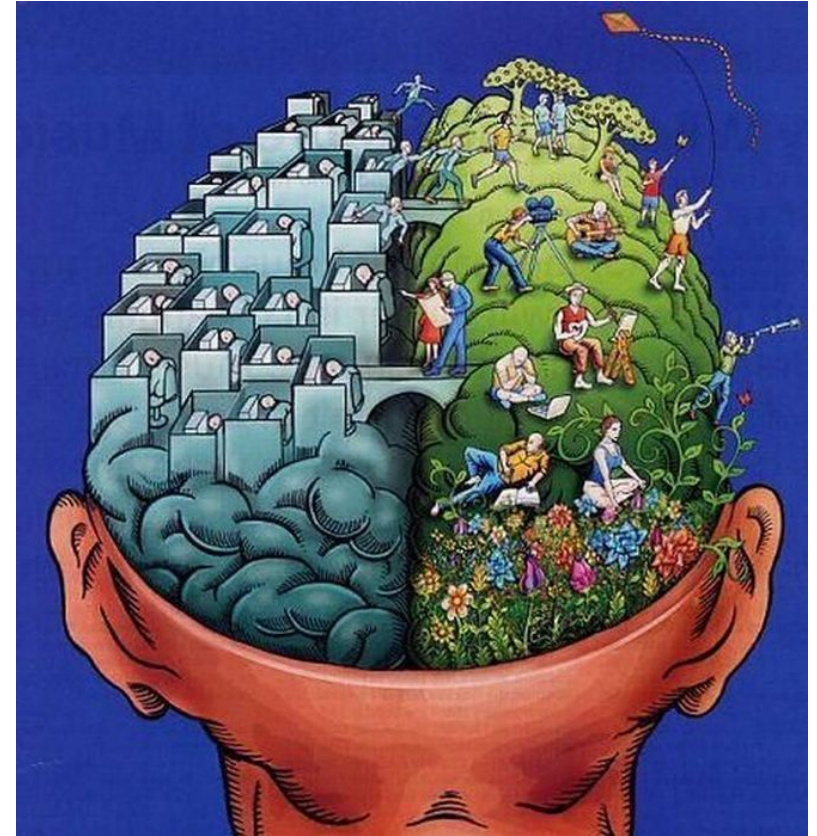


# Mesterséges Intelligencia MI

Komplex  
döntések

Dobrowiecki Tadeusz  
Eredics Péter, és mások

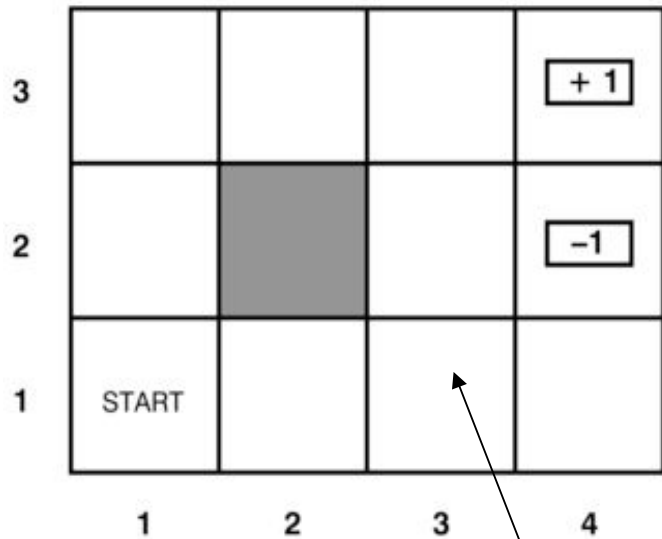


BME I.E. 437, 463-28-99

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

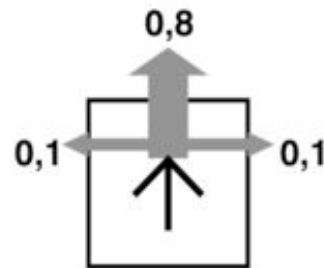
<http://www.mit.bme.hu/general/staff/tade>

