

Lec. 1

State / State space: $\{S_i\}_{i=1}^n = S$ Action / Action space: $\{a_i, a_{i+1}\} = A(s_i)$ State Transition: $P(s'|s, a) = P(s_{t+1} = s' | s_t = s, a_t = a)$

$$\sum_{s' \in S} P(s'|s, a) = 1$$

Policy: $\pi(a|s) = P(A=a | S=s)$

Return: 看到 trajectory 和 Reward 和

Discounted Return: $G_t = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots$

Lec 2. Bellman Equation:

State value: $V_{\pi}(s) = V(s) = E[\gamma G_t | S_t = s]$

$$\Rightarrow G_t = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots = G_t = V_{\pi}(s) + \gamma G_{t+1}$$

$$\therefore G_t(s) = E[V_{\pi}(s') | S_{t+1} = s'] = E[V_{\pi}(s') | S_t = s, A_t = a]$$

$$= E[V_{\pi}(s') | S_t = s] = \sum_a \pi(a|s) E[V_{\pi}(s') | S_t = s, A_t = a]$$

$$= \sum_a \pi(a|s) \sum_{s'} P(s'|s, a) \cdot r$$

$$E[G_t(s)] = \sum_{s'} P(s'|s) \cdot E[G_t(s') | S_t = s'] = \sum_{s'} P(s'|s) \cdot \frac{E[G_t(s') | S_t = s']}{V_{\pi}(s')}$$

$$= \sum_{s'} \sum_a \pi(a|s) \cdot P(s'|s, a) \cdot V_{\pi}(s') = \sum_a \pi(a|s) \cdot \sum_{s'} P(s'|s, a) \cdot V_{\pi}(s')$$

$$\therefore V_{\pi}(s) = \sum_a \pi(a|s) \left[\sum_{s'} P(s'|s, a) \cdot r + \gamma \sum_{s'} P(s'|s, a) \cdot V_{\pi}(s') \right] \quad (\text{BE的 Element-wise form})$$

矩阵形式 (Matrix-Vector Form):

$$V_{\pi} = \begin{bmatrix} V_{\pi}(s_1) \\ V_{\pi}(s_2) \\ \vdots \\ V_{\pi}(s_n) \end{bmatrix} \quad r_{\pi} = \begin{bmatrix} r_{\pi}(s_1) \\ r_{\pi}(s_2) \\ \vdots \\ r_{\pi}(s_n) \end{bmatrix} = \begin{bmatrix} \vdots \\ \sum_a \pi(a|s_1) \cdot \sum_{s'} P(s'|s_1, a) \cdot r \\ \vdots \\ \sum_a \pi(a|s_n) \cdot \sum_{s'} P(s'|s_n, a) \cdot r \end{bmatrix}$$

$$P_{\pi} \in R^{n \times n}: [P_{\pi}]_{ij} = P(s_j | s_i) = \sum_a \pi(a|s_i) \cdot P(s_j | s_i, a)$$

$$\Rightarrow V_{\pi} = r_{\pi} + \gamma P_{\pi} \cdot V_{\pi}$$

Lec 3. Bellman Optimal Equation

Optimal Policy: $V_{\pi^*}(s) \geq V_{\pi}(s) \quad \forall s \in S, \forall \pi$

$$\text{BE: } V_{\pi}(s) = \sum_a \pi(a|s) \left[\underbrace{\sum_{s'} P(s'|s, a) \cdot r + \gamma \sum_{s'} P(s'|s, a) \cdot V_{\pi}(s')}_{q_{\pi}(s, a)} \right]$$

Bob: $V(s) = \max_a q_{\pi}(s, a)$ 把 π 代入这个式子算入≈ π 把价值量最大化 $q_{\pi}(s, a)$

$$\Rightarrow \max_a \sum_a q_{\pi}(s, a) = \max_a q_{\pi}(s, a)$$

$$\Rightarrow V(s) = \max_a \left[\sum_a P(s|s, a) \cdot r + \gamma \sum_{s'} P(s'|s, a) \cdot V(s') \right]$$

Matrix-Vector Form: $V = \max_a (r_{\pi} + \gamma P_{\pi} V)$

$$f(V) = \max_a (r_{\pi} + \gamma P_{\pi} V) \quad V = f(W)$$

由压缩映射定理: 有解 $\|f(W) - f(V)\| \leq \epsilon \|\omega - V\| \leq \gamma \epsilon \|V\| \leq \gamma \epsilon \|V^*\|$ ∴ 存在唯一不动点 $V^* = f(V^*)$ ∵ 算法收敛

$$\text{值迭代求解 BE: } \text{①: } q_{\pi}(s, a) = \sum_{s'} q(s', a) + \gamma \sum_{s'} P(s'|s, a) \cdot V(s')$$

$$\text{②: } V_{\pi} = \max_a q_{\pi}(s, a)$$

优化结果:

$$\pi^*(s, a) = \begin{cases} 1 = \arg \max_a q^*(s, a) \\ 0 \text{ else} \end{cases}$$

奖励函数不变性: 对奖励进行线性变换 $R' = \alpha R + \beta$ ($\alpha > 0$)

