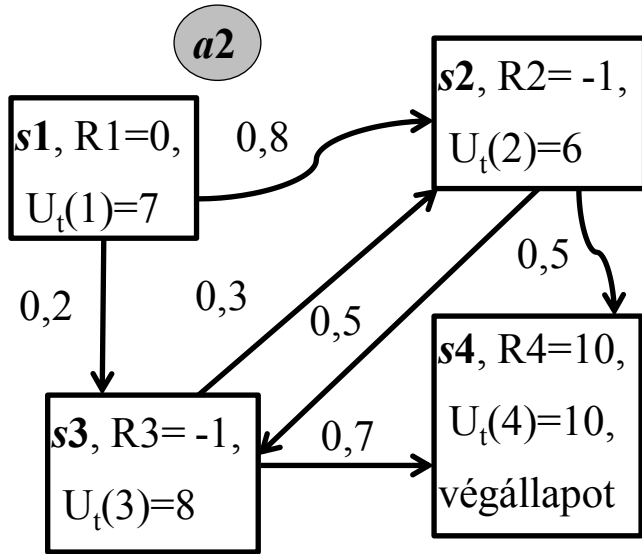
# A 2. zh megerősítéses tanulási feladata

Aktív megerősítéses tanulást végzünk a 4 állapottal rendelkező térben. Minden állapotban két cselekvés közt választhatunk:  $a_1$ , illetve  $a_2$ . A végállapot  $s_4$ , a kiinduló állapot mindig  $s_1$ . Az általunk használt felfedezési függvény (a szokásos jelölésekkel):

$$f(u, n) = \begin{cases} 20 & \text{ha } n < 4 \\ u & \text{különben} \end{cases}$$

Az eddigi futások során nyert információinkat az alábbi két ábrán, illetve az ábrák alatti táblázatban foglaltuk össze. Az egyes szövegdobozokban az állapot, az állapotban kapható jutalom, illetve az állapotra vonatkozó jelenlegi ( $t$  iterációs lépés utáni) optimista hasznosságbecslésünk van. A baloldali ábrán az  $a_1$ , a jobboldalin az  $a_2$  cselekvés esetén tapasztalt állapotátmenetvalószínűségek láthatók a dobozok közt az átmenetet mutató nyilak mellett. A táblázatban az látható, hogy az egyes állapotokban hányszor választottuk az egyik, illetve a másik cselekvést.

# A 2. zh megerősítéses tanulási feladata



$$f(u, n) = \begin{cases} 20 & \text{ha } n < 4 \\ u & \text{különben} \end{cases}$$

<b>N(sk,am)</b>	<b>s1</b>	<b>s2</b>	<b>s3</b>	<b>s4</b>
<b>a1</b>	9	2	3	
<b>a2</b>	7	4	6	

Jelenleg az s2 állapotban vagyunk, milyen cselekvést fog választani az aktív tanuló ágens? (Természetesen röviden indokolja a választ!)

## A 2. zh Q-tanulási feladata

Aktív megerősítéses Q-tanulást végzünk időbeli különbség SARSA módszerrel. Az alábbi lépéssorozat (állapot-cselekvés párok) mentén módosítunk, a tanulás bátorsági faktora 0,1; a leszámítolási tényező 0,9. Jutalmat csak az  $s_4$  végállapotban kapunk,  $R(s_4)=1$ .

$s_1, a_1 \Rightarrow s_2, a_2 \Rightarrow s_2, a_1 \Rightarrow s_1, a_1 \Rightarrow s_3, a_1 \Rightarrow s_1, a_2 \Rightarrow s_4$ , végállapot

Mi lesz az  $s_3$ -as állapot Q-értékeinek ( $Q(a_1, s_3)$  és  $Q(a_2, s_3)$ ) új becslése a lépéssorozat után, ha a lépéssorozatot megelőzően az egyes állapot-cselekvés párok Q értékeinek becslése a következő volt.

$Q(a_m, s_k)$	$s_1$	$s_2$	$s_3$	$s_4$
$a_1$	0,3	0,22	0,3	
$a_2$	0,8	0,37	0,6	

# Passzív megerősítéses tanulás

Az alábbi útvesztőben passzív megerősítéses tanulást végzünk. Az egyes mezőkben a felső szám az állapot sorszáma, az alsó, zárójelben lévő szám – ha van – az állapothoz rendelt jutalom értéke. Ahol nincs feltüntetve a jutalom értéke, ott nulla. Nem minden mezőben tüntettük fel az állapot sorszáma, értelemszerűen sorban növekednek a számok. A végállapot a szürke színnel is megjelölt 36-os. A tanulást rekurzív átlagolással végezzük, a leszámítolási tényező 1.

