



中国科学技术大学  
University of Science and Technology of China

# 马尔可夫决策过程

吉建民

[jianmin@ustc.edu.cn](mailto:jianmin@ustc.edu.cn)

**01** 马尔可夫过程 (Markov Process)

02 马尔可夫奖励过程 (Markov Reward Process)

03 马尔可夫决策过程 (Markov Decision Process)

04 部分可观察马尔可夫决策过程 (Partially observable MDPs)

# 目录

- **随机过程 (Stochastic Process)**：随时间演化的一连串随机事件，每个时间点都对应一个可能出现多种取值的随机变量
  - 概率论的研究对象是静态的随机现象，而随机过程的研究对象是随时间演变的随机现象
    - 例如，天气随时间的变化、城市交通随时间的变化
  - 在随机过程中，随机现象在某时刻  $t$  的取值是一个随机变量，用  $s_t$  表示
  - 随机过程可以表示为： $\{S_t : t \in T\}$
  - 已知历史信息  $(S_1, \dots, S_t)$  时，下一个时刻状态为  $S_{t+1}$  的概率为

$$P(S_{t+1} | S_1, \dots, S_t)$$

# 马尔可夫性质

- 一个随机过程被称为具有**马尔可夫性质 (Markov property)**，当且仅当任意时刻的状态只取决于上一时刻的状态，用公式表示为

$$P(S_{t+1}|S_t) = P(S_{t+1}|S_1, \dots, S_t)$$

- 当前状态是未来的充分统计量，即下一个状态只取决于当前状态，而不会受到过去状态的影响
- 马尔可夫性可以简化运算，只要当前状态可知，所有的历史信息都不再需要了，利用当前状态信息就可以决定未来

# 状态转移矩阵

- 对一个马尔科夫状态  $s$  及其后继状态  $s'$ , 状态转移概率 (State Transition Probability) 定义为  $P_{ss'} = P(S_{t+1} = s' | S_t = s)$
- **状态转移矩阵 (State Transition Matrix)**  $\mathcal{P}$  则定义了从所有状态  $s$  到其所有后继状态  $s'$  的转移概率, 即

$$\mathcal{P} = \begin{array}{c} \text{to} \\ \left[ \begin{array}{ccc} \mathcal{P}_{11} & \dots & \mathcal{P}_{1n} \\ \vdots & & \\ \mathcal{P}_{n1} & \dots & \mathcal{P}_{nn} \end{array} \right] \\ \text{from} \end{array}$$

其中矩阵的每一行元素之和都等于 1

- 马尔可夫过程是具有马尔可夫性质的随机过程

## 定义

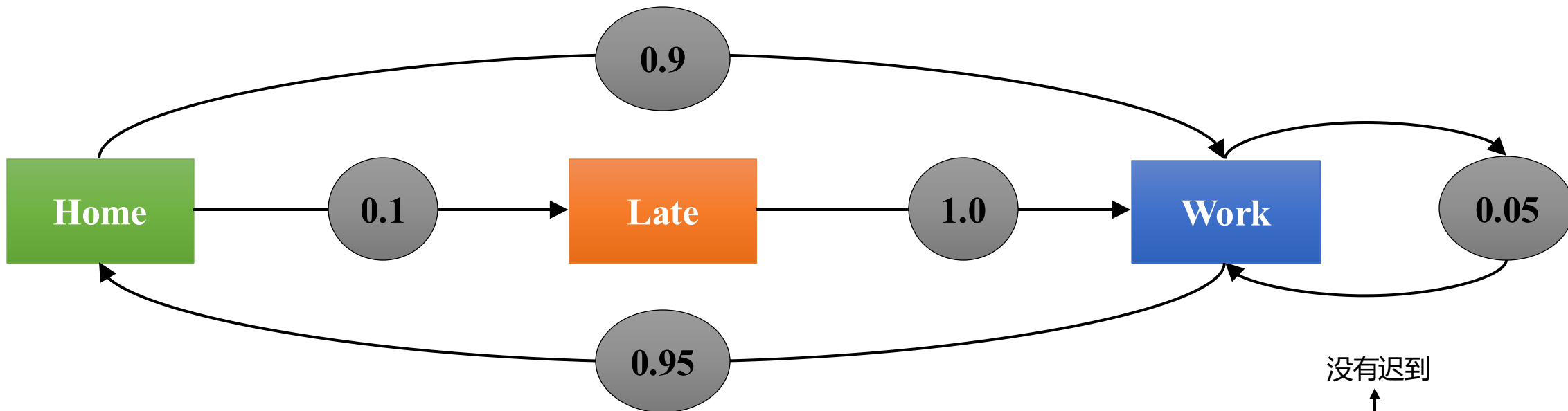
一个**马尔可夫过程 (Markov Process)** 或**马尔可夫链 (Markov Chain)**

可以表示为一个二元组  $\langle S, \mathcal{P} \rangle$ , 其中:

- $S$  是状态的 (有限) 集合
- $\mathcal{P}$  是状态转移矩阵,  $P_{ss'} = P(S_{t+1} = s' | S_t = s)$

