

<p><b>强化学习1:</b>强化学习是一种优化智能体在环境中行为的一种方法,根据环境反馈的奖励,调整智能体的行为策略,提升智能体实现目标的能力。</p> <p>○强化学习考虑的是<b>序贯决策过程</b>:</p> <p>智能体处在产生的环境中产生一系列的动作,而环境能够根据这些动作改变智能体的当前状态。</p> <p>举例</p> <p>遥控直升飞机的特技表演;打败围棋世界冠军;管理股票证券;发电厂调控;控制人型机器人双足行走;视频游戏上超越人类</p> <p>目标:选择一组动作使未来奖励和最大化</p> <p>动作可能在未来很久才会产生影响;奖励可能是延时的;有可能会牺牲短期利益从而获得长期的回报</p> <p>举例:</p> <p>投资(要几个月才会收益);直升机加油(防止几小时后没油坠机)</p> <p><b>强化学习:</b></p> <p>①产生的结果(动作) 能够改变数据的分布(状态)</p> <p>②最终的目标可能要很长时间才能观察到(下棋)</p> <p>③没有明确的标签(label) 数据</p> <p>④根据当前的奖励,解决实现长远的目标</p> <p><b>监督学习(Supervised Learning, SL)/非监督学习:</b></p> <p>①产生的结果(输出) 不会改变数据的分布②结果是瞬时的③要有明确的标签数据(SL)④要完全没有任何标签数据(USL)</p> <p>●<b>马尔可夫性 Markov property</b></p> <p>智能体未来的状态只与当前时刻的状态<math>S_t</math>有关,而与过去的状态<math>\{S_1, \dots, S_{t-1}\}</math>无关, 那么称智能体的模型具有马尔可夫性。</p> <p><math>P[S_{t+1} S_t] = P[S_{t+1} S_1, \dots, S_t]</math></p> <p>未来只与当前有关,与历史无关<math>S_{1:t} \rightarrow S_t \rightarrow S_{t+1:\infty}</math></p> <p>一旦当前状态确定了,历史状态都可以弃掉了,也就是说当前状态足以决定未来状态是什么样的</p> <p><b>强化学习</b>主要研究的是具有马尔可夫性的问题。</p> <p><b>状态转移矩阵:</b> <b>内在状态</b></p> <p>对于一个马尔可夫状态<math>s</math>和后继状态<math>s'</math>, 状态转移概率定义为</p> <p><math>P_{ss'} = P(S_{t+1} = s'   S_t = s)</math> <b>o.t.智能体观测量</b></p> <p>状态转移矩阵<math>P</math> 定义从所有状态<math>s</math>到所有后继状态<math>s'</math>的转移概率</p> $P = \text{from} \begin{bmatrix} P_{11} & \dots & P_{1n} \\ \vdots & \ddots & \vdots \\ P_{n1} & \dots & P_{nn} \end{bmatrix} \text{ 其中矩阵的每一行和为 1}$ <p>●<b>马尔可夫过程 Markov Process</b></p> <p>一个马尔可夫过程是一个无记忆的随机过程, 即一组具有马尔可夫性的随机状态序列 <math>S_1, S_2, \dots</math></p> <p>定义:</p> <p>马尔可夫过程 (马尔可夫链) 可以用一组<math>(S, P)</math>表示</p> <p><math>S</math> 是 (有限) 状态集</p> <p><math>P</math> 是状态转移概率矩阵 <math>P_{ss'} = P(S_{t+1} = s'   S_t = s)</math></p> <p>●<b>马尔可夫奖励过程 Markov Reward Process</b></p> <p>一个马尔可夫奖励过程是一个马尔可夫链加上奖励定义:</p> <p>一个马尔可夫奖励过程由一组<math>(S, P, R, \gamma)</math>构成</p> <p><math>S</math> 是一组有限状态集</p> <p><math>P</math> 是状态转移概率矩阵 <math>P_{ss'} = P(S_{t+1} = s'   S_t = s)</math></p> <p><math>R</math> 是奖励函数, <math>R_t = E(R_{t+1}   S_t = s)</math></p> <p><math>\gamma</math> 是折扣因子, <math>\gamma \in [0, 1]</math></p> <p><b>回报 Return:</b> 回报<math>G_t</math> 代表从<math>t</math>时刻往后所有的折扣奖励:</p> $G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$ <p>折扣因子<math>\gamma</math>代表未来的奖励在当前时刻贡献的价值</p> <p><math>k+1</math>时刻后的奖励<math>R</math>对当前回报的贡献只有<math>\gamma^k R</math></p> <p>0-目光短浅; 1=目光长远 <b>回报对应某一具体的轨迹</b></p> <p>如果想要调整奖励的重要性, 数学上很方便, 在循环马尔可夫过程中能避免无穷回报。但是有可能会忽视未来奖励。如果奖励代表金钱, 近期的奖励会比远期产生更多的收益。自然界的人类或动物行为模式更倾向于近期奖励。有时候也会使用无折扣的马尔可夫奖励过程 (即1)。例如所有的事件序列都有终止状态</p> <p><b>价值函数<math>v(s)</math></b> 代表智能体在状态<math>s</math>下的长期价值</p> <p>一个马尔可夫奖励过程的状态价值函数等于从状态<math>s</math>出发的期望回报 <b>价值代表所有轨迹的期望</b></p> $v(s) = E[G_t   S_t = s]$ <p><b>MRPs 的贝尔曼方程 Bellman equation</b></p> <p>价值函数拆分成两部分: 瞬间奖励<math>R_{t+1}</math>; 后继状态的折扣价值<math>\gamma v(S_{t+1})</math></p> $v(s) = E[G_t   S_t = s] = E[R_{t+1} + \gamma R_{t+2} + \dots   S_t = s] = E[R_{t+1} + \gamma v(S_{t+1})   S_t = s] = E[R_{t+1} + \gamma v(S_{t+1})   S_t = s] = E[R_{t+1} + \gamma v(S_{t+1})   S_t = s]$ $v(s) = E[R_{t+1} + \gamma v(S_{t+1})   S_t = s] = R_s + \gamma \sum_{s' \in S} P_{ss'} v(s')$ <p>矩阵: <math>V = R + \gamma P V</math></p> $\begin{bmatrix} v(1) \\ \vdots \\ v(n) \end{bmatrix} = \begin{bmatrix} R(1) \\ \vdots \\ R(n) \end{bmatrix} + \gamma \begin{bmatrix} P_{11} & \dots & P_{1n} \\ \vdots & \ddots & \vdots \\ P_{n1} & \dots & P_{nn} \end{bmatrix} \begin{bmatrix} v(1) \\ \vdots \\ v(n) \end{bmatrix}$ <p><math>V = (I - \gamma P)^{-1} R</math>: <math>n</math>个状态下的计算复杂度 <math>O(n^3)</math></p> <p>●<b>马尔可夫决策过程 Markov Decision Process, MDP</b></p> <p>马尔可夫决策过程是马尔可夫奖励过程加上决策。问题的所有状态都具有马尔可夫性。</p> <p>一个马尔可夫决策过程由一组<math>(S, \underline{A}, P, R, \gamma)</math>构成</p> <p><math>S</math> 是一组有限状态集</p> <p><math>A</math> 是有限动作集</p> <p><math>P</math> 是状态转移概率矩阵</p> $P_{ss'}^a = P(S_{t+1} = s'   S_t = s, A_t = a)$ $R$ 是奖励函数, $R_t^a = E(R_{t+1}   S_t = s, A_t = a)$ <p><math>\gamma</math> 是折扣因子, <math>\gamma \in [0, 1]</math></p> <p>●<b>强化学习的主要组成元素:</b></p> <p><b>策略:</b> 智能体的行为</p> <p><b>价值函数</b> (值函数、性能指标函数): 智能体在某一状态和/或某一动作时是好还是坏</p> <p><b>模型:</b> 智能体对真实环境的估计</p> <p><b>策略:</b> 代表了智能体是如何行为的, 是从状态到动作的映射, 是状态到动作的一种分布</p> <p>确定性策略: <math>a = \pi(s)</math></p> <p>随机性策略: <math>\pi(a s) = P(A_t = a   S_t = s)</math></p> <p>一个策略定义了一个智能体的行为, MDP 问题的策略取决于当前时刻的状态(与历史无关), 即策略是静态的(时不变性), 迷宮中的箭头 (上下左右) 代表策略<math>\pi(s)</math>在状态<math>s</math>的动作</p> <p>给定一个 MDP 的 <math>M = (S, \underline{A}, P, R, \gamma)</math> 和策略<math>\pi</math>:</p> <p>状态序列<math>S_1, S_2, \dots</math> 是一个马尔可夫过程(<math>S, P^\pi</math>)</p> <p>状态和奖励序列<math>S_1, R_1, S_2, \dots</math> 是一个马尔可夫奖励过程(<math>S, P^\pi, R, \gamma</math>)</p> $P_{ss'}^\pi = \sum_{a \in A} \pi(a s) P_{ss'}^a \quad R_s^\pi = \sum_{a \in A} \pi(a s) R_t^a$ <p><b>价值函数 Value Function</b></p> <p>价值函数是对未来奖励的预测, 评估智能体在某一状态下是好还是坏, 因而可以用来选择对智能体最有利的动作</p> <p>① <b>状态-价值函数</b> <math>v_\pi(s) = E_\pi[G_t   S_t = s]</math></p> <p>② <b>动作-价值函数</b> <math>q_\pi(s, a) = E_\pi[G_t   S_t = s, A_t = a]</math></p> $q_\pi(s, a) = E_\pi[R_{t+1} + \gamma q_\pi(S_{t+1}, A_{t+1})   S_t = s, A_t = a]$ $v_\pi(s) \leftarrow s$ $q_\pi(s, a) \leftarrow a$	<p><math>V_\pi(s)</math> 的轨迹</p> $s_t = s, a_t \sim \pi, r_{t+1} \sim \mathcal{R}, s_{t+1} \sim \mathcal{P}, a_{t+1} \sim \pi, \dots$ <p><math>Q_\pi(s, a)</math> 的轨迹</p> $s_t = s, a_t = a, r_{t+1} \sim \mathcal{R}, s_{t+1} \sim \mathcal{P}, a_{t+1} \sim \pi, \dots$ <p><b>贝尔曼期望方程</b></p> $v_\pi = \sum_{a \in A} \pi(a s) \left( R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a v_\pi(s') \right)$ $q_\pi(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a \sum_{a' \in A} \pi(a' s') q_\pi(s', a')$ $v_\pi = (I - \gamma P^\pi)^{-1} R^\pi$ <p><b>最优价值函数:</b></p> <p><b>最优状态-价值函数:</b> 在所有策略中价值函数最大的</p> $v_*(s) = \max_{\pi} v_\pi(s)$ <p><b>最优动作-价值函数:</b> 在所有策略中动作价值函数最大</p> $q_*(s, a) = \max_{\pi} q_\pi(s, a)$ <p>最优价值函数代表了智能体在该 MDP 问题上最好的性能; 如果得到了最优价值函数, 那么 MDP 问题就已经解决了</p> <p><b>最优策略:</b> <math>\pi \geq \pi'</math> if <math>v_\pi(s) \geq v_{\pi'}(s), \forall s</math></p> <p>定理</p> <p>对任意马尔可夫决策过程</p> <p>总是存在一个最优策略<math>\pi_*</math>比其它所有策略都不差</p> $\pi_* \geq \pi \quad \forall \pi$ <p>所有最优策略的价值函数都相等, 且等于最优价值函数, <math>V_*(s) = v_*(s)</math></p> <p>所有最优策略的动作-价值函数都相等, 且等于最优动作-价值函数, <math>q_*(s, a) = q_*(s, a)</math></p> <p><b>寻找最优策略</b></p> <p>一个最优策略可以通过最大化<math>q_*(s, a)</math>来确定</p> $\pi_*(a s) = \begin{cases} 1 & \text{if } a = \arg \max_{a'} q_*(s, a) \\ 0 & \text{otherwise} \end{cases}$ <p>对任何 MDP 都存在一个确定性的最优策略</p> <p>如果<math>q_*(s, a)</math>已知, 即可得到最优策略</p> <p><b>最优化原理:</b></p> <p>强化学习目标是找到一组时间序列的动作<math>\{A_0, A_1, \dots\}</math>使得智能体从<math>S_0</math>出发得到的期望累加奖励最大化</p> <p><b>最优价值函数:</b></p> $v^*(S_0) = E \left[ \max_{a_0} (R_1 + \gamma R_2 + \gamma^2 R_3 + \dots) \right]$ <p>最优策略具有如下性质: 不论初始状态和初始决策(第一步决策)如何, 以第一步决策所形成的阶段和状态作为初始条件来考虑时, 余下的决策对余下的问题而言也必构成最优策略</p> $v^*(S_0) = E \left[ \max_{a_0} (R_1 + \gamma R_2 + \gamma^2 R_3 + \dots) \right]$ $= \max_{a_0} E \left[ R_1 + \gamma \max_{a_1} (R_2 + \gamma R_3 + \dots) \right]$ $= \max_{a_0} E \left[ R_1 + \gamma E[v^*   S_1 = s_1] \right]$ <p><math>v_\pi(s) \leftarrow s</math></p> <p><math>q_\pi(s, a) \leftarrow a</math></p> <p><math>v_*(s) = \max_a \sum_{s'} P_{ss'}^a v_*(s')</math></p> $q_*(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a v_*(s')$ <p><math>v_\pi(s) \leftarrow s</math></p> <p><math>q_\pi(s, a) \leftarrow a</math></p> <p><math>q_\pi(s, a') \leftarrow a'</math></p> <p><math>q_\pi(s', a') \leftarrow a'</math></p> <p>关于<math>v_*</math>的 <b>Bellman 最优方程</b>:</p> $v_*(s) = \max_a \left( R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a v_*(s') \right)$ <p>关于<math>q_*</math>的: <math>q_*(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a \max_{a'} q_*(s', a')</math></p> <p>确定性的最优策略:</p> $\pi_*(s) = \arg \max_{a \in A} q_*(s, a)$ <p>求解贝尔曼最优方程需要:</p> <ol style="list-style-type: none"> <li>1 求解非线性算子 <math>\max</math></li> <li>2 模型已知</li> <li>3 足够的计算空间</li> </ol> <p><b>动态规划 DP</b></p> <p>通过把原问题分解为相对简单的子问题来求解复杂问题的方法, 适用于:</p> <p><b>最优子结构性质:</b> 问题的最优解所包含的子问题的解也是最优的, 也就是满足最优化原理, 将问题分成子问题。