

Homework 02: Review

Saturday, March 26, 2022 10:41 AM

- ② since $f_R(s_t) \rightarrow$ the reward has nothing to do with the action done or the previous state.

$$T_1 = s_1 \xrightarrow{1} s_2 \xrightarrow{-1} s_3 \xrightarrow{2} s_1 \xrightarrow{1} s_2 \xrightarrow{2} s_1$$

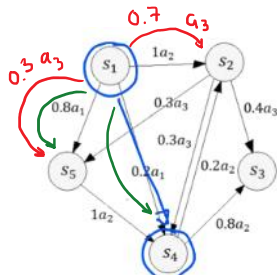
even if there were two actions that $s_1 \rightarrow s_2$, the reward is the same
thus, one path

$$f_{R_1} = 1 + \gamma(-1) + \gamma^2(2) + \gamma^3(1) + \gamma^4(2)$$

if they asked

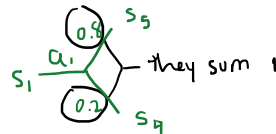
$$f_{R_1}(s_1) = V(s_1) = \text{average acc reward of all trajectories starting on } s_1$$

4.



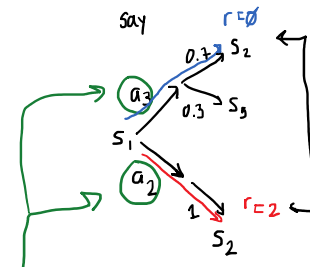
→ coincidentally, on every state there is one action that takes you to another state.

The probabilities of transition are not used since we are calculating the f_{R_1} of just one trajectory.



Let's suppose what's on red now we have two actions that take us from s_1 to s_2 (a_3 and a_2)

Thus the decision of which action to take can be decided based on a distribution. Let's say

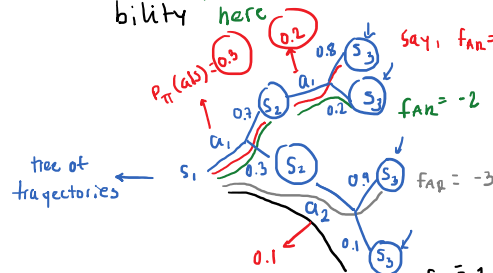


$T_1 = s_1, s_2 \rightarrow$ we have two possible rewards

Thus the policy must be stochastic

$$P_{\pi}(a|s)$$

this defines the probability here



say, $f_{R_1} = 3 \rightarrow$ this will be multiplied by the probabilities $(0.3)(0.2)$ so that it can be summed with the others, instead of doing the average like before.

$$\overline{f_{R_1}} = \underbrace{(0.3)(0.7)(0.2)(0.2)}_{\text{prob of happening}} (3) + \dots$$

Let's say

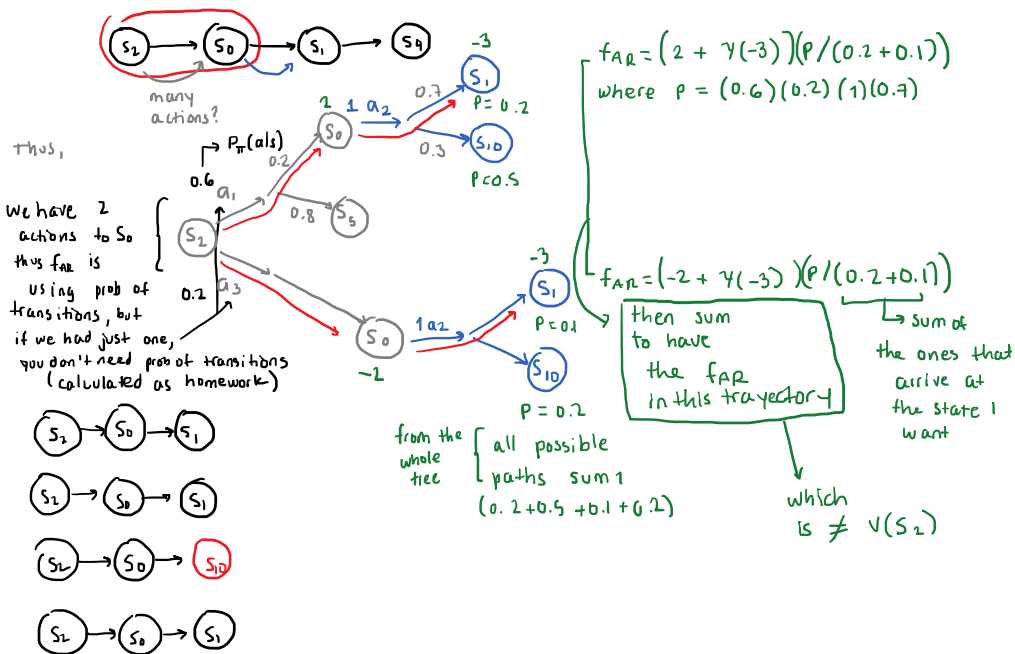
$$T_1 = s_2, s_0, s_1, s_4$$

based on state with the policy $P_{\pi}(a|s)$

Let's say

$$\mathcal{T} = \{S_2, S_0, S_1, S_4\}$$

with the policy $\pi(a|s)$



Class:

Let's use M_T non deterministic Bellman's Equation

M_T NO determinista

$$V(s) = \overline{f_{RA}(\tau)}_{s(0)=s}$$

this can be written as this

$$= \overline{f_R(s, a, s_f) + \gamma V(s_f)}$$

also like this

$$= \sum_{s_f \in \mathcal{S}} P_{M_T}(s_f | s, a) [f_R(s, a, s_f) + \gamma V(s_f)]$$

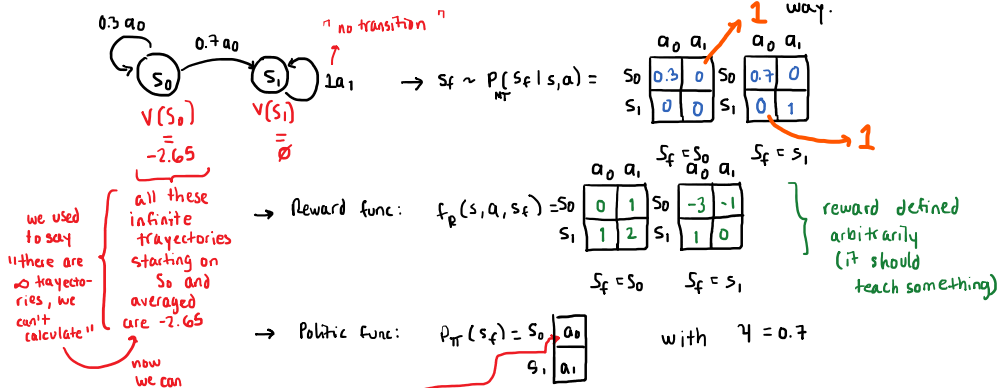
plug this

con: $a = f_\pi(s)$

which result in

$$V(s) = \sum_{s_f \in \mathcal{S}} P_{M_T}(s_f | s, f_\pi(s)) [f_R(s, f_\pi(s), s_f) + \gamma V(s_f)]$$

Bellman's equation in a general way.



Thus, there must be two $V(S) = V(S_0)$ and $V(S_1)$ (one for each state)

or, all final states (all states are final states)

$$V(S_0) = \sum_{i=0}^{\infty} P_{M_T}(S_i | S_0, a_0) [f_R(S_0, a_0, S_i) + \gamma V(S_i)]$$

$$= \underbrace{P_{M_T}(S_0 | S_0, a_0)}_{0.3} \underbrace{[f_R(S_0, a_0, S_0)]}_{0} + \underbrace{P_{M_T}(S_1 | S_0, a_0)}_{0.7} \underbrace{[f_R(S_0, a_0, S_1) + \gamma V(S_1)]}_{-3}$$

from P_{M_T} from f_R from P_{M_T}

plugging known values:

$$V(S) = 0.3 [0 + 0.7 V(S_0)] + 0.7 [(-3) + \gamma V(S_1)]$$

$$= 0.21 V(s_0) - 2.1 + 0.49 V(s_1)$$

Thus,

$$V(s_0) - 0.21 V(s_0) = -2.1 + 0.49 V(s_1)$$

Giving out the first equation of a 2x2 system

$$\boxed{0.79 V(s_0) = -2.1 + 0.49 V(s_1)} \quad (1)$$

Let's calculate eq 2

$$\begin{aligned} V(s_1) &= \sum_{i=0}^1 P_{MT}(s_i | s_1, a_1) [f_R(s_1, a_1, s_i) + \gamma V(s_i)] \\ &= \underbrace{P_{MT}(s_0 | s_1, a_1)}_{\emptyset} \underbrace{[f_R(s_1, a_1, s_0) + \gamma V(s_0)]}_{\emptyset} + \underbrace{P_{MT}(s_1 | s_1, a_1)}_1 \underbrace{[f_R(s_1, a_1, s_1) + \gamma V(s_1)]}_{\emptyset} \end{aligned}$$

$$V(s_1) = 1 [0.7 V(s_1)]$$

$$V(s_1) - 0.7 V(s_1) = 0$$

$$0.3 V(s_1) = 0$$

$$V(s_1) = \frac{0}{0.3}$$

$$\boxed{V(s_1) = 0} \quad (2) \quad \text{2nd bellman's eq}$$

Thus, the system is:

$$0.79 V(s_0) = -2.1 + 0.49 \cancel{V(s_1)}$$

$$V(s_0) = \frac{-2.1}{0.79} = -2.65$$

which we now go to the graph and complete.