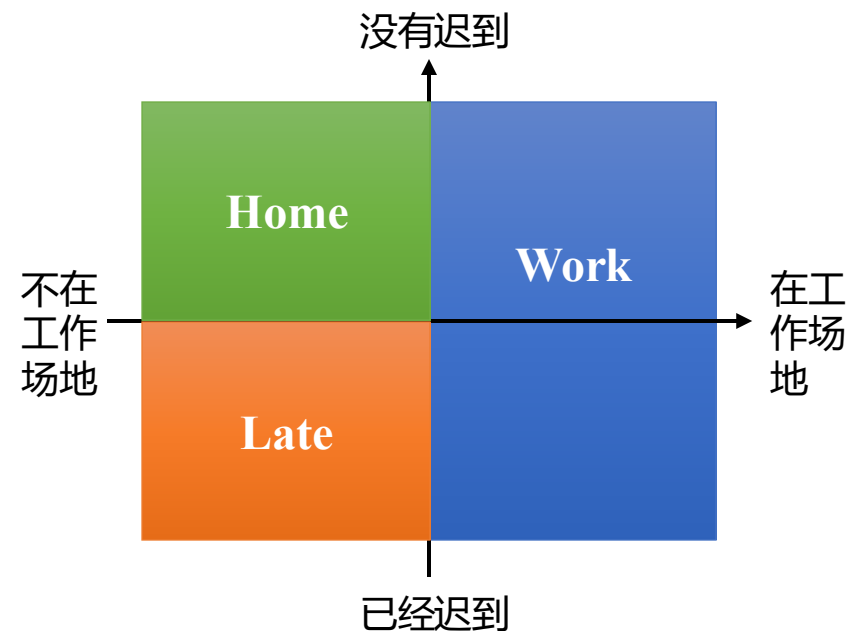
# 示例：通勤问题

□ 通勤问题中三个状态，Home, Late, Work, 构成马尔可夫链



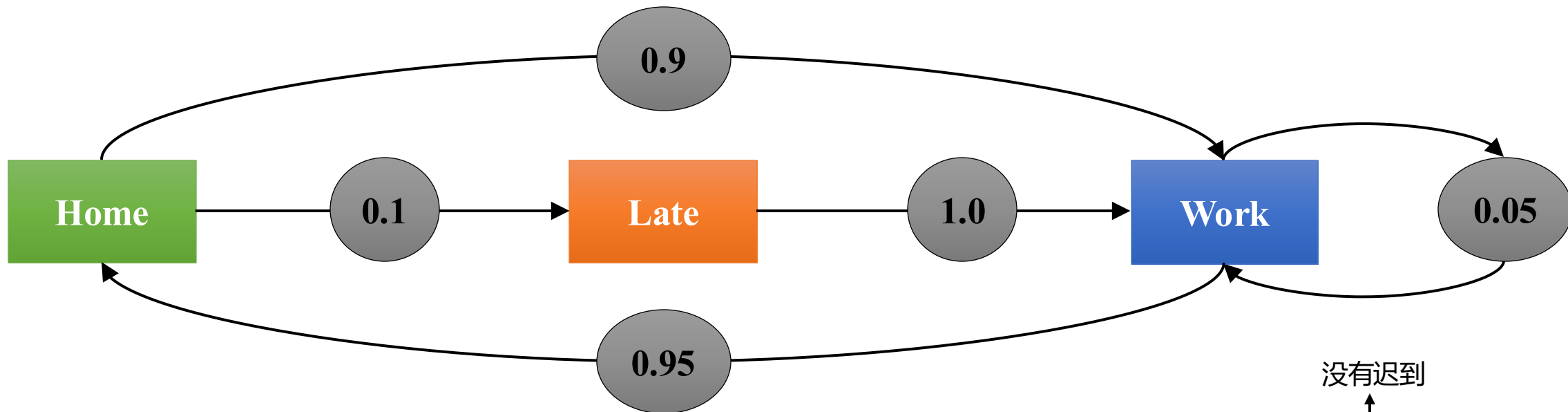
针对 (Home, Late, Work) 的状态转移矩阵为：

$$P_{ss'} = \begin{bmatrix} 0.0 & 0.1 & 0.9 \\ 0.0 & 0.0 & 1.0 \\ 0.95 & 0.0 & 0.05 \end{bmatrix}$$



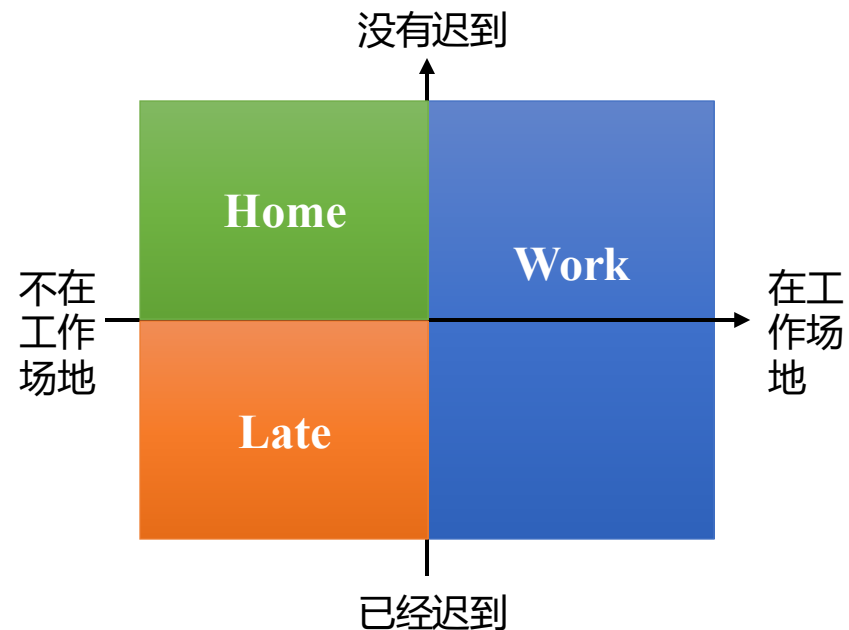
# 示例：通勤问题

□ 通勤问题中三个状态，Home, Late, Work, 构成马尔可夫链



□ 初始状态为 Home 的马尔可夫链样本序列示例：

- Home, Late, Work, Home, ...
- Home, Late, Work, Work, ...
- Home, Work, Work, Home, ...
- Home, Work, Work, Work, ...
- Home, Work, Home, Late, ...
- Home, Work, Home, Work, ...



01

马尔可夫过程 (Markov Process)

02

马尔可夫奖励过程 (Markov Reward Process)

03

马尔可夫决策过程 (Markov Decision Process)

04

部分可观察马尔可夫决策过程 (Partially observable MDPs)

# 目录

# 马尔可夫奖励过程

- 马尔可夫奖励过程是在马尔可夫链基础上增加了“价值（奖励）”

## 定义

一个**马尔可夫奖励过程** (Markov Reward Process, MRP) 可以表示为一个

四元组  $\langle \mathcal{S}, \mathcal{P}, \mathcal{R}, \gamma \rangle$ , 其中:

- $\mathcal{S}$  是状态的 (有限) 集合
- $\mathcal{P}$  是状态转移矩阵,  $P_{ss'} = P(S_{t+1} = s' | S_t = s)$
- $\mathcal{R}$  是奖励函数,  $R_s = E[R_{t+1} | S_t = s]$
- $\gamma$  是折扣因子,  $\gamma \in [0, 1]$

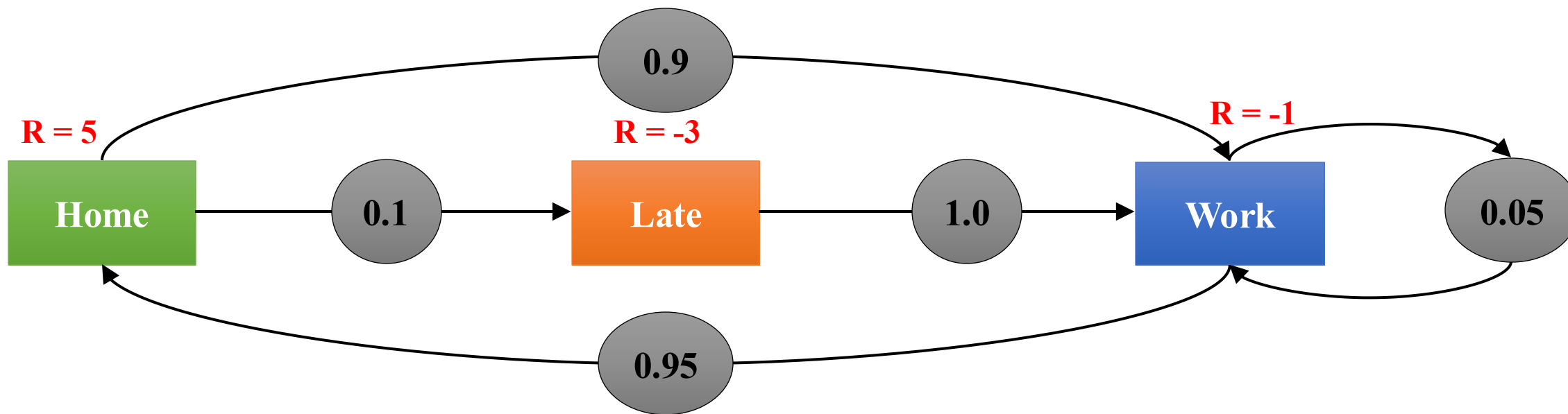
- 引入折扣因子是因为远期收益具有一定不确定性, 有时我们更希望能够尽快获得一些奖励, 所以我们需要对远期收益打一些折扣
  - 接近 1 的  $\gamma$  更关注长期的累计奖励, 接近 0 的  $\gamma$  更考虑短期奖励

# 示例：通勤问题

□ 通勤问题中三个状态, Home, Late, Work, 构成马尔可夫链

□ 通勤问题扩展为马尔可夫奖励过程

➤  $R(\text{Home}) = 5$ ,  $R(\text{Late}) = -3$ ,  $R(\text{Work}) = -1$



## 定义

在一个马尔可夫奖励过程 (MRP) 中, 在时间步  $t$  开始计算的**回报 (Return)**  $G_t$  指的是从时间步  $t$  开始累积的折扣奖励总和:

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

其中, 折扣因子  $\gamma \in [0, 1]$ ,  $R_t$  表示在时刻  $t$  获得的奖励。

在通勤问题中, 当  $\gamma = 0.5$ , 一个马尔可夫链样本序列 Home, Late, Work, Home 的回报为:

$$G_1 = 5 + 0.5 \times (-3) + 0.5^2 \times (-1) + 0.5^3 \times 5 = 5 - 1.5 - 0.25 + 0.625 = 3.875$$