# Szekvenciális döntési probléma

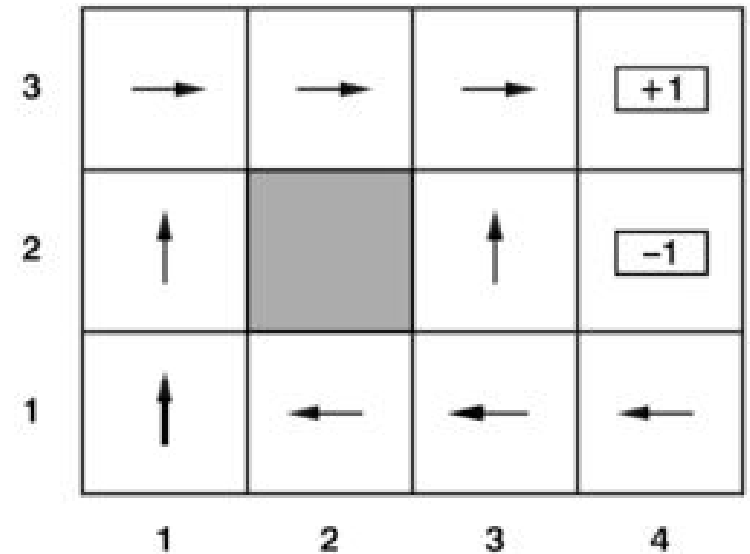


(a)

- 0.04



(b)



# Szekvenciális döntési probléma

## Markov döntési folyamat

Kezdőállapot:	$S_0$
Állapotátmenet-modell:	$T(s, a, s')$
Jutalomfüggvény:	$R(s)$ , vagy $R(s, a, s')$

**Optimális eljárás mód** = optimális mozgás, döntés cselekvés megválasztására, de nem elég egyszer, folyamatosan kell, amíg nincs a probléma vége.

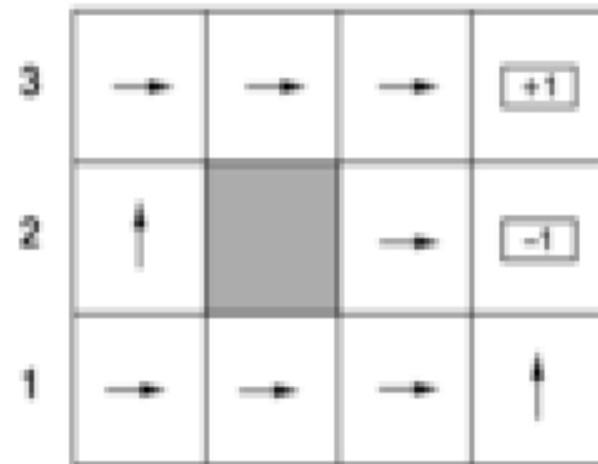
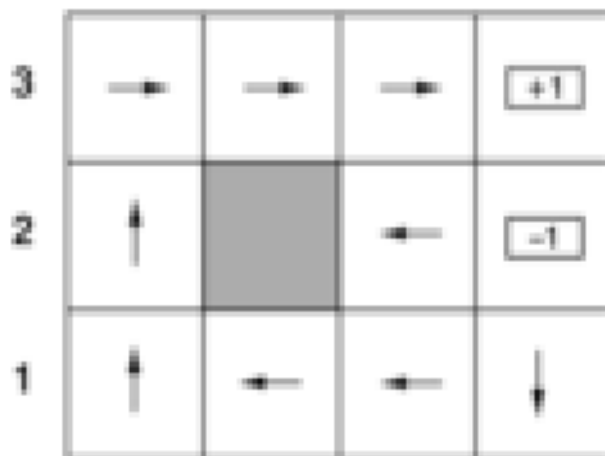
# Szekvenciális döntési probléma

