## Q-Learning Algorithm

1: Initialize $Q(s, a)$ to 0 for all $\alpha \in A$ in each $s \in S$
2: Initialize learning rate $\alpha \in (0, 1]$
3: Initialize discount factor $\gamma \in [0, 1]$
4: Initialize exploration rate $\epsilon \in [0, 1]$
5: **while** not converged **do**
6:      $s \leftarrow s_0$
7:      **while** $s$ not terminal **do**
8:          Observe current state $s$
9:          **if** $explore()$ **then**
10:             $a \leftarrow$ random action
11:          **else**
12:             $a \leftarrow \arg\max_a Q(s, a)$
13:          **end if**
14:          Take action $a$, observe reward $R(s, a)$ and next state $s'$
15:          Update Q-value:

$$Q(s, a) \leftarrow (1 - a) \cdot Q(s, a) + \alpha \left[ R(s, a) + \gamma \cdot \max_a Q(s', a) \right]$$

16:          $s \leftarrow s'$
17:      **end while**
18: **end while**