此序列发生的概率为 (初始状态为 Home) :

$$\begin{aligned} P_1 &= P(\text{Home}) \times P(\text{Late}|\text{Home}) \times P(\text{Work}|\text{Late}) \times P(\text{Home}|\text{Work}) \\ &= 1 \times 0.1 \times 1 \times 0.95 = 0.095 \end{aligned}$$

## 定义

在一个马尔可夫奖励过程（MRP）中，**状态价值函数（State Value Function）**  $v(s)$  定义为从状态  $s$  开始时的期望回报：

$$v(s) = E[G_t | S_t = s]$$

价值函数  $v(s)$  表示状态  $s$  的长期价值，是从状态  $s$  开始的所有可能的马尔可夫链样本序列的回报**期望**，即，这些序列的回报乘以相应发生的概率，再进行累加。

若从状态  $s$  开始的所有可能的马尔可夫链样本序列的集合为  $T$ ，则

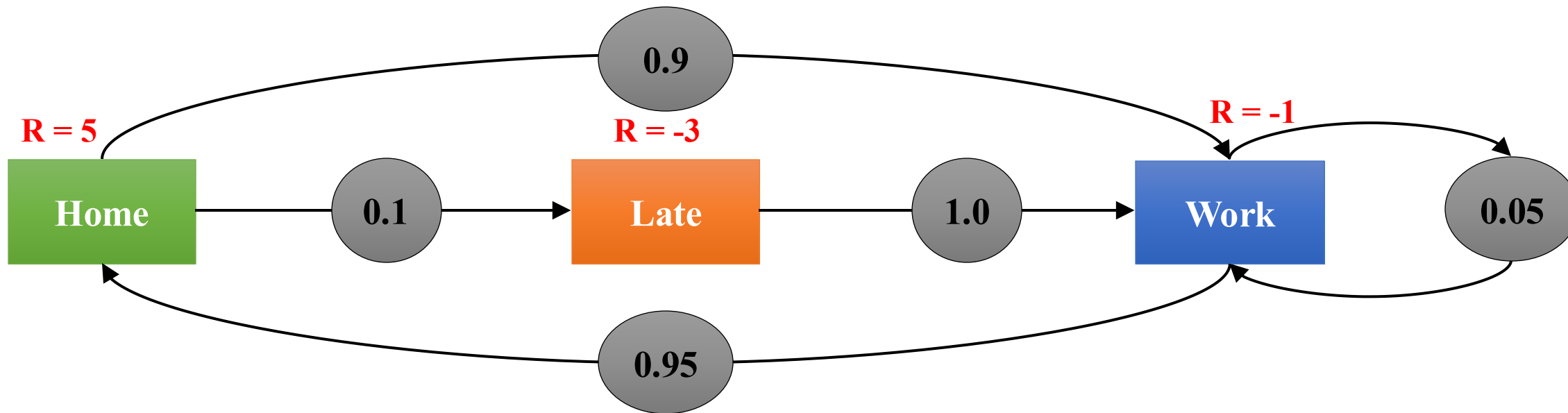
$$v(s) = E[G_t^\tau | \tau \in T] = \sum_{\tau \in T} G_t^\tau \times P(\tau)$$

# 示例：通勤问题

□ 当  $\gamma = 0$ ，通勤问题对应的 MRP 中价值函数为

➤  $v(\text{Home}) = 5, v(\text{Late}) = -3, v(\text{Work}) = -1$

□ 当  $\gamma \neq 0$ ，如何计算价值函数？



# MRP 的贝尔曼方程 (Bellman Equation)

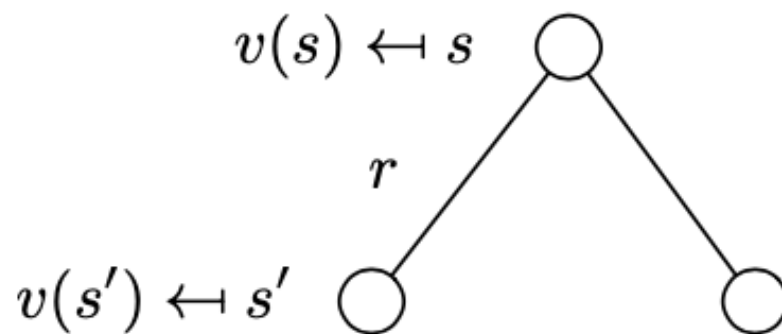
□ 价值函数可以分解为两部分：

- 即时奖励  $R_{t+1}$
- 后继状态折扣价值  $\gamma v(S_{t+1})$

$$\begin{aligned}v(s) &= \mathbb{E}[G_t \mid S_t = s] \\&= \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \mid S_t = s] \\&= \mathbb{E}[R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \dots) \mid S_t = s] \\&= \mathbb{E}[R_{t+1} + \gamma G_{t+1} \mid S_t = s] \\&= \mathbb{E}[R_{t+1} + \gamma v(S_{t+1}) \mid S_t = s]\end{aligned}$$

# MRP 的贝尔曼方程 (Bellman Equation)

$$v(s) = \mathbb{E} [R_{t+1} + \gamma v(S_{t+1}) \mid S_t = s]$$



$$v(s) = \mathcal{R}_s + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'} v(s')$$

# 矩阵形式的贝尔曼方程

贝尔曼方程可以使用矩阵形式简洁地表示为

$$v = \mathcal{R} + \gamma \mathcal{P} v$$

其中,  $v$  是一个列向量, 每个状态对应向量中的一个元素

$$\begin{bmatrix} v(1) \\ \vdots \\ v(n) \end{bmatrix} = \begin{bmatrix} \mathcal{R}_1 \\ \vdots \\ \mathcal{R}_n \end{bmatrix} + \gamma \begin{bmatrix} \mathcal{P}_{11} & \dots & \mathcal{P}_{1n} \\ \vdots & & \\ \mathcal{P}_{n1} & \dots & \mathcal{P}_{nn} \end{bmatrix} \begin{bmatrix} v(1) \\ \vdots \\ v(n) \end{bmatrix}$$

# 求解贝尔曼方程

贝尔曼方程是一个线性方程组，可直接求解

$$v = \mathcal{R} + \gamma \mathcal{P}v$$

$$(I - \gamma \mathcal{P})v = \mathcal{R}$$

$$v = (I - \gamma \mathcal{P})^{-1} \mathcal{R}$$

- 当状态数为  $n$  时，直接求解的计算复杂度为  $O(n^3)$ ，因此只适用于较小规模的 MRP
- 对于规模较大的 MRP，则通常采用迭代方法，例如：
  - 动态规划 (Dynamic Programming)
  - 蒙特卡洛评估 (Monte-Carlo Evaluation)
  - 时序差分学习 (Temporal-Difference Learning)

# 示例：通勤问题

□ 当  $\gamma = 0.5$ ，通勤问题对应的 MRP 中价值函数计算

$$\begin{bmatrix} v(\text{Home}) \\ v(\text{Late}) \\ v(\text{Work}) \end{bmatrix} = \begin{bmatrix} 5 \\ -3 \\ -1 \end{bmatrix} + 0.5 \begin{bmatrix} 0.0 & 0.1 & 0.9 \\ 0.0 & 0.0 & 1.0 \\ 0.95 & 0.0 & 0.05 \end{bmatrix} \begin{bmatrix} v(\text{Home}) \\ v(\text{Late}) \\ v(\text{Work}) \end{bmatrix}$$

$$\begin{bmatrix} v(\text{Home}) \\ v(\text{Late}) \\ v(\text{Work}) \end{bmatrix} = \left( \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} - 0.5 \begin{bmatrix} 0.0 & 0.1 & 0.9 \\ 0.0 & 0.0 & 1.0 \\ 0.95 & 0.0 & 0.05 \end{bmatrix} \right)^{-1} \begin{bmatrix} 5 \\ -3 \\ -1 \end{bmatrix}$$

$$\begin{bmatrix} v(\text{Home}) \\ v(\text{Late}) \\ v(\text{Work}) \end{bmatrix} = \begin{bmatrix} \frac{6806}{1199} \\ -\frac{2554}{1199} \\ \frac{2086}{1199} \end{bmatrix} \approx \begin{bmatrix} 5.6772 \\ -2.1301 \\ 1.7397 \end{bmatrix}$$