</p> <p><b>子问题重叠性质:</b> 使用递归算法自顶向下对问题进行求解, 每次产生的子问题并不总是新问题。(子问题重复多次出现, 保存第一次的解)</p> <p>马尔可夫决策过程满足以上两个属性: 贝尔曼方程具有递归性: 价值函数保存和重复利用。</p> <p><b>广告位招租</b></p>	<p><b>贝尔曼最优方程:</b></p> $v_*(s) = [T(v_*)](s) = \max_a \left( R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a v_*(s') \right)$ <p>价值迭代收敛性: <math>\gamma &lt; 1</math>(此时价值迭代算子是收缩算子)</p> <p>使用 <math>\infty</math>-范数比较两个价值函数的不同, 即状态集上两个函数最大的差值 <math>\ u - v\ _\infty = \max_i  u(s_i) - v(s_i) </math></p> <p>收缩算子: 对于任意两个函数<math>f, g</math>, 如果两个函数的距离经过一个算子<math>T</math>后能够被缩小, 那么这个算子被称为收缩算子</p> $\ T(f - g)\ _\infty \leq \ f - g\ _\infty$ <p><math>[T(v)](s) = [T(v)](s), \quad \forall s \in S</math></p> <p><b>是否需要证明过程?</b></p> $= \max_{a_1} \left( R^{a_1} + \gamma \sum_{s' \in S} P_{ss'}^{a_1} u(s') \right) - \max_{a_2} \left( R^{a_2} + \gamma \sum_{s' \in S} P_{ss'}^{a_2} v(s') \right)$ $\leq \max_{a_1} \left  \left( R^{a_1} + \gamma \sum_{s' \in S} P_{ss'}^{a_1} u(s') \right) - \left( R^{a_1} + \gamma \sum_{s' \in S} P_{ss'}^{a_1} v(s') \right) \right $ $= \gamma \max_{a_1} \left  \sum_{s' \in S} P_{ss'}^{a_1} (u(s') - v(s')) \right $ $\leq \gamma \ u - v\ _\infty \quad (\gamma < 1)$ <p><math>\Rightarrow \ T(u - v)\ _\infty \leq \ u - v\ _\infty</math></p> <p>从任意初始价值函数 <math>v_1</math> 出发, 经过价值迭代计算的新函数 <math>v_{k+1}</math> 与最优价值函数之间的距离</p> $\ v_{k+1} - v_*\ _\infty = \ T v_k - T v_*\ _\infty$ $\leq \gamma \ v_k - v_*\ _\infty$ $\vdots$ $\leq \gamma^k \ v_1 - v_*\ _\infty$ <p>当<math>k \rightarrow \infty, v_k \rightarrow v_*</math></p> <p><b>策略</b></p> <p>智能体可根据给定的策略决定在当前状态下如何采取下一步动作</p> <ol style="list-style-type: none"> <li>1 确定性策略: <math>a_t = \pi(s_t)</math></li> <li>2 随机性策略: <math>a_t \sim \pi(s_t)</math></li> </ol> <p>给定一个策略, 就可以对它的好坏进行评估</p> <p>智能体按照给定的策略执行动作, 获得的期望回报称为该策略的价值函数 <math>v_\pi(s) = E_\pi[R_{t+1} + \gamma R_{t+2} + \dots   S_t = s, A_t \sim \pi(S_t)]</math></p> <p>在给定策略<math>\pi</math>的情况下, 智能体有限个状态之间的转移可以用矩阵形式表示</p> $P^\pi = \begin{bmatrix} P_{s_1, s_1}^\pi & \dots & P_{s_1, s_n}^\pi \\ \vdots & \ddots & \vdots \\ P_{s_n, s_1}^\pi & \dots & P_{s_n, s_n}^\pi \end{bmatrix}$ <p>其中<math>P_{s_i, s_j}^\pi</math>代表了在策略<math>\pi</math>下, 从状态<math>s_i</math>转移到<math>s_j</math>的概率</p> <p><b>贝尔曼期望方程:</b> <math>v_\pi = (I - \gamma P^\pi)^{-1} R^\pi</math> 矩阵很大, 矩阵稀疏</p> $v_\pi = \begin{bmatrix} v_\pi(s_1) \\ \vdots \\ v_\pi(s_n) \end{bmatrix}, \quad R_\pi = \begin{bmatrix} R_{s_1} \\ \vdots \\ R_{s_n} \end{bmatrix}$ <ol style="list-style-type: none"> <li>1: 给定一个初始策略 <math>\pi_1, k = 1</math></li> <li>2: <b>loop</b></li> <li>3: <b>策略评估</b> policy evaluation: 对当前策略 <math>\pi_k</math> 计算它的价值函数 <math>V_{\pi_k}</math></li> </ol> $V_{\pi_k}(s) = E[(r_{t+1} + \gamma r_{t+2} + \dots)   s_t = s, a_t \sim \pi_k(s_t)]$ $= \sum_a \pi_k(a s) \left( R_s^a + \gamma \sum_{s'} P_{ss'}^a V_{\pi_k}(s') \right)$ <ol style="list-style-type: none"> <li>4: <b>策略提升</b> policy improvement: 根据 <math>V_{\pi_k}</math> 提取贪心策略</li> </ol> $\pi_{k+1}(s) = \arg \max_a \left( R_s^a + \gamma \sum_{s'} P_{ss'}^a V_{\pi_k}(s') \right)$ <ol style="list-style-type: none"> <li>5: <math>k \leftarrow k + 1</math></li> <li>6: <b>end loop</b></li> </ol> <p>■ 策略迭代收敛时策略不再变化</p> <p>■ 迭代的次数不超过 MDP 的最大 (确定性) 策略数量</p> $\pi_1 \rightarrow \pi_2 \rightarrow \dots \rightarrow \pi_k = \pi_*$ $k \leq  A ^S$ <p>■ e.g. <math> A  = 2,  S  = 10, k \leq 2^{10} = 1024</math></p> <p>■ <b>思考: 策略迭代过程会出现死循环吗? 即</b>  <math>\pi_i \rightarrow \pi_{i+1} \rightarrow \dots \rightarrow \pi_j = \pi_i</math>  <b>不会, 策略的价值只会越来越小</b></p> $\ V_k - V_*\ _\infty \leq \gamma^{k-1} \ V_1 - V_*\ _\infty$ <p>■ 假设初始 <math>V_1(s) = 0, \forall s \in S</math></p> <p>■ 收敛时间<math>\epsilon</math>, 即 <math>\ V_1 - V_*\ _\infty \leq \epsilon</math> 时认为收敛到 <math>V_*</math></p> $\Rightarrow \begin{cases} \gamma^{k-1} \ V_1 - V_*\ _\infty \leq \epsilon \\ k \geq 1 + \log_{1/\gamma} \frac{\epsilon}{\ V_1 - V_*\ _\infty} \end{cases}$ <p>■ e.g. <math>\max V_*(s) = 10, \gamma = 0.95, \epsilon = 0.001, k \geq 180.56</math></p>
--	--	--

