

# Week 6: Homework

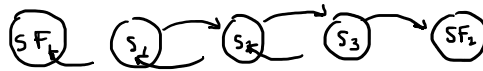
Saturday, March 12, 2022 10:40 AM

## ⑤ 1 dimensional World

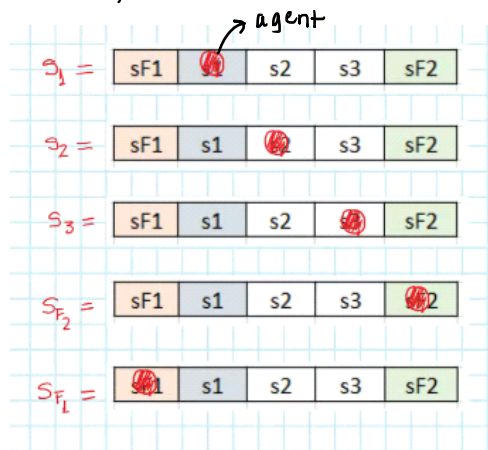
$$S = \{s_1, s_2, s_3, sF_1, sF_2\}$$

initial state  $s_1$  : we always begin here.

final states : whenever you reach a final state, the program ends.  
You cannot go out of it.



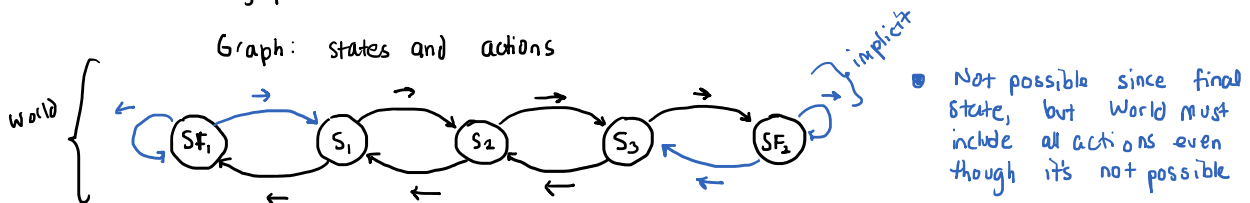
When we say "state  $s_i$ ", we mean the agent is in that position in the world, and so on:



The same happens with the reward diagram

↳ does not depend on which state you came from, only on the  $s_f$  (final state after an action)  
i.e. no matter from where I arrive to  $s_2$ , the reward is always  $-0.4$   $s_f$

a) Build the graph of the world



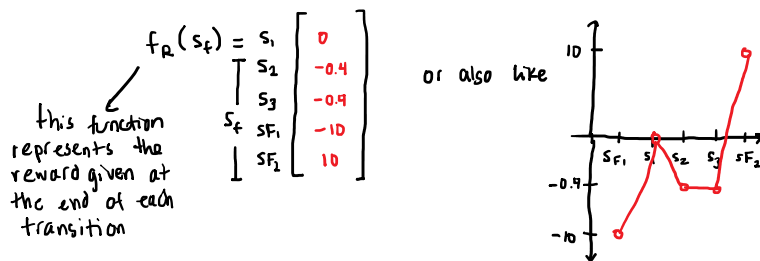
Transition model: deterministic, since given an action  $a$ , the agent always moves 1 unit to the same state.

b) Write the transition function  $f_{MT}(s, a)$

$$f_{MT}(s, a) = \begin{matrix} s_1 \\ s_2 \\ s_3 \\ sF_1 \\ sF_2 \end{matrix} \begin{matrix} \leftarrow & \rightarrow \\ \begin{matrix} sF_1 & s_2 \\ s_1 & s_3 \\ s_2 & sF_2 \\ sF_1 & s_1 \\ s_3 & sF_2 \end{matrix} \end{matrix}$$

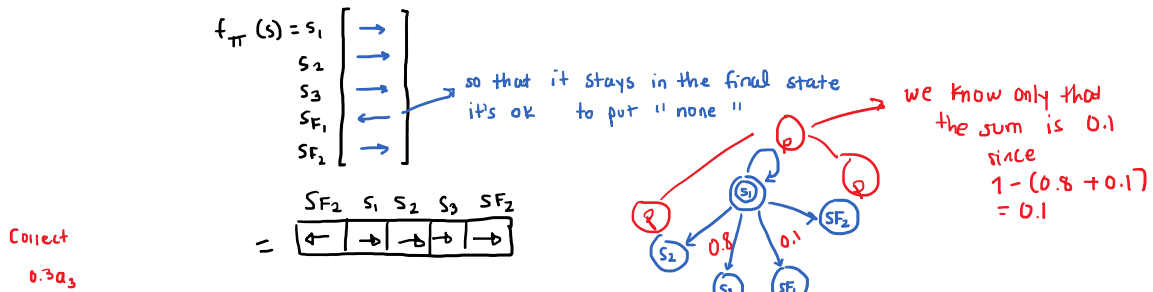
if you didn't put transition out of  $sF_1$ , this should be not defined and it's ok.

c) Write the reward function (function of one variable)

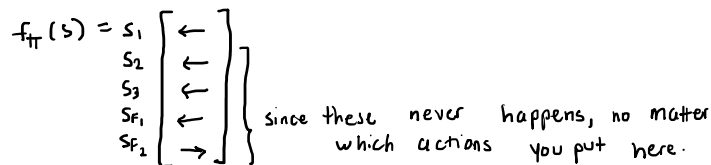


d) 2 functions of action  $f_\pi(s)$  : make two policies

I. from  $s_1$  to final state  $s_{F2}$

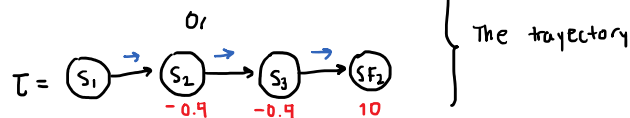


II. from  $s_1$  to final state  $s_{F1}$



Let's calculate the accumulated reward function:

$\mathcal{T} = s_1, s_2, s_3, s_{F2}$



(gamma)  $\gamma = 0.8$ , let's say.

