

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 $V(s) = \overline{f_{RA}(\bar{T})}_{s(0)=s}$ } The average of the acc reward of all trajectories that start in state s

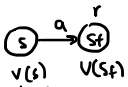
→ Why average? In non deterministic worlds, we need to calculate the average reward. From this case, we can derive the deterministic case: the average in a deterministic world will always be the same.

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

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

Bellman

MT Deterministic

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

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:  $f_\pi(s) = s_2 \rightarrow$
 $s_3 \rightarrow$
 $s_f \rightarrow$

Considering $V(s) = f_R(s, f_\pi(s), f_{MT}(s, a)) + \gamma V(s_f)$
 $\quad \quad \quad \text{↳ immediate reward after the transition (given by policy)}$

$V(s_1) = -0.4 + \gamma V(s_2)$

Now we write $V(s_2)$ as reward function from s_2 on, according to the policy:

$$\left. \begin{array}{l} (2) \quad V(s_2) = -0.4 + \gamma V(s_3) \\ (3) \quad V(s_3) = 10 + \gamma V(s_{F_2}) \\ (4) \quad V(s_{F_1}) = -10 + \gamma V(s_{F_2}) \\ (5) \quad V(s_{F_2}) = 10 + \gamma V(s_{F_2}) \end{array} \right\} \begin{array}{l} \text{all these are calculated} \\ \text{for getting } f_{\pi} \end{array}$$

with $\gamma = 0.9$

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

$$\begin{aligned} (5) \quad V(s_{F_2}) &= 10 + \gamma V(s_{F_2}) \\ V(s_{F_2}) - \gamma V(s_{F_2}) &= 10 \\ (1 - \gamma) V(s_{F_2}) &= 10 \\ V(s_{F_2}) &= 10 \\ V(s_{F_2}) &= \frac{10}{1 - \gamma} = \frac{10}{1 - 0.9} = \frac{10}{0.1} = 100 \end{aligned}$$

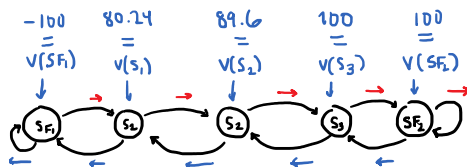
$$\begin{aligned} (4) \quad V(s_{F_1}) &= -10 + \gamma V(s_{F_2}) \\ [1 - \gamma] V(s_{F_1}) &= -10 \\ V(s_{F_1}) &= \frac{-10}{1 - \gamma} = \frac{-10}{1 - 0.9} = \frac{-10}{0.1} = -100 \end{aligned}$$

$$\begin{aligned} (3) \quad V(s_3) &= 10 + \gamma V(s_{F_2}) \\ V(s_3) &= 10 + (0.9)(100) \\ V(s_3) &= 10 + 90 \\ V(s_3) &= 100 \rightarrow \text{you can use the geometric series to solve this and get the same 100} \end{aligned}$$

$$\begin{aligned} (2) \quad V(s_2) &= -0.4 + \gamma V(s_3) \\ V(s_2) &= -0.4 + (0.9)(100) \\ V(s_2) &= -0.4 + 90 \\ V(s_2) &= 89.6 \end{aligned} \quad \begin{aligned} (1) \quad V(s_1) &= -0.4 + \gamma V(s_2) \\ V(s_1) &= -0.4 + (0.9)(89.6) \\ V(s_1) &= 80.24 \end{aligned}$$

exam: April 2nd

Now, we can say,



The policy originally is to begin in s_1 and move until s_{F_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

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

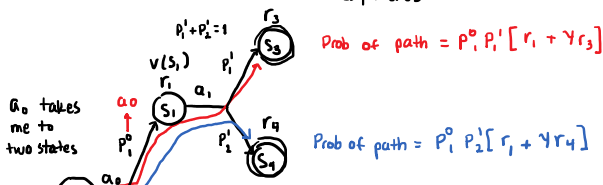
Using the definition of the sum this sum is $= f_{\pi} + \gamma V(s')$

$$= \overline{f_{\pi}(s, a, s')} + \gamma V(s') = \sum_{s' \in S} p_{\pi}(s' | s, a) [f_{\pi}(s, a, s') + \gamma V(s')]$$

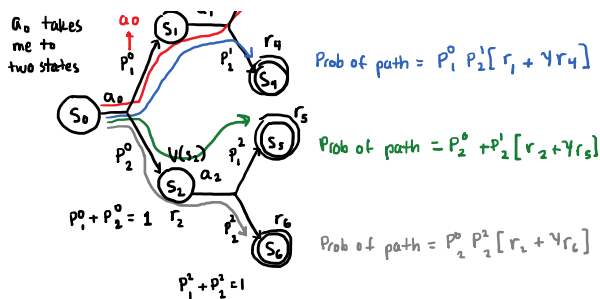
with $a = f_{\pi}(s)$

Example

Let's imagine we start in s_0
Assume in final states, it ends



We got a total of 100



We got
a total of
4 trajectories
to final states

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

Let's write this in Bellman's form:

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

$$= p_1^0 r_1 + \gamma p_1^0 V(S_1) + p_2^0 r_2 + \gamma p_2^0 V(S_2)$$

$$= p_1^0 [r_1 + \gamma V(S_1)] + p_2^0 [r_2 + \gamma V(S_2)]$$

$$= p_{MT}(S_1 | S_0, a_0) [r_1 + \gamma V(S_1)] + p_{MT}(S_2 | S_0, a_0) [r_2 + \gamma V(S_2)]$$

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

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

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

keep in mind:

- $S_f \sim p_{MT}(S_f | s, a)$
- $p_1^0 = p_{MT}(S_1 | S_0, a_0)$
- $p_2^0 = p_{MT}(S_2 | S_0, a_0)$

for this world, we would need 6 eqs at this

Two unknowns,

V and a