

Reinforcement learning

① Deterministic / Non-deterministic Transition

$$\begin{aligned} \downarrow & \\ f(s, a) = s' & \quad f(s, a, s') = P_{ss'}^a \quad r(s, a) = \sum_{s'} P_{ss'}^a p(s, a, s') \end{aligned}$$

② policy \rightarrow $\begin{cases} \text{deterministic} & \pi(s) = a \\ \text{non-deterministic} & \pi(a|s) \rightarrow \underline{\text{probability distribution}} \end{cases}$

Key Question: How to make sense that $\pi^*(s) := \arg\max_a Q^*(s, a)$

③ **Main Topic** \rightarrow Idea: determine **Best Policy** (optimal) $\Downarrow Q^*(a|s) = \max_{\pi} Q^{\pi}(a|s)$

{ State-Value Function $V^{\pi}(s) = \mathbb{E}^{\pi}[R_t | S_t = s]$

{ Action-Value Function $Q^{\pi}(s, a) = \mathbb{E}^{\pi}[R_t | S_t = s, A_t = a]$

$R_t = r_{t+1} + \gamma r_{t+2} + \dots$ **Reward (long term)**

\rightarrow The Relationship Between $V^{\pi}(s)$ & $Q^{\pi}(s, a)$

$$\begin{aligned} \textcircled{1} \quad V^{\pi}(s) &= \mathbb{E}^{\pi}[R_t | S_t = s] \\ &= \mathbb{E}^{\pi}\left[\mathbb{E}^{\pi}[R_t | S_t = s, A_t]\right] \end{aligned}$$

$$= \sum_a \pi(a|s) \mathbb{E}^{\pi}[R_t | S_t = s, A_t = a]$$

$$= \sum_a \pi(a|s) \cdot Q^{\pi}(s, a)$$

\rightarrow when $\pi(\cdot)$ is deterministic,

then $V^{\pi}(s)$ degenerate to $Q^{\pi}(s, \pi(s))$

$$\text{i.e., } V^{\pi}(s) = Q^{\pi}(s, \pi(s))$$

Generally, $V^{\pi}(s) = \sum_a \pi(a|s) \cdot Q^{\pi}(s, a)$

②

$$\begin{aligned} Q^{\pi}(s, a) &= \mathbb{E}^{\pi}[R_t | S_t = s, A_t = a] \\ &= \sum_s P_{ss'}^a \cdot \mathbb{E}^{\pi}[p(s, a, s') + r R_{t+1} | S_{t+1} = s'] \end{aligned}$$

$$= \sum_{s'} P_{ss'}^a (p(s, a, s') + r V^\pi(s'))$$

$$= r(s, a) + \gamma \sum_{s'} P_{ss'}^a V^\pi(s')$$

When $\pi(\cdot)$ is deterministic, $V^\pi(s')$ degenerate to $\boxed{Q^\pi(s', \pi(s'))}$

$$\Rightarrow Q^\pi(s, a) = r(s, a) + \gamma \sum_{s'} P_{ss'}^a Q^\pi(s', \pi(s'))$$

$$\text{Generally } Q^\pi(s, a) = r(s, a) + \gamma \sum_{s'} P_{ss'}^a V^\pi(s')$$

These 2 are Bellman Equation !

Then, Back to our goal \rightarrow find optimal policy $\underline{\underline{\pi^*(\cdot)}}$

Firstly, consider deterministic case $\underline{\pi(\cdot)}$ \Rightarrow [easy part]

$$\rightarrow \begin{cases} V^\pi(s) = Q^\pi(s, \pi(s)) \\ Q^\pi(s, a) = r(s, a) + \gamma \sum_{s'} P_{ss'}^a V^\pi(s') \end{cases}$$

\rightarrow Intuitively, if we have $\underline{\underline{\pi^*(\cdot)}}$, then we want

$$\underbrace{Q^{\pi^*}(s, a)}_{\text{something like 'optimal'}} \geq Q^\pi(s, a) \quad \text{for } \forall s, a$$

