

强化学习

- MDP 马尔可夫决策过程
- 值函数
- 值迭代
- 策略迭代

监督 → 训练集

学习理论

非监督 → 未标记数据

水平监督, 比

强化学习

奖励信号

→ 动态下棋 → 难以分配到具体步

$MDP(S, A, \{P_{sa}\}, \gamma, R)$

S - 状态集

A - 动作集

P_{sa} - 状态转移分布

$$\sum_{s'} P_{sa}(s') = 1 \quad P_{sa}(s') \geq 0$$

γ - discount factor ($0 \leq \gamma \leq 1$)

R 奖励函数 $R: S \rightarrow R$

例 1

	1	2	3	4
3				+1
2				-1
1				

$S = 11$ 个

$A = \{\uparrow, \downarrow, \leftarrow, \rightarrow\}$

↑ 指定 80%
← 10% → 10%

$$P_{(3,1), \uparrow}[(3,2)] = 0.8$$

$$P_{(3,1), \uparrow}[(4,1)] = 0.1$$

$$P_{(3,1), \uparrow}[(2,1)] = 0.1$$

$$P_{(3,1), \uparrow}(\dots) = 0$$

$$R(4,3) = +1$$

$$R(4,2) = -1$$

$$R(s) = -0.02 \text{ other}$$

$$S_0 \xrightarrow{a_0} S_1 \sim P_{S_0 a_0}$$

$$\xrightarrow{a_1} S_2 \sim P_{S_1 a_1}$$

$$R(S_0) + \gamma R(S_1) + \gamma^2 R(S_2) + \gamma^3 R(S_3)$$

$$0 \leq \gamma < 1$$

$$\sum_{t=0}^{\infty} \gamma^t R(S_t) \text{ 记 } E(\downarrow) \text{ 期望}$$

policy $\pi: S \rightarrow A$

\rightarrow	\rightarrow	\rightarrow	\neq
\uparrow	\swarrow	\uparrow	-1
\uparrow	\leftarrow	\leftarrow	\leftarrow

10% 概率
巧女计中, 避免 -1

定义 V^π, V^*, π^*
最佳 π^*

$\forall \pi$, 定义值函数 $V^\pi: S \rightarrow \mathbb{R}$

$V^\pi(s) =$ 预期执行 π 总收益

$$= E[R(s_0) + \gamma V^\pi(s_1) + \gamma^2 R(s_2) + \dots \mid \pi, s_0 = s]$$

$$= E[R(s) + \gamma (R(s_1) + \gamma R(s_2) + \dots) \mid \pi, s = s_0]$$

$V^\pi(s_1)$

$$V^\pi(s) = R(s) + \gamma \sum_{s_1} P_{s_0 \pi(s)}(s_1) V^\pi(s_1)$$

贝尔曼方程

固定 $\pi \rightarrow V^\pi = ?$

$P:$

→	→	→	+
↓	///	→	-
→	→	↑	↑

 $V^\pi((3,1)) = R((3,1)) + \gamma [0.8 V^\pi((3,2)) + 0.1 V^\pi((4,1)) + 0.1 V^\pi((2,1))]$

每点

↓ 11个等式约束 → 解出 V^π

$$V^*(s) = \max_{\pi} V^\pi(s)$$

V^* 贝尔曼最优

$$V^*(s) = R(s) + \max_a \gamma \sum_{s_1} P_{sa}(s_1) V^*(s_1)$$

$$\pi^*(s) = \arg \max_a \sum_{s_1} P_{sa}(s_1) V^*(s_1)$$

↓ 最优

$\therefore P^*$

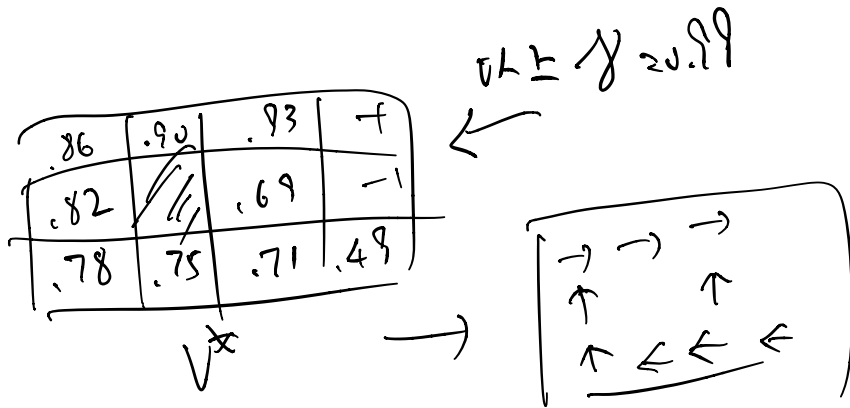
值迭代 用 Bellman 方程迭代最优

$$V(s) = 0$$

$$\forall s \quad V(s) := R(s) + \max_a \gamma \sum_{s_1} P_{sa}(s_1) V(s_1)$$

$$V(s) \xrightarrow{\text{收敛}} V^*(s) \quad // \quad V := B(V)$$

每次只更新1个, 异步 (2+2?)



Policy iteration

random π

Repeat: $V := V^\pi$ ($V = B(V)$)

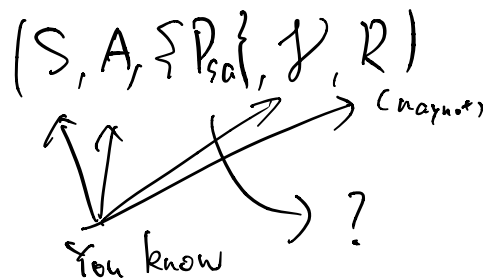
if 成本 < 1000
 至少 Policy 迭代
 至少 1000 次迭代

let $\pi(s) := \arg\max_a \sum_{s_1} P_{sa}(s_1) V(s_1)$

$V \rightarrow V^*$ $\pi \rightarrow \pi^*$

凸函数吗? (... 老师说凸函数不知道
 但可以证收敛, 证明不难但长, 不写)

如果已知转移概率?



$$P_{sa}(s_i) = \frac{\text{到 } s_i \text{ by } s \text{ and action } a}{s \text{ and action } a}$$

有很多混合处理方法

Repeat {

在MDP 中 用 π action 得到结果

update P_{sa} 估计

解 Bellman 方程 得 V (估计)

$$\pi_i = \arg \max_{\pi} \sum_{s_i} P_{sa}(s_i) V(s_i)$$

}