$$R(s) \leq -1,6284$$

$$-0,4278 \leq R(s) \leq -0,0850$$

$$-0,0221 < R(s) < 0$$

$$R(s) > 0$$

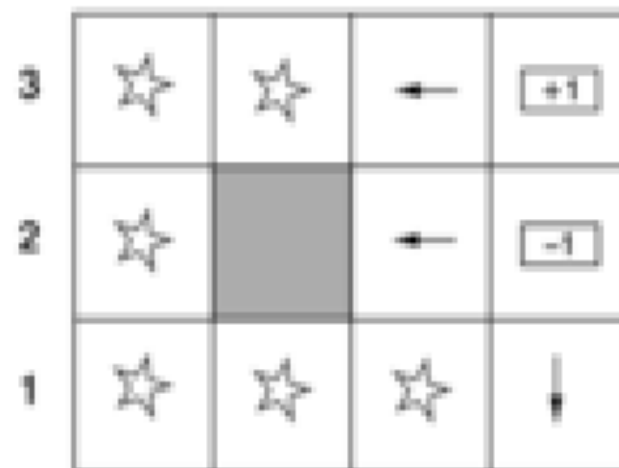
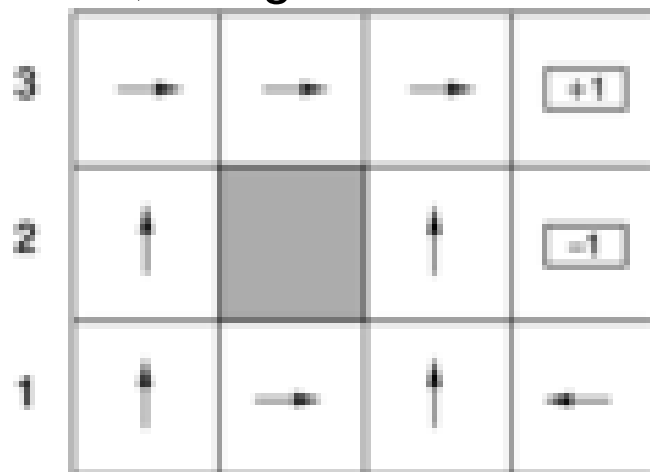


élet elviselhetetlen, ki!

élet kellemetlen; +1 állapot, -1 kockázattal

az élet csak kevéssé bánatos, ne legyen kockázat!

élet kifejezetten élvezhető, az ágens benn akar maradni



# Szekvenciális döntési probléma

## Optimalis szekvenciális döntési probléma

végtelen horizont

véges horizont

optimális eljárás mód nem-stacionárius

stacionárius

## Többattribútumú hasznosságelmélet:

ágens preferenciái az állapotsorozatok között stacionáriusok

additív jutalmak  $U_h([s_0, s_1, s_2, \dots]) = R(s_0) + R(s_1) + R(s_2) + \dots$

leszámított jutalmak

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

# Szekvenciális döntési probléma

Leszámított jutalmak, egy végtelen sorozat hasznossága

$$U_h([s_0, s_1, s_2, \dots]) = \sum_{t=0 \dots \infty} \gamma^t R(s_t) = < \sum_{t=0 \dots \infty} \gamma^t R_{\max} = R_{\max} / (1 - \gamma)$$

Ha van végállapot, ha garantált, hogy az ágens végül bele kerül, akkor nincs szükség végtelen sorozatok összehasonlítására.

Egy eljárás mód, ami garantáltan végállapotba juttat, véges eljárás mód,  $\gamma = 1$

Végtelen sorozatok összehasonlítása:

az időegységenkénti átlagjutalom

Optimális  
eljárás mód

$$\pi^* = \arg \max_{\pi} E[\sum_{t=0 \dots \infty} R(s_t) | \pi]$$

## Értékiteráció

Egy **állapot hasznossága** – a belőle kiinduló állapotsorozatok várható hasznossága

Az állapotsorozatok függenek a végrehajtott eljárásmodtól, így elsőként egy adott  $\pi$  eljárásmodra definiáljuk a hasznosságot:

$$U^\pi(s) = E \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t) \mid \pi, s_0 = s \right] \quad \pi^*(s) = \arg \max_a \sum_{s'} T(s, a, s') U(s')$$

## Optimális eljárásmod

Az **állapot hasznossága** - az állapotban tartózkodás közvetlen jutalmának és a következő állapot várható leszámított hasznosságának az összege, feltéve, hogy az ágens az optimális cselekvést választja (**Bellman** egyensúlyi **egyenlet**)

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

# Szekvenciális döntési probléma

Legyen  $\gamma = 1$  és a nem végállapotoknál  $R(s) = -0,04$

Nézzük meg a  $4 \times 3$ -as világ Bellman-egyenleteinek egyikét.