\curvearrowright something like 'optimal'

Can we achieve this?

$$\text{Since } Q^\pi(s, a) = r(s, a) + \gamma \sum_{s'} P_{ss'}^a Q^\pi(s', \pi(s'))$$

if $\pi^*(\cdot)$ attains the maximum for $\forall (s, a)$ pairs,

then it is necessary that: $F(\pi, s') := Q^\pi(s'; \pi(s'))$

Necessary Condition

to attain the maxima at π^*

\Leftrightarrow consider $\pi(s') = a'$

$$\max Q^\pi(s', \pi(s'))$$

$$\text{then } \max_{a' \in \Pi} Q^\pi(s', a')$$



Actually, what we do is to decompose $\pi \rightarrow \begin{cases} \pi(s') \rightarrow a' \\ \pi(s \setminus \{s'\}) \rightarrow A \end{cases}$

can determine $\pi^*(s') \rightarrow a'$

do for every s'

$\pi^*: S \rightarrow A$ optimal
exists.

⇒ the \star existence of optimal policy $\pi^*(\cdot)$

$$\text{s.t. } Q^{\pi^*}(s, a) \geq Q^\pi(s, a) \quad \forall s, a, \pi$$

[Answer]: We can achieve this! By picking π^* as follows:

$$\pi^*: S \mapsto \arg \max_a \left\{ \max_{\pi} Q^{\pi}(s, a) \right\}$$

a Naive Method!

Greedy policy → actually is optimal policy

Since we know there exists $Q^*(s, a) \geq Q^\pi(s, a) \quad \forall (s, a, \pi)$

then we turn to check the optimality condition!

$$\begin{aligned} Q^*(s, a) &= \max_{\pi} Q^{\pi}(s, a) \\ &= \max_{\pi} \left\{ r(s, a) + \gamma \sum_{s'} P_{ss'}^a Q^{\pi}(s', \pi(s')) \right\} \\ &= r(s, a) + \gamma \sum_{s'} P_{ss'}^a \max_{\pi} \left\{ Q^{\pi}(s', \pi(s')) \right\} \\ &\xrightarrow{\text{decompose}} = r(s, a) + \gamma \sum_{s'} P_{ss'}^a \max_{a'} \max_{\tilde{\pi}} \left\{ Q^{\tilde{\pi}}(s', a') \right\} \end{aligned}$$

$$= r(s, a) + \gamma \sum_{s'} P_{ss'}^a \max_{a'} Q^*(s', a')$$

↓
determine $\pi^*(s') = \arg \max_{a'} Q^*(s', a')$

from Optimality Condition

(necessary)

→ For Non-deterministic policy $\pi(a|s)$

Decompose

 $\pi: s \rightarrow a$
 $\tilde{\pi}(s|s_0) \rightarrow A$

$$\begin{cases} V^\pi(s) = \sum_a \pi(a|s) \cdot Q^\pi(s, a) \\ Q^\pi(s, a) = r(a, s) + \gamma \sum_{s'} P_{ss'}^a V^\pi(s') \end{cases}$$

$\Rightarrow V^\pi(s) = \sum_a \pi(a|s) \left(r(a, s) + \gamma \sum_{s'} P_{ss'}^a V^\pi(s') \right)$

Dynamic Program

optimality Bellman Equation:

$$\begin{cases} V^*(s) = \max_\pi V^\pi(s) = \max_a Q^*(s, a) = \max_a \{ r(a, s) + \gamma \sum_{s'} P_{ss'}^a V^*(s') \} \\ Q^*(s, a) = r(a, s) + \gamma \sum_{s'} P_{ss'}^a V^*(s') \\ = r(a, s) + \gamma \sum_{s'} P_{ss'}^a \max_{a'} Q^*(s', a') \end{cases}$$

policy $\pi^*(a|s) = \begin{cases} 1 & a \in \arg\max Q^*(s, a) \\ 0, \text{ else} & \end{cases}$