# 示例：通勤问题

□ 当  $\gamma = 0.5$ ，通勤问题对应的 MRP 中价值函数计算

□ 对计算结果进行验证

$$\begin{bmatrix} v(\text{Home}) \\ v(\text{Late}) \\ v(\text{Work}) \end{bmatrix} \approx \begin{bmatrix} 5 \\ -3 \\ -1 \end{bmatrix} + 0.5 \begin{bmatrix} 0.0 & 0.1 & 0.9 \\ 0.0 & 0.0 & 1.0 \\ 0.95 & 0.0 & 0.05 \end{bmatrix} \begin{bmatrix} 5.6772 \\ -2.1301 \\ 1.7397 \end{bmatrix}$$
$$\approx \begin{bmatrix} 5.67636 \\ -2.13015 \\ 1.74016 \end{bmatrix}$$

$$\begin{bmatrix} v(\text{Home}) \\ v(\text{Late}) \\ v(\text{Work}) \end{bmatrix} = \begin{bmatrix} 5 \\ -3 \\ -1 \end{bmatrix} + 0.5 \begin{bmatrix} 0.0 & 0.1 & 0.9 \\ 0.0 & 0.0 & 1.0 \\ 0.95 & 0.0 & 0.05 \end{bmatrix} \begin{bmatrix} \frac{6806}{1199} \\ -\frac{2554}{1199} \\ \frac{2086}{1199} \end{bmatrix}$$
$$= \begin{bmatrix} \frac{6806}{1199} \\ -\frac{2554}{1199} \\ \frac{2086}{1199} \end{bmatrix}$$

01

马尔可夫过程 (Markov Process)

02

马尔可夫奖励过程 (Markov Reward Process)

03

马尔可夫决策过程 (Markov Decision Process)

04

MDP 的扩展 (Extensions to MDPs)

# 目录

- 马尔可夫决策过程是在马尔可夫奖励过程基础上增加了动作决策

## 定义

一个**马尔可夫决策过程 (Markov Decision Process, MDP)** 可以表示为一个五元组  $\langle S, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$ , 其中:

- $S$  是状态的 (有限) 集合
- $\mathcal{A}$  是动作的 (有限) 集合
- $\mathcal{P}$  是状态转移矩阵,  $P_{ss'}^a = P(S_{t+1} = s' | S_t = s, A_t = a)$
- $\mathcal{R}$  是奖励函数,  $R_s^a = E[R_{t+1} | S_t = s, A_t = a]$
- $\gamma$  是折扣因子,  $\gamma \in [0, 1]$

# 示例：通勤问题

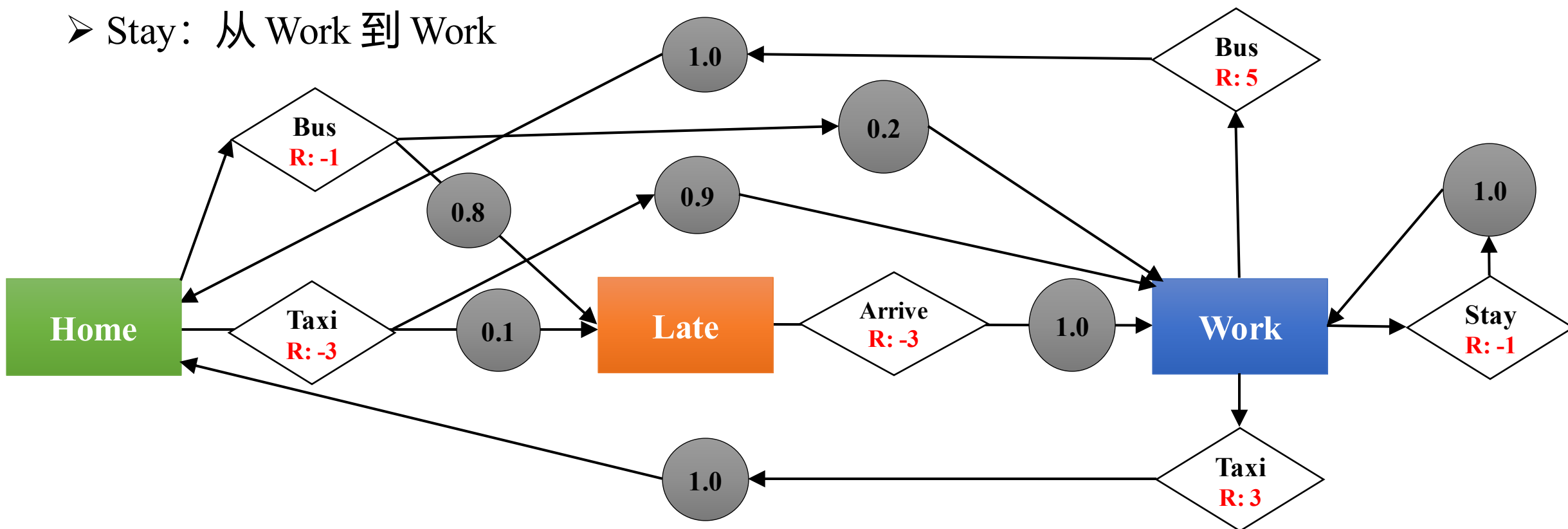
□ 通勤问题在状态 Home, Late, Work 的基础上，增加四个行动

➤ Bus: 乘坐巴士, 从 Home 到 Late 或 Work, 或者从 Work 到 Home

➤ Taxi: 乘坐出租车, 从 Home 到 Late 或 Work, 或者从 Work 到 Home

➤ Arrive: 从 Late 到 Work

➤ Stay: 从 Work 到 Work



# 示例：通勤问题

## □ 通勤问题的 MDP 描述

- 状态集合：Home, Late, Work
- 动作集合：Bus, Taxi, Arrive, Stay
- 状态转移矩阵：

$$P_{ss'}^{\text{Stay}} = \begin{bmatrix} 1.0 & 0.0 & 0.0 \\ 0.0 & 1.0 & 0.0 \\ 0.0 & 0.0 & 1.0 \end{bmatrix}$$

$$P_{ss'}^{\text{Bus}} = \begin{bmatrix} 0.0 & 0.8 & 0.2 \\ 0.0 & 1.0 & 0.0 \\ 1.0 & 0.0 & 0.0 \end{bmatrix}$$

$$P_{ss'}^{\text{Taxi}} = \begin{bmatrix} 0.0 & 0.1 & 0.9 \\ 0.0 & 1.0 & 0.0 \\ 1.0 & 0.0 & 0.0 \end{bmatrix}$$

$$P_{ss'}^{\text{Arrive}} = \begin{bmatrix} 1.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 1.0 \\ 0.0 & 0.0 & 1.0 \end{bmatrix}$$

- 奖励函数：

$$R_s^{\text{Bus}} = \begin{bmatrix} -1 \\ -\infty \\ 5 \end{bmatrix}$$

$$R_s^{\text{Taxi}} = \begin{bmatrix} -3 \\ -\infty \\ 3 \end{bmatrix}$$

$$R_s^{\text{Arrive}} = \begin{bmatrix} -\infty \\ -3 \\ -\infty \end{bmatrix}$$

$$R_s^{\text{Stay}} = \begin{bmatrix} -\infty \\ -\infty \\ -1 \end{bmatrix}$$

## 定义

**策略 (Policy)**  $\pi$  是在给定状态  $s$  时, 采取各动作  $a$  的概率分布:

$$\pi(a|s) = P[A_t = a | S_t = s]$$

- 策略完全决定了智能体的行为
- 在 MDP 中, 策略仅依赖于当前状态 (与历史无关)
  - 换言之, 策略是平稳的 (与时间无关): 对于所有  $t > 0$ , 有

$$A_t \sim \pi(\cdot | S_t)$$

# MDP 与 MP 和 MRP 的关系

给定一个 MDP  $\mathcal{M} = \langle S, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$  及策略  $\pi$  :

