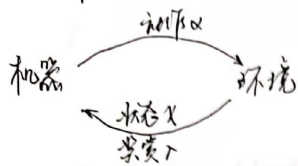


## 第16章 强化学习



马尔可夫决策过程

环境  $E$ , 状态空间  $X$ ,  $x \in X$

动作空间  $A$ ,  $a \in A$

$a$  作用在状态  $x$  上, 潜在的转移函数  $p$  将使

环境状态发生变化

一定概率

同时会以潜在的“奖励”函数反馈给机器

一个奖励

四元组  $E = \langle X, A, p, R \rangle$

$p: X \times A \times X \rightarrow R$  指定状态转移概率

$R: X \times A \times X \rightarrow R$  奖励

or  $X \times X \rightarrow R$

注意机器与环境的界限

状态转移, 奖励返回不受机器控制

机器只能通过要执行的动作影响环境

通过观察状态和奖励来感知环境

策略  $\pi: X \rightarrow A$  确定性

:  $X \times A \rightarrow R$  概率性

目标: 最优之策略

长期累积奖励最大化

### 16.2 探索与利用

① 探索: 获取每个动作的期望奖励估计

利用: 采用目前最优的动作 (但实际可能错)

需要折中

### ② $\epsilon$ -贪心

以  $\epsilon$  概率探索,  $1-\epsilon$  概率利用

随机选择随机 选取当前平均  
奖励最高

$$Q(k) = \frac{1}{n} \sum_{i=1}^n v_i \quad \text{摇臂 } k \text{ 的平均奖励}$$

$$Q_{n+1}(k) = \frac{1}{n} ((n-1) \times Q_n(k) + v_n) \\ = Q_n(k) + \frac{1}{n} (v_n - Q_n(k))$$

### ③ Softmax

基于当前已知的摇臂奖励来对探索和利用进行折中。

若各摇臂的平均奖励相当, 选取概率也相当。

若某臂更高, 则选取概率也更高

$$p(k) = \frac{e^{\frac{Q(k)}{\tau}}}{\sum_{i=1}^K e^{\frac{Q(i)}{\tau}}} \quad \tau > 0 \text{ (温度)}$$

$\tau$  越小则平均奖励越高的摇臂选取概率越高

高,  $\tau \rightarrow 0 \Rightarrow$  只利用

$\tau \rightarrow \infty \Rightarrow$  只探索

将  $K$ -摇臂赌博问题算法用于每个状态上动作的选择, 用强化学习的累积奖励来代替  $K$ -摇臂奖励函数

### 16.3 有模型学习

$E = \langle X, A, p, R \rangle$  已知

$V^\pi(x)$  状态值函数

$Q^\pi(x, a)$  状态-动作值函数

$$V_T^\pi(x) = E\pi\left[\frac{1}{T} \sum_{t=0}^{T-1} r_t \mid x_0 = x\right] \quad T \text{ 步}$$

$$V_T^\pi(x) = E\pi\left[\sum_{t=0}^{\infty} \gamma^t r_{t+1} \mid x_0 = x\right] \quad \gamma \text{ 折扣}$$

$$Q_T^\pi(x, a) = E\pi\left[\frac{1}{T} \sum_{t=0}^{T-1} r_t \mid x_0 = x, a_0 = a\right]$$

$$Q_T^\pi(x, a) = E\pi\left[\sum_{t=0}^{\infty} \gamma^t r_{t+1} \mid x_0 = x, a_0 = a\right]$$

MDP 具有马尔可夫性

系统任一时刻状态仅由当前时刻状态决定，不依赖于任何其它状态

$$V_T^*(x) = E\pi\left[\frac{1}{T} \sum_{t=1}^T r_t \mid x_0 = x\right]$$

$$= E\pi\left[\frac{1}{T} r_1 + \frac{T-1}{T} \frac{1}{T-1} \sum_{t=2}^T r_t \mid x_0 = x\right]$$

动作-价值展开

$$= \sum_{a \in A} \pi(x, a) \sum_{x' \in X} P_{x \rightarrow x'}^a \left( \frac{1}{T} R_{x \rightarrow x'}^a + \frac{T-1}{T} E\pi\left[\frac{1}{T-1} \sum_{t=2}^T r_t \mid x_0 = x'\right] \right)$$

$$= \sum_{a \in A} \pi(x, a) \sum_{x' \in X} P_{x \rightarrow x'}^a \left( \frac{1}{T} R_{x \rightarrow x'}^a + \frac{T-1}{T} V_T^*(x') \right)$$

$$V_T^*(x) = \sum_{a \in A} \pi(x, a) \sum_{x' \in X} P_{x \rightarrow x'}^a \left( R_{x \rightarrow x'}^a + \gamma V_T^*(x') \right)$$

