

# Week 7: Bellman's Equations

Saturday, March 19, 2022 10:41 AM

Trajectory:

→ time function

→ graph

↳ includes  $r(t)$  and  $a(t)$  (reward and actions as time functions)

The graph of the world  $\neq$  the graph of a trajectory

→ if  $\gamma = 0$  we are ignoring all future rewards, and we are only considering the current transition's reward.

→  $\gamma < 1$  so that it converges

For deterministic worlds:  $r_t + \gamma [f_R(\bar{T})]$

Average Accumulated Reward  $V(s) = \overline{f_{RA}(\bar{T})}_{s(s)=s}$  } The average of the acc reward of all trajectories that start in state  $s$

→ Why average?

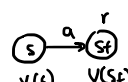
$\overline{f_{RA}(\bar{T})}$

In non-deterministic worlds, we need to calculate the average reward. From this case, we can derive the deterministic case: the average in a deterministic world will always be the same.

Since in non-deterministic worlds, trajectories can be infinite. → how do we calculate  $f_{RA}(\bar{T})$

Bellman

MT Deterministic

$V(s) = \overline{f_{RA}(\bar{T})}$    $\left. \begin{aligned} V(s) &= r + V(S_f) \\ &\quad \text{reward from that on} \\ &\quad \text{↳ immediate reward} \end{aligned} \right\} V(s) \text{ is now written in terms of } V(S_f) \text{ (next state)}$   
 $\Rightarrow f_R(s, a, S_f) + \gamma V(S_f)$

Since it is the same average because it is deterministic

$= f_R(s, a, f_{MT}(s, a)) + \gamma V(f_{MT}(s, a))$   
 $f_{MT} = \text{Transition Model function: } S_f = f_{MT}(s, a)$

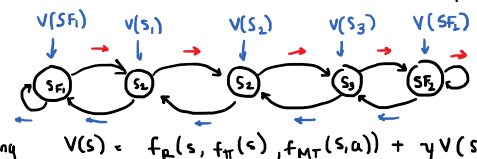
If we consider  $a = f_\pi(s)$

$V(s) = f_R(s, f_\pi(s), f_{MT}(s, f_\pi(s))) + \gamma V(f_{MT}(s, f_\pi(s)))$   
 Which is now a function only depending  $S_f$  on  $s$ .

We have, in a world with  $N$  states, we have  $N$  equations in the form:

$V(s)$  = written in  $V(S_f) \rightarrow$  Bellman's eq

Thus we have a linear system of equations with  $V(S_f)$  as the unknown.

Example:   $f_\pi(s) = s_1$   $f_R(S_f) = \begin{bmatrix} Sf_1 & -10 \\ S_1 & 0 \\ S_2 & -0.4 \\ S_3 & -0.4 \\ Sf_2 & 10 \end{bmatrix}$   
 Considering  $V(s) = f_R(s, f_\pi(s), f_{MT}(s, a)) + \gamma V(S_f)$   
 ↳ immediate reward after the transition

(given by politic)

$$V(S_1) = -0.4 + \gamma V(S_2)$$

Now we write  $V(S_2)$  as reward function from  $S_2$  on, according to the politic:

$$\begin{aligned} (2) \quad V(S_2) &= -0.4 + \gamma V(S_3) \\ (3) \quad V(S_3) &= 10 + \gamma V(SF_2) \\ (4) \quad (SF_1) \quad V(SF_1) &= -10 + \gamma V(SF_1) \\ (5) \quad V(SF_2) &= 10 + \gamma V(SF_2) \end{aligned} \quad \left. \begin{array}{l} \\ \\ \\ \end{array} \right\} \begin{array}{l} \text{all these are calculated} \\ \text{for getting } f_{AA} \end{array}$$

with  $\gamma = 0.9$

Thus we have a linear system of equations (5 eq)

$$\begin{aligned} (5) \quad V(SF_2) &= 10 + \gamma V(SF_2) \\ V(SF_2) - \gamma V(SF_2) &= 10 \\ (1 - \gamma) V(SF_2) &= 10 \\ V(SF_2) &= 10 \\ V(SF_2) &= \frac{10}{1 - \gamma} = \frac{10}{1 - 0.9} = \frac{10}{0.1} = \boxed{100} \end{aligned}$$

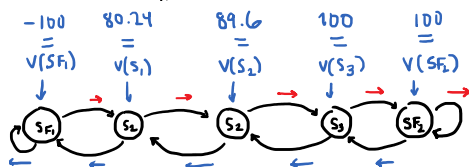
$$\begin{aligned} (4) \quad V(SF_1) &= -10 + \gamma V(SF_1) \\ [1 - \gamma] V(SF_1) &= -10 \\ V(SF_1) &= \frac{-10}{1 - \gamma} = \frac{-10}{1 - 0.9} = \frac{-10}{0.1} = \boxed{-100} \end{aligned}$$

$$\begin{aligned} (3) \quad V(S_3) &= 10 + \gamma V(SF_2) \\ V(S_3) &= 10 + (0.9)(100) \\ V(S_3) &= 10 + 90 \\ \boxed{V(S_3) = 100} &\rightarrow \text{you can use the geometric series to solve this and get the same 100} \end{aligned}$$

$$\begin{aligned} (2) \quad V(S_2) &= -0.4 + \gamma V(S_3) & (1) \quad V(S_1) &= -0.4 + \gamma V(S_2) \\ V(S_2) &= -0.4 + (0.9)(100) & V(S_1) &= -0.4 + (0.9)(89.6) \\ V(S_2) &= -0.4 + 90 & V(S_1) &= \boxed{80.24} \\ \boxed{V(S_2) = 89.6} & & & \end{aligned}$$

exam: April 2nd

Now, we can say,



The politic originally is to begin in  $S_1$  and move until  $SF_2$ , and with the results given, the agent coincidentally has to move to the right (biggest reward)

Now, for the nondeterministic world, where the trajectories can be infinite. Their average value converges, but how to find it?

