
Algorithm 1: The DQN algorithm with experience replay

Input: *workflows*; *Amazon EC2 instances*

Output: Q-values Q , action profile a , reward r

```
1 Initialize replay memory  $D$ , action-value function  $Q$  with random weights;
2 observe initial state  $S$ ;
3 while not at max_episode do
4   select an action  $a$ ;
5   if with probability  $\varepsilon$  then
6     select a random action;
7   else
8     select replace  $a = \operatorname{argmax}_{a'} Q(s, a')$  with  $r_m$ ;
9   carry out action  $a$ ;
10  observe reward  $r$  and new state  $s'$ ;
11  store experience  $\langle s, a, r, s' \rangle$  in replay memory  $D$ ;
12  sample random transitions  $\langle ss, aa, rr, ss' \rangle$  from replay memory  $D$ ;
13  calculate target for each minibatch transition;
14  if  $ss'$  is terminal state then
15     $tt = rr$ ;
16  else
17     $tt = rr + \gamma \max_{a'} Q(ss', aa')$ ;
18  train the Q-network using  $(tt - Q(ss, aa))^2$  as loss;
19   $s = s'$ 
20 return Q-values  $Q$ , action profile  $a$ , reward  $r$ 
```

Algorithm 2: DECENTRALIZED(Γ, f, g, α, i)

Input: game Γ , selection mechanism f , decay schedule g , learning rate α , DQN-based agent i

Output: values V , Q-values Q , joint policy π^{i*}

```
1 initialize Q-values  $Q$ , state  $s$ , action profile  $a$ ;
2 while not at max_episode do
3   simulate action  $a_i$  in state  $s$ ;
4   observe action profile  $a_{-i}$ , rewards  $R(s, a)$ , and next state  $s'$ ;
5   select  $\pi_s^{i*} \in f(Q(s'))$ ;
6   for all DQN-based agent  $j$  do
7     update  $V_j(s')$ ;
8     update  $Q_j(s, a)$ ;
9   choose action  $a'_i$ ;
10  update  $s = s', a = a'$ ;
11  decay  $\alpha$  via  $g$ ;
12 return values  $V$ , Q-values  $Q$ , joint policy  $\pi^{i*}$ ;
```