注意：P, R 均已知

16.3.2 策略改进

$$\pi^* = \arg \max_{\pi} \sum_{x \in X} V^*(x)$$

$$V_T^*(x) = \max_{a \in A} \sum_{x' \in X} P_{x \rightarrow x'}^a \left( \frac{1}{T} R_{x \rightarrow x'}^a + \frac{T-1}{T} V_T^*(x') \right)$$

$$V_T^*(x) = \max_{a \in A} \sum_{x' \in X} P_{x \rightarrow x'}^a \left( R_{x \rightarrow x'}^a + \gamma V_T^*(x') \right)$$

$$V^*(x) = \max_{a \in A} Q^*(x, a)$$

$$Q_T^*(x, a) = \sum_{x' \in X} P_{x \rightarrow x'}^a \left( \frac{1}{T} R_{x \rightarrow x'}^a + \frac{T-1}{T} \max_{a' \in A} Q_T^*(x', a') \right)$$

$$Q_T^*(x, a) = \sum_{x' \in X} P_{x \rightarrow x'}^a \left( R_{x \rightarrow x'}^a + \gamma \max_{a' \in A} Q_T^*(x', a') \right)$$

$$V^*(x) \leq Q^*(x, \pi(x))$$

$$= \sum_{x' \in X} P_{x \rightarrow x'}^{\pi(x)} \left( R_{x \rightarrow x'}^{\pi(x)} + \gamma V^*(x') \right)$$

$$\leq \sum_{x' \in X} P_{x \rightarrow x'}^{\pi(x)} \left( R_{x \rightarrow x'}^{\pi(x)} + \gamma Q^*(x', \pi(x')) \right)$$

$$= V^*(x)$$

$$\pi^*(x) = \arg \max_{a \in A} Q^*(x, a)$$

16.3.3 策略迭代 值迭代

16.4 免模型学习

环境状态, 状态转移概率, 奖励函数往往很难得知

用策略蒙特卡罗强化学习

策略评估和改进是-7策略

策略：

f 在 p 分布下期望：

$$E[f] = \int_{\mathcal{X}} p(x) f(x) dx$$

从 p 上采样  $x_1, x_2, \dots, x_m$

$$\hat{E}[f] = \frac{1}{m} \sum_{i=1}^m f(x_i)$$

若引入 g 分布 q, 则 f 在 p 下的期望

$$E[f] = \int_{\mathcal{X}} g(x) \frac{p(x)}{q(x)} f(x) dx$$

看作  $\frac{p(x)}{q(x)} f(x)$  在分布 g 下的期望

$$\hat{E}[f] = \frac{1}{m} \sum_{i=1}^m \frac{p(x_i)}{q(x_i)} f(x_i)$$

策略  $\pi$  的采样轨迹来评估  $\pi$  (累积奖励估计期望)

$$Q(x, a) = \frac{1}{m} \sum_{i=1}^m r_i$$

若用  $\pi$  的采样来评估  $\pi$ , 则仅需对累积奖励加权

$$Q(x, a) = \frac{1}{m} \sum_{i=1}^m \frac{p_i^{\pi}}{p_i^{\pi}} r_i$$

$$p^{\pi} = \prod_{i=0}^{T-1} \pi(x_i, a_i) p_{x_{i+1} | x_i, a_i}^{\pi}$$

$$\frac{p^{\pi}}{p^{x^*}} = \prod_{i=0}^{T-1} \frac{\pi(x_i, a_i)}{\pi^*(x_i, a_i)}$$

10.4.2  
时序差分学习

蒙特卡罗强化学习局限：多次尝试后求平均  
来作为期望累积奖励的近似

"批处理式"：在一个完整的采样轨迹完成后  
再对所有的状态-动作对更新

增量式进行：

$$Q_t^{\pi}(x, a) = \frac{1}{t} \sum_{i=1}^t r_i$$

$$Q_{t+1}^{\pi}(x, a) = Q_t^{\pi}(x, a) + \underbrace{\frac{1}{t+1}}_{\text{增量}} (r_{t+1} - Q_t^{\pi}(x, a))$$

$$\Downarrow$$
$$\frac{\alpha_{t+1}}{\alpha} (r_{t+1} - Q_t^{\pi}(x, a))$$

$\alpha$ : 更新步长

$\alpha$  越大，越容易遗忘旧数据