Bellman's Equation, nondeterministic

$$V(s) = \overline{f_{AA}(T)} \Big|_{s(0)=s} \quad \left. \begin{array}{l} \text{given that} \\ \text{the average } f_{AA} \text{ of all trajectories} \\ \text{that start in start } s \end{array} \right\}$$

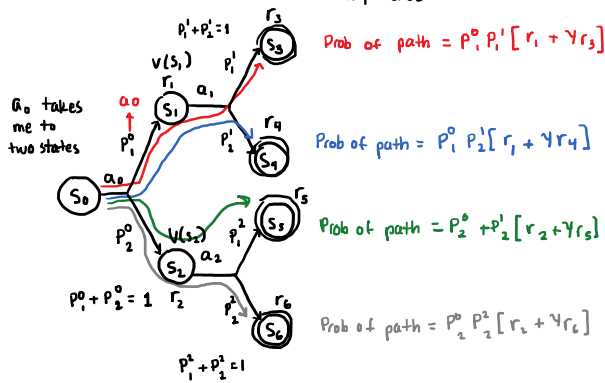
$$\begin{aligned} \text{Using the definition of the sum} \quad & \text{this sum is } = \overline{f_R + \gamma V(S_t)} \\ & = \overline{f_R(s, a, S_t) + \gamma V(S_t)} = \sum_{S_t \in S} p_{NT}(S_t | s, a) [f_R(s, a, S_t) + \gamma V(S_t)] \end{aligned}$$

with  $a = f_{\pi}(s)$

Example

Let's imagine we start in  $S_0$   
Assume in final states, it ends  
r.

Let's imagine we start in  $S_0$   
Assume in final states, it ends



We got  
a total of  
4 trajectories  
to final states

$$V(S_0) = p_1^0 p_1^1 [r_1 + \gamma r_3] + p_1^0 p_2^1 [r_1 + \gamma r_4] + p_2^0 p_1^2 [r_2 + \gamma r_5] + p_2^0 p_2^2 [r_2 + \gamma r_6]$$

Let's write this in Bellman's form:

$$\begin{aligned} &= p_1^0 [p_1^1 [r_1 + \gamma r_3] + p_2^1 [r_1 + \gamma r_4]] + p_2^0 [p_1^2 [r_2 + \gamma r_5] + p_2^2 [r_2 + \gamma r_6]] \\ &= p_1^0 [\underbrace{(p_1^1 + p_2^1)}_1 r_1 + \gamma (p_1^1 r_3 + p_2^1 r_4)] + p_2^0 [\underbrace{(p_1^2 + p_2^2)}_1 r_2 + \gamma (p_1^2 r_5 + p_2^2 r_6)] \\ &\quad \underbrace{V(S_1)}_{\text{all reward from } S_1 \text{ till the end}} \\ &= p_1^0 r_1 + \gamma p_1^0 (p_1^1 r_3 + p_2^1 r_4) + p_2^0 r_2 + \gamma p_2^0 (p_1^2 r_5 + p_2^2 r_6) \\ &\quad \underbrace{V(S_2)}_{\text{all reward from } S_2 \text{ till the end}} \end{aligned}$$

$$\begin{aligned} &= p_1^0 r_1 + \gamma p_1^0 V(S_1) + p_2^0 r_2 + \gamma p_2^0 V(S_2) \\ &= p_1^0 [r_1 + \gamma V(S_1)] + p_2^0 [r_2 + \gamma V(S_2)] \end{aligned}$$

Keep in mind:

$$\begin{aligned} S_f &\sim P_{\text{MT}}(S_f | S_i, a) \\ p_i^a &= P_{\text{MT}}(S_i | S_0, a_0) \\ p_2^a &= P_{\text{MT}}(S_2 | S_0, a_0) \end{aligned}$$

$$\begin{aligned} &= P_{\text{MT}}(S_1 | S_0, a_0) [r_1 + \gamma V(S_1)] + P_{\text{MT}}(S_2 | S_0, a_0) [r_2 + \gamma V(S_2)] \end{aligned}$$

In sum form, we arrive at a Bellman's form,

$$V(S_0) = \sum_{i=1}^2 P_{\text{MT}}(S_i | S_0, a_0) [r_i + \gamma V(S_i)]$$

where  $a$ 's are given by a policy  
We have an equation per state, and this forms a linear system of equations where  $V$  is the unknown

Two unknowns,

$a = \pi(s)$  is the policy we want the agent to learn, and thus in real cases we need to solve this through approximations (Fixed Point Iteration)