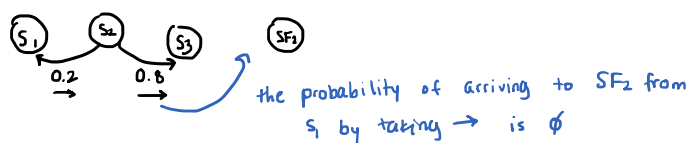
Thus, for  $\mathcal{T}$ ,

$$\begin{aligned}
 f_{AR}(\mathcal{T}) &= -0.4 + \gamma(-0.4) + \gamma^2(10) \\
 &= -0.4 + (0.8)(-0.4) + (0.8)^2(10) \\
 &= -0.4 + (-0.32) + 6.4 \\
 &= \boxed{5.68}
 \end{aligned}$$

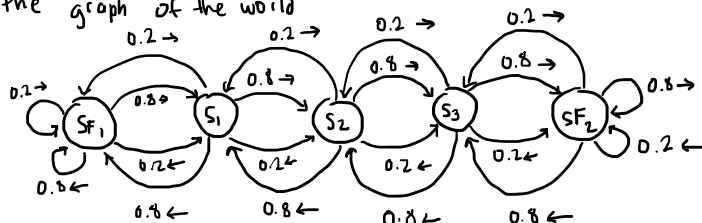
⑥

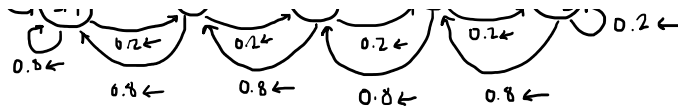
$\rightarrow$  = to the right with  $p=0.8$   
 to the left with  $p=0.2$

8 out of 10 times,  $\rightarrow$  moves to the right



a) Build the graph of the world





b) Transition function  $P_{MT}(s_f | s_i, a)$  ↖ stochastic  
↘ World  
3D matrix

$$P_{MT}(s_f | s_i, a) = \begin{matrix} & \leftarrow & \rightarrow \\ s_i & \begin{bmatrix} 0 & 0 \\ 0.8 & 0.2 \\ 0 & 0 \\ 0.2 & 0.8 \\ 0 & 0 \end{bmatrix} & s_i & \begin{bmatrix} 0.2 & 0.8 \\ 0 & 0 \\ 0.8 & 0.2 \\ 0 & 0 \\ 0 & 0 \end{bmatrix} & s_i & \begin{bmatrix} 0 & 0 \\ 0.2 & 0.8 \\ 0 & 0 \\ 0 & 0 \\ 0.8 & 0.2 \end{bmatrix} & s_i & \begin{bmatrix} 0.8 & 0.2 \\ 0 & 0 \\ 0 & 0 \\ 0.8 & 0.2 \\ 0 & 0 \end{bmatrix} & s_i & \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 0.2 & 0.8 \\ 0 & 0 \\ 0.2 & 0.8 \end{bmatrix} \\ s_f = s_1 & & s_f = s_2 & & s_f = s_3 & & s_f = SF_1 & & s_f = SF_2 \end{matrix}$$

c) Write the reward function

$$f_R(s_f) = \begin{matrix} s_1 \\ s_2 \\ s_3 \\ s_f \\ SF_1 \\ SF_2 \end{matrix} \begin{bmatrix} 0 \\ -0.4 \\ -0.9 \\ -10 \\ 10 \end{bmatrix} \quad (\text{same as 5.c})$$

d) 2 action functions  $f_{\pi}(s)$

- look for a policy that most of the time arrives to  $SF_2$   
 ↳ take the case of 0.8 prob (→)

I.  $f_{\pi}(s) = \begin{matrix} s_1 \\ s_2 \\ s_3 \\ SF_1 \\ SF_2 \end{matrix} \begin{bmatrix} \rightarrow \\ \rightarrow \\ \rightarrow \\ \rightarrow \\ \rightarrow \end{bmatrix}$  ↳ so that it stays in  $SF_1$

II.  $f_{\pi}(s) = \begin{matrix} s_1 \\ s_2 \\ s_3 \\ SF_1 \\ SF_2 \end{matrix} \begin{bmatrix} \leftarrow \\ \leftarrow \\ \leftarrow \\ \leftarrow \\ \rightarrow \end{bmatrix}$  ↳ in case of going to  $s_2$ , we go back and so on

e) Probability of arriving to  $SF_2$  and  $SF_1$  using the last 4. policies.  
 all probabilities are independent

↳ all variables are independent (decision trees)

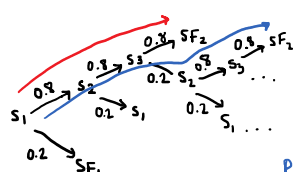
I.

$$P = (0.8)(0.8)(0.8) = 0.512$$

probability of this path only

↳ half of times this happens

the other half is all the other options in the tree



Since they are independent from history (indep. events) and so we multiply to get their 'together' distribution

$$P = (0.8)(0.8)(0.2)(0.8)(0.8) = 0.08$$

If we do not find a pattern to the convergence, just 512