

GRIDWORLD

Jarod DURET
Jonathan HENO

19 janvier 2021

1 Architecture

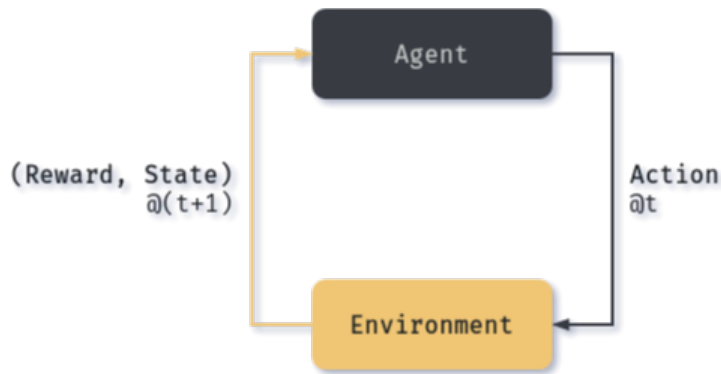


FIGURE 1 – Interactions between the decision maker and the environment

2 ϵ -greedy policy

There is a chance (valued at $\epsilon \in [0; 1]$) that an action initially chosen by the agent will not take effect. In our code we call it the **disobey factor**.

$$p_{\pi}(s, a) \triangleq \begin{cases} 1 - \epsilon & \text{if } a \in \pi(s) \\ \epsilon / (\text{Card}(A) - 1) & \text{otherwise} \end{cases} \quad (2.1)$$

It helps weighting the decision process given the fact that an action can fail and lead to another random decision, and thus to an unplanned reward. It also helps the agent to widen its field of view, to explore a wider space.

3 Markov Decision Process

The rewards progressively propagate accross the world and the optimal policy is questioned at each time step, by gauging every possible actions for each state.

INPUT : $S(\{i|i = 0..N^2\})$, $A(\{N, E, S, W\})$, $T(S, A) \in S$, p_π , γ , ϵ , δ
 OUTPUT : $\pi^*(S)$, $v_\pi^*(S)$

Initialization

$\pi_0^*(S) \leftarrow (rand(A))_{s \in S}$

$v_{\pi^*,0}^*(S) \leftarrow \mathbf{0}$

FOR $t = 1$ TO t_{obs} DO

Policy evaluation

$v_{\pi^*,t}(S) \leftarrow \left(\sum_{s' \in T(s, A(s))} p_{\pi_{t-1}^*}(s, a) \left(r(s') + \gamma v_{\pi^*,t-1}^*(s') \right) \right)_{a \in \pi_{t-1}^*(S)}$

Policy improvement

We value all actions available for each state $s \in S$

$\pi_t^*(S) \leftarrow \left(\underset{a \in A}{argmax} v_{a,t}(s) \right)_{s \in S}$

$v_{\pi^*,t}^*(S) \leftarrow \left(\underset{a \in A}{max} v_{a,t}(s) \right)_{s \in S}$

IF $\max_{s \in S} (v_{\pi^*,t}^*(s) - v_{\pi^*,t-1}^*(s)) \leq \delta$ THEN

RETURN $\pi_t^*(S)$, $v_{\pi^*,t}^*(S)$

END IF

END FOR

RETURN $\pi_{t_{obs}}^*(S)$, $v_{\pi^*,t_{obs}}^*(S)$

The other methods rely more on the policy than the previous one. At each time step, for a given epoch, each action taken, with respect to current policy, is weighted and learned through environment feedbacks.

4 State Action Reward State Action (SARSA)

The key formulae of this algorithm is the following equation :

$$Q_{s_t, a_t} = Q_{s_t, a_t} + \alpha(R_{s_t, a_t} + \gamma Q_{s_{t+1}, a_{t+1}} - Q_{s_t, a_t}) \quad (4.1)$$

With this in mind, in every iteration, we should generate two actions, with respect to the ongoing policy, and get the arrival state of both actions. We then update the Q matrix with the reward associated with the first action taken in the process : the optimal state is built upon the action value function, by mapping the optimal action given a fixed state.

INPUT : $S(\{i|i = 0..N^2\})$, $s_0 \in S$, $A(\{N, E, S, W\})$, $T(S, A) \in S$, p_π , γ , ϵ , α , $n_{u, max}$
 OUTPUT : $\pi^*(S), Q(S)$

Initialization

$Q^{(0)}, \pi_0^*(S) \leftarrow (0)_{a \in A, s \in S}, (rand(A))_{s \in S}$

$s_0, n_{unchanged} \leftarrow s_0, 0$

FOR $t = 0$ TO $(t_{obs} - 1)$ DO

Generating first action with respect to policy π_t^* at state s_t

$a_t, s' \leftarrow \text{gen_move}(s_t, \pi_t^*)$

Generating second action with respect to ongoing policy π_t^* from state s_{t+1}

$a_{t+1}, s'' \leftarrow \text{gen_move}(s', \pi_t^*)$

Correcting model given the actions taken

$Q_{s_t, a_t}^{(t+1)} \leftarrow Q_{s_t, a_t}^{(t)} + \alpha(R_{s_t, a_t} + \gamma Q_{s', a_{t+1}}^{(t)} - Q_{s_t, a_t}^{(t)})$

Improving policy given the update of the Q -values for state s_t

$\pi_{t+1}^*(s_t) \leftarrow \underset{a \in A(s_t)}{\text{argmax}} Q_{s_t, a}(s)$

Updating next state for next iteration

$s_{t+1} \leftarrow s'$

IF $\pi_{t+1}^* = \pi_t^*$ THEN

$n_{unchanged} \leftarrow n_{unchanged} + 1$

IF $n_{unchanged} = n_{u, max}$ THEN

RETURN $\pi_{t+1}^*(S), Q^{(t+1)}$

END IF

ELSE

$n_{unchanged} \leftarrow 0$

END IF

END FOR

RETURN $\pi_{t_{obs}}^*(S), Q^{(t_{obs})}$

And voilà !

5 Q-learning

The logic is more or less the same than the previous, except the fact that we are valuing the best action value function in the landing state to optimize the model :

$$Q_{s_t, a_t} = Q_{s_t, a_t} + \alpha \left(R_{s_t, a_t} + \gamma \left(\max_{a \in A(s_{t+1})} Q_{s_{t+1}, a} \right) - Q_{s_t, a_t} \right) \quad (5.1)$$

INPUT : $S(\{i|i = 0..N^2\})$, $s_0 \in S$, $A(\{N, E, S, W\})$, $T(S, A) \in S$, p_π , γ , ϵ , α , $n_{u, max}$
 OUTPUT : $\pi^*(S), Q(S)$

Initialization

$Q^{(0)}, \pi_0^*(S) \leftarrow (0)_{a \in A, s \in S}, (rand(A))_{s \in S}$

$s_0, n_{unchanged} \leftarrow s_0, 0$

FOR $t = 0$ TO $(t_{obs} - 1)$ DO

Generating first action with respect to policy π_t^* at state s_t

$a_t, s' \leftarrow \text{gen_move}(s_t, \pi_t^*)$

Correcting model given the actions taken

$Q_{s_t, a_t}^{(t+1)} \leftarrow Q_{s_t, a_t}^{(t)} + \alpha \left(R_{s_t, a_t} + \gamma \left(\max_{a \in A(s')} Q_{s', a}^{(t)} \right) - Q_{s_t, a_t}^{(t)} \right)$

Improving policy given the update of the Q -values for state s_t

$\pi_{t+1}^*(s_t) \leftarrow \underset{a \in A(s_t)}{\text{argmax}} Q_{s_t, a}(s)$

Updating next state for next iteration

$s_{t+1} \leftarrow s'$

IF $\pi_{t+1}^* = \pi_t^*$ THEN

$n_{unchanged} \leftarrow n_{unchanged} + 1$

IF $n_{unchanged} = n_{u, max}$ THEN

RETURN $\pi_{t+1}^*(S), Q^{(t+1)}$

END IF

ELSE

$n_{unchanged} \leftarrow 0$

END IF

END FOR

RETURN $\pi_{t_{obs}}^*(S), Q^{(t_{obs})}$