最优策略不变

∴ 解析解: $V_{\pi} = (I - \gamma P_{\pi})^{-1} r_{\pi}$ 迭代解: $V_{\pi+1} = r_{\pi} + \gamma P_{\pi} V_{\pi}$ 收敛性证明: $\left\{ \begin{array}{l} \text{① } V_{\pi+1} = r_{\pi} + \gamma P_{\pi} V_{\pi} \\ \text{② } V_{\pi} = r_{\pi} + \gamma P_{\pi} V_{\pi} \end{array} \right.$

$$\text{①-②: } V_{\pi+1} - V_{\pi} = \gamma P_{\pi} (V_{\pi+1} - V_{\pi}) \quad \delta_{\pi+1} = \gamma P_{\pi} \delta_{\pi}$$

$$\therefore \delta_{\pi+1} = \gamma P_{\pi}^T P_{\pi} \delta_{\pi} \quad P_{\pi} \text{ 为随机矩阵 } V_{\pi} \in [0, 1]^n$$

矩阵 ϵ

$$\therefore \delta_{\pi+1} \rightarrow 0$$

Action Value (动作价值): $q_{\pi}(s, a) = E[G_t | S_t = s, A_t = a]$

$$V_{\pi}(s) = \sum_a q_{\pi}(s, a) \cdot q_{\pi}(s, a)$$

如果知道所有 Action Value → 能得到这个状态的 State Value

如果知道所有状态的 State Value → 能求出所有 Action Value

Lec 4. 值迭代 / 值迭代

值迭代: $V_{\pi+1} = f(V_{\pi}) = \max_{\pi} (r_{\pi} + \gamma P_{\pi} V_{\pi})$

$$\text{Step 1: } \pi_{\text{optimal}} = \text{argmax}_{\pi} (r_{\pi} + \gamma P_{\pi} V_{\pi}) \quad (\text{PI})$$

$$\text{Step 2: } V_{\pi+1} = r_{\pi_{\text{optimal}}} + \gamma P_{\pi_{\text{optimal}}} \cdot V_{\pi} \quad (\text{PE})$$

策略迭代: $\pi_0 \xrightarrow{\text{PE}} V_{\pi_0} \xrightarrow{\text{PI}} \pi_1 \xrightarrow{\text{PE}} V_{\pi_1} \xrightarrow{\text{PI}} \dots$ Step 1. (PI): 对应当前策略 π_0 对应真实 V_{π_0}

$$\text{Bellman Equation: } V_{\pi_0} = r_{\pi_0} + \gamma P_{\pi_0} V_{\pi_0}$$

Step 2 (PI): 根据 V_{π_0} 寻找更好策略 π_{optimal}

$$\pi_{\text{optimal}} = \text{argmax}_{\pi} (r_{\pi} + \gamma P_{\pi} V_{\pi_0})$$

PE 和值迭代 (固定策略的值迭代)

$$V_{\pi_{\text{optimal}}}^{(1)} = r_{\pi_0} + \gamma P_{\pi_0} V_{\pi_0}^{(0)}$$

$$\vdots \rightarrow V_{\pi_{\text{optimal}}}^{(1)} \rightarrow V_{\pi_{\text{optimal}}}$$

可证: $V_{\pi_{\text{optimal}}} \geq V_{\pi_0}$

$$\text{Step 1: } V_{\pi_{\text{optimal}}} \geq V_{\pi_0} \Rightarrow V_{\pi_{\text{optimal}}} + \gamma P_{\pi_{\text{optimal}}} V_{\pi_0} \geq V_{\pi_0} + \gamma P_{\pi_0} V_{\pi_0} = V_{\pi_0} \quad 0$$

$$\Delta_0 = V_{\pi_{\text{optimal}}} - V_{\pi_0} \Rightarrow \Delta_0 = (V_{\pi_{\text{optimal}}} + \gamma P_{\pi_{\text{optimal}}} V_{\pi_0}) - V_{\pi_0}$$

$$= V_{\pi_{\text{optimal}}} + \gamma P_{\pi_{\text{optimal}}} V_{\pi_0} - \gamma P_{\pi_0} V_{\pi_0} = V_{\pi_0} - V_{\pi_0}$$

$$= \gamma P_{\pi_{\text{optimal}}} (V_{\pi_{\text{optimal}}} - V_{\pi_0}) + (V_{\pi_{\text{optimal}}} - \gamma P_{\pi_{\text{optimal}}} V_{\pi_0} - V_{\pi_0})$$

$$\therefore \Delta_0 = V_{\pi_{\text{optimal}}} - \gamma P_{\pi_{\text{optimal}}} V_{\pi_0} - V_{\pi_0} \quad \Delta_0 \geq 0$$

$$\therefore \Delta_0 = \gamma P_{\pi_{\text{optimal}}} \Delta_0 + \delta = \gamma (P_{\pi_{\text{optimal}}} \Delta_0 + \delta) = \delta + \gamma P_{\pi_{\text{optimal}}} \delta + (P_{\pi_{\text{optimal}}} \delta) \delta + \dots$$

$$\therefore \delta \geq 0 \quad \gamma > 0 \quad \Delta_0 \geq 0 \quad \therefore V_{\pi_{\text{optimal}}} \geq V_{\pi_0}$$

截断策略迭代:

PE 阶段迭代固定次, 不像值迭代一样 PE 阶段迭代一次

也不像策略迭代一样 PE 阶段迭代无数次