

Agentes Cognitivos e Adaptativos 2020.1
Alunos: Diana Marcela e Vinícius Andrade

Exercício Decisões Sequenciais

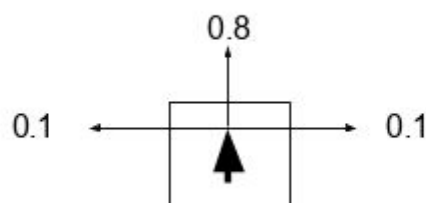
Implemente o algoritmo de MDP com as equações de Bellman e aplique no exercício em anexo. O exercício é similar ao ambiente 4 x 3 visto nas aulas. No entanto, a posição não é mais um obstáculo e sim uma casa onde há uma recompensa igual a -0.5. Além disso, existe uma casa com recompensa positiva igual +0.2.

● **Exemplo: Ambiente 4x3**

- Recompensas $R(s)$:
- Dois estados finais $R(s) = +1$ ou $R(s) = -1$
- $R(s) = -0.5$ na posição (2,2)
- $R(s) = +0.2$ na posição (4,1)
- $R(s) = r$ (valor fixo) para todos estados não terminais

3	<div>r</div>	<div>r</div>	<div>r</div>	<div>+1</div>
2	<div>r</div>	<div>-0.5</div>	<div>r</div>	<div>-1</div>
1	<div>r</div>	<div>r</div>	<div>r</div>	<div>+0.2</div>
	1	2	3	4

- Implemente o algoritmo de MDP com as equações de Bellman e retorne a política retornada para as seguintes situações:
 - $r = -0.4$
 - $r = -0.04$
 - $r = -0.0004$



- Comente as diferenças entre as políticas retornadas
 - As políticas observadas:

- $r = -0.4$

Como apresentado a imagem abaixo com uma recompensa muito baixa chegando a se tornar um custo, o algoritmo procura sempre ir na direção do estado que maximiza sua recompensa, e como o estado (1,4) possui uma recompensa ativa de **0.2** as políticas tende a levar o caminho para a mesma com exceção da terceira linha que tende a levar para o estado terminal de utilidade 1 como demonstrado na **Figura 1**.



Figura 1

Curiosamente se iniciarmos pelo estado (1,4) ele tenderá a ficar nesse estado indo sempre para baixo e voltando, já o estado (1,3) o redireciona para o estado (1,4). E isto é algo que repete em todos os ciclos abaixo que possuem recompensas menores.

Porém, o mesmo não é visto quando a recompensa do estado (1,4) é negativa, no qual ele sempre tenta fugir desse estado como na **Figura 2**.

-0.7213	-0.1591	0.4170	1.0000
-1.1682	-0.7744	-0.2427	-1.0000
-1.6297	-1.2326	-0.8224	-1.0350
R= -0.4 POLICY:			
RIGHT	RIGHT	RIGHT	*
UP	UP	UP	*
UP	RIGHT	UP	LEFT

Figura 2

Já se a recompensa geral é de **-0.6** for menor que a do estado (1,4), de **-0.2**, ao chegar nesse estado o agente procura terminar logo com o caminho e se mata no estado terminal de (2,4), pois o custo para qualquer outro estado vizinho seria maior que o do terminal negativo.

-1.5108	-0.6881	0.1423	1.0000
-2.0923	-1.3317	-0.7193	-1.0000
-2.7590	-2.0925	-1.5126	-1.2792
R= -0.6 POLICY:			
RIGHT	RIGHT	RIGHT	*
RIGHT	UP	UP	*
UP	UP	UP	UP

Figura 3

- **r = -0.04**

Já para a recompensa de **-0.04** as primeira interações produzem políticas que tendem ao caminho que leva para o estado terminal (3,4) que possui a utilidade positiva de **1** demonstrado na **Figura 4**, porém após algumas dezenas interações sem muita alteração, a política muda começa a encaminhar para o estado (1,4) que possui uma recompensa de **0.2** e eventualmente se torna superior à utilidade do estado terminal apresentado da **Figura 5**.

0.6107	0.7524	0.9058	1.0000
0.3716	0.1588	0.5799	-1.0000
0.1552	0.3680	0.6458	0.9023
R= -0.04 POLICY:			
RIGHT	RIGHT	RIGHT	*
UP	UP	UP	*
UP	RIGHT	RIGHT	DOWN

Figura 4: 10 interações

1.9030	1.8128	2.0111	1.0000
2.1655	1.9788	2.3768	-1.0000
2.4439	2.7338	3.0830	3.4261
R= -0.04 POLICY:			
DOWN	RIGHT	DOWN	*
DOWN	DOWN	DOWN	*
RIGHT	RIGHT	RIGHT	DOWN

Figura 5: 23 interações

- $r = -0.0004$

De modo geral, a recompensa **-0.0004** produz política que leva a encontrar o caminho do estado terminal positivo (3,4), de valor **1**, mais rápido que os demais por ter um custo muito baixo. Apesar de, estando nas posições (1,2) e (1,3), o agente tende a escolher o caminho que leva ao estado (1,4) com a recompensa **0.2**, como pode ser visualizado na **Figura 6**.

0.7567	0.8531	0.9595	1.0000
0.5429	0.2605	0.6709	-1.0000
0.3435	0.4827	0.7151	0.9249
R= -0.0004 POLICY:			
RIGHT	RIGHT	RIGHT	*
UP	UP	UP	*
UP	RIGHT	RIGHT	DOWN

Figura 6

Neste caso, sendo a recompensa geral de **-0.6**, observa-se que o agente tem uma maior tendência a escolher caminhos que o levem a recompensa **-0.2**, na posição (1,4), com exceção se o início se der na posição (3,1), que leva o agente direto ao estado terminal positivo de posição (3,4) e valor **1**.

-1.3977	-0.6004	0.1936	1.0000
-1.5392	-0.7774	-0.1608	-1.0000
-1.0119	-0.0783	0.8797	1.8828
R= -0.6 POLICY:			
RIGHT	RIGHT	RIGHT	*
RIGHT	RIGHT	DOWN	*
RIGHT	RIGHT	RIGHT	DOWN

Figura 7