

# ML\_强化学习

## 填空题

- 1, 最优动作价值函数 $Q_*$ 依赖于\_\_\_\_\_。
- 2, DQN是对\_\_\_\_\_的近似。
- 3, 驾车按照“甲, 乙, 丙”行驶, 从甲地出发, 模型预计需要行驶20小时, 实际行驶6小时到达乙地, 模型预计还需12个小时到达丙地, 如果我们用TD算法更新模型, 那么TD目标 $\hat{y} =$ \_\_\_\_\_小时, TD绝对误差值 $|\delta|$ \_\_\_\_\_小时;

## 选择题

- 1, 设 $A = \{\text{上, 下, 左, 右}\}$ 为动作空间,  $s_t$ 为当前状态,  $\pi$ 为策略函数, 策略函数的输出:

$$\begin{aligned}\pi(\text{上}|s_t) &= 0.2, \\ \pi(\text{下}|s_t) &= 0.05, \\ \pi(\text{左}|s_t) &= 0.7, \\ \pi(\text{右}|s_t) &= 0.15.\end{aligned}\tag{1}$$

请问, 哪个动作会成为 $a_t$ ?

- A, 下
  - B, 左
  - C, 4种动作都有可能
- 2, 设随机变量 $U_t$ 为 $t$ 时刻的回报, 请问 $U_t$ 依赖于哪些变量?
- A,  $t$ 时刻的状态 $S_t$
  - B,  $t$ 时刻的动作 $A_t$
  - C,  $S_t$ 和 $A_t$
  - D,  $S_t, S_{t+1}, S_{t+2}, \dots$ 和 $A_t, A_{t+1}, A_{t+2}, \dots$

- 3, 动作价值函数是什么的期望?

- A, 奖励
- B, 回报
- C, 状态
- D, 动作

- 4, 设 $A = \{\text{上, 下, 左, 右}\}$ 为动作空间,  $s_t$ 为当前状态,  $Q_*$ 为最优动作价值函数, 策略函数的输出: