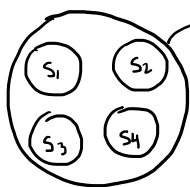


# Value Iteration: Frozen Lake

Saturday, April 30, 2022 9:58 AM



trajectory  
if we got  $s = \{s_1, s_2, s_3, s_4\}$  in our world,  
all possible trajectories are  $v(s_1), v(s_2)$   
 $v(s_3)$  and  $v(s_4)$ . There can't be more  
than that.

In stochastic and deterministic action functions, the optimal trajectory is the same, just one. In stochastic

the 100 tries build  $P_{MT}$  and  $f_R$ , then in each  $i$ th try,  
represent each Value Iteration for solving Bellman's equations.

0.85  $\rightarrow$  85% percent of playing, it reaches the goal.

## Value Iteration

We saw previously how to solve Bellman's Optimality Equations using the numerical method of Value Iteration. These equations were defined either for  $V(s)$  or  $Q(s, a)$ :

$$V(s) \leftarrow \max_a \left[ \sum_{s_f \in S} P_{MT}(s_f | s, a) \underbrace{[f_R(s, a, s_f) + \gamma V(s_f)]}_{\text{known}} \right]$$

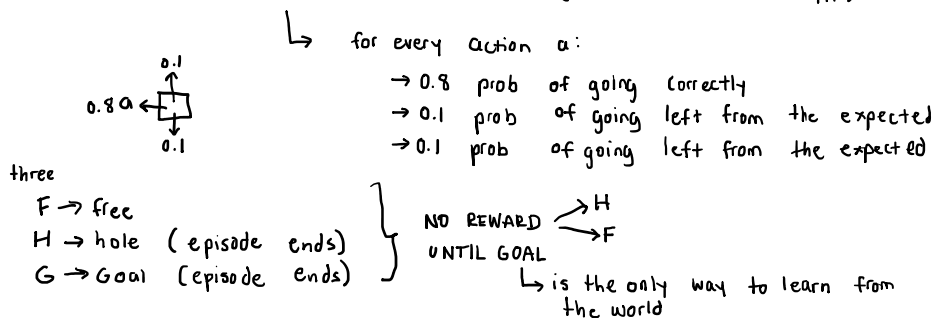
$$Q(s, a) \leftarrow \sum_{s_f \in S} P_{MT}(s_f | s, a) \underbrace{[f_R(s, a, s_f) + \gamma \max_{a_f} Q(s_f, a_f)]}_{\text{known}}$$

By 'known' we mean that either those were given explicitly or we got them by making the agent experience the world randomly and calculating them.

Therefore, for solving the equations using this method, we need to know the Transition Model and the Reward Function.

To define them, we sent the agent to experience the world randomly 100 times in order to LEARN THE TRANSITION MODEL and REWARD FUNCTION, since we need to know them for solving the Eqs using Value Iteration.

Frozen Lake World: Stochastic result of actions (transition model is  $P_{MT}$ )

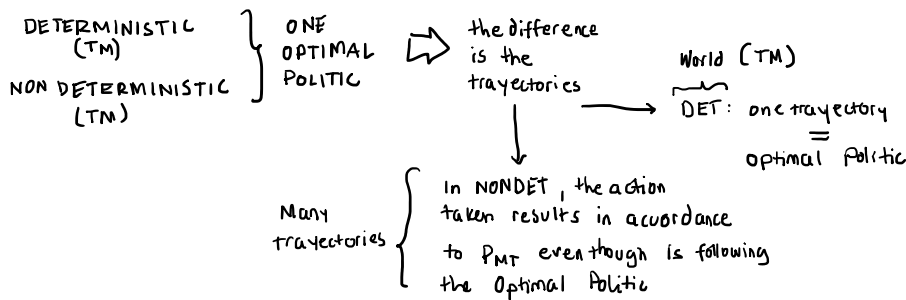


The whole point is to iterate to solve Bellman's Equations

$\rightarrow$  the info of the reward is implicit in the Bellman Equations we are trying to solve with Value Iteration

$\rightarrow$  If the transition model function is deterministic, the actions in the Optimal Policy are always the same. That is, once we solve Bellman's System of Eqs and thus we get the optimal policy, if we test

→ If the transition model function is deterministic, the actions in the Optimal Policy are always the same. That is, once we solve Bellman's System of Eqs and thus we get the optimal policy, if we test the agent using that policy, every test we will see the same trajectory. But, since frozen lake is stochastic, every test using the optimal Policy will be always different: the Policy is ONE, but transition may diverge.



→ In frozen lake program: what is the agent truly experimenting?