□ 状态序列  $S_1, S_2, \dots$  构成马尔可夫过程 (Markov Process)  $\langle S, P^\pi \rangle$

□ 状态-奖励序列  $S_1, R_2, S_2, \dots$  构成马尔可夫奖励过程 (MRP)  
 $\langle S, P^\pi, R^\pi, \gamma \rangle$

其中:

$$P_{SS'}^\pi = \sum_{a \in \mathcal{A}} \pi(a|s) P_{SS'}^a$$

$$R_S^\pi = \sum_a \pi(a|s) R_S^a$$

# 价值函数

## 定义

在给定策略  $\pi$  的 MDP 中，**状态价值函数 (state-value function)**  $v_\pi(s)$  表示从状态  $s$  开始，随后按照策略  $\pi$  行动时的期望回报：

$$v_\pi(s) = E_\pi[G_t | S_t = s]$$

## 定义

在给定策略  $\pi$  的 MDP 中，**动作价值函数 (action-value function)**  $q_\pi(s, a)$  表示从状态  $s$  开始，先执行动作  $a$ ，然后按照策略  $\pi$  行动时的期望回报：

$$q_\pi(s, a) = E_\pi[G_t | S_t = s, A_t = a]$$

# 贝尔曼期望方程 (Bellman Expectation Equation)

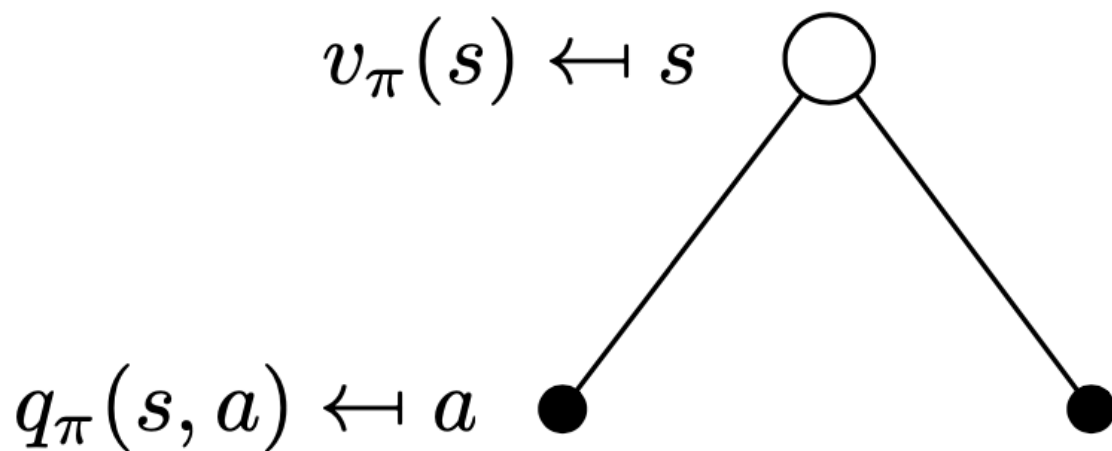
□ 对于状态价值函数, 有以下分解:

$$v_{\pi}(s) = E_{\pi}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) \mid S_t = s]$$

□ 对于动作价值函数, 同样可以分解:

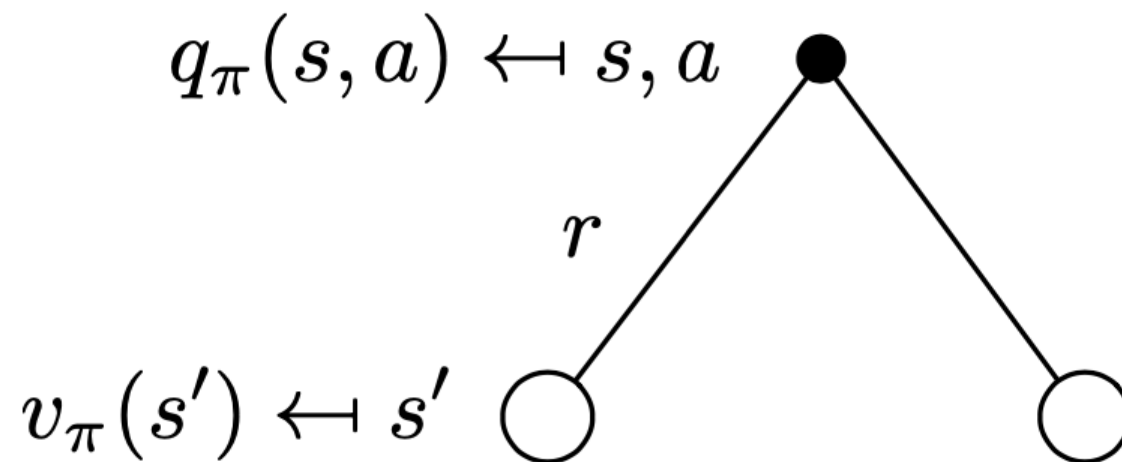
$$q_{\pi}(s, a) = E_{\pi}[R_{t+1} + \gamma q_{\pi}(S_{t+1}, A_{t+1}) \mid S_t = s, A_t = a]$$

# $v_\pi$ 的贝尔曼期望方程



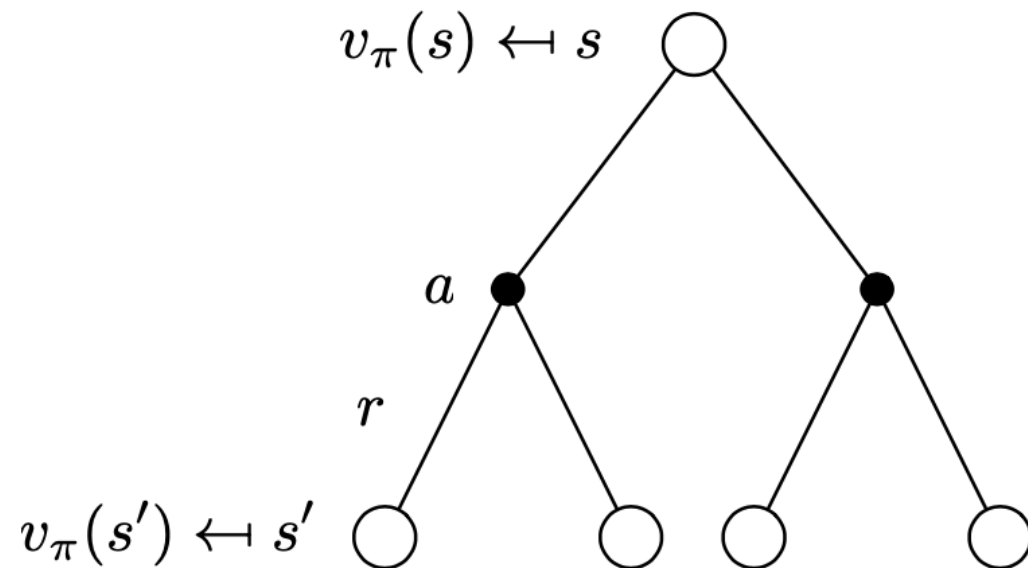
$$v_\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) q_\pi(s, a)$$