<b>U1 = -4</b>	<b>U2 = -3</b>	<b>U3 = -2</b>	<b>U4 = -1</b>	<b>U5 = 0</b>
<b>U6 = 1</b>	<b>U7 = 1</b>	<b>U8 = 0</b>	<b>U9 = -6</b>	<b>U10 = -5</b>
<b>U11 = 4</b>	<b>U12 = -8</b>	<b>U13 = -9</b>	<b>U14 = -4</b>	<b>U15 = -2</b>
<b>U16 = -8</b>	<b>U17 = -1</b>	<b>U18 = 0</b>	<b>U19 = 2</b>	<b>U20 = 3</b>
<b>U21 = -8</b>	<b>U22 = -7</b>	<b>U23 = -6</b>	<b>U24 = 15</b>	<b>U25 = 20</b>
<b>U26 = 25</b>	<b>U27 = 32</b>	<b>U28 = 36</b>	<b>U29 = 30</b>	<b>U30 = 25</b>
<b>U31 = 55</b>	<b>U32 = 30</b>	<b>U33 = 32</b>	<b>U34 = 36</b>	<b>U35 = 47</b>
<b>U36 = 50</b>	<b>U37 = 48</b>	<b>U38 = 47</b>	<b>U39 = 66</b>	<b>U40 = 65</b>

1	2	3						8
9				13 (-10)			16 (-5)	17
18					21 (-5)			24
25			27 (-5)					31
32 (-10)				36 (+50)			39 (+20)	40

Közvetlen (naiv) módosítás eljárással tanítunk az alábbi lépéssorozat mentén:

$$1 \Rightarrow 9 \Rightarrow 10 \Rightarrow 11 \Rightarrow 19 \Rightarrow 20 \Rightarrow 27 \Rightarrow 28 \Rightarrow 36$$

Mi lesz a 11-es állapot hasznosságértékének új becslése, ha a tanulást megelőzően a hasznosságbecslések a táblázatban láthatók voltak, és előtte már 7-szer jártunk a 11-ben?

Egy **szekvenciális döntési problémában** 4 állapot alkotja az állapotteret,  $S_3$  a végállapot. Minden állapotban kétféle cselekvés ( $a_1$  és  $a_2$ ) közt választhatunk. Az alábbi baloldali táblázatban láthatók az  $a_1$ -hez tartozó állapotátmenet-valószínűségek, a jobboldali táblázatban az  $a_2$  cselekvéshez tartozók. A leszámítolási tényező  $\gamma=0,5$ .

$a_1$ esetén		$s'$			
		S1	S2	S3	S4
$T(s \rightarrow s')$					
s	S1	0,5	0	0,5	0
	S2	0	0,5	0	0,5
	S3	0	0	0	0
	S4	0	0	0,5	0,5

$a_2$ esetén		$s'$			
		S1	S2	S3	S4
$T(s \rightarrow s')$					
s	S1	0	0,5	0	0,5
	S2	0,2	0	0,8	0
	S3	0	0	0	0
	S4	1	0	0	0

Eljárás móditerációs algoritmust alkalmazunk, az eddigi iterációk eredményeképp a  $t$ . lépésben a becsült eljárás mód:  $\pi_t(S_1) = a_1; \pi_t(S_2) = a_2; \pi_t(S_4) = a_1$ . Ez alapján a  $t$ . lépésben a becsült hasznosságok  $U_t(S_1) = 0; U_t(S_2) = 1; U_t(S_3) = R(S_3) = 1,2; U_t(S_4) = 1$ .

Adja meg a  $t+1$ -dik iterációs lépésre előálló eljárás módot!

# Szekvenciális döntési probléma

$$U(s) = R(s) + \gamma \cdot \max_a \sum_{s'} T(s, a, s') \cdot U(s')$$

$$\pi^*(s) = \arg \max_a \sum_{s'} T(s, a, s') \cdot U(s')$$

## Eljárásmód-iteráció:

1.  $t=0$  - valamilyen kiinduló eljárás mód  $\pi_0(s)$
2.  $U_t(s)$  meghatározása az  $U_t(s) = R(s) + \gamma \sum_{s'} T(s, \pi_t(s), s') \cdot U_t(s')$  lineáris egyenletrendszerből ( $n$  állapot,  $n$  egyenlet)
3. Amelyik  $s$ -re  $\max_a \sum_{s'} T(s, a, s') \cdot U_t(s') > \sum_{s'} T(s, \pi_t(s), s') \cdot U_t(s')$  arra  $\pi_{t+1}(s) = \arg \max_a \sum_{s'} T(s, a, s') \cdot U_t(s')$  a többire  $\pi_{t+1}(s) = \pi_t(s)$
4.  $t \leftarrow t+1$
5. Ha már egyik  $s$ -nél sincs változás - KÉSZ, ha volt változás, akkor folytassuk újra a 2.-nél

Egy **szekvenciális döntési problémában** 3 állapot alkotja az állapotteret,  $S_2$  a végállapot. Minden állapotban kétféle cselekvés ( $a_1$  és  $a_2$ ) közt választhatunk. Az alábbi baloldali táblázatban láthatók az  $a_1$ -hez tartozó állapotátmenet-valószínűségek, a jobboldali táblázatban az  $a_2$  cselekvéshez tartozók. A leszámítolási tényező  $\gamma=0,5$ .

