彪哥带你学强化学习

8.深入理解 SARSA 算法及使用场景

DEEPLY UNDERSTAND REINFORCEMENT LEARNING

讲师: 韩路彪





SARSA: S_t A_t R S_{t+1} A_{t+1}

SARSA算法:学习的是Qm,进而可以根据状态选出获得最大价值的动作

 S_1

a_1	2
a ₂	4
a_3	3

	S ₁	S ₂	S ₃	S ₄	0 0 0
a_1					
a_2					
a ₃					

SARSA算法和Q算法的区别是什么?既然已经有Q算法了,为什么还需要SARSA算法?



时序差分算法:temporal difference

训练好之后的结果没训好时候不成立

理论依据:哈夫曼方程

$$Q^{\pi}(s_t, a_t) = \mathbb{E}^{\pi}_{S_{t+1}, A_{t+1}}[R(s_t, a_t, S_{t+1}) + \gamma Q^{\pi}(S_{t+1}, A_{t+1})]$$

$$Q^{\pi}(s_t, a_t) \sim r(s_t, a_t, s_{t+1}) + \gamma Q^{\pi}(s_{t+1}, a_{t+1})$$

Α

3

tderror = B - A

 $A_{new} = A + \alpha*tderror$

其中α∈(0,1], 学习率

蒙特卡洛采样

SARSA算法

$$Q^{\pi}(s_t, a_t) \sim r(s_t, a_t, s_{t+1}) + \gamma Q^{\pi}(s_{t+1}, a_{t+1})$$

$$Q_t^\pi \sim r + \gamma Q_{t+1}^\pi \qquad \qquad tderror = B - A$$

$$tderror = r + \gamma Q_{t+1}^{\pi} - Q_{t}^{\pi}$$

$$Q^\pi_{t_{new}} = Q^\pi_t + lpha(r + \gamma Q^\pi_{t+1} - Q^\pi_t)$$



探索性强,没有利用经验 容易多次踩坑

适用于开始训练阶段

采样策略

随机策略:在可能的动作里边随机采样

贪婪策略:哪个动作的价值高就用哪个

ε-贪婪策略: ε的概率随机, (1-ε)的概率贪婪 ε ∈ [0,1]

经验利用,探索不足 容易错过更优动作 适用于训练后期或推理阶段

SARSA算法

SARSA学习步骤

1、创建Q表格,并全部置零

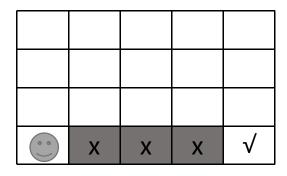
	S ₁	S ₂	S ₃	S ₄	0 0 0
a_1	0	0	0	0	0
a ₂	0	0	0	0	0
a ₃	0	0	0	0	0

- 2、采样S_t A_t
- 3、在 S_t 执行动作 A_t ,得的奖励 R 和状态 S_{t+1}
- 4、使用策略π在状态S_{t+1}选择动作A_{t+1}
- 5、根据公式更新Q(S_t, A_t)

$$Q^\pi_{t_{new}} = Q^\pi_t + lpha(r + \gamma Q^\pi_{t+1} - Q^\pi_t)$$

6、循环3~5步,直到所有Q值收敛





悬崖漫步

- ▶ 从左下角出发,每走一步奖励-1
- ▶ 如果在边缘,下个位置不存在,位置不动,奖励-1
- ▶ 掉入悬崖奖励-100 , 结束
- ▶ 走到右下角奖励100 , 结束

	S ₁₁	S ₁₂	S ₁₃	S ₁₄	S ₁₅	S ₂₁	S ₂₂	S ₂₃	S ₂₄	S ₂₅	S ₃₁	S ₃₂	S ₃₃	S ₃₄	S ₃₅	S ₄₁	S ₄₂	S ₄₃	S ₄₄	S ₄₅
上																				
下																				
左																				
右																				