# $q_\pi$ 的贝尔曼期望方程



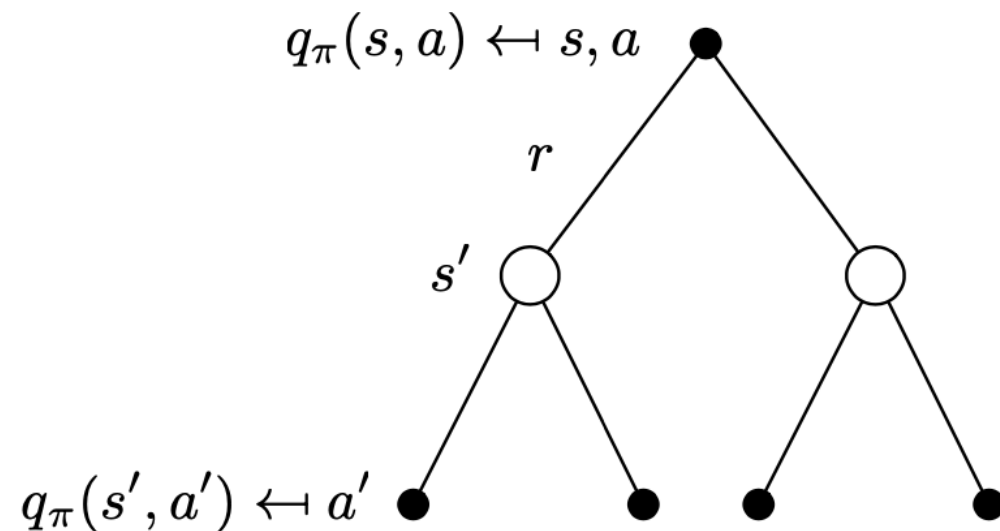
$$q_\pi(s, a) = \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_\pi(s')$$

# $v_\pi$ 的贝尔曼期望方程



$$v_\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \left( \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_\pi(s') \right)$$

# $q_\pi$ 的贝尔曼期望方程



$$q_\pi(s, a) = \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a \sum_{a' \in \mathcal{A}} \pi(a'|s') q_\pi(s', a')$$

# 矩阵形式的贝尔曼期望方程

借鉴 MDP 与 MRP 之间的关系，贝尔曼期望方程可以使用矩阵形式简洁地表示为：

$$V_{\pi} = \mathcal{R}^{\pi} + \gamma \mathcal{P}^{\pi} V_{\pi}$$

直接结果为：

$$V_{\pi} = (I - \gamma \mathcal{P}^{\pi})^{-1} \mathcal{R}^{\pi}$$

$$P_{SS'}^{\pi} = \sum_{a \in A} \pi(a|s) P_{SS'}^a$$

$$R_s^{\pi} = \sum_a \pi(a|s) R_s^a$$

# 最优价值函数

## 定义

**最优状态价值函数 (optimal state-value function)**  $v_*(s)$  是对所有策略的状态价值函数的最大值:

$$v_*(s) = \max_{\pi} v_{\pi}(s)$$

**最优动作价值函数 (optimal action-value function)**  $q_*(s, a)$  是对所有策略的动作价值函数的最大值:

$$q_*(s, a) = \max_{\pi} q_{\pi}(s, a)$$

- 最优价值函数刻画在该 MDP 中能够获得的最佳可能表现
- 当我们知道了最优价值函数时, 就意味着该 MDP 已被“解决”

# 最优策略

- 在策略之间定义偏序关系:

$$\pi \geq \pi' \text{ 当且仅当 } v_{\pi}(s) \geq v_{\pi'}(s), \text{ 对任意状态 } s$$

## 定理

对于任意马尔可夫决策过程 (MDP) :

- 存在一个**最优策略 (optimal policy)** , 它不劣于任何其他策略, 即

$$\pi_* \geq \pi \text{ 对所有策略 } \pi \text{ 都成立}$$

- 所有最优策略都能实现最优状态价值函数, 即

$$\text{对所有状态 } s, v_{\pi_*}(s) = v_*(s)$$

- 所有最优策略都能实现最优动作价值函数, 即

$$\text{对所有状态 } s \text{ 和动作 } a, q_{\pi_*}(s, a) = q_*(s, a)$$

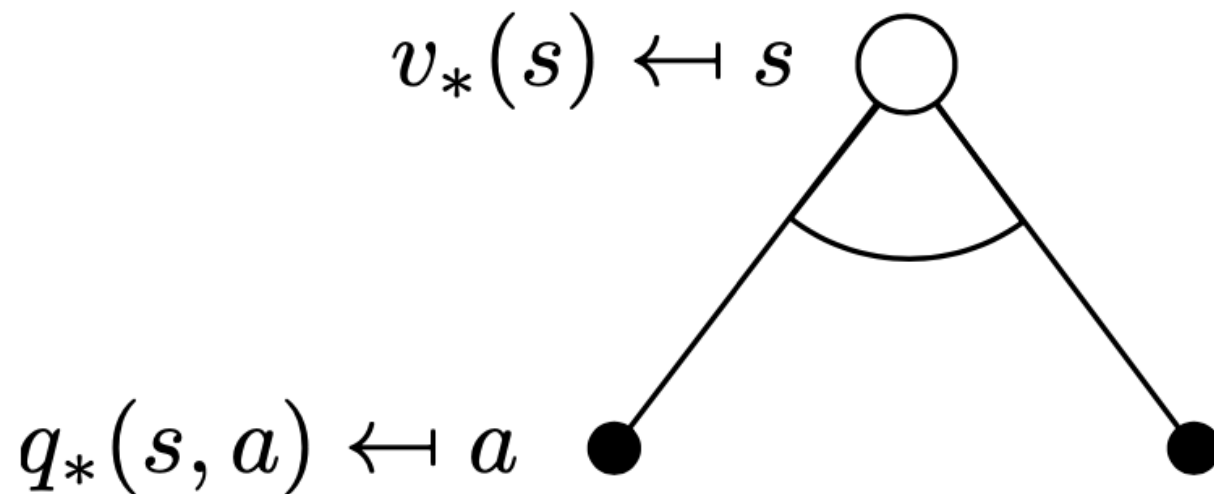
# 寻找最优策略

- 可以通过对  $q_*(s, a)$  取最大值来构造最优策略:

$$\pi_*(a|s) = \begin{cases} 1 & \text{if } a = \operatorname{argmax}_{a \in \mathcal{A}} q_*(s, a) \\ 0 & \text{otherwise} \end{cases}$$

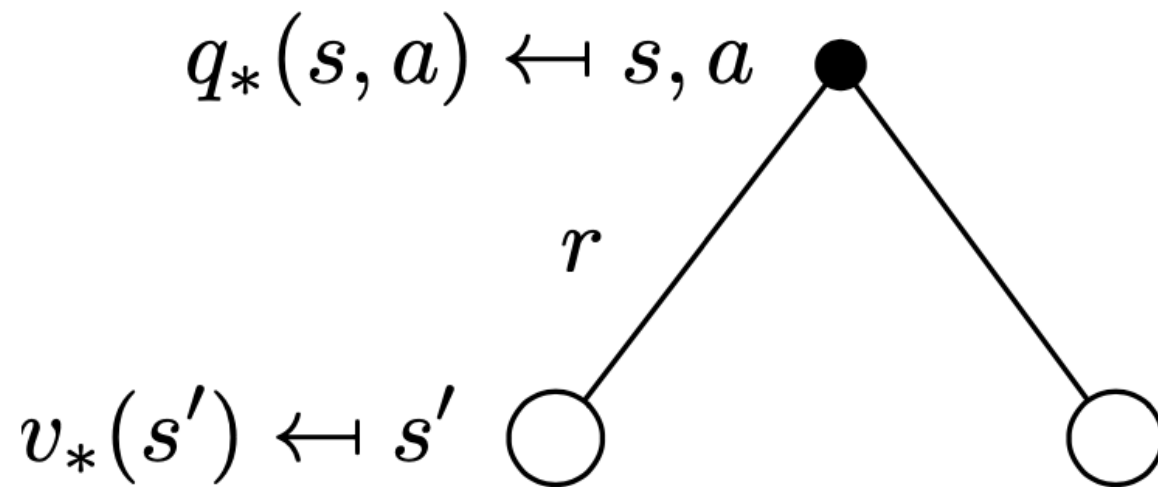
- 对于任意 MDP, 总存在一个确定性的最优策略
- 一旦我们知道了  $q_*(s, a)$ , 便可立刻得到该最优策略

