

Aprendizagem por Reforço - Exercícios

- Para cada um dos agentes indicados especifique o ambiente, agente, ações e recompensas:
 - Robot aspirador.
 - Campanha de marketing para venda de um produto com base no perfil do cliente.
 - Controlo de temperatura de uma sala (com sistema de ar condicionado).
- Considere um agente num mundo em grelha 3x2, conforme mostrado na figura 1. Começamos no estado “1” e terminamos no canto superior direito (estado 6). Ao atingir o estado 6, recebe uma recompensa de +10 e inicia um novo episódio. Em todas as outras ações que não levam ao estado “6”, a recompensa é -1. O agente inicia na célula A e pretende alcançar F. Quando alcançar F recebe uma recompensa de +10 e o episódio termina. Para qualquer outro movimento que não conduza a F recebe a recompensa de “-1”.

A start	B	C
D	E	F finish

Figura 1 – Ambiente.

Em cada estado temos quatro ações possíveis: cima, baixo, esquerda e direita. Para cada ação, o agente move-se de forma determinista na direção pretendida. Não são possíveis movimentos para fora da grelha.

As estimativas atuais para os valores de $Q(s,a)$ são apresentadas na tabela abaixo:

S/A				
	Up	Down	Left	Right
A	4	--	--	3
B	6	--	3	8
C	9		7	
D	--	2		5
E	--	6	5	8

- a. Considerando que existe o conhecimento completo do ambiente, atualize o valor de $Q(C, \text{left})$ com base na equação de Bellman, para uma política sôfrega e taxa de desconto de 0,9.
 - b. Assuma agora que não tem o conhecimento completo do ambiente. Atualize a Q-table com o algoritmo “SARSA”. A partir de B, seguiu-se a trajetória: “B – baixo - E - direita – F”, terminando o episódio.
Atualize a tabela considerando coeficiente de aprendizagem de 0,2 e desconto de 0,8.
3. Considere um sistema com dois estados (S1,S2) e duas ações (a1,a2). Um agente executa ações e observa as recompensas e transições de acordo com as iterações (Estado atual; recompensa; ação; transição resultante):

$$\begin{aligned} It1: S1 \ r=-10 \ a1: S1 \rightarrow S1 \\ It2: S1 \ r=-10 \ a2: S1 \rightarrow S2 \\ It3: S2 \ r=+20 \ a1: S2 \rightarrow S1 \\ It4: S1 \ r=-10 \ a2: S1 \rightarrow S2 \end{aligned}$$
 - a. Represente a Q-table, com entradas inicializadas a zero
 - b. Atualize a tabela com o algoritmo de Q-learning para as quatro iterações com coeficiente de aprendizagem de 0.5 e taxa de desconto de 0.5.
4. Considere o ambiente do openAI Gym, “Taxi-v3”. Aplique o método “Q-Learning” para ensinar o agente (táxi) a apanhar e largar passageiros nos locais certos.
 - a. Especifique o ambiente, estado, ações do agente e recompensas.
 - b. Implemente o algoritmo e treine o agente para 5000 episódios.
 - i. Avalie o número de épocas e penalizações por episódio e comente o resultado.
 - ii. Qual a melhor ação para o estado 328 representado na Figura?
 - iii. Compare com um algoritmo de “força bruta”, isto é com escolha aleatória das ações.
 - c. Com base numa política “epsilon-greedy” analise o efeito dos hiper-parâmetros: coeficiente de aprendizagem, taxa de desconto e taxa de exploração.
Inicie com $\alpha=0.1$; $\gamma=0.7$ e $\epsilon=0.2$.
5. Considere o ambiente do openAI Gym, “CartPole-v1”. Implemente um agente SARSA com base na biblioteca keras-rl.
 - a. Identifique o ambiente, o agente, ações o estado e recompensas.
 - b. Implemente uma gente com base no algoritmo Deep Q-Learning.

Soluções:

2a.

A equação de Bellman específica indica:

$$Q(s, a) = \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma \sum_{a'} \pi(s', a') Q(s', a')]$$

Uma vez que o ambiente é determinístico, apenas consideramos um estado seguinte s' (Probabilidade (C, left) \rightarrow B = 1) e sendo uma política "greedy", consideramos apenas a melhor ação a'

$$Q(\text{C, left}) = 1 * [-1 + 0.9 * (1 * 8)] = 6.2$$

2a.

Para um agente SARSA - $Q(s_t, a_t) = Q(s_t, a_t) + \alpha(r_t + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t))$

	SARSA policy	Q-learning policy
Escolher A'	ϵ -greedy ($\epsilon > 0$ com exploration)	ϵ -greedy ($\epsilon > 0$ com exploration)
Atualizar Q	ϵ -greedy ($\epsilon > 0$ com exploration)	greedy policy ($\epsilon=0$, sem exploration)

Assim:

$$Q(\text{B, Baixo}) = 6 + 0,2 * (-1 + 0,8 * 8 - 6) = 6 + 0,2 * (-0,6) = 6 - 0,12 = 5,88$$

$$Q(\text{E, Direita}) = 8 + 0,2 * (10 + 0,8 * 0 - 8) = 8,4.$$

3.a

Q	a1	a2
S1	0	0
S2	0	0

3.b

Q-learning: $Q(\text{state}, \text{action}) \leftarrow (1 - \alpha)Q(\text{state}, \text{action}) + \alpha(\text{reward} + \gamma \max_a Q(\text{next state}, \text{all actions}))$

- Iteração 1: S1 $r = -10$ a2: S1 \rightarrow S1

$$Q(s1, a1) = (1 - 0,5) * 0 + 0,5 * (-10 + 0,5 * \max[0, 0]) = 0 - 5 = -5$$

Q	a1	a2
S1	-5	0
S2	0	0