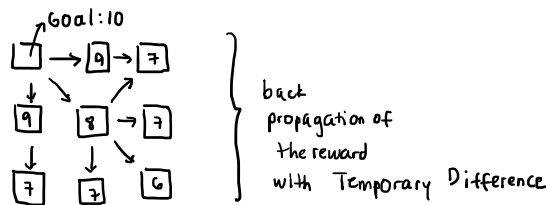
↳ the agent moves randomly to know/learn the Transition Model (1) and Reward Function (2)

here, the agent is not really experiencing the world to do this, it's just iterating

(1) and (2) are needed for solving the Bellman's Eq System using value iteration

implicit info about reward in Goal

Temporary Difference Method: in value iteration, the agent doesn't arrive at the goal and thus solves the policy: the agent knows  $f_R(s, a, s')$  and with that solves Bellman's Equations which output the policy. Meanwhile, in temporary Difference, the agent does need to arrive to the goal to know there is reward and it has to propagate it backwards to the states near the goal as if it was a ripple.



In [1]: runfile('C:/Users/José Ab (1).py', wdir='C:/Users/José Abdó')

each i is one value iteration

with the update of  $V(s)$ 's we are getting avg reward of 0.0 from the 20 episodes

In this case, the Bellman Eqs are solved in 19 iterations is like one excel iteration which updates all  $V(s)$ 's

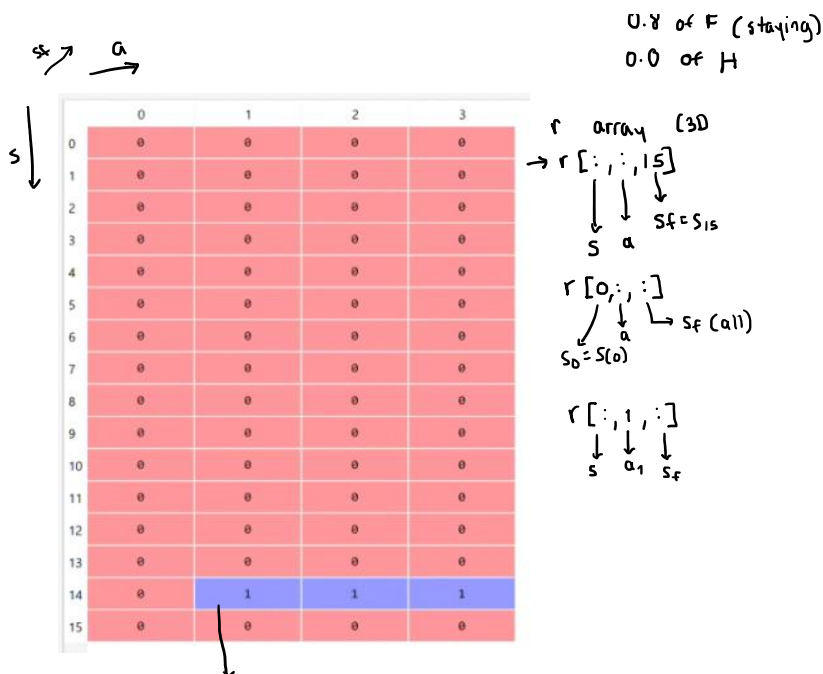
from the 20 episodes, we have avg reward = 0.1

the value iteration is getting close to the solution of the Equations.

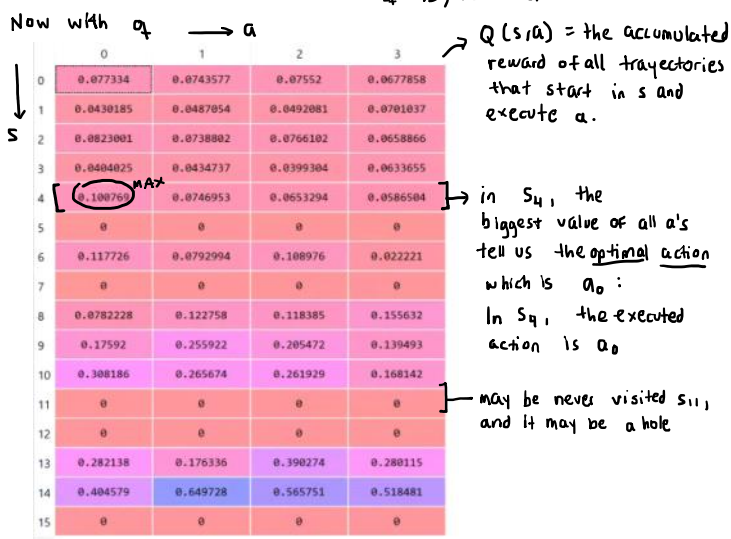
→ the 20 episodes: check whether or not  $V(s)$  is the Final solution to Eqs by computing average reward.

→ the 100 random actions: experience world to calculate and adjust PMT and  $f_R$ : needed for Eqs.

notes Page 3



$s = 14$  with  $a_1, a_2$  or  $a_3$   
 and we arrive at  $s_f = 15$ , reward of 1



Thus:

- $\rightarrow V(s)$ : less memory (vector) but slower since it needs to compute  $a$ .
- $\rightarrow Q(s, a)$ : more memory (matrix) but faster since it already has  $a$ .