# $v_*$ 的贝尔曼最优方程 (Bellman Optimality Equation)



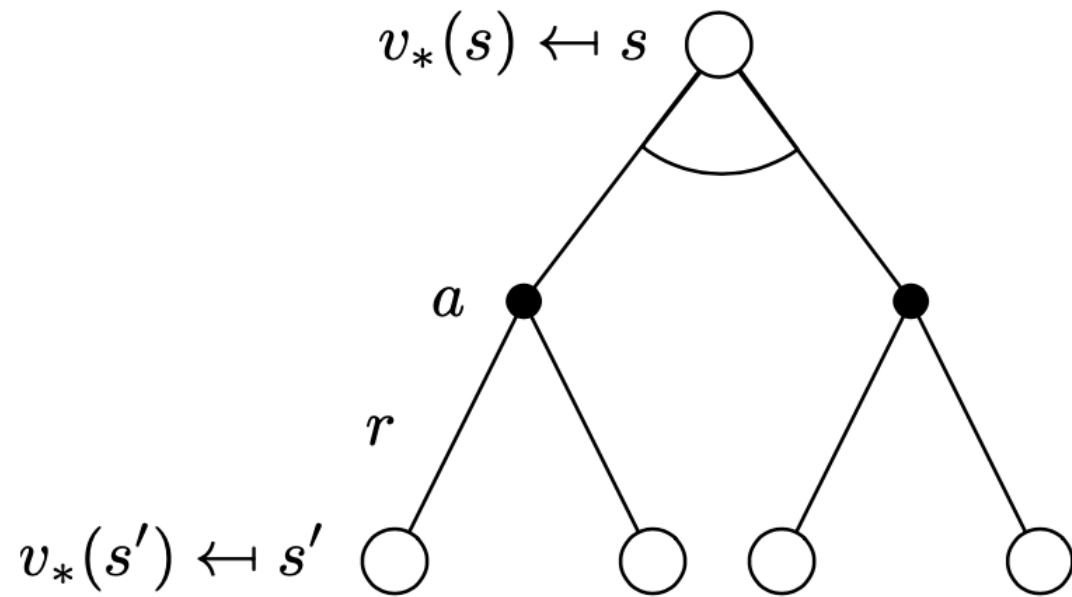
$$v_*(s) = \max_a q_*(s, a)$$

# $q_*$ 的贝尔曼最优方程



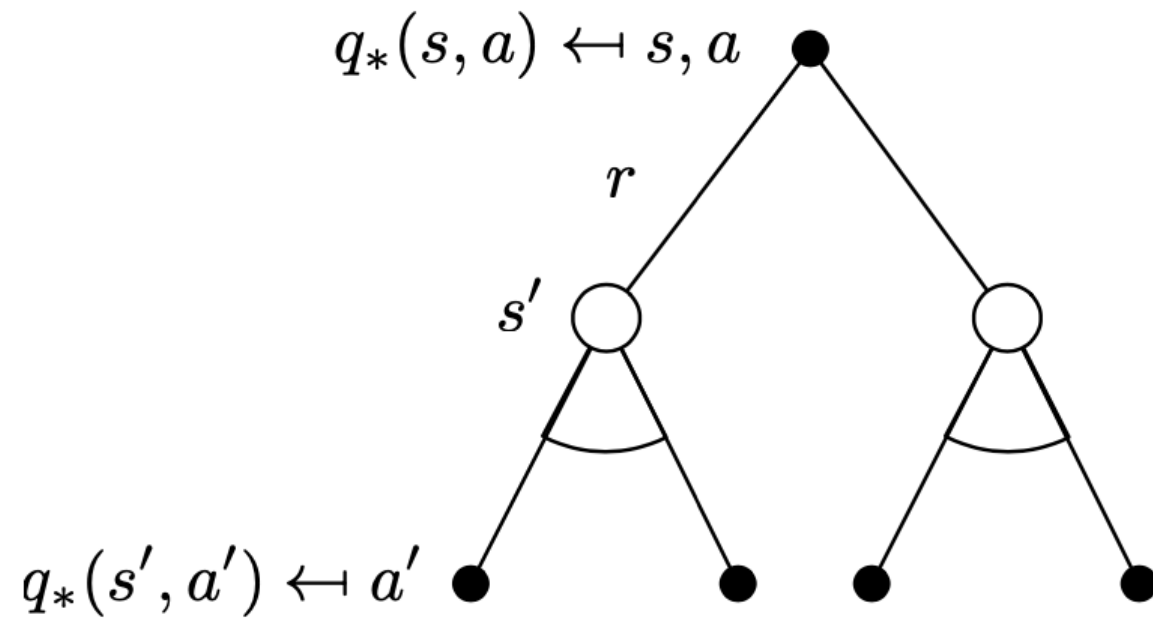
$$q_*(s, a) = \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_*(s')$$

# $v_*$ 的贝尔曼最优方程



$$v_*(s) = \max_a \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_*(s')$$

# $q_*$ 的贝尔曼最优方程



$$q_*(s, a) = \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a \max_{a'} q_*(s', a')$$

# 贝尔曼期望方程和贝尔曼最优方程

$$V_{\pi}(s) = \sum_{a \in A} \pi(a | s) \left( R(s, a) + \gamma \sum_{s' \in S} P(s' | s, a) V_{\pi}(s') \right),$$

$$V^*(s) = \max_{a \in A} \left( R(s, a) + \gamma \sum_{s' \in S} P(s' | s, a) V^*(s') \right),$$

$$Q_{\pi}(s, a) = R(s, a) + \gamma \sum_{s' \in S} P(s' | s, a) \sum_{a' \in A} \pi(a' | s') Q_{\pi}(s', a'),$$

$$Q^*(s, a) = R(s, a) + \gamma \sum_{s' \in S} P(s' | s, a) \max_{a' \in A} Q^*(s', a').$$

# 求解贝尔曼最优方程

- 贝尔曼最优方程是非线性的
- (一般) 不存在封闭形式解
- 常见的迭代求解方法:
  - 值迭代 (Value Iteration)
  - 策略迭代 (Policy Iteration)
  - Q 学习 (Q-Learning)
  - Sarsa

# 值迭代算法

值迭代 (value iteration) 算法求  $V_t^*$  序列, 借助于辅助  $Q_t^a(s)$ , 其直观含义为: 在  $s$  执行  $a$ , 然后执行  $t-1$  步的最优策略所产生的预期回报

对所有  $s \in S$ ,  $V_0(s) := 0$ ;  $t := 0$ ;

loop

$t := t + 1$ ;

loop 对所有  $s \in S$

loop 对所有  $a \in A$

$$Q_t^a(s) := R(s, a) + \gamma \sum_{s'} T(s, a, s') V_{t-1}(s')$$

end loop

$$V_t(s) := \max_a Q_t^a(s) \text{ (Bellman equation)}$$

end loop

until  $|V_t(s) - V_{t-1}(s)| < \epsilon$  对所有  $s \in S$  成立

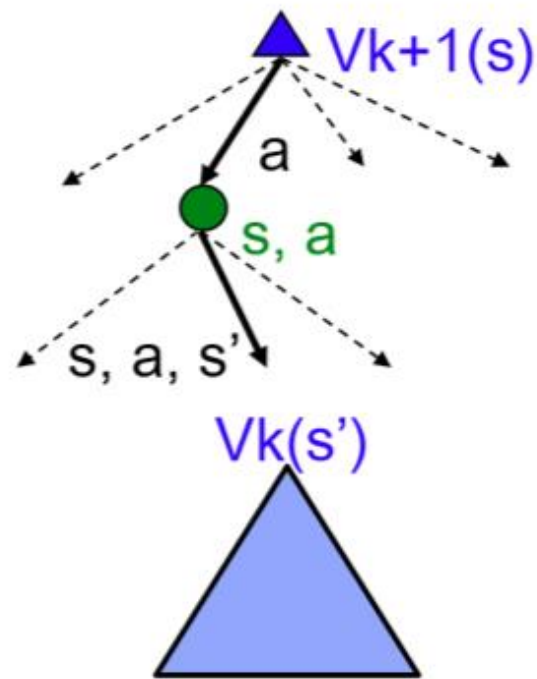
# 值迭代算法

1. 对所有状态 $s$ ，初始化  $V_0(s) = 0$
2. 给定一组 $V_k(s)$ 的值，对所有的状态进行下列迭代更新：
$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k(s')]$$
3. 不断重复步骤2，直到收敛

