



## MAS 2022-23 HW 5 MDP Solutions v2

Multi-agent systems (Vrije Universiteit Amsterdam)

HW5

HW5

1

Bellman equations for deterministic policy iteration

General

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

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

### ① Deterministic policy

Under the policy  $\pi$ , each state is mapped to unique action  $a_s$

$$\pi: s \mapsto a_s$$

Hence, summation over action collapses in singleton.

$$V_{\pi}(s) = \sum_{s'} P(s'|s,a_s) [r(s,a_s,s') + \gamma V_{\pi}(s')]$$

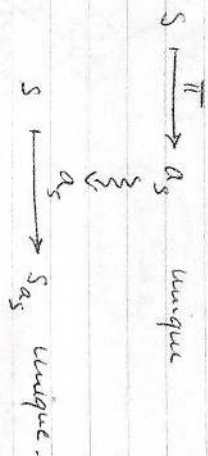
$q_{\pi}(s,a)$  = Value when taking action  $a$  in state  $s$   
(action  $a$  is arbitrary, not necessarily dictated by policy!) and THEN following policy  $\pi$ :

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

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

② Deterministic policy and transition.

We now have the following deterministic mapping:

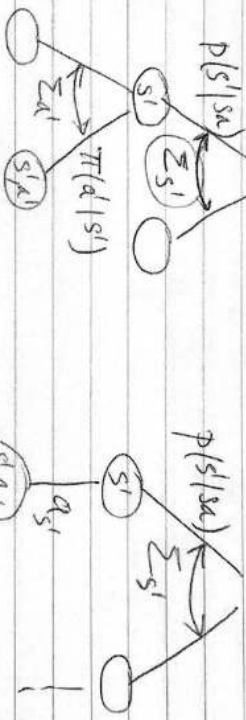
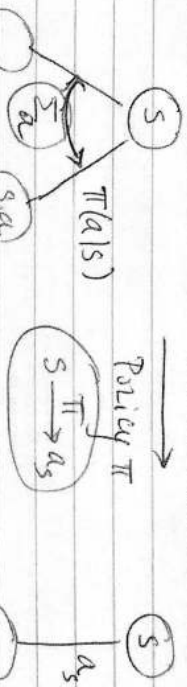


$$V_{\pi}(s) = V(s, a_s, s_{a_s}) + \gamma V_{\pi}(s_{a_s})$$

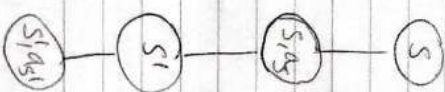
$$Q_{\pi}(s, a) = V(s, a, s_a) + \gamma Q_{\pi}(s_a, a_{s_a})$$

# Backup backups

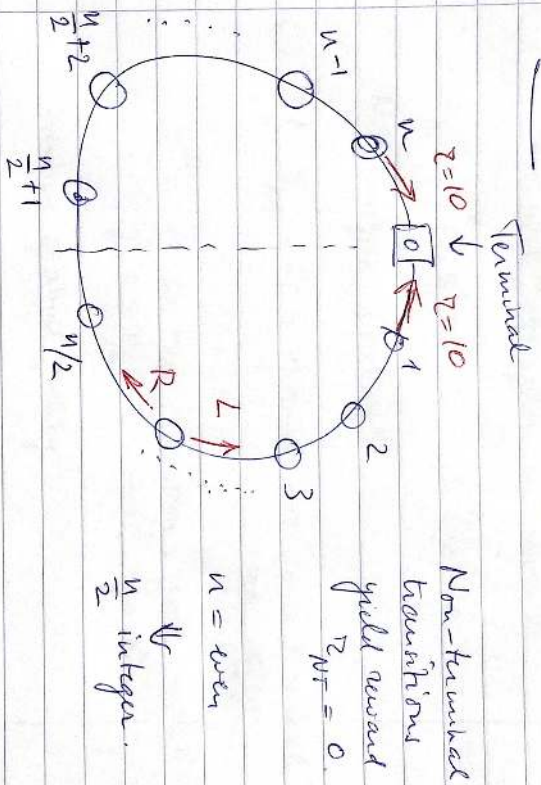
DETERMINISTIC



DETERMINISTIC  
TRANSITION



## MDP 1 Circular state space



Note 0 = absorbing: transition yields  $r=10$ .

$\gamma=1$  (no discounting)

$\pi$  = equiprobable policy: each action has prob  $\frac{1}{2}$   
two actions in each node  
move clockwise (R)  
move counterclockwise (L)

- (1) Since transitions b/w non-terminal states  
have no cost and the agent will eventually  
end up in absorbing state 0 we conclude:

$$V_{\pi}^0(s) = 10 \quad \forall s.$$

$$Q_{\pi}^0(s,a) = 10 \quad \forall s,a.$$



② Optimal policy: any policy that ensures eventual absorption in state 0.

$$q^*(s, a) = 10 \quad v^*(s) = 10.$$

Not unique. (policy).

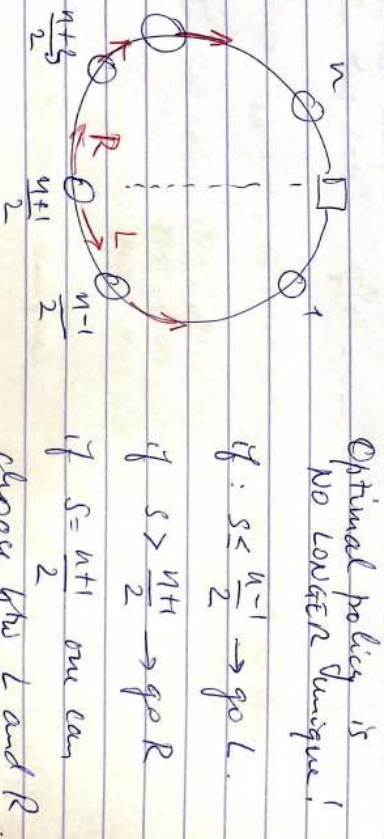
③ If  $v_{0T} = -1$ : optimal  $\rightarrow$  go to terminal state asap.

