

1. MDP:

① Markov-chain = $\{S, T\}$

S : State-space

T : State-transform-operator

② $\mu_{t,i} = p(S_t = i)$

$T_{i,j} = p(S_{t+1} = i | S_t = j)$

$\mu_{t+1} = T \mu_t$

$$\mu_t = \begin{bmatrix} p(S_t = i) \\ p(S_t = j) \\ p(S_t = k) \end{bmatrix} \quad T = \begin{bmatrix} p(S_{t+1} = i | S_t = i) & \vdots \\ p(S_{t+1} = i | S_t = j) & \vdots \\ p(S_{t+1} = i | S_t = k) & \vdots \end{bmatrix}$$

③ Markov-Decision-Process:

$\mathcal{M} = \{S, A, T, r\}$

S : ...

A : action-space

T : Matrix \rightarrow Tensor

Match the model
in Reinforce-Learning

$T_{i,j,k} = p(S_{t+1} = i | S_t = j, a_t = k)$

$\mu_{t,j} = p(S_t = j), \xi_{t,k} = p(a_t = k)$

$$\mu_{t+1,i} = \sum_{j,k} \gamma_{i,j,k} \cdot \mu_{t,j} \cdot \xi_{t,k}$$

$$\gamma = \begin{bmatrix} \begin{bmatrix} p(s_{t+1}=i_0 | s_t=j_0, a_t=k_0) & \dots \\ \vdots \\ p(s_{t+1}=i_n | s_t=j_0, a_t=k_0) & \dots \end{bmatrix} \\ \vdots \\ \begin{bmatrix} p(s_{t+1}=i_n | s_t=j_0, a_t=k_0) & \dots \\ \vdots \\ p(s_{t+1}=i_n | s_t=j_n, a_t=k_0) & \dots \end{bmatrix} \end{bmatrix}$$

$$r \rightarrow r(s,a) \rightarrow S \times A \rightarrow R$$

④ Partially-Observed-Markov-Decision-Process

$$\mathcal{M} = \{S, A, O, T, \xi, r\}$$

$\xi \rightarrow$ emission-probability.

\rightarrow decide $p(o_t | s_t)$

⑤ trick: $\mu \cdot \xi \rightarrow (s, a)$

将 state-action 视为一个联合状态, 简化了分析

$$P(s_{t+1}, a_{t+1} | (s_t, a_t)) = P(s_{t+1} | s_t, a_t) \cdot \pi_\theta(a_{t+1} | s_{t+1})$$


⑥ 平稳分布: 经过一次转移后不发生改变的分布.

$\boxed{\mu = T\mu} \rightarrow \mu$ 是 T 特征值为 1 的特征向量

2. Expectation of Reinforce-Learning

① Original Reward / Cost Function 是不平滑的, 无法直接求 gradient, 然后 backward

② 平稳分布下的 Expectation-of-Reward 是光滑的

$$\boxed{E_{(s,a) \sim p_\psi(s,a)} r(s,a)}$$


* θ 收敛是 μ 达到平稳分布的必要条件.