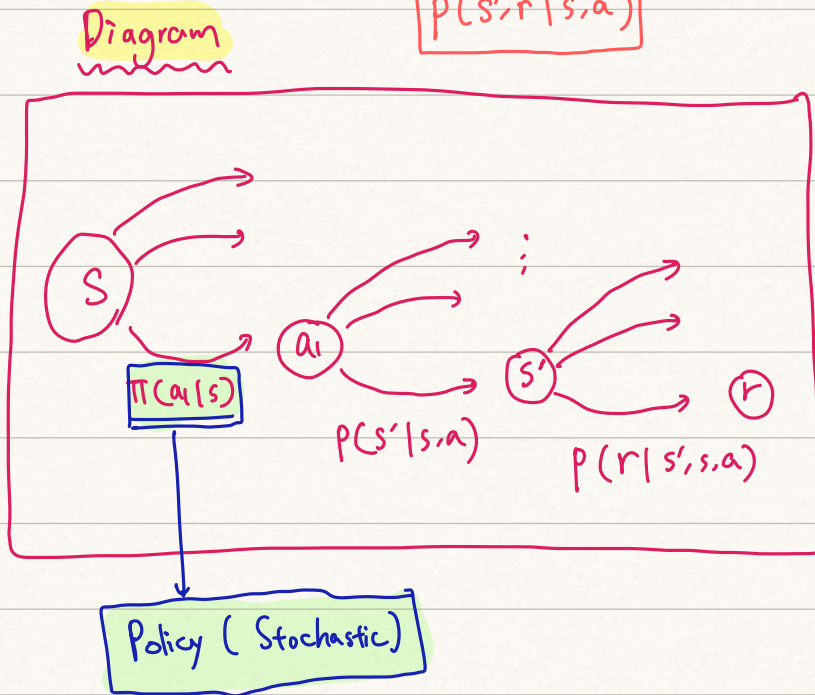
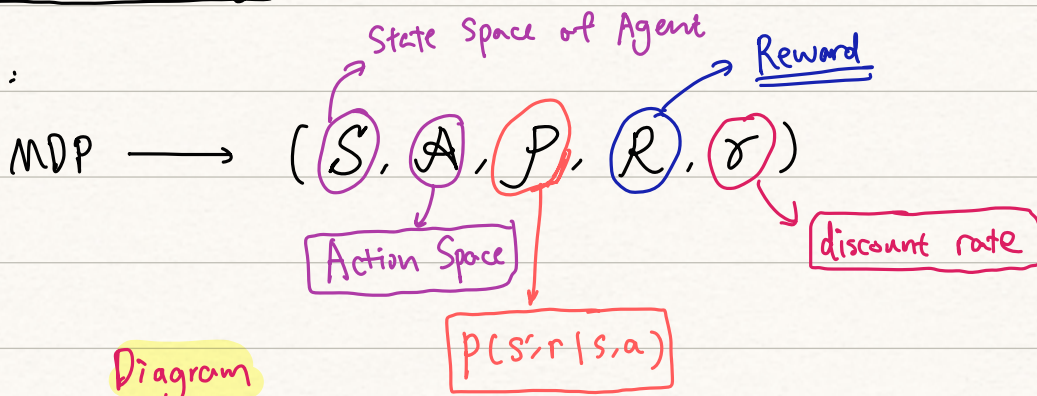


Recap:



DP \longrightarrow RL Model-Based Model

\longrightarrow State-value Function

$$V_{\pi}(s) = \mathbb{E}^{\pi} [G_t | S_t = s]$$

$$= \mathbb{E}^{\pi} \left[\sum_{l=0}^{\infty} \gamma^l R_{t+l+1} \mid S_t = s \right]$$

We consider a time-homogeneous MDP

action-value Function

$$Q_{\pi}(s, a) = \mathbb{E}^{\pi} [G_t \mid S_t = s, A_t = a]$$

$$= \mathbb{E}^{\pi} \left[\sum_{l=0}^{\infty} \gamma^l R_{t+l+1} \mid S_t = s \right]$$

→ Goal: derive recursive relation with $q_\pi(s, a) / V_\pi(s)$

Bellman Equation

(Bellman Equation)