⇒ Value Iteration:

$$V_{k+1}(s) \leftarrow \max_a \{ r(a, s) + \gamma \sum_{s'} P_{ss'}^a V_k(s') \}$$

$Q_k(\cdot)$ represents under π_k

↓
optimal Bellman Equation

KEY: $Q_k(s, \pi_k(s))$
 $\geq V_k(s)$
 $= Q_k(s, \pi_{k+1}(s))$

⇒ Policy Iteration

it can be proved:

$$V_{k+1}(s) \geq V_k(s) \text{ with greedy policy } \pi_k$$

$$V_{k+1}(s) \leftarrow \sum_a \pi_k(a|s) (r(a, s) + \gamma \sum_{s'} P_{ss'}^a V_k(s'))$$

↓
converge to

V_{π_k}

where $\pi_k(a|s) = \begin{cases} 1, & a \in \arg\max_{a'} Q_k(s, a') \\ & = r(s, a') + \gamma \sum_{s'} P_{ss'}^a V_k(s') \\ 0, & \text{else} \end{cases}$

Last Part Discussion About DP \leftrightarrow RL (State-Value $V(s)$
& optimality condition)

Bellman Eqn.

$$\textcircled{1} \quad V^\pi(s) = \sum_a \pi(a|s) \cdot Q^\pi(s, a)$$

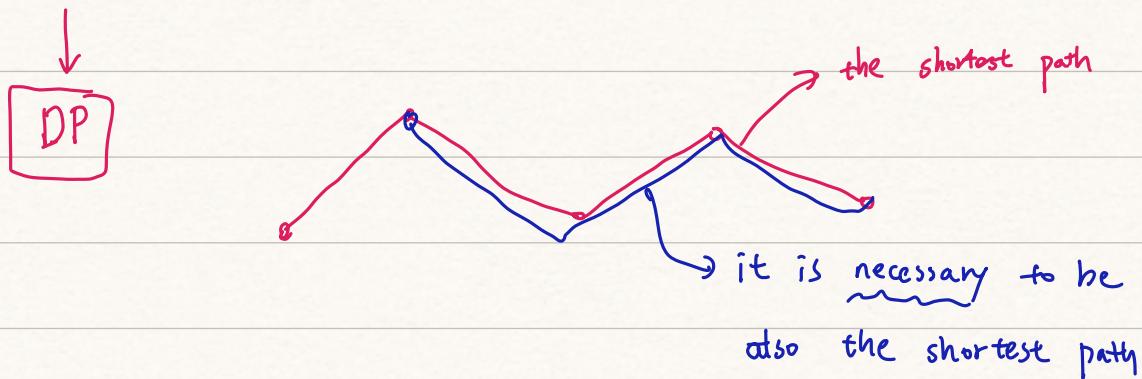
$$= \sum \pi(a|s) \cdot [r(a, s) + \sum_{s'} P_{ss'}^a V^\pi(s')]$$



If $V^\pi(s) \rightarrow$ maxima at state s (optimal)

then for the SUBSEQUENT $V^\pi(s')$,

it still needs to be optimal
(necessary condition)



$$\textcircled{2} \quad Q^\pi(s, a) = r(s, a) + \gamma \sum_{s'} P_{ss'}^a \underbrace{Q^\pi(s', \pi(s'))}_{\text{subsequently optimal}}$$

Slide

Standard

$$1. \quad V^\pi(s) = \sum_a \pi(a|s) \cdot Q^\pi(s, a)$$

$$= \sum_a \pi(a|s) \sum_{s'} P_{ss'}^a (r(s, a, s') + \gamma V^\pi(s'))$$

$$= \sum_a \pi(a|s) \left(R_s^a + \gamma \sum_{s'} P_{ss'}^a V^\pi(s') \right)$$

$$= R_\pi(s) + \gamma \sum_{s'} P_\pi(s, s') V^\pi(s')$$

→ 2. Define Bellman Operator

$$B^\pi(\underline{V}) = R_\pi + \gamma P_\pi \cdot \underline{V}_\pi$$

From (1.), we know $B^\pi(\underline{V}_\pi) = R_\pi + \gamma P_\pi \cdot \underline{V}_\pi$

\downarrow
Bellman State-Value Equation

\downarrow

\underline{V}_π is fixed point of operator B^π !

