

Q-learning

2019年7月26日 9:43

一、算法思想

初始化 $Q(s, a), \forall s \in S, a \in A(s)$, 任意的数值, 并且 $Q(\text{terminal} - \text{state}, \cdot) = 0$

重复 (对每一节 episode) :

初始化 状态 S

重复 (对 episode 中的每一步) :

使用某一个 policy 比如 ($\epsilon - greedy$) 根据状态 S 选取一个动作执行

执行完动作后, 观察 reward 和新的状态 S'

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha (R_{t+1} + \lambda \max_a Q(S_{t+1}, a) - Q(S_t, A_t))$$

$$S \leftarrow S'$$

循环直到 S 终止

特点 :

Q-learning 是基于离散结点的。属于 Value-based 算法, 输出的是 action 的 value。

可以把 Q 值当成矩阵

单步更新

二、问题

1、维度灾难 : 状态太多

2、用函数去拟合 Q 表 $Q(s, a) = f(s, a, w)$

3、高维输入低维输出

因此选择采用神经网络的方式