MC 和 TD 对比

**MC 是高方差，零偏差：**

需要完整的事件序列计算出回报后学习

好的收敛性(即使是使用逼近器也能保证收敛性);

与价值函数初始值无关;

原理简单, 使用方便;

要求事件达到终止状态或序列足够长

能够利用马尔可夫性, 因此在马尔可夫环境下更有效

**TD 是低方差, 有偏差:**

可以在智能体运行过程中的每一步在线学习

可以从不完整的事件序列中学习

通常情况是比 MC 更有效

TD(0) 能够收敛到  $v_\pi(s)$

(但是与逼近器结合后并不一定收敛)

受价值函数初始值影响

不能利用马尔可夫性, 因此在非马尔可夫环境下更有效

偏差/方差的权衡:

回报  $G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-t} R_t = v_\pi(S_t)$  的无偏估计

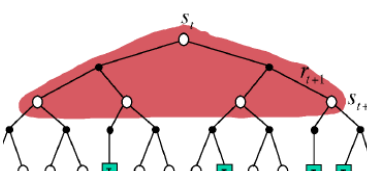
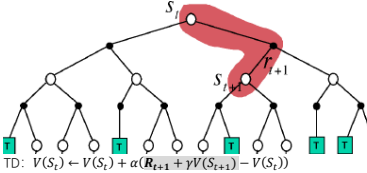
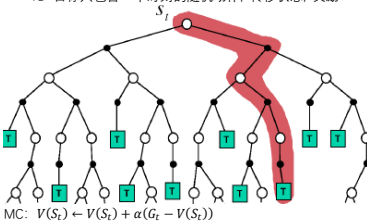
真实的 TD 目标  $R_{t+1} + \gamma v_\pi(S_{t+1}) = v_\pi(S_t)$  的无偏估计