注: 下考慮 Optimal

we define

$$\begin{aligned} V^*(s) &= \max_\pi V_\pi(s) \\ Q^*(s, a) &= \max_\pi Q_\pi(s, a) \end{aligned}$$

the existence of Optimal Policy \longleftrightarrow Greedy Policy

$G: \underline{V} \rightarrow$ greedy policy ($G(\underline{V})$ is a policy)

$$G(\underline{V})(s) := \underset{a}{\operatorname{argmax}} \left\{ R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V(s') \right\}$$

Construction !!!

→ 考慮 policy:

$$\pi^*: S \rightarrow A$$

$$Q^*(s, a)$$

$$s \mapsto \underset{a}{\operatorname{argmax}} \left\{ \max_\pi Q^\pi(s, a) \right\}$$

Punchline:

the existence of optimal policy

Therefore, $Q^*(s, \pi^*(s))$

$$= \max_a Q^*(s, a)$$

$$= V^*(s)$$

Optimal Bellman Equation

⇒ π^* can yield optimal

State-value function

$$Q^*(s, a) = R(s, a) + \sum_{s'} P_{ss'}^a V^*(s')$$

⇒ $\pi^*(\cdot)$ can lead to $V_{\pi^*} = V^*$ for all $s \in S$

also, $Q_{\pi^*}(s, a) = Q^*(s, a)$ for all

$$3. V_\pi(s) = R_\pi(s) + \gamma \sum_{s'} P_\pi(s, s') V_\pi(s')$$

Define $V^*(s) = \max_\pi V_\pi(s)$ → do not have corresponding policy explicitly

$$\begin{aligned} V_\pi(s) &= \sum_a \pi(a|s) Q_\pi(s, a) = \max_\pi \left\{ \sum_a \pi(a|s) Q_\pi(s, a) \right\} \\ &\leq Q_\pi(s, a) = Q_{\pi^*}(s, a) = \max_\pi \left\{ \sum_a \pi(a|s) (R_s^a + \gamma \sum_{s'} P_{ss'}^a V_\pi(s')) \right\} \\ &\Rightarrow V_\pi(s) \geq Q_{\pi^*}(s, a) = Q^*(s, a) = \max_\pi Q_\pi(s, a) = \max_\pi \left\{ R_s^a + \gamma \sum_{s'} P_{ss'}^a V_\pi(s') \right\} \end{aligned}$$

3. Bellman Optimality Equation

Punchline.

$$\textcircled{1} \quad V^*(s) = \max_\pi V_\pi(s)$$

4. Define Bellman Optimal Operator

$$B_\pi[\underline{V}(s)] = \max_a \{ R_s^a + \gamma \sum_{s'} P_{ss'}^a \underline{V}(s') \}$$

$$\rightarrow \underline{V}^*(s) = \max_\pi \underline{V}_\pi(s) \in \text{fixed point of } B_\pi(\cdot)$$

i.e. $B_\pi(V_\pi(s)) = V_\pi(s)$

$$Q^*(s, a) = \max_\pi Q_\pi(s, a)$$

$$\textcircled{2} \quad \underline{V}^*(s) = \max_{a \in A} Q^*(s, a)$$

$$Q^*(s, a) = R(s, a) + \sum_{s'} P_{ss'}^a \underline{V}^*(s')$$

$$\underline{V}^*(s) = \max_\pi \underline{V}_\pi(s) = \max_\pi Q_\pi(s, a)$$

* the maxima π must have the property that $\pi(a|s)=1$ for some a

from defn

$$\{\underline{V}(s)\} \rightarrow \{Q(s, a)\}$$

解説 \Rightarrow greedy algorithm (policy)