①

Derivation:

$$V_\pi(s) = \mathbb{E}^\pi [G_t | S_t = s] = \sum_a \pi(a|s) \cdot q_\pi(s, a) \\ = \mathbb{E}^\pi [R_{t+1} | S_t = s] + \gamma \mathbb{E}^\pi \left[\sum_{\tau=0}^{\infty} \gamma^\tau R_{t+\tau+2} | S_t = s \right]$$

$$= b(\pi)_s + \gamma \sum_{s'} P(\pi)_{ss'} \cdot V_\pi(s')$$

$$\left\{ \begin{aligned} b(\pi)_s &:= \mathbb{E}^\pi [R_{t+1} | S_t = s] = \sum_a \pi(a|s) \cdot \sum_{r, s'} p(s', r | s, a) \cdot r \\ P(\pi)_{ss'} &= \mathbb{P}^\pi (S_{t+1} = s' | S_t = s) \\ &= \sum_a \pi(a|s) \mathbb{P} (S_{t+1} = s' | S_t = s, A_t = a) \\ &= \sum_a \pi(a|s) \cdot \sum_r p(s', r | s, a) \end{aligned} \right.$$

give the first-step transition

probability



easy to compute

② $q_\pi(s, a) = \mathbb{E}^\pi [G_t | \underline{S_t = s, A_t = a}]$

$$= \sum_{s', r} r \cdot p(s', r | s, a) + \gamma \mathbb{E}^\pi \left[\sum_{\tau=0}^{\infty} \gamma^\tau R_{t+\tau+2} | s, a \right]$$

$$= \sum_{s', r} p(s', r | s, a) \cdot r + \gamma \sum_{s'} p(s' | s, a) \cdot V_\pi(s')$$

$$= \sum_{s', r} p(s', r | s, a) (r + \gamma V_\pi(s'))$$

Rmk: $V_\pi = b(\pi) + \gamma \cdot P(\pi) \cdot V_\pi$

$$\Rightarrow V_\pi = (I - \gamma P(\pi))^{-1} b(\pi)$$

⇒ idea: when we have $\begin{cases} P(s', r | s, a) \\ \text{policy } \pi(a|s) \end{cases}$, then we can

compute for $V_\pi(s)$ & $q_\pi(s, a)$ $\forall s \in S$ and $\forall a \in A$

Policy Comparison

Defn:

$$\pi' \succcurlyeq \pi$$

$$\Leftrightarrow V_{\pi'}(s) \geq V_\pi(s) \quad \forall s \in S$$

(POSET)
partial order

① not necessary happen

Property

$$\textcircled{2} \begin{cases} a \geq a \\ a \geq b \quad b \leq a \Rightarrow a = b \\ a \geq b \quad b \geq c \Rightarrow a \geq c \end{cases}$$

A trivial question: optimal policy exists or not!

Bellman Optimality Introduced

uniques or not!

Outline: ① Policy Improvement

→ if $\sum_a \pi'(a|s) \cdot q_\pi(s, a) \geq \sum_a \pi(a|s) q_\pi(s, a)$
then $\pi' \succcurlyeq \pi$ ($V_{\pi'}(s) \geq V_\pi(s) \quad \forall s \in S$)

Rmk: ① Give us the motivation of policy iteration

② Necessary Condition for Optimal Policy

③ a natural choice to improve policy is:

$$\pi'(a|s) = \begin{cases} 1 & a = \operatorname{argmax}_{a \in A} q_{\pi}(s, a) \\ 0 & \text{o/w} \end{cases}$$

↓
deterministic policy

(necessary part)

② Bellman Optimality Condition

$\pi^*(a|s)$ is an optimal policy
 \Rightarrow ① if $\pi^*(a|s) > 0$

then $a \in \operatorname{argmax}_{a'} q_{\pi^*}(a', s)$

$$\Rightarrow$$
 ② $V_{\pi^*}(s) = \max_a q_{\pi^*}(a, s)$

$$= \max_a \sum_{s', r} p(s', r | s, a) (r + \gamma V_{\pi^*}(s'))$$

$$q_{\pi^*}(s, a) = \sum_{s', r} p(s', r | s, a) (r + \gamma V_{\pi^*}(s'))$$

requires Contraction Mapping Thm to show
the uniqueness of Bellman
Optimality Equation

$$= \sum_{s', r} p(s', r | s, a) (r + \gamma \max_{a'} q_{\pi^*}(s', a'))$$

(Sufficient Part)

given a policy $\pi^*(a|s) \rightarrow \begin{cases} V_{\pi^*}(s) \\ q_{\pi^*}(s, a) \end{cases}$

$$\rightarrow \text{if } V_{\pi^*}(s) = \max_a \sum_{s', r} p(s', r | s, a) (r + \gamma V_{\pi^*}(s'))$$

then π^* is the optimal policy

$$\rightarrow \text{if } \pi^*(a|s) > 0 \Rightarrow a \in \operatorname{argmax}_{a'} q_{\pi^*}(a', s)$$

then π^* is the optimal policy