TD 目标  $R_{t+1} + \gamma V(S_{t+1}) = v_\pi(S_t)$  的有偏估计

TD 目标与回报的方差差小很多

回报的计算包含整个序列中的随机动作, 转移状态, 奖励

TD 目标只包含一个时刻的随机动作, 转移状态, 奖励



动态规划更新:  $V(S_t) \leftarrow E_t[R_{t+1} + \gamma V(S_{t+1})]$

自举法: 更新时包含一个预测量

MC 不使用自举 DP 使用自举 TD 使用自举

采样法: 使用采样的数据计算期望

MC 使用采样 DP 不使用采样 TD 使用采样

**等价模型:**

MC 收敛结果对应最小二乘误差, 最佳匹配观测的回报

$$\sum_{k=1}^T \sum_{i=1}^K (G_t^{(k)} - V(S_t^{(k)}))^2$$

TD(0) 收敛结果对应最小似然马尔可夫模型  $(S, A, \hat{P}, \hat{R}, \gamma)$

$$\hat{P}_{s,s'} = \frac{1}{N(s,a)} \sum_{k=1}^K 1(s_t^k, a_t^k, s_{t+1}^k = s, a, s')$$
$$\hat{R}_{s,s'} = \frac{1}{N(s,a)} \sum_{k=1}^K 1(s_t^k, a_t^k = s, a) r_t^k$$

$n=1$  (TD)  $G_t^{(1)} = R_{t+1} + \gamma V(S_{t+1})$

$n=2$   $G_t^{(2)} = R_{t+1} + \gamma R_{t+2} + \gamma^2 V(S_{t+2})$

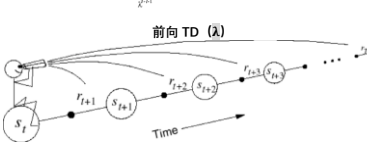
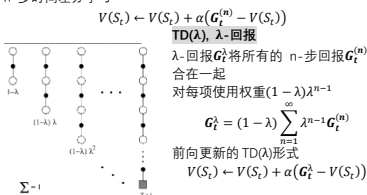
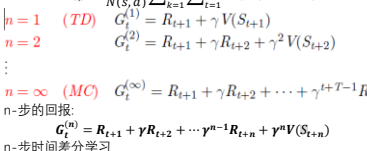
$\vdots$

$n=\infty$  (MC)  $G_t^{(\infty)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{t+T-1} R_T$

n-步的回报:

$$G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V(S_{t+n})$$

n-步时间差分学习



更新价值函数向  $\lambda$ -回报逼近, 面向未来时刻计算  $G_t^{(\lambda)}$ , 与 MC 一样要求完整的事件序列

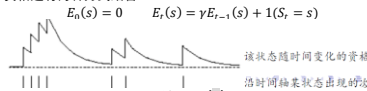
**资格迹 Eligibility Traces**

可信度: 警铃和闪光是否和高电压直接相关?

受频率启发: 对发生次数最多的状态给予最高的可信度

受近因启发: 对时间上最近发生的状态给予最高的可信度

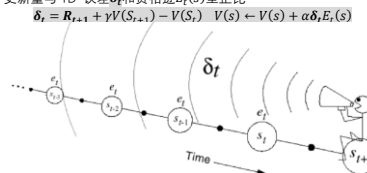
资格迹将两种方式结合

$$E_t(s) = 0 \quad E_t(s) = \gamma E_{t-1}(s) + 1(S_t = s)$$


后向 TD ( $\lambda$ )

为每个状态  $s$  记录它的资格迹, 更新  $s$  的价值  $V(s)$

更新量与 TD 误差  $\delta_t$  和资格迹  $E_t(s)$  成正比

$$\delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t) \quad V(s) \leftarrow V(s) + \alpha \delta_t E_t(s)$$


**TD(A) 和 TD(0)**

当  $\lambda=0$  时, 只有当前状态被更新

$$E_t(s) = 1(S_t = s) \quad V(s) \leftarrow V(s) + \alpha \delta_t E_t(s)$$

等同于 TD(0) 方法  $V(s) \leftarrow V(s) + \alpha \delta_t$

**TD(A) 和 MC**

当  $\lambda=1$  时, 可信度被推迟到事件结束, 考虑每个事件都到达终止状态的环境, 使用离线更新, 对每个事件序列, TD(1) 的所有更新和 MC 的所有更新是一致的

当使用离线更新时, 前向和后向 TD(A) 所有的更新是相同的

$$\sum_{t=1}^T \alpha \delta_t E_t(s) = \sum_{t=1}^T \alpha (G_t^{(\lambda)} - V(S_t)) 1(S_t = s)$$

**MC 和 TD(1)**

考虑在一次事件中  $s$  在  $k$  时刻被唯一访问一次

该状态的 TD(1) 资格迹会在被访问之后随时间衰减

$$E_t(s) = \gamma E_{t-1}(s) + 1(S_t = s) = \begin{cases} 0 & \text{if } t < k \\ (\gamma)^{t-k} & \text{if } t \geq k \end{cases}$$

TD(1) 在更新时会在在线累计误差

$$\sum_{t=1}^{T-1} \alpha \delta_t E_t(s) = \alpha \sum_{t=k}^{T-1} (\gamma)^{t-k} \delta_t = (G_k - V(S_k))$$

在事件结束时累计误差等于

$$\delta_k + \gamma \delta_{k+1} + \gamma^2 \delta_{k+2} + \dots + \gamma^{T-1-k} \delta_{T-1}$$

当  $\lambda=1$  时, TD 误差可以转换成 MC 误差

$$\delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$$
$$\delta_k + \gamma \delta_{k+1} + \gamma^2 \delta_{k+2} + \dots + \gamma^{T-1-k} \delta_{T-1} = G_t - V(S_t)$$

**TD(A) 和 TD(1)**

TD(1) 大致等价于每次经过的 MC 方法, 误差会随在线运行逐步累积, 如果价值函数只在事件结束时离线更新, 那么所有的更新量等同于 MC 方法

**前向和后向 TD(A)**

考虑在一次事件中  $s$  在  $k$  时刻被唯一访问一次

该状态的 TD(A) 资格迹会在被访问之后随时间衰减

$$E_t(s) = \gamma \lambda E_{t-1}(s) + 1(S_t = s) = \begin{cases} 0 & \text{if } t < k \\ (\lambda \gamma)^{t-k} & \text{if } t \geq k \end{cases}$$

后向 TD(A) 在更新时会在在线累计误差

$$\sum_{t=1}^{T-1} \alpha \delta_t E_t(s) = \alpha \sum_{t=k}^{T-1} (\lambda \gamma)^{t-k} \delta_t = (G_k^{(\lambda)} - V(S_k))$$

在事件结束时总共累计误差等于  $\lambda$ -回报

如果  $s$  被多次访问,  $E_t(s)$  也会累积多次的误差

前向提供理论原理, 后向给出算法实现: 在线学习; 每一时刻更新; 可以适用事件序列中的一小段, 非完整序列

目标: 获得最优策略, 基于模型: 动态规划、价值迭代、策略迭代

目标: 获得策略的价值函数, 基于样本: 迭代策略评估, 蒙特卡洛方法、时间差分学习 TD(A)

我们想获得最优策略, 但是 MDP 的模型未知, 或者问题空间太大只能考虑和智能体最相关的状态: 无模型方法; 使用智能体的经验样本: 最优价值函数

上节课: 无模型预测, 估计一个模型未知 MDP 问题的价值函数

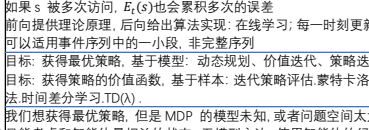
这节课: 无模型控制, 最优价值函数 MDP 问题的价值函数

在策略(on-policy) 和离策略(off-policy) 学习

在策略学习: 根据策略  $\pi$  产生的样本来学习关于  $\pi$  的相关知识

离策略学习: 根据另一策略  $\mu$  产生的样本来学习关于  $\pi$  的相关知识

策略迭代:



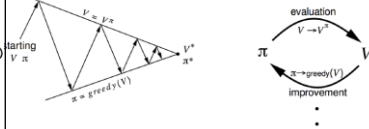
策略评估: 估计  $v_\pi$

- 例如迭代策略评估

策略提升: 生成新的  $\pi' \geq \pi$

- 例如贪心策略提升

基于蒙特卡洛的策略迭代:



策略评估: MC 策略评估,  $Q = q_\pi$

策略提升: 贪心策略提升

策略提升: MC 策略提升,  $Q = q_\pi$

策略提升: 贪心策略提升

动作-价值函数实现无模型的策略迭代:

基于  $V(s)$  进行贪心策略提升需要 MDP 的模型信息

$$\pi^*(s) = \arg \max_{a \in A} (R(s,a) + \gamma P_{s,a}^T V(S_t))$$

基于  $Q(s,a)$  进行贪心策略提升是无模型的

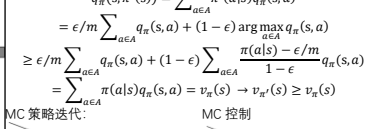
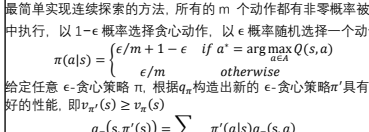
$$\pi^*(s) = \arg \max_{a \in A} Q(s,a)$$

**$\epsilon$ -贪心探索**

最简单实现连续探索的方法, 所有的  $m$  个动作都有非零概率被选中执行, 以  $1-\epsilon$  概率选择贪心动作, 以  $\epsilon$  概率随机选择一个动作

$$\pi(a|s) = \begin{cases} \epsilon/m + 1-\epsilon & \text{if } a = \arg \max_{a \in A} Q(s,a) \\ \epsilon/m & \text{otherwise} \end{cases}$$

给定任意  $\epsilon$ -贪心策略  $\pi$ , 根据  $q_\pi$  构造出新的  $\epsilon$ -贪心策略  $\pi'$  具有更好的性能, 即  $v_{\pi'}(s) \geq v_\pi(s)$

$$q_\pi(s, \pi'(s)) = \sum_{a \in A} \pi'(a|s) q_\pi(s,a)$$
$$= \epsilon/m \sum_{a \in A} q_\pi(s,a) + (1-\epsilon) \arg \max_{a \in A} q_\pi(s,a)$$
$$\geq \epsilon/m \sum_{a \in A} q_\pi(s,a) + (1-\epsilon) \sum_{a \in A} \frac{\pi(a|s) - \epsilon/m}{1-\epsilon} q_\pi(s,a)$$
$$= \sum_{a \in A} \pi(a|s) q_\pi(s,a) = v_\pi(s) \rightarrow v_{\pi'}(s) \geq v_\pi(s)$$


MC 策略迭代:

MC 控制

策略评估: MC 策略评估,  $Q = q_\pi$

策略提升: 贪心策略提升

策略提升: 贪心策略提升

策略提升: 贪心策略提升

无限探索, 无穷时收敛为贪心策略 (Greedy in the Limit with Infinite Exploration, GLIE)

智能体能无限次数地探索所有的状态-动作对:  $\lim_{k \rightarrow \infty} N_k(s,a) = \infty$

策略在无穷时收敛到贪心策略

$$\lim_{k \rightarrow \infty} \pi_k(a|s) = 1(a = \arg \max_{a \in A} Q_k(s,a))$$

$\epsilon$ -贪心策略的探索率  $\epsilon$  能随时间衰减到零  $\epsilon_k = \frac{1}{k}$ , 则它就满足 GLIE

GLIE 蒙特卡洛控制

使用策略  $\mu$  采第  $k$  次事件  $\{S_1, A_1, R_1, \dots, S_T\} \sim \pi$  对事件中的每个状态  $s$  和动作  $A_t$

$$N(S_t, A_t) \leftarrow N(S_t, A_t) + 1$$
$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{1}{N(S_t, A_t)} (G_t - Q(S_t, A_t))$$

基于新得到的动作-价值函数对策略进行提升

$$\epsilon \leftarrow \frac{1}{k} \quad \pi \leftarrow \text{greedy}(Q)$$

GLIE 蒙特卡洛控制能够收敛到最优动作-价值函数

$$Q(s,a) \leftarrow q_*(s,a)$$

**MC vs TD 控制**

TD 学习和 MC 方法相比有如下优势:

低方差, 在线学习, 无完整事件序列

在控制方法中使用 TD 替代 MC:

基于 TD 方法训练  $Q(S,A)$ , 使用  $\epsilon$ -贪心方法进行策略提升, 每个时刻都对应状态更新

**基于 Sarsa 方法更新动作-价值函数**

动作-价值函数的贝尔曼期望方程:

$$q_\pi(s,a) = R_s^a + \gamma \sum_{s',a'} P_{ss'}^{aa'} \pi(a'|s') q_\pi(s',a')$$

在某一时刻智能体观察到的一段事件  $(S,A,R,S',A')$  基于样本的 TD 更新:

$$Q(S,A) \leftarrow Q(S,A) + \alpha (R + \gamma Q(S',A') - Q(S,A))$$

基于 Sarsa 的在策略控制:

在每一时刻:

策略评估:

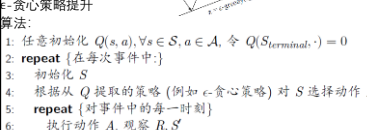
使用 Sarsa 方法,  $Q \approx q_\pi$

策略提升:

$\epsilon$ -贪心策略提升

算法:

1. 任意初始化  $Q(s,a), \forall s \in S, a \in A$ , 令  $Q(S_{\text{terminal}}, \cdot) = 0$
2. repeat {在每次事件中:}
3. 初始化  $S$
4. 根据从  $Q$  提取的策略 (例如  $\epsilon$ -贪心策略) 对  $S$  选择动作  $A$
5. repeat {对事件中的每一时刻:}
6. 执行动作  $A$ , 观察  $R, S'$
7. 根据从  $Q$  提取的策略 (例如  $\epsilon$ -贪心策略) 对  $S'$  选择动作  $A'$
8.  $Q(S,A) \leftarrow Q(S,A) + \alpha [R + \gamma Q(S',A') - Q(S,A)]$
9.  $S \leftarrow S', A \leftarrow A'$
10. until  $S$  是终止状态
11. until  $Sarsa$  的收敛性:



Sarsa 的收敛性:

如下条件满足时 Sarsa 算法是收敛到最优动作-价值函数:

策略序列  $\pi_t(s)$  是 GLIE 的

更新步长  $\alpha_t$  满足 Robbins-Monro 序列要求

$$\sum_{t=1}^{\infty} \alpha_t = \infty \quad \sum_{t=1}^{\infty} \alpha_t^2 < \infty$$

**Sarsa vs MC**

注意: 蒙特卡洛方法并不能直接应用在这个问题上

并不是所有的策略都能到达终止状态, 如果学到的某一策略使智能体停留在原地不动, MC 在基于该策略对事件采样时永远都无法结束

Sarsa 算法没有这样的问题

每一步都在学习, 很快就会发现这样的策略性能很差, 然后转向其它策略

**n-步 Sarsa**

考虑 n-步的回报,  $n = 1, 2, \dots$ :

$n=1$  (Sarsa)  $q_t^{(1)} = R_{t+1} + \gamma Q(S_{t+1}, A_{t+1})$

$n=2$   $q_t^{(2)} = R_{t+1} + \gamma R_{t+2} + \gamma^2 Q(S_{t+2}, A_{t+2})$

$n=\infty$  (MC)  $q_t^{(\infty)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{t+T-1} R_T$

定义 n-步的 Q-回报

$$q_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n Q(S_{t+n}, A_{t+n})$$

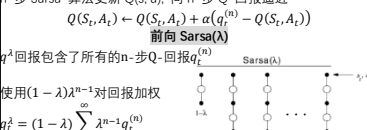
n-步 Sarsa 算法更新  $Q(s,a)$ , 向 n-步 Q-回报逼近

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha (q_t^{(n)} - Q(S_t, A_t))$$

**前向 Sarsa(X)**

$q_t^{(n)}$  回报包含了所有的 n-步 Q-回报  $q_t^{(n)}$

使用  $(1-\lambda)\lambda^{n-1}$  对回报加权

$$q_t^{(n)} = (1-\lambda) \sum_{i=1}^n \lambda^{i-1} q_t^{(i)}$$


前向 Sarsa(X) 算法

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha (q_t^{(n)} - Q(S_t, A_t))$$

和 TD(X) 算法一样, 我们将策略迹引入到在线的控制算法中

但是 Sarsa(X) 算法不同的地方在于为每个状态-动作对都记录资格迹-之前只是对状态记录

$$E_t(s,a) = 0 \quad E_t(s,a) = \gamma \lambda E_{t-1}(s,a) + 1(S_t = s, A_t = a)$$

在每一时刻会对每个状态-动作对更新  $Q(s,a)$  更新

更新量与 TD 误差  $\delta_t$  和资格迹  $E_t(s,a)$  成正比

$$\delta_t = R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)$$

**后向 Sarsa(X)**

一步 Sarsa 只对最终导致高奖励的最后一步动作强化它的价值

资格迹方法能够对事件中的多个动作强化它们的价值

强化的幅度(箭头大小) 随着离奖励步数的增加而衰减

衰减率等于  $\lambda$

这里  $\lambda = 1, \lambda = 0.9$

**高策略学习**

想要对目标策略  $\mu(a|s)$  进行评估, 计算  $v_\mu(s)$  或  $q_\mu(s,a)$ , 但是智能体实际的行为策略是  $\mu(a|s)$ :  $\{S_1, A_1, R_1, \dots, S_T\} \sim \mu$

为什么要考虑这种情况?

