

Solution of Reinforcement Learning: An Introduction by Sutton and Barto

Chapter - 3 Finite Markov Decision Process

Exercise 3.1

- ① Mars Rover, it can have multiple objective but for the time being, say it has objective to go from point A. to point B.

Rewards

Small +ve reward for each step in the ^{correct} direction and
large +ve " " reaching point B.
Small -ve " " moving in wrong direction.

States

The charge of its battery.
The other type of state may be the type of terrain it is on. (like sandy or rocky).

Action

Whether it has to stand and charge itself.
The direction in which it has to go.

Note: S, A, R are not ~~finite~~ ^{discrete} here.

- ② A system that could sense the presence of a person in a room and then adjust the voltage of light and fans, so as to minimize the total power consumption.

Reward

~~0~~ -1 for high electricity consumed.
1 " low " "
0 " normal " "

States

(A, B), A is the prob. of person in room.
B " current power supplied

Action

Whether to inc. or dec. power supplied.

Ex 3.2

No, there are situations where $p(s', r | s, a)$ can not be defined or is not discrete probability distribution. Example is the Mars Rover in Ex 3.1 ①.

Ex 3.3

The general rule is that anything that can not be changed arbitrarily by the agent is considered a part of its environment.

So, the correct choice for the line b/w agent and env. is accelerator, steering wheel and brakes.

Ex 3.4

~~Ex 3.4~~

$$E[R_{t+1} | S_t] = \sum_{a \in A} \pi(a | S_t) \sum_{s' \in S} r \cdot p(s', r | S_t, a)$$

Ex 3.6

If $s \neq s_T$

$$\sum_{s' \in S} \sum_{r \in R} p(s', r | s, a) = 1$$

If $s = s_T$

$$\sum_{s' \in S} \sum_{r \in R} p(s', r | s, a) = 0 \quad (\text{As the state } s_T \text{ is final state and no transition can be made from it})$$

Ex 3.7

$$R_t = -\gamma^k \quad [k \text{ is number of time steps until failure}]$$

Assuming that once it has failed, it will be in failure state without external intervention and the task to be non episodic the return would be

$$\begin{aligned} & -\gamma^k - \gamma^{k+1} - \dots \\ & = -\gamma^k (1 + \gamma + \gamma^2 + \dots) = \frac{-\gamma^k}{1-\gamma} \quad (|\gamma| < 1) \end{aligned}$$

Ex 3.8. While formulating a MDP for practical purpose one of the things to keep in mind is that the agent should receive reward in regular intervals for effective learning.

A much nicer way to formulate this problem is to give a reward -1 for each time step the agent is inside the maze and a big +ve reward when the agent has escaped the maze.

Ex 3.9. Assuming R_6 & G_6 to be 0.

$$G_5 = R_6 + \gamma G_6 \\ = R_6 + \frac{G_6}{2} = 0$$

$$G_4 = R_5 + \gamma G_5 = 2$$

$$G_4 = 2 + \frac{1}{2} G_5 = 2$$

$$G_3 = 3 + \frac{G_4}{2} = 4$$

$$G_2 = 6 + \frac{G_3}{2} = 8$$

$$G_1 = 2 + \frac{G_2}{2} = 6$$

$$G_0 = -1 + \frac{G_1}{2} = 2.$$

Ex 3.10

$$G_0 = R_1 + \gamma R_2 + \gamma^2 R_3 + \dots$$

$$= 2 + \gamma [7 + \gamma 7 + \gamma^2 7 + \dots]$$

$$= 2 + \frac{0.9 \times 7}{1 - 0.9}$$

$$= 2 + 63$$

$$= 65$$

$$G_1 = 7 + \gamma 7 + \gamma^2 7 + \dots$$

$$= \frac{63}{0.9} = 70$$

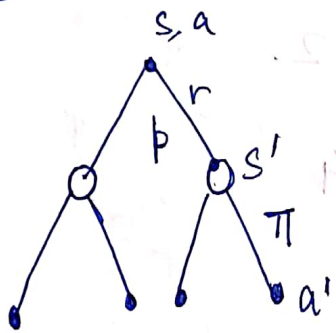
Ex 3.11

$$\begin{aligned}
 G_t &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \\
 &= 1 + \gamma + \gamma^2 + \gamma^3 + \dots \\
 &= \frac{1}{1-\gamma} \quad (|\gamma| < 1)
 \end{aligned}$$