Az  $(1, 1)$  állapothoz tartozó egyenlet:

$$U(1, 1) = -0,04 +$$

$$\gamma \max\{0,8 U(1, 2) + 0,1 U(2, 1) + 0,1 U(1, 1) \text{ (Fel)},$$
$$0,8 U(2, 1) + 0,1 U(1, 2) + 0,1 U(1, 1) \text{ (Jobbra)},$$
$$0,9 U(1, 1) + 0,1 U(1, 2) \text{ (Balra)},$$
$$0,9 U(1, 1) + 0,1 U(2,1)\} \text{ (Le)}$$

Sajnos nemlineáris –  
iteráció!

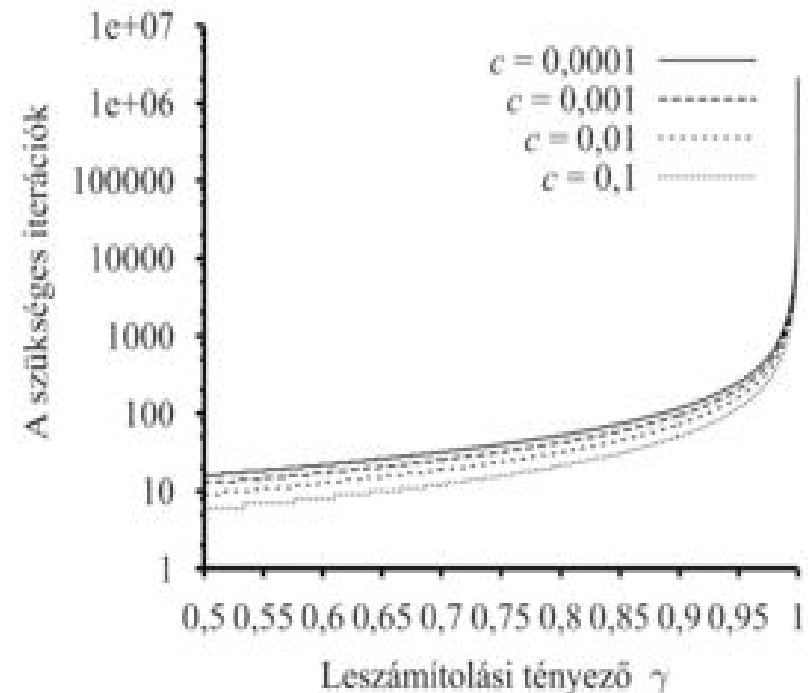
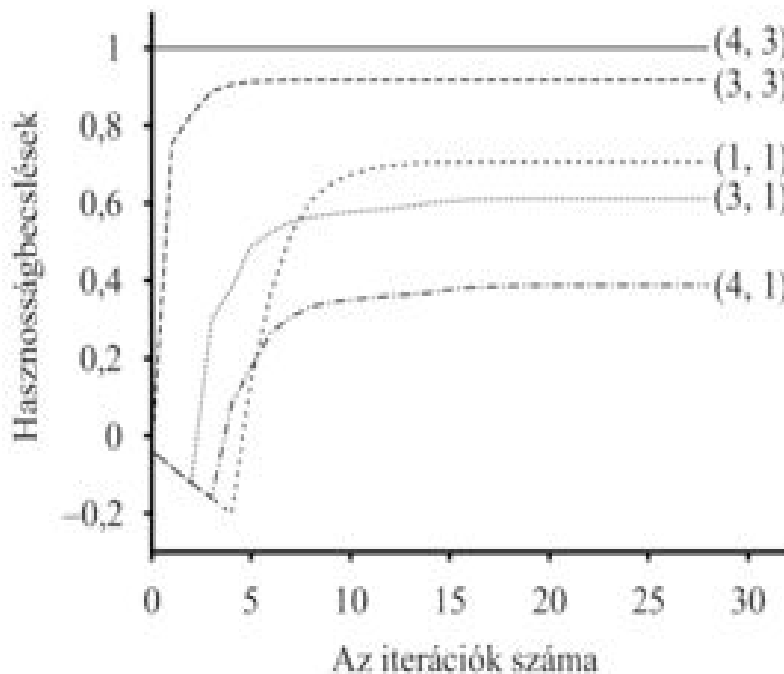
3	0,812	0,868	0,918	+ 1
2	0,762		0,660	-1
1	0,705	0,655	0,611	0,388
	1	2	3	4