always yielding a better policy ($\pi(a|s)=1$)

$$\textcircled{3} \quad \underline{V}^*(s) = \max_{a \in A} \{ R(s, a) + \sum_{s'} P_{ss'}^a \underline{V}^*(s') \}$$

$$\textcircled{4} \quad Q^*(s, a) = R(s, a) + \sum_{s'} P_{ss'}^a \max_{a' \in A} \{ Q^*(s', a') \}$$

$$V_\pi = B_\pi(V_\pi) \leq B_\star(V_\pi)$$

$$= B_{G(V_\pi)}(V_\pi)$$

$$:= B_\pi(V_\pi)$$

$$\leq B_\pi^N(V_\pi)$$

$$\approx V_\pi$$

→ Define Bellman Optimal Operator:

$$B^*(\underline{V}) := \max_{a \in A} \{ R(s, a) + \sum_{s'} P_{ss'}^a \underline{V}(s') \}$$

\textcircled{3} $\Rightarrow \underline{V}^* \in$ the fixed point of B^* operator

$$\text{i.e., } B^*(\underline{V}^*) = \underline{V}^*$$

4. Greedy policy for Value Function \underline{V}

$$G(\underline{V})(s) := \underset{a}{\operatorname{argmax}} \{ R_s^a + \gamma \sum_{s'} P_{ss'}^a \underline{V}(s') \}$$

相当于筛选

$$\Rightarrow B_{G(\underline{V})}\underline{V} = B^*\underline{V} \quad \text{from definition}$$

policy operator

and $G(\underline{V})$ is greedy policy for value function \underline{V}

Key point:

$$\textcircled{1} \quad B_\pi(V_\pi) = V_\pi$$

$$\textcircled{2} \quad B_\star(V_\star) = V_\star$$

$$\textcircled{3} \quad B_{G(\underline{V})}(\underline{V}) = B^*(\underline{V})$$

γ -contradiction

\Rightarrow fixed point

$$\Rightarrow \lim_{N \rightarrow \infty} B_\pi^N(V_0) = V_\pi$$

→ Recall

Policy Iteration

→ ① Policy Evaluation

given policy $\pi_k \rightarrow$ calculate V_{π_k}

$$\left. \begin{array}{l} \text{Naive} \rightarrow V_{\pi_k} = R_{\pi_k} + P_{\pi_k} \cdot V_{\pi_k} \rightarrow \boxed{\text{equation}} \\ \text{Fixed point} \rightsquigarrow \lim_{N \rightarrow \infty} B_{\pi_k}^N (v_0) = V_{\pi_k} \end{array} \right.$$

$B_{\pi_k}^N \rightarrow$ linear operator

$$\rightarrow \boxed{R_{\pi_k} + P_{\pi_k} \cdot \cdot \cdot}$$

→ ② Policy improvement → Greedy Policy

$$\pi_{k+1} = G(V_{\pi_k}) \rightarrow \boxed{\text{greedy policy}}$$

theoretically, we have

$$V_{\pi_k} = B_{V_{\pi_k}}(V_{\pi_k}) \leq B_x(V_{\pi_k}) = B_{G(V_{\pi_k})}(V_{\pi_k})$$

$$= B_{\pi_{k+1}}(V_{\pi_k})$$

$$\text{Therefore, considering } \underbrace{V_{\pi_{k+1}}}_{\geq V_{\pi_k}} = \lim_{N \rightarrow \infty} B_{\pi_{k+1}}^N(\pi_k)$$

$$\geq V_{\pi_k} !$$

从 operator 角度分析！ (妙) !