

Ex 02 RL Samuel Pitch, Jakob Grosz

1) chess $S = \{\text{every possible position}\}$

64 squares \times 13 figures (6 black, 6 white, empty) \Rightarrow 2 dim discrete

$A = \{\text{every legal move}\}$

16 figures \times different moves with each figure \Rightarrow 2 dim discrete

$R = \{\text{game ended in win/lose}\} \Rightarrow$ 1 dim discrete

b) pick & place $S = \{\text{position endeffector, position object, item grabbed}\}$

6 dim cont \times 6 dim cont \times 1 dim discrete

$A = \{\text{change position endeffector, grabbing}\}$

6 dim cont \times 1 dim discrete

$R = \{\text{MSE of object and goal}\}$

6 dim cont.

c) drone $S = \{\text{position, angles and derivatives}\}$

12 dim cont

$A = \{\text{voltage of motors}\}$

4 dim cont

$R = \{\text{MSE of current position, angles and derivatives to goalstate (0)}\}$

1 dim cont

d) vacuum robot $S = \{x, y, \varphi, \text{on/off}\}$

3 dim cont \times 1 dim discrete

$A = \{\text{acceleration, turning, switching on/off}\}$

2 dim cont \times 1 dim discrete

$R = \{\text{sum of vacuumed area + derivative of vacuumed area}\}$

1 dim cont.

2, a) We don't need to consider future rewards in the bandit setting since each action and reward are independent of earlier and later actions, so the bandits don't need a state.

$$b) \quad v_{\pi}(s) \stackrel{\text{slide 31}}{=} \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma v_{\pi}(s')] \quad \forall s \in S \quad (1)$$

$$q_{\pi}(s) \stackrel{\text{slide 33}}{=} \sum_{s', r} p(s', r | s, a) [r + \gamma v_{\pi}(s')]$$

Since this right hand side of the equation is also present in equation (1)

$$\Rightarrow v_{\pi}(s) = \sum_a \pi(a|s) q_{\pi}(s)$$

$$c) \quad v_{\pi}(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma v_{\pi}(s')]$$

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

$$\stackrel{\text{slide 27}}{=} \sum_a \pi(a|s) \sum_{s'} p(s' | s, a) r(s, a, s') + \sum_r p(s', r | s, a) \gamma v_{\pi}(s')$$

$$= \sum_a \pi(a|s) \sum_{s'} p(s' | s, a) (r(s, a, s') + \gamma v_{\pi}(s'))$$

3, a) Be $\mathcal{P} = \{\pi_i\}$
 then $|\mathcal{P}| = |A|^{|S|}$

$$b) \quad v_{\pi} = R + \gamma P_{\pi} v_{\pi}$$

$$\Leftrightarrow v_{\pi} - \gamma P_{\pi} v_{\pi} = R$$

$$\Leftrightarrow v_{\pi} (I - \gamma P_{\pi}) = R$$

$$\Leftrightarrow v_{\pi} = (I - \gamma P_{\pi})^{-1} R$$

$$3, c) \quad V^* = \begin{bmatrix} 0,498 & 0,831 & 1,311 \\ 0,536 & 0,977 & 2,235 \\ 0,306 & 0 & 5 \end{bmatrix}$$

There exist 2 optimal policies

→	→	→
↑	↑	→
←	—	—

- d) The method doesn't work well, since computation time gets very large. This solution method need small discrete state and action spaces.