# Szekvenciális döntési probléma

Bellman-frissítés:

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



A hasznosságok fejlődése

A szükséges értékiterációk száma,  
hogy a hiba garantáltan legfeljebb  
 $\varepsilon = c R_{\max}$  legyen

# Szekvenciális döntési probléma

Az értékiteráció konvergenciája – **kontrakció**  $U_{i+1} = B U_i$

$$\|U\| = \max_s |U(s)|$$

$$\|BU_i - U^*\| \leq \gamma \|U_i - U^*\|, \quad U^* \text{ az igazi: } B(U^*) = U^*$$

Az összes állapot hasznossága korlátos  $\pm R_{\max} / (1-\gamma)$  értékkel A maximális kezdeti hiba  $\|U_0 - U^*\| \leq 2R_{\max} / (1-\gamma)$

Ha  $\|U_{i+1} - U_i\| < \varepsilon(1-\gamma)/\gamma$ , akkor  $\|U_{i+1} - U^*\| < \varepsilon$

# Szekvenciális döntési probléma

## Eljárás mód-iteráció

(optimális eljárás módot is kaphatunk, ha a hasznosság függvény becslése pontatlan - ha egy cselekvés egyértelműen jobb, mint a többi, akkor a releváns állapotok pontos hasznosságát nem szükséges precízen tudnunk)

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

## Eljárás mód-értékelés

Egy adott  $\pi_i$  eljárás módnál számítsuk ki  $U_i = U^{\pi_i}$ -t, az egyes állapotok hasznosságát mintha  $\pi_i$  volna végrehajtva.

lineáris!

$$U_i(s) = R(s) + \gamma \sum_{s'} T(s, \pi_i(s), s') U_i(s')$$

$$U_{i+1}(s) \leftarrow R(s) + \gamma \sum_{s'} T(s, \pi_i(s), s') U_i(s')$$

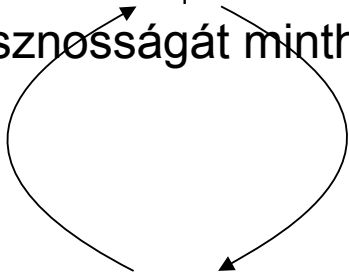
## Eljárás mód-javítás

for each  $s$  állapotra in  $S$  do

if  $\max_a \sum_{s'} T(s, a, s') U[s'] > \sum_{s'} T(s, \pi[s], s') U[s']$  then

$$\pi[s] \leftarrow \operatorname{argmax}_a \sum_{s'} T(s, a, s') U[s']$$

Módosított eljárás mód-iteráció



# Szekvenciális döntési probléma

## Részlegesen megfigyelhető Markov döntési folyamat

Kezdőállapot:  $S_0$

Állapotátmenet-modell:  $T(s, a, s')$

Jutalomfüggvény:  $R(s)$ , v.  $R(s, a, s')$

Megfigyelési modell,

az  $s$  állapotban az  $o$  megfigyelés érzékelésének a valószínűsége

$O(s, o)$

**Hiedelmi állapot** =  $b(s)$  = eloszlás állapotok felett

0,111	0,111	0,111	0,000
0,111		0,111	0,000
0,111	0,111	0,111	0,111

$$\left\langle \frac{1}{9}, \frac{1}{9}, \frac{1}{9}, \frac{1}{9}, \frac{1}{9}, \frac{1}{9}, \frac{1}{9}, \frac{1}{9}, \frac{1}{9}, 0, 0 \right\rangle$$

$$b'(s') = \alpha O(s', o) \sum_s T(s, a, s') b(s)$$

(szűrés)

