

目录

0.1	<i>K</i>-摇臂赌博机	2
0.1.1	ϵ -贪心	2
0.1.2	Softmax 算法	3
0.1.3	UCB 方法	3
0.2	有模型学习	3
0.2.1	策略评估	3
0.2.2	策略改进	4
0.2.3	策略迭代与值迭代	4
0.3	免模型学习	5
0.3.1	蒙特卡罗强化学习	5
0.3.2	时序差分学习	6
0.4	值函数近似	6
0.5	模仿学习	6
0.5.1	直接模仿学习	6
0.5.2	逆强化学习	6
0.6	其他	7

第十六章 强化学习

强化学习任务可以用 $E = \langle X, A, P, R \rangle$ 的马尔科夫决策过程 (Markov Decision Process, 简称 MDP) 表示, 其中 E, X, A, P, R 分别表示当前环境、状态空间、动作空间、指定状态转移概率和奖赏函数。在强化学习中可以把状态、动作和策略分别对应监督学习中的示例、标记和分类器或回归器; 但是强化学习中并没有监督学习中的有标记样本 (示例-标记对), 只有等到策略执行完, 才能通过最终的学习结果才能评估动作执行是否正确。强化学习可以看做具有“延迟标记信息”的监督学习问题。

0.1 K -摇臂赌博机

0.1.1 ϵ - 贪心

每个臂被等概率探索。令 $Q(k)$ 记录摇臂 k 的平均奖赏, 若摇臂 k 被尝试了 n 次, 得到的奖赏为 v_1, v_2, \dots, v_n , 则平均奖赏为:

$$\begin{aligned} Q(k) &= \frac{1}{n} \sum_{i=1}^n v_i \\ &= \frac{1}{n} ((n-1) \times Q_{n-1}(k) + v_n) \\ &= Q_{n-1}(k) + \frac{1}{n} (v_n - Q_{n-1}(k)) \end{aligned}$$

ϵ - 贪心算法

输入: 摇臂数 K ; 奖赏函数 R ; 尝试次数 T ; 探索概率 ϵ 。

```
1:  $r = 0$ 
2:  $\forall i = 1, 2, \dots, K : Q(i) = 0, count(i) = 0$ ;  $count(i)$  记录臂  $i$  被选中的次数
3: for  $i = 1, 2, \dots, T$  do
4:   if  $\text{rand}() < \epsilon$  then
5:      $k = \text{从 } 1, 2, \dots, K \text{ 中以均匀分布随机选取}$ 
6:   else
7:      $k = \arg \max_i Q(i)$ 
8:   end if
9:    $v = R(k)$ ;
10:   $r = r + v$ ;
11:   $Q(k) = \frac{Q(k) \times count(k) + v}{count(k) + 1}$ 
12:   $count(k) = count(k) + 1$ 
13: end for
```

输出: 累积奖赏 r

0.1.2 Softmax 算法

平均奖赏越高的臂被选中的概率越大。Softmax 算法中摇臂概率的分配是基于 Boltzmann 分布的：

$$P(k) = \frac{e^{\frac{Q(k)}{\tau}}}{\sum_{i=1}^K e^{\frac{Q(i)}{\tau}}} \quad (1)$$

其中 $Q(i)$ 记录第 i 个臂的平均奖赏； $\tau > 0$ 称为“温度”； τ 越小则平均奖赏越高的摇臂被选取的概率越高。

Softmax 算法

输入：摇臂数 K ；奖赏函数 R ；尝试次数 T ；探索概率 ϵ 。

```

1:  $r = 0$ 
2:  $\forall i = 1, 2, \dots, K : Q(i) = 0, count(i) = 0$ ;  $count(i)$  记录臂  $i$  被选中的次数
3: for  $i = 1, 2, \dots, T$  do
4:    $k =$  从  $1, 2, \dots, K$  中以式子 (1) 选取
5:    $v = R(k)$ ;
6:    $r = r + v$ ;
7:    $Q(k) = \frac{Q(k) \times count(k) + v}{count(k) + 1}$ 
8:    $count(k) = count(k) + 1$ 
9: end for