$P(s \rightarrow s' | a_1)$

$s \setminus s'$	1	2	3
1	0	0,3	0,7
2	0	0	0
3	0,4	0,6	0

$P(s \rightarrow s' | a_2)$

$s \setminus s'$	1	2	3
1	0	1	0
2	0	0	0
3	0,6	0,4	0

Az állapotokban elnyerhető jutalmak, és a kiinduló hasznosságbecslések:

$s$	1	2	3
$R(s)$	+1	+20	-4
$U_0(s)$	0	+20	0

Értékiterációs algoritmussal adja meg a következő lépésben kapott hasznosságbecsléseket!



# Szekvenciális döntési probléma

**Bellman egyenlet:**  $U(s) = R(s) + \gamma \cdot \max_a \sum_{s'} T(s, a, s') \cdot U(s')$

**Értékiteráció** ( $t$  változó az iterációs lépések számlálója):

1.  $t=0$  - valamilyen kiinduló hasznosságfüggvény  $U_0(s)$
2.  $U_{t+1}(s)$  meghatározása (nem oldjuk meg az egyenletrendszer, csak kiszámítunk mindegyikből egy-egy új  $U(s)$  értéket, a max-ot persze használjuk!)

$$U_{t+1}(s) \leftarrow R(s) + \gamma \max_a \sum_{s'} T(s, a, s') \cdot U_t(s')$$

3.  $t \leftarrow t+1$
4. Ha már egyik  $U(s)$ -nél sincs változás (vagy egy adott értéknél kisebb) - KÉSZ, ha volt változás, folytassuk újra a 2.-nél

**Aktív megerősítéses Q-tanulást** végzünk felfedezési függvényt használó IK módszerrel. A tanulás bátorsági faktora 0,1; a leszámítolási tényező 0,4. Jutalmat csak az  $s_4$  végállapotban kapunk,  $R(s_4)=1$ .

Mi lesz az  $s_3$ -as állapot Q-értékeinek ( $Q(a_1, s_3)$  és  $Q(a_2, s_3)$ ) új becslése a lépéssorozat után, ha a lépéssorozatot megelőzően az egyes állapot-cselekvés párok Q értékeinek becslése a következő volt.

$N(s_k, a_m)$	$s_1$	$s_2$	$s_3$	$s_4$	$Q(a_m, s_k)$	$s_1$	$s_2$	$s_3$	$s_4$
$a_1$	4	6	7		$a_1$	0,3	0,2	0,3	
$a_2$	7	5	6		$a_2$	0,8	0,3	0,6	

$$f(u, n) = \begin{cases} 10 & \text{ha } n < 5 \\ u & \text{egyébként} \end{cases}$$

A következő átmenet után mi lesz a  $Q(s_3, a_1)$  új becslése? Milyen cselekvést javasol most a  $s_3$  állapotban az eljárás?

**$s_3, a_1 \Rightarrow s_1$**

**function** Q-Tanuló-Ágens(*észlelés*) **returns** egy cselekvés

**inputs:** *észlelés* , egy észlelés, amely a pillanatnyi  $s'$  állapotot és az  $r'$  jutalmat tartalmazza

**static:**  $Q$  , egy cselekvésérték-tábla; állapottal és cselekvéssel indexelünk

$N_{sa}$  , az állapot-cselekvés párok gyakorisági táblája  
 $s, a, r$  az előző állapot, előző cselekvés, jutalom, kezdeti értékük nulla

**if**  $s$  nem nulla **then do**  
 inkrementáljuk  $N_{sa}[s,a]$ -t  
 $Q[a,s] \leftarrow Q[a,s] + \alpha(r + \gamma \max_{a'} Q[a',s'] - Q[a,s])$   
**if** Végállapot? $[s']$  **then**  $s, a, r \leftarrow$  nulla  
**else**  $s, a, r \leftarrow s', \operatorname{argmax}_{a'} f(Q[a',s'], N_{sa}[a',s']), r'$

**return**  $a$

(aztán a meghívó függvény a következő lépésben végrehajtja  $s'$ -ben és a visszakapott  $a$  cselekvést, eljut  $s''$ -be és frissíti  $Q[a,s']$ -t stb.)