

初始化 $Q(s, a)$

循环(对每个episode)

初始化 $s$

循环(对每个episode的每个step)

按照 $\epsilon$ 贪婪策略基于 $Q$ 选择 $s$ 状态下要执行的动作 $a$

执行动作 $a$ ，观测得到 $r, s'$

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left[ r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right]$$
$$s \leftarrow s'$$

直到 $s$ 是结束状态