```

输出：累积奖赏 r

0.1.3 UCB 方法

UCB 方法考虑了每个臂被探索的次数。UCB(Upper Confidence Bound, 上置信界) 方法每次选择 $Q(k) + UC(k)$ 最大的摇臂, $UC(k)$ 定义为置信区间, 例如:

$$Q(k) + \sqrt{\frac{2 \ln n}{n_k}}$$

其中 n 为已执行所有摇臂的总次数, n_k 为已执行摇臂 k 的次数。

0.2 有模型学习

有模型学习 (model-based learning) 即马尔科夫决策过程四元组 $E = \langle X, A, P, R \rangle$ 均为已知。

0.2.1 策略评估

用 $V^\pi(x)$ 表示从状态 x 出发, 使用策略 π 所带来的累积奖赏, 可分为 T 步累积奖赏和 $\gamma \in (0, 1)$ 折扣累计奖赏, 可得状态值函数:

$$\begin{cases} V_T^\pi(x) = \mathbb{E}_\pi \left[\frac{1}{T} \sum_{t=1}^T r_t | x_0 = x \right] = \sum_{a \in A} \pi(x, a) \sum_{x' \in X} P_{x \rightarrow x'}^a \left(\frac{1}{T} R_{x \rightarrow x'}^a + \frac{T-1}{T} V_{T-1}^T(x') \right) \\ V_\gamma^\pi(x) = \mathbb{E}_\pi \left[\sum_{t=0}^{+\infty} \gamma^t r_{t+1} | x_0 = x \right] = \sum_{a \in A} \pi(x, a) \sum_{x' \in X} P_{x \rightarrow x'}^a (R_{x \rightarrow x'}^a + \gamma V_\gamma^T(x')) \end{cases}$$

函数 $Q^\pi(x, a)$ 表示从状态 x 出发, 执行动作 a 后在使用策略 π 带来的累积奖赏, 用 x_0 表示起始状态, a_0 表示起始状态上采取的第一个动作, 可得状态-动作值函数:

$$\begin{cases} Q_T^\pi(x, a) = \mathbb{E}_\pi \left[\frac{1}{T} \sum_{t=1}^T r_t | x_0 = x, a_0 = a \right] = \sum_{x' \in X} P_{x \rightarrow x'}^a \left(\frac{1}{T} R_{x \rightarrow x'}^a + \frac{T-1}{T} V_{T-1}^T(x') \right) \\ V_\gamma^\pi(x, a) = \mathbb{E}_\pi \left[\sum_{t=0}^{+\infty} \gamma^t r_{t+1} | x_0 = x, a_0 = a \right] = \sum_{x' \in X} P_{x \rightarrow x'}^a (R_{x \rightarrow x'}^a + \gamma V_\gamma^T(x')) \end{cases}$$

0.2.2 策略改进

理想的策略应能最大化累积奖赏：

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

无论是状态值函数还是状态-动作值函数对于策略的每一点改进都是单调递增的。

0.2.3 策略迭代与值迭代

策略迭代 (policy iteration) 是按照策略评估-策略迭代-策略评估-策略迭代。。。以此迭代。

基于 T 步累积奖赏的策略迭代算法

输入: MDP 四元组 $E = \langle X, A, P, R \rangle$; 累积奖赏参数 T

```

1:  $\forall x \in X : V(x) = 0, \pi(x, a) = \frac{1}{|A(x)|}$ 
2: loop
3:   for  $t = 1, 2, \dots$  do
4:      $\forall x \in X : V'(x) = \sum_{a \in A} \pi(x, a) \sum_{x' \in X} P_{x \rightarrow x'}^a (\frac{1}{t} R_{x \rightarrow x'}^a + \frac{t-1}{t} V_{T-1}^T(x'))$ 
5:     if  $t = T + 1$  then
6:       break
7:     else
8:        $V = V'$ 
9:     end if
10:  end for
11:   $\forall x \in X : \pi'(x) = \arg \max_{a \in A} Q(x, a);$ 
12:  if  $\forall x : \pi'(x) = \pi(x)$  then
13:    break
14:  else
15:     $\pi = \pi'$ 
16:  end if
17: end loop
输出: 最优策略  $\pi$ 

