

### Question 1

$$1) \quad G_0 = R_1 + \gamma R_2 + \gamma^2 R_3 + \gamma^3 R_4$$

$$R_1 : s_1 \rightarrow s_2 \quad -1$$

$$R_2 : s_2 \rightarrow s_1 \quad -1 \quad \gamma = 0.9$$

$$R_3 : s_1 \rightarrow s_2 \quad -1$$

$$R_4 : s_2 \rightarrow s_{\text{Term}} \quad 10$$

$$\therefore G_0 = -1 + 0.9(-1) + 0.81(-1) + 10(0.729)$$

$$= \underline{\underline{4.58}}$$

$$2) \quad V_\alpha(s_2) = 0.5(-1 + 0.9 V_\alpha(s_1)) + 0.5(10)$$

$$= \underline{\underline{4.5 + 0.45 V_\alpha(s_1)}}$$

### Question 2:

Agent will learn to suck dust and spit it out again and it will keep repeating this.

Room keeps getting dirtier but agent gets infinite +1s. So what we can do instead is we can reward based on final room cleanliness.

### Question 3: PART A

We mathematically need  $\gamma < 1$  for infinite horizon tasks to ensure that the value function converges. If the task never ends with rewards always +1 and  $\gamma = 1$ .

$$V_\alpha(s) = \sum_{t=0}^{\infty} 1^t x 1 = \infty$$

### PART B:

$\gamma = 0$  (Impulsive Agent): This agent only cares about immediate rewards since future rewards are discounted to 0.

$\gamma = 0.99$  (Strategic Agent): The agent values future rewards as much as immediate ones. It makes long term strategic decisions.

### Question 4:

Yes, the optimal policy changes.

In the original setup with  $R = -1$ , the return is  $G = -n$  where  $n$  is the number of steps.

The agent maximises return by minimizing steps.

with the modified reward  $R_{\text{new}} = +1$ , return becomes  $G_t = +n$ . Now, the agent maximizes return by maximizing steps. The new optimal policy will ~~not~~ take the longest possible path.

### Question 5:

$$V_\alpha(s) = E_\alpha [G_t | s_t = s]$$

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

$$V_\alpha(s) = E_\alpha [R_{t+1} + \gamma G_{t+1} | s_t = s]$$

$$V_\alpha(s) = E_\alpha [R_{t+1} | s_t = s] + \gamma E_\alpha [G_{t+1} | s_t = s]$$

$$V_\alpha(s) = \sum_a \alpha(a|s) E[R_{t+1} + \gamma G_{t+1} | s_t = s, A_t = a]$$

$$V_\alpha(s) = \sum_a \alpha(a|s) \sum_{s', r} p(s', r | s, a)$$

$$V_\alpha(s) = \sum_a \alpha(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma E[G_{t+1} | s_{t+1} = s']]$$

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

### Question 6:

$$1) V_\alpha = r_\alpha + \gamma P^\alpha V_\alpha$$

$$V_\alpha = (I - \gamma P^\alpha)^{-1} r_\alpha$$

$$2) O(N^3) = O((10^{20})^3) = O(10^{60})$$

$$\text{Supercomputer speed : } 10^{18} \text{ FLOPs/second}$$

$$\text{Time required} = \frac{10^{60}}{10^{18}} = 10^{42} \text{ seconds}$$

3) This is why model-free  $\approx 10^{34}$  years Monte-Carlo are important methods such as tractable for realistic problems. They are

### Question 7:

$$1) \Pi'(s) = \arg \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V^*(s')]$$

$$2) \Pi'(s) = \arg \max_a q^*(s, a)$$

3) In model-free environments, we don't know  $p(s', r | s, a)$  beforehand. Using only  $V^*(s)$ , we cannot construct an optimal policy because we cannot compute which action leads to the best next states without the transition model. This is why Monte Carlo ~~can't be used~~ was action value functions - they enable optimal decision-making in model-free settings.