有时智能体需要观察人类或别的智能体的行为去学习

有时需要重复利用旧策略  $\pi_1, \pi_2, \dots, \pi_{t-1}$  产生的经验去学习

执行探索性的策略去学习最优策略

执行单一的策略去学习多个策略

**Q-学习**

考虑基于动作-价值  $Q(s,a)$  的离策略学习, 不再使用重要性采样

智能体下一时刻执行的动作是由行为策略产生  $A_{t+1} \sim \mu(S_t)$

但是学习算法考虑的是由另一个目标策略产生的后继动作

更新  $Q(S_t, A_t)$  向另一个后继动作的价值逼近

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha (R_{t+1} + \gamma Q(S_{t+1}, A') - Q(S_t, A_t))$$

Q-学习实现离策略控制:

现在允许行为策略和目标策略都能提升

目标策略  $\pi$  是由  $Q(s,a)$  生成的贪心策略

$$\pi(S_{t+1}) = \arg \max_{a \in A} Q(S_{t+1}, a)$$

行为策略  $\mu$  是由  $Q(s,a)$  生成的例如  $\epsilon$ -贪心策略

**那么 Q-学习的目标就是**

$$R_{t+1} + \gamma Q(S_{t+1}, A') = R_{t+1} + \gamma \arg \max_{a \in A} Q(S_{t+1}, a)$$
$$= R_{t+1} + \gamma \max_{a \in A} Q(S_{t+1}, a)$$

Q-学习控制算法

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha (R + \gamma \max_{a \in A} Q(S', a') - Q(S_t, A_t))$$

Q-学习控制能够收敛到最优动作-价值函数。

Q-学习算法实现高策略控制:

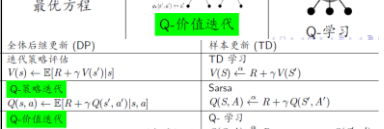
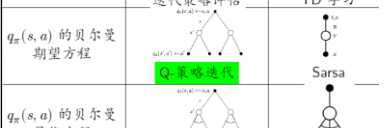
1. 任意初始化  $Q(s,a), \forall s \in S, a \in A$ , 令  $Q(S_{\text{terminal}}, \cdot) = 0$
2. repeat {在每次事件中:}
3. 初始化  $S$
4. repeat {对事件中的每一时刻:}
5. 根据从  $Q$  提取的策略 (例如  $\epsilon$ -贪心策略) 对  $S$  选择动作  $A$
6. 执行动作  $A$ , 观察  $R, S'$
7.  $Q(S,A) \leftarrow Q(S,A) + \alpha [R + \gamma \max_{a \in A} Q(S', a) - Q(S,A)]$
8.  $S \leftarrow S'$
9. until  $S$  是终止状态
10. until

**Sarsa vs Q-学习**

Q-学习正确地学到了最优路径, 即沿着悬崖的边缘行走, 但是由于使用的  $\epsilon$ -贪心策略会采取一定概率的随机动作, 有时智能体跌入悬崖

Sarsa 学到的是安全路径, 在学习过程中考虑了随机探索动作的影响, 即使 Sarsa 学到的安全路径比 Q-学习的最优路径行走步数要多, 但是每次获得的奖励和对比 Q-学习的高

动态规划和时间差分学习的关系



探索(exploration) 与利用(exploitation)

在线决策过程存在的一个关键问题是:

利用: 根据当前的信息做出最佳的决策

探索: 采样更多的信息

想要做出长期的最佳决策有时需要牺牲当前短期的部分利益

收集更多的信息从而做出整体最佳的决策

Q-学习

如果 Q-学习没有探索, 即  $\epsilon = 0$ , 学习结果容易陷入局部最优解, 无法找到真正最优策略, 如果  $\epsilon$  很大也不好

1. 状态、动作空间很大时, 想要探索整个空间是费时费力的, 对在线学习过程不利

2. 随机的动作会降低智能体的实际回报

非均匀探索概率

$\epsilon$ -贪心策略除了最优动作, 其它动作以等概率随机选择

假如: 有两个动作看起来不错, 而其它的动作完全不可行

那么比较合理的探索方式是选择那些看起来不错的动作, 从不错的动作中选出最优的动作, 而不是将精力浪费在看起来不可行的动作, 根据动作-价值决定动作被选中的概率

波尔兹曼探索, Boltzmann exploration

$$\pi(a|s) = \frac{e^{Q(s,a)/T}}{\sum_{a \in A} e^{Q(s,a)/T}}$$

$\alpha$  被选中的概率与  $e^{Q(s,a)/T}$  成正比 ( $Q(s,a) \geq 0, \forall s, a$ )

温度系数  $T > 0$  决定了策略的随机性能

如果  $T$  很大, 所有动作几乎以等概率选择 (探索)

如果  $T$  很小, 高价值的动作更容易被选中 (利用)

极端情况  $T \rightarrow 0$ , 只选择最优动作

求解非线性算子 max

动态规划: 价值迭代/策略迭代

模型已知

预测(prediction): 蒙特卡洛方法, 时间差分学习, TD(A)

控制(control): Sarsa, Sarsa(X), Q-学习

算法	查表法	线性逼近	非线性逼近
MC 控制	✓	(✓)	×
Sarsa	✓	(✓)	×
Q 学习	✓	×	×
LSPI	✓	(✓)	-

(✓) 代表算法会在最优价值函数附近振荡

最小二乘策略迭代:

策略评估: 使用最小二乘策略的 Q 函数

策略提升: 贪心策略提升

在/离策略 算法 查表法 线性逼近 非线性逼近

在策略	算法	查表法	线性逼近	非线性逼近
MC	✓	✓	✓	✓
LSMC	✓	✓	✓	✓
TD	✓	✓	✓	×
LSTD	✓	✓	✓	×
TD	✓	✓	×	×
LSTD	✓	✓	×	×

**批处理强化学习, Batch RL:**

梯度下降法简单, 易实现, 但是对样本的利用率不足

批处理的方法能够基于智能体的经验(训练数据), 试着找到最合适价值函数

**使用逼近器的好处**

使用较少数量的参数就能表达复杂的函数 (计算复杂度)

对于一个权重的调整就可以影响到很多的点 (泛化能力)

多种特征表示和逼近器结构(多样性)

可微的函数逼近器, 例如: 特征函数的线性组合, 神经网络, 决策树, 最近邻, 傅里叶/小波基函数