$\pi^*$ : if  $s \leq \frac{n}{2}$ : action: go L with prob = 1  
if  $s > \frac{n}{2} + 1$ : go R

unique optimal policy.

④ If  $\gamma < 1$ : go to terminal asap to be optimal.  
(similar to 3).

⑤ If  $n = \text{odd}$  ( $v_{0T} = -1, \gamma = 1$ )



HW 5

Markov decision process (MDP) MDP2



① Equi-prob. policy  $\pi$

$$V_{\pi}(2) = V_{\pi}(5) = 10 \text{ because of symmetry.}$$

Equal prob to end up in A (reward 20) and B (reward 0).

The other values can be computed using the Bellman eq:

$$V(5) = \sum_a \pi(a|5) \sum_{s'} p(s'|s,a) [v(s,a,s') + \gamma V(s')]$$

$$V(1) = \frac{1}{2}(20 + 10) = 15; \quad V(3) = \frac{1}{2}(10 + 0) = 5$$

$$V(6) = V(1), \quad V(5) = V(2), \quad V(4) = V(3)$$

② An optimal policy is any policy that avoids absorption by B.

So not unique.

In this case,  $V_{\pi}^*(1) = \dots = V_{\pi}^*(6) = 20$ .

③ Since  $V_{\text{abs}} = -1$ , optimal policy = "go to A as fast as you can".

$$\text{In that case, } V_{\pi}^*(1) = V_{\pi}^*(6) = 20$$

$$V_{\pi}^*(2) = V_{\pi}^*(5) = 19$$

$$V_{\pi}^*(3) = V_{\pi}^*(4) = 18$$

This policy is unique.

(6.3 continued)

④  $V_{NT} = -10.$

Optimal values

$$\begin{aligned} v^*(1) &= v^*(6) = 20 \\ v^*(2) &= v^*(5) = 10 \\ v^*(3) &= v^*(4) = 0 \end{aligned}$$

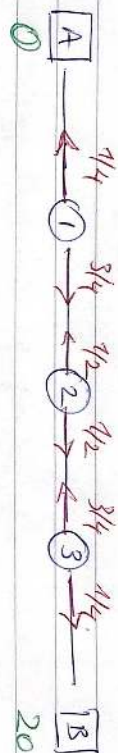
Policy is NOT unique since in 3 and 4 it does not matter which direction you choose.

7/8



# HW5, question 4: MDP3

- ① State value  $V_{\pi}(s)$  under given policy  $\pi$



Since the transitions are deterministic we can simplify the Bellman eq:

$$s \xrightarrow{a} s_a$$

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

Therefore,  $V_{\pi}(1) = V_1$ ,  $V_{\pi}(2) = V_2$ ,  $V_{\pi}(3) = V_3$

Notice:  $V_{\pi}(A) = 0 = V_{\pi}(B)$

$$V_1 = \frac{1}{4}(0+0) + \frac{3}{4}(-2+V_2) = -\frac{3}{2} + \frac{3}{4}V_2$$

$$V_2 = \frac{1}{2}(-2+V_1) + \frac{1}{2}(-2+V_3) = \frac{V_1+V_3}{2} - 2$$

$$V_3 = \frac{3}{4}(-2+V_2) + \frac{1}{4}(20+0) = \frac{3}{4}V_2 + 5 - \frac{3}{2}$$

$$= \frac{3}{4}V_2 + \frac{7}{2}$$

Summing  $V_1$  and  $V_3$ :

$$\begin{aligned} V_1 + V_3 &= \left(-\frac{3}{2} + \frac{3}{4}V_2\right) + \left(\frac{3}{4}V_2 + \frac{7}{2}\right) \\ &= \frac{3}{2}V_2 + 2 \end{aligned}$$

Substituting this into eq. for  $V_2$ :

$$\begin{aligned} V_2 &= \frac{1}{2}(V_1 + V_3) - 2 = \frac{1}{2}\left(\frac{3}{2}V_2 + 2\right) - 2 \\ &= \frac{3}{4}V_2 - 1 \end{aligned}$$

$$\Rightarrow \boxed{V_2 = -4}$$

$$\Rightarrow \boxed{V_1 = -\frac{3}{2} + \frac{3}{4}V_2 = -\frac{3}{2} + \frac{3}{4}(-4) = -\frac{9}{2}}$$

$$\Rightarrow \boxed{V_3 = \frac{3}{4}V_2 + \frac{7}{2} = \frac{3}{4}(-4) + \frac{7}{2} = \frac{1}{2}}$$

② Compute state-action value  $q_{\pi}(2,R)$  and  $q_{\pi}(3,L)$

$$q_{\pi}(2,R) = -2 + V_{\pi}(3) = -2 + \frac{1}{2} = -\frac{3}{2}$$

$$q_{\pi}(3,L) = -2 + V_{\pi}(2) = -2 + \left(-\frac{9}{2}\right) = -\frac{13}{2}$$

③ Optimal policy: go R in each state.  
unique!