```

值迭代 (value iteration) 改进就是先更新值函数，最后在形成策略。

基于 T 步累积奖赏的值迭代算法

输入: MDP 四元组 $E = \langle X, A, P, R \rangle$; 累积奖赏参数 T ; 收敛阈值 θ

```

1:  $\forall x \in X : V(x) = 0$ 
2: for  $t = 1, 2, \dots$  do
3:    $\forall x \in X : V'(x) = \max_{a \in A} \sum_{x' \in X} P_{x \rightarrow x'}^a (\frac{1}{t} R_{x \rightarrow x'}^a + \frac{t-1}{t} V_{T-1}^T(x'))$ 
4:   if  $\max_{x \in X} |V(x) - V'(x)| < \theta$  then
5:     break
6:   else
7:      $V = V'$ 
8:   end if
9: end for
输出: 最优策略  $\pi(x) = \arg \max_{a \in A} Q(x, a)$ 

```

0.3 免模型学习

免模型学习 (model-free learning) 指算法不依赖与环境建模, 有两个困难: **策略无法评估和奖赏函数未知**。也称为无模型学习。

0.3.1 蒙特卡罗强化学习

多次采样, 求取平均累积奖赏来作为期望累积奖赏的近似。因为采样次数必须有限, 所以这个方法更适合于 T 步累积奖赏的强化学习任务。采样时对策略的选取采用 ϵ -贪心法, 将确定性的策略 π 称为原始策略:

$$\pi^\epsilon(x) = \begin{cases} \pi(x), & \text{以概率 } 1 - \epsilon \\ A \text{中以均匀概率选取的动作}, & \text{以概率 } \epsilon \end{cases}$$

因为这里被评估和被改进的是同一个策略, 因此被称为“同策略”(one-policy) 蒙特卡罗强化学习算法。

同策略蒙特卡罗强化学习算法

输入: 环境 E ; 动作空间 A ; 起始状态 x_0 ; 策略执行步数 T

```

1:  $Q(x, a) = 0, count(x, a) = 0, \pi(x, a) = \frac{1}{|A(x)|}$ 
2: for  $s = 1, 2, \dots$  do
3:   在  $E$  中执行策略  $\pi$  产生轨迹:  $\langle x_0, a_0, r_1, x_1, a_1, r_2, \dots, x_{T-1}, a_{T-1}, r_T, x_T \rangle$ 
4:   for  $t = 0, 1, \dots, T-1$  do
5:      $R = \frac{1}{T-t} \sum_{i=t+1}^T r_i$ ;
6:      $Q(x_t, a_t) = \frac{Q(x_t, a_t) \times count(x_t, a_t) + R}{count(x_t, a_t) + 1}$ 
7:      $count(x_t, a_t) = count(x_t, a_t) + 1$ 
8:   end for
9:   对所有已见状态  $x$ :
```

$$\pi(x) = \begin{cases} \arg \max_{a'} Q(x, a'), & \text{以概率 } 1 - \epsilon \\ \text{以均匀概率从 } A \text{ 中选取动作}, & \text{以概率 } \epsilon \end{cases}$$

```

10: end for
```

