彪哥带你学强化学习

7.深入理解[]学习算法及优化思路

DEEPLY UNDERSTAND REINFORCEMENT LEARNING

讲师: 韩路彪





Q:最大动作状态价值函数 $Q^*(S_t, A_t)$ ,指在状态 $S_t$ 采取动作 $A_t$ 的最大价值

Q-learning: 学习的是Q\*, 进而可以根据状态选出获得最大价值的动作

 $\mathsf{S}_1$ 

$a_1$	2
a <sub>2</sub>	4
<b>a</b> <sub>3</sub>	3

	S <sub>1</sub>	S <sub>2</sub>	<b>S</b> <sub>3</sub>	S <sub>4</sub>	0 0 0
$a_1$					
$a_2$					
<b>a</b> <sub>3</sub>					



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

$$Q^{\star}(s_t, a_t) = \mathbb{E}_{S_{t+1} \sim p(\cdot | s_t, a_t)}[R(s_t, a_t, S_{t+1}) + \gamma \max_{A_{t+1} \in \mathcal{A}} Q^{\star}(S_{t+1}, A_{t+1})]$$

$$Q^{\star}(s_t, a_t) \sim r(s_t, a_t, s_{t+1}) + \gamma \max_{A_{t+1} \in \mathcal{A}} Q^{\star}(s_{t+1}, A_{t+1})]$$

Α

В

蒙特卡洛采样

$$td = B - A$$

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

其中α∈(0,1], 学习率

## Q-learning

$$Q^{\star}(s_t, a_t) \sim r(s_t, a_t, s_{t+1}) + \gamma \max_{A_{t+1} \in \mathcal{A}} Q^{\star}(s_{t+1}, A_{t+1})]$$

td = B - A

$$Q_t^\star \sim r + \gamma \max Q_{t+1}^\star$$

$$td = r + \gamma \max Q_{t+1}^\star - Q_t^\star$$

$$Q_{t_{new}}^\star = Q_t^\star + lpha(r + \gamma \max Q_{t+1}^\star - Q_t^\star)$$



采样策略

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

随机策略:在可能的动作里边随机采样

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

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]



## Q-learning学习步骤

## 1、创建Q表格,并全部置零

	S <sub>1</sub>	s <sub>2</sub>	<b>S</b> <sub>3</sub>	S <sub>4</sub>	0 0 0
$a_1$	0	0	0	0	0
a <sub>2</sub>	0	0	0	0	0
<b>a</b> <sub>3</sub>	0	0	0	0	0

- 2、采样S<sub>t</sub> A<sub>t</sub>
- 3、在St执行动作At,得的奖励r和状态St+1
- 4、根据公式更新Q(S<sub>t</sub>, A<sub>t</sub>)

$$Q_{t_{new}}^\star = Q_t^\star + lpha(r + \gamma \max Q_{t+1}^\star - Q_t^\star)$$

5、循环2~4步,直到所有Q变化不大



Х	Х	Х	√

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