## 值迭代算法的一些结论

- 每一步值迭代的计算复杂度是  $O(|A||S|^2)$
- 根据**不动点理论**，算法将最终收敛到最优值
  - 收敛到最优值所需的迭代次数为：

$$\text{poly}(|S|, |A|, 1/(1-\gamma))$$



本质上是一种**动态规划**的方法

# 策略迭代算法

## ▶ 策略迭代算法

1. 先给定初始策略  $\pi$ ，求解线性方程，计算出  $\pi$  相应的值函数  $V$

$$V(s) = R(s, \pi(s)) + \gamma \sum_{s'} T(s, a, s') V(s').$$

2. 基于计算出的值函数  $V$ ，按下式更新策略

$$\pi(s) \leftarrow \operatorname{argmax}_a \left\{ R(s, a) + \gamma \sum_{s'} T(s, a, s') V(s') \right\}$$

3. 不断重复步骤 1, 2，直到收敛

## ▶ 策略迭代算法最终收敛到最优值

- ▶ 收敛到最优值所需的迭代不大于  $|A|^{|S|}$ （确定性策略的总数目）

01

马尔可夫过程 (Markov Process)

02

马尔可夫奖励过程 (Markov Reward Process)

03

马尔可夫决策过程 (Markov Decision Process)

04

部分可观察马尔可夫决策过程 (Partially observable MDPs)

# 目录

# 部分可观察马尔可夫决策过程

- 当状态不可完全观察时，通过引入观察和观察函数来扩展 MDP

## 定义

一个**部分可观察马尔可夫决策过程** (Partially Observable Markov Decision Process, POMDP) 可以表示为一个七元组  $\langle \mathcal{S}, \mathcal{A}, \mathcal{O}, \mathcal{P}, \mathcal{R}, \mathcal{Z}, \gamma \rangle$ , 其中:

- $\mathcal{S}$  是状态的 (有限) 集合
- $\mathcal{A}$  是动作的 (有限) 集合
- $\mathcal{O}$  是观察的 (有限) 集合
- $\mathcal{P}$  是状态转移矩阵,  $P_{SS'}^a = P(S_{t+1} = s' | S_t = s, A_t = a)$
- $\mathcal{R}$  是奖励函数,  $R_s^a = E[R_{t+1} | S_t = s, A_t = a]$
- $\mathcal{Z}$  是观察函数,  $Z_{s'o}^a = P[O_{t+1} = o | S_{t+1} = s', A_t = a]$
- $\gamma$  是折扣因子,  $\gamma \in [0, 1]$

## 定义

**历史 (history)**  $H_t$  是动作、观测和奖励构成的序列：

$$H_t = A_0, O_1, R_1, \dots, A_{t-1}, O_t, R_t$$

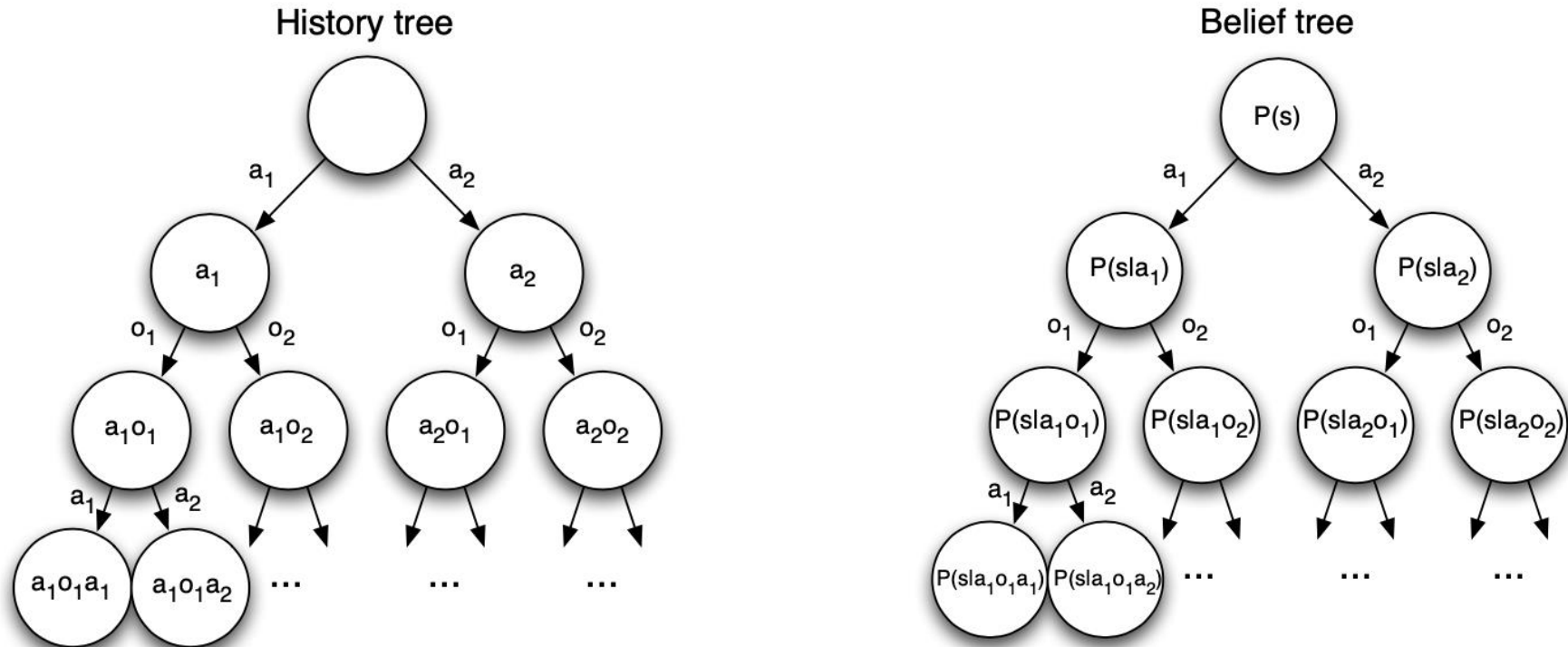
## 定义

**信念状态 (belief state)**  $b(h)$  是在给定历史  $h$  的条件下，对所有可能状态的概率分布：

$$b(h) = ( P[S_t = s^1 \mid H_t = h], \dots, P[S_t = s^n \mid H_t = h] )$$

# POMDP 的简化

- 历史  $H_t$  满足马尔可夫性质, POMDP 可以简化为一个 (无限的) 历史树
  - 给定完整的历史序列  $H_t$ , 未来状态只依赖于该历史
- 信念状态  $b(H_t)$  满足马尔可夫性质, POMDP 可以简化为一个 (无限的) 信念状态树
  - 给定当前的信念分布  $b_t = b(H_t)$ , 未来只依赖于该分布, 而与过往历史无关





中国科学技术大学  
University of Science and Technology of China

**谢谢!**