输出: 策略 π

异策略 (off-policy) 蒙特卡罗强化学习算法是策略评估时引用 ϵ -贪心, 策略改进时改进原始策略。

Algorithm 1 *

异策略蒙特卡罗强化学习算法

输入: 环境 E ; 动作空间 A ; 起始状态 x_0 ; 策略执行步数 T

```

1:  $Q(x, a) = 0, count(x, a) = 0, \pi(x, a) = \frac{1}{|A(x)|}$ 
2: for  $s = 1, 2, \dots$  do
3:   在  $E$  中执行  $\pi$  的  $\epsilon$ -贪心策略  $\pi$  产生轨迹:  $\langle x_0, a_0, r_1, x_1, a_1, r_2, \dots, x_{T-1}, a_{T-1}, r_T, x_T \rangle$ 
4:
```

$$p_i = \begin{cases} 1 - \epsilon + \epsilon/|A|, & a_i = \pi(x_i); \\ \epsilon/|A|, & a_i \neq \pi(x_i) \end{cases}$$

```

5:   for  $t = 0, 1, \dots, T-1$  do
6:      $R = \frac{1}{T-t} (\sum_{i=t+1}^T r_i) \prod_{i=t+1}^{T-1} \max(1, \frac{\mathbb{I}(a_i = \pi(x_i))}{p_i})$ ;
7:      $Q(x_t, a_t) = \frac{Q(x_t, a_t) \times count(x_t, a_t) + R}{count(x_t, a_t) + 1}$ 
8:      $count(x_t, a_t) = count(x_t, a_t) + 1$ 
9:   end for
10:  对所有已见状态  $\pi(x) = \arg \max_{a'} Q(x, a')$ 
```

11: **end for**
输出: 策略 π

0.3.2 时序差分学习

蒙特卡罗强化学习算法通过完成一个采样估计后再更新策略的值估计，解决了模型未知问题，但是效率很低。时序差分学习 (Temporal Difference Learning) 则结合了动态规划和蒙特卡罗方法的思想，充分利用了强化学习任务重的 MDP 结构：每执行一步策略就更新一次值函数估计。

分别可得到同策略的 Sarsa 算法和异策略的 Q-学习算法 (Q-learning)。

0.4 值函数近似

前面的强化学习算法都是在有限状态空间上进行的，如果状态空间无限，那么就需要将连续状态空间的每一个状态表示成一个状态向量 \mathbf{x} ，然后结合参数向量 $\boldsymbol{\theta}$ 来表示状态的线性函数 (也是值函数)：

$$V_{\boldsymbol{\theta}}(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x}$$

线性值函数近似 Sarsa 算法

输入: 环境 E ；动作空间 A ；起始状态 \mathbf{x}_0 ；奖赏折扣 γ ；更新步长 α

```

1:  $\boldsymbol{\theta} = \mathbf{0}, \mathbf{x} = \mathbf{x}_0, a = \pi(\mathbf{x}) = \arg \max_{a''} \boldsymbol{\theta}^T(\mathbf{x}; a'')$ 
2: for  $s = 1, 2, \dots$  do
3:    $r, \mathbf{x}' =$  在  $E$  中执行动作  $a$  产生的奖赏与转移的状态；
4:    $a' = \pi^{\epsilon}(\mathbf{x}')$ 
5:    $\boldsymbol{\theta} = \boldsymbol{\theta} + \alpha(r + \gamma \boldsymbol{\theta}^T(\mathbf{x}'; a') - \boldsymbol{\theta}^T(\mathbf{x}; a))(\mathbf{x}; a)$ ；
6:    $\pi(\mathbf{x}) = \arg \max_{a''} \boldsymbol{\theta}^T(\mathbf{x}; a'')$ ；
7:    $\mathbf{x} = \mathbf{x}', a = a'$ 
8: end for
```

输出: 策略 π

0.5 模仿学习

0.5.1 直接模仿学习

直接模仿人类专家的“状态-动作对”。

0.5.2 逆强化学习

逆强化学习 (inverse reinforcement learning) 就是从人类专家提供的范例数据中反推出奖赏函数。

假设奖赏函数可以表达为状态特征的线性函数即： $R(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ ，那么策略 π 的累积奖励可以写成状态向量的加权求和的期望与系数 \mathbf{w} 的内积：

$$\rho^{\pi} = \mathbb{E} \left[\sum_{i=0}^{+\infty} \gamma^i R(\mathbf{x}_i) | \pi \right] = \mathbb{E} \left[\sum_{i=0}^{+\infty} \gamma^i \mathbf{w}^T \mathbf{x}_i | \pi \right] = \mathbf{w}^T \mathbb{E} \left[\sum_{i=0}^{+\infty} \gamma^i \mathbf{x}_i | \pi \right]$$

把 $\mathbb{E} \left[\sum_{i=0}^{+\infty} \gamma^i \mathbf{x}_i | \pi \right]$ 简写为 $\bar{\mathbf{x}}^{\pi}$ ，把每条范例轨迹上的状态加权求和再平均记为 $\bar{\mathbf{x}}^*$ ，对于最优奖赏函数 $R(\mathbf{x}) = \mathbf{w}^{*T} \mathbf{x}$ 和其他任意策略产生的 $\bar{\mathbf{x}}^{\pi}$ ，有：

$$\mathbf{w}^{*T} \bar{\mathbf{x}}^* - \mathbf{w}^{*T} \bar{\mathbf{x}}^{\pi} = \mathbf{w}^{*T} (\bar{\mathbf{x}}^* - \bar{\mathbf{x}}^{\pi}) \geq 0$$

所以有最优策略：

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \min_{\pi} \mathbf{w}^T (\bar{\mathbf{x}}^* - \bar{\mathbf{x}}^{\pi}) \quad \text{s.t.} \quad \|\mathbf{w}\| \leq 1$$

迭代式逆强化学习算法

输入： 环境 E ；动作空间 A ；状态空间 X ；范例轨迹数据集 $D = \{\tau_1, \tau_2, \dots, \tau_m\}$

- 1: $\bar{\mathbf{x}}^* =$ 从范例轨迹中算出状态加权求和的均值向量；
- 2: $\pi =$ 随机策略
- 3: **for** $t = 1, 2, \dots$ **do**
- 4: $\bar{\mathbf{x}}^{\pi} =$ 从 π 的采样轨迹算出状态加权求和的均值向量；
- 5: 求解 $\mathbf{w}^* = \arg \max_{\mathbf{w}} \min_{\pi} \mathbf{w}^T (\bar{\mathbf{x}}^* - \bar{\mathbf{x}}^{\pi}) \quad \text{s.t.} \quad \|\mathbf{w}\| \leq 1$
- 6: $\pi =$ 在环境 $\langle X, A, R(\mathbf{x}) = \mathbf{w}^{*T} \mathbf{x} \rangle$ 中的求解最优策略
- 7: **end for**

输出： 奖赏函数 $R(\mathbf{x}) = \mathbf{w}^{*T} \mathbf{x}$ 与策略 π

0.6 其他

RL 的缺点：

- 需要大量有效的样本
- 最终表现很多时候并不好；表现好的时候可能是过拟合
- 需要奖励函数，但是通常难以获得
- 因不当探索和利用 (exploration and exploitation) 导致的局部最优 (local optima) 难以避免
- 结果不稳定，而且很难重现

RL 表现好的条件：

- 容易获得大量样本
- 问题可以化简成一个简单的形式
- 可以进行左右互搏 (self-play)
- 奖励函数可以被定义；奖励信号丰富，反馈及时

RL 未来的方向和发展：

- 局部最优或许已经足够好
- 靠硬件加速
- 增加更多的学习信号
- 用基于模型的 RL 来提高样本利用率
- 把 RL 当做一种微调手段 (fine-tuning)
- 迁移学习和强化学习结合
- 自学习奖励函数
- 好的先验知识可以大大降低学习时间
- 复杂的问题可能更容易用 RL 解决