

- 进一步把公式 (4.2) 中的 Q_* 近似成 \tilde{Q} , 得到

$$\hat{y}_t \triangleq r_t + \gamma \cdot \max_{a \in \mathcal{A}} \tilde{Q}(s_{t+1}, a).$$

把它称作 TD 目标。它是表格在 $t+1$ 时刻对 $Q_*(s_t, a_t)$ 做出的估计。

$\tilde{Q}(s_t, a_t)$ 和 \hat{y}_t 都是对最优动作价值 $Q_*(s_t, a_t)$ 的估计。由于 \hat{y}_t 部分基于真实观测到的奖励 r_t , 我们认为 \hat{y}_t 是更可靠的估计, 所以鼓励 $\tilde{Q}(s_t, a_t)$ 更接近 \hat{y}_t 。更新表格 \tilde{Q} 中 (s_t, a_t) 位置上的元素:

$$\tilde{Q}(s_t, a_t) \leftarrow (1 - \alpha) \cdot \tilde{Q}(s_t, a_t) + \alpha \cdot \hat{y}_t.$$

这样可以使得 $\tilde{Q}(s_t, a_t)$ 更接近 \hat{y}_t 。Q 学习的目的是让 \tilde{Q} 逐渐趋近于 Q_* 。

收集训练数据: Q 学习更新 \tilde{Q} 的公式不依赖于具体的策略。我们可以用任意策略控制智能体, 与环境交互, 把得到的轨迹划分成 (s_t, a_t, r_t, s_{t+1}) 这样的四元组, 存入经验回放数组。这个控制智能体的策略叫做行为策略 (Behavior Policy), 比较常用的行为策略是 ϵ -greedy:

$$a_t = \begin{cases} \operatorname{argmax}_a \tilde{Q}(s_t, a), & \text{以概率 } (1 - \epsilon); \\ \text{均匀抽取 } \mathcal{A} \text{ 中的一个动作,} & \text{以概率 } \epsilon. \end{cases}$$

事后用经验回放更新表格 \tilde{Q} , 可以重复利用收集到的四元组。

经验回放更新表格 \tilde{Q} : 随机从经验回放数组中抽取一个四元组, 记作 (s_j, a_j, r_j, s_{j+1}) 。设当前表格为 \tilde{Q}_{now} 。更新表格中 (s_j, a_j) 位置上的元素, 把更新之后的表格记作 \tilde{Q}_{new} 。

1. 把表格 \tilde{Q}_{now} 中第 (s_j, a_j) 位置上的元素记作:

$$\hat{q}_j = \tilde{Q}_{\text{now}}(s_j, a_j).$$

2. 查看表格 \tilde{Q}_{now} 的第 s_{j+1} 行, 将该行的最大值记作:

$$\hat{q}_{j+1} = \max_a \tilde{Q}_{\text{now}}(s_{j+1}, a).$$

3. 计算 TD 目标和 TD 误差:

$$\hat{y}_j = r_j + \gamma \cdot \hat{q}_{j+1}, \quad \delta_j = \hat{q}_j - \hat{y}_j.$$

4. 更新表格中 (s_j, a_j) 位置上的元素:

$$\tilde{Q}_{\text{new}}(s_j, a_j) \leftarrow (1 - \alpha) \cdot \tilde{Q}_{\text{now}}(s_j, a_j) - \alpha \cdot \delta_j.$$

收集经验与更新表格 \tilde{Q} 可以同时进行。每当智能体执行一次动作, 我们可以用经验回放在 \tilde{Q} 做几次更新。也可以当完成一局游戏, 对 \tilde{Q} 做几次更新。