Ex 3.12

$$\begin{aligned}
 0.7 &= \frac{1}{4} [0 + 0.9 \times 0.7] + \frac{1}{4} [0 + 0.9 \times 2.3] \\
 &\quad + \frac{1}{4} [0 + 0.9 \times 0.4] + \frac{1}{4} [0 - 0.9 \times 0.4] \\
 &= \frac{0.9}{4} [3] = \frac{2.7}{4} = 0.68 \approx 0.7
 \end{aligned}$$

Ex 3.13



$$\begin{aligned}
 q_{\pi}(s, a) &= \sum_{s'} p(s' | s, a) \left\{ E[r | s', s, a] + \gamma V_{\pi}(s') \right\} \\
 &= \sum_{s'} p(s' | s, a) \left\{ E[r | s', s, a] + \gamma \sum_a \pi(a | s') q^{\pi}(s', a) \right\}
 \end{aligned}$$

Ex 3.14

$$G_{t+1} = R_{t+1} + \gamma G_{t+1}$$

$$G'_t = (R_{t+1} + c) + \gamma (R_{t+2} + c) + \gamma^2 (R_{t+2} + c) + \dots$$

$$= G_t + c + c\gamma + c\gamma^2 + \dots$$

$$G'_t = G_t + \frac{c}{1-\gamma}$$

$$V_c = \frac{c}{1-\gamma}$$

Ex 3.15 Yes, it would have effect as we get
 $V_c = \frac{c}{1-\gamma}$ for all states only if the task is
 on episodic, when the task is episodic V_c will be
 different for all the states.

Ex 3.16

$$V_{\pi}(s) = \sum_a \pi(a|s) q_{\pi}(s, a)$$

Ex 3.17

$$q_{\pi}(s, a) = \sum_{s'} p(s'|s, a) [E[r|s', s, a] + \gamma V_{\pi}(s')]$$

$$\text{or } q_{\pi}(s, a) = \sum_{s', r} p(s', r|s, a) [r + \gamma V_{\pi}(s')]$$

Ex 3.20

$$q_*(s, a) = \sum_{s', r} p(s', r|s, a) [r + \gamma \max_{a'} q_*(s', a')]$$

$$q_*(\text{high}, \text{wait}) = r_{\text{wait}} + \gamma \max(q_*(\text{high}, \text{wait}), q_*(\text{high}, \text{search}))$$

$$q_*(\text{high}, \text{search}) = \alpha [r_{\text{search}} + \gamma \max(q_*(\text{high}, \text{wait}), q_*(\text{high}, \text{search}))]$$

$$+ (1-\alpha) [r_{\text{search}} + \gamma \max(q_*(\text{low}, \text{wait}),$$

$$q_*(\text{low}, \text{recharge}),$$

$$q_*(\text{low}, \text{search}))]$$

Ex 3.22

$$\gamma = 0$$

$$q(q_0, \text{left}) = 1 + 0 = 1$$

$$q(q_0, \text{right}) = 0$$

$$1 > 0$$

correct choice left.

$$\gamma = 0.9$$

$$q(q_0, \text{left}) = 1 \times (1 + 0.9 \times 0) = 1$$

$$q(q_0, \text{right}) = 0 + 0.9 \times 2 = 1.8$$

$$\gamma = 0.5 \quad 1.8 > 1 \rightarrow \text{right}$$

$$q(q_0, \text{left}) = 1 \times (1 + 0.5 \times 0) = 1$$

$$q(q_0, \text{right}) = 1 \times (0 + 0.5 \times 2) = 1$$

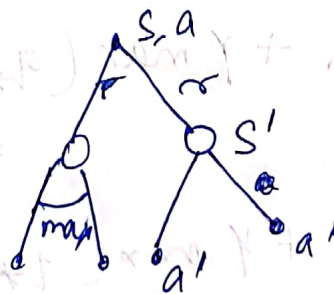
$= 1$ (any direction)

Ex 3.23

$$v_*(s) = \max_a (q_*(s, a))$$

Ex 3.24

$$q_*(s, a) = \sum_{s', r} p(s', r | s, a) [\gamma + \gamma v_*(s')]$$



Ex 3.25

$$a_* = \arg \max_a q_*(s, a)$$

$$Ex 3.26 \quad a_* = \arg \max_a \sum_{s', r} p(s', r | s, a) [\gamma + v_*(s')]$$