

Week 9:

Saturday, April 2, 2022 9:57 AM

Solución: Iteración de Valor	$V(s) \leftarrow f_R(s, f_R(s), s_f) + \gamma V(s_f)$	$V(s) \leftarrow \sum_{s_f \in S} P_{s_f s} (s_f s, f_R(s)) [f_R(s, f_R(s), s_f) + \gamma V(s_f)]$
------------------------------	---	--

$$V(s) = \overline{f_R(T)} \Big|_{s(0)=s}$$

we said we can write this function in a nondeterministic way, and also in a det form. Both cases lead us to a linear system of equations

↳ we can solve this either by an analytical way or a numerical way, which is the method: Value Iteration

We can apply this num method to either deterministic and non det.

We said $V(s)$ is the average of all reward of all trajectories starting on s .
↳ all this is done assuming we have a defined politic
but how to define the politic?

↳ the principle:

Obtain the maximum average reward



Bellman's idea:

Bellman's Optimality Equation

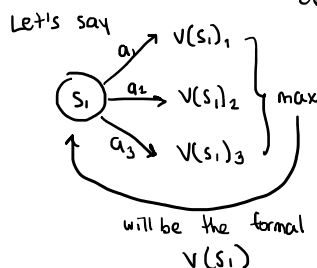
How?

Instead of $V(s) = r + \gamma V(s_f)$ we thus

$$V(s) = \max_a [f_R(s, a, s_f) + \gamma V(s_f)]$$

new det of V max with the action this a iterates and becomes each and every action

} the maximum average of the accumulated reward (previous $V(s)$) for all available actions (s in $V(s)$ for a would be the state that has the available actions)



Thus, by obtaining the $\max V(s)$, the agent discovers the optimal politic.
↳ for all states

This $V(s) = \max_a [\text{Bellman}]$ creates a different equation system: it becomes

a non-linear system of equations, since $\max()$ is not continuous. Therefore, to solve this nonlinear system we use a numerical method.

In literature, to avoid this, people defined:

$$V(s) = \max_a [f_R(s, a, s_f) + \gamma V(s_f)] \quad (1) \Rightarrow \begin{array}{c} \text{Q}(s, a) \quad \text{Q}(s_f, a_f) \\ \text{Q}(s, a) \quad \text{Q}(s_f, a_f) \\ \text{Q}(s, a) \quad \text{Q}(s_f, a_f) \end{array}$$

considering $V(s_f) = \max_a [f_R(s_f, a_f, s_{f2}) + \gamma V(s_{f2})]$ state after s_f

so, by defining $V(s) = \max_a [Q(s, a)]$, thus we rewrite $V(s)$ as $Q(s, a)$

Considering $V(s_f) = \max_a [f_R(s_f, a_f, s_{f2}) + \gamma V(s_{f2})]$
state after s_f

so, by defining $V(s) = \max_a [Q(s,a)]$, thus we rewrite $V(s)$ as $Q(s,a)$

$$Q(s,a) = [f_R(s,a,s_f) + \gamma \max_a [Q(s_f,a_f)]] \quad (2)$$

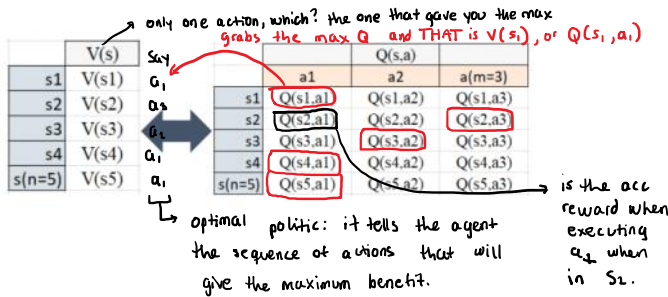
(1) and (2) are two ways of visualizing the same thing

The relationship between V and Q is

$$V(s) = \max_a [Q(s,a)]$$

one two variables

In theory, $V(s,a)$, but we don't write it, since to $V, 'a'$ doesn't matter



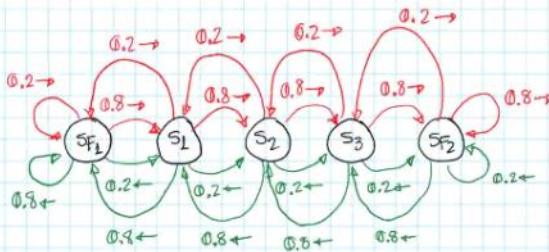
in order to know this politic, we need to calculate the whole matrix for all a 's in order to know the optimal a .

Thus, the problem is basically solving the nonlinear systems $V(s)$ or $Q(s,a)$ and the num method is once again the value iteration.

Example

s_f	s_1	s_2	s_3	s_f
-------	-------	-------	-------	-------

a. Construya el grafo del mundo



b. Escriba la función de transición $P_{MT}(s_f | s, a)$

$P_{MT}(s_f | s, a) =$

	$s_f = s_1$	$s_f = s_2$	$s_f = s_3$	$s_f = s_{f1}$	$s_f = s_{f2}$
s1	0.2	0.8	0	0.2	0
s2	0.8	0.2	0	0	0
s3	0	0	0.2	0	0.8
s_f1	0.2	0.8	0	0.8	0.2
s_f2	0	0	0	0	0.2

$$V(s) = \max_a [f_R(s, a, s_f) + \gamma V(s_f)] \quad \text{Bellman's Optimality Eq.}$$

First Bellman's Eq for optimal politic:

$$V(s_1) = \max_a [f_R(s_1, \leftarrow, s_f) + \gamma V(s_f), f_R(s_1, \rightarrow, s_f) + \gamma V(s_f)]$$

stochastic: s_f is either one that is s_f iterate each possible action

$$= \max_a \left[\sum_{s_f \in S} P_{MT}(s_f | s_1, \leftarrow) [f_R(s_1, \leftarrow, s_f) + \gamma V(s_f)], \right. \\ \left. \sum_{s_f \in S} P_{MT}(s_f | s_1, \rightarrow) [f_R(s_1, \rightarrow, s_f) + \gamma V(s_f)] \right]$$

action: \leftarrow (5 terms in the sum)

$$= \max_a \left[P_{MT}(s_2 | s_1, \leftarrow) [f_R(s_1, \leftarrow, s_2) + \gamma V(s_2)] + P_{MT}(s_{f1} | s_1, \leftarrow) [f_R(s_1, \leftarrow, s_{f1}) + \gamma V(s_{f1})], \right. \\ \left. P_{MT}(s_2 | s_1, \rightarrow) [f_R(s_1, \rightarrow, s_2) + \gamma V(s_2)] + P_{MT}(s_{f1} | s_1, \rightarrow) [f_R(s_1, \rightarrow, s_{f1}) + \gamma V(s_{f1})] \right]$$

$$= \max_a [0.2[0.4 + 4V(s_2)] + 0.8[-10 + 4V(s_{F_1})], 0.8[-0.4 + 4V(s_2)] + 0.2[-10 + 4V(s_{F_2})]] \quad (1)$$

↳ we will have 5 equations like this

Apart from the numerical method (value iteration), there is another method: Temporary Difference