# Szekvenciális döntési probléma

$$\begin{aligned} P(o | a, b) &= \sum_{s'} P(o | a, s', b) P(s' | a, b) \\ &= \sum_{s'} O(s', o) P(s' | a, b) = \sum_{s'} O(s', o) \sum_s T(s, a, s') b(s) \end{aligned}$$

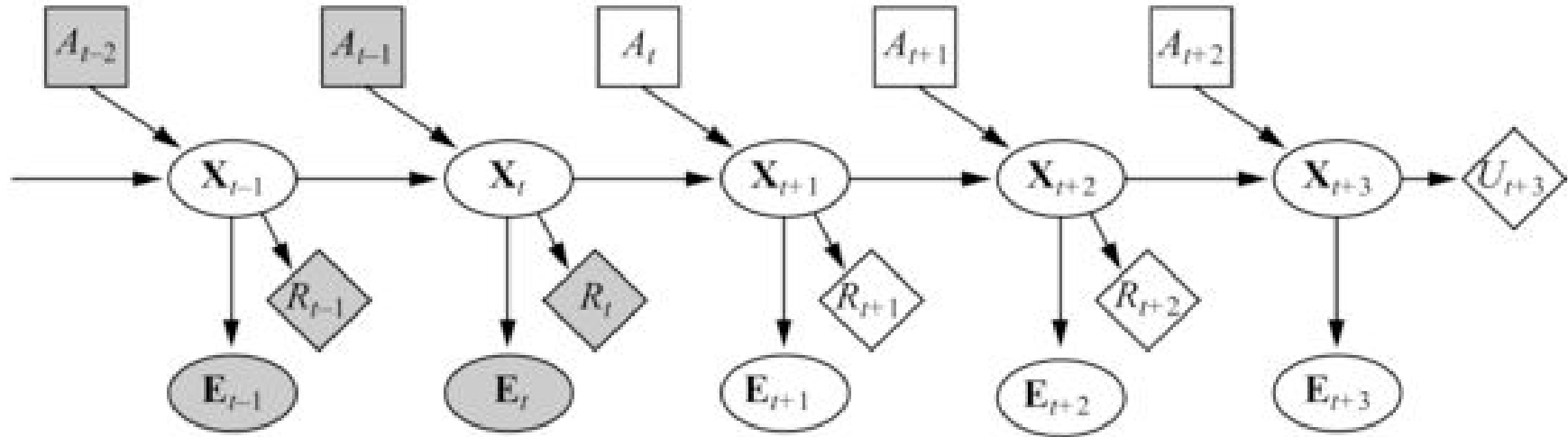
$$\begin{aligned} \tau(b, a, b') &= P(b' | a, b) = \sum_o P(b' | o, a, b) P(o | a, b) \\ &= \sum_o P(b' | o, a, b) \sum_{s'} O(s', o) \sum_s T(s, a, s') b(s) \end{aligned}$$

$$\rho(b) = \sum_s b(s) R(s)$$

RMMDF megoldása a fizikai (véges) állapottérben redukálható egy MDF megoldására a hozzá tartozó hiedelmi állapot térben (val. eloszlások folytonos terében)

# Szekvenciális döntési probléma

Dinamikus Döntési Hálók (nem foglalkozunk vele)



## Játékelmélet

(vele sem)

	<i>Aliz:tanúskodik</i>	<i>Aliz:tagad</i>
<i>Bendegúz:tanúskodik</i>	$A = -5, B = -5$	$A = -10, B = 0$
<i>Bendegúz:tagad</i>	$A = 0, B = -10$	$A = -1, B = -1$

## Működési mód tervezés (vele sem)

(egy olyan játék tervezése, aminek a megoldása az egyes ágensek által követett saját racionális stratégiáik együttese, és ez egy globális hasznosságfüggvény maximálását eredményezi, pl. árverés)