

1 Markov Decision Process, MDP

Def [MDP]: $M = (S, A, P, r, \gamma, \mu)$

- ① state space: S , finite / countably infinite
- ② action space: A , finite / infinite
- ③ transition function: $P: S \times A \rightarrow \Delta(S)$ distributions over S
- ④ reward function: $r: S \times A \rightarrow [0, 1]$ or generally $\Delta([0, 1])$
- ⑤ discount factor: $\gamma \in [0, 1)$
- ⑥ initial state distribution $\mu \in \Delta(S)$ to generate S_0 .

Note: (i) for convenience, we assume that:

- ★1: A is finite;
- ★2: r is deterministic;
- ★3: μ is one-point distribution at S_0 .

Def [trajectory]: a trajectory τ_t at time t is the interaction record at time t ,

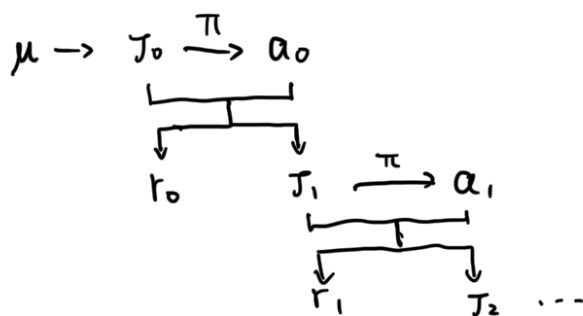
i.e. $\tau_t = (S_0, a_0, r_0, \dots, S_{t-1}, a_{t-1}, r_{t-1}, S_t)$

① $H := \{ \tau_t : t \geq 0 \}$

Def [policy]: a policy is such a mapping: $\pi: H \rightarrow \Delta(A)$

- ① a stationary policy: $\pi: S \rightarrow \Delta(A)$, only depends on S_t .
- ② a deterministic stationary policy: $\pi: S \rightarrow A$.

Note: MDP with policy π runs in such a flow:



Def [value function] $V_M^\pi : S \rightarrow \mathbb{R}$, π : policy, M : MDP.

$$V_M^\pi(s) = E \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid \pi, s_0 = s \right]$$

Note: ① since $\gamma \in [0, 1)$ and $r \in [0, 1]$, $V_M^\pi \in [0, \frac{1}{1-\gamma}]$

Def [action-value function]: $Q_M^\pi : S \times A \rightarrow \mathbb{R}$

$$Q_M^\pi(s, a) = E \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid \pi, s_0 = s, a_0 = a \right]$$

Note: ① Q_M^π is also bounded by $\frac{1}{1-\gamma}$.

② we ignore M when it is clear from context.

The goal of a MDP problem is to find an optimal policy π^* s.t.

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

Assertion: there exists an optimal deterministic and stationary policy π^* .

1.2. Bellman Consistency Equations for stationary policies.

Lemma 1.4 Suppose π : stationary policy, then we have

Bellman
consistency
equations \Rightarrow

$$V^\pi(s) = Q^\pi(s, \pi(s)).$$

$$Q^\pi(s, a) = r(s, a) + \gamma E_{s' \sim p(\cdot | s, a)} [V^\pi(s')].$$

$$\text{pf: } ① Q^\pi(s, \pi(s)) = E_{a \sim \pi(\cdot | s)} [Q^\pi(s, a)]$$

$$= E_{a \sim \pi(\cdot | s)} [E[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid \pi, s_0 = s, a_0 = a]]$$

$$\text{rule of conditional expectation} = E[E[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid \pi, s_0 = s, a_0 = a] \mid a \sim \pi(s)]$$

$$= E[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, \pi]$$

$$= V^\pi(s)$$

$$② Q^\pi(s, a) = E[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid \pi, s_0 = s, a_0 = a]$$

$$\begin{aligned}
&= E[r(s_0, a_0) + \sum_{t=1}^{\infty} \gamma^t r(s_t, a_t) \mid \pi, s_0=s, a_0=a] \\
&= r(s, a) + \gamma E[\sum_{t=1}^{\infty} \gamma^{t-1} r(s_t, a_t) \mid \pi, s_0=s, a_0=a] \\
&= r(s, a) + \gamma E[E[\sum_{t=1}^{\infty} \gamma^{t-1} r(s_t, a_t) \mid s_1=s'] \mid s' \sim p(s, a)] \\
&= r(s, a) + \gamma E_{s' \sim p(s, a)}[V^{\pi}(s')]
\end{aligned}$$

It is easy to see that if a is substituted by $\pi(s)$,

$$\text{we have } Q^{\pi}(s, \pi(s)) = r(s, \pi(s)) + \gamma E_{\substack{a \sim \pi(s) \\ s' \sim p(s, a)}}[V^{\pi}(s')] \quad \square$$

Notation: $P_{(s, a), s'} := P(s' \mid s, a)$

Def [transition matrix on (s, a) with a stationary policy π]

$$P_{(s, a), (s', a')}^{\pi} := P(s' \mid s, a) \pi(a' \mid s')$$

Note: ① for deterministic policies:

$$P_{(s, a), (s', a')}^{\pi} := \begin{cases} P(s' \mid s, a), & \text{if } a' = \pi(s') \\ 0 & , \text{ else.} \end{cases}$$

With this notation, we have the following equations:

$$\begin{aligned}
\text{(H)} \quad Q^{\pi} &= r + \gamma P V^{\pi}; \\
Q^{\pi} &= r + \gamma P^{\pi} Q^{\pi}.
\end{aligned}$$

Pf: for a stationary policy π ,

according to Lemma 1.4, we have

$$\begin{aligned}
Q^{\pi}(s, a) &= r(s, a) + \gamma E_{s' \sim p(s, a)}[V^{\pi}(s')] \\
&= r(s, a) + \gamma \sum_{s'} P_{(s, a), s'} V^{\pi}(s') \quad (1)
\end{aligned}$$

we can rewrite it as

$$Q^{\pi} = r + \gamma P V^{\pi}.$$

since Lemma 1.4, $V^\pi(s') = Q^\pi(s', \pi(s'))$

$$= \sum_{a'} \pi(a'|s') Q^\pi(s', a') \quad (2)$$

Introduce (2) into (1):

$$\begin{aligned} Q^\pi(s, a) &= r(s, a) + \gamma \sum_{s'} \sum_{a'} P_{(s, a), s'} \pi(a'|s') Q^\pi(s', a') \\ &= r(s, a) + \gamma \sum_{(s', a')} P_{(s, a), (s', a')}^\pi Q^\pi(s', a'). \end{aligned}$$

We can rewrite it as

$$Q^\pi = r + \gamma P^\pi Q^\pi. \quad \square$$

Corollary 1.5: Suppose π is stationary, we have

$$Q^\pi = (I - \gamma P^\pi)^{-1} r.$$

Pf: only need to show $I - \gamma P^\pi$ is invertible.

for any non-zero vector $x \in \mathbb{R}^{|S| \cdot |A|}$,

$$\begin{aligned} \|(I - \gamma P^\pi)x\|_\infty &= \|x - \gamma P^\pi x\|_\infty \\ &\geq \|x\|_\infty - \|\gamma P^\pi x\|_\infty \\ &\geq \|x\|_\infty - \gamma \|x\|_\infty > 0 \quad (\text{since } \gamma < 1 \text{ \& } \|x\|_\infty > 0) \end{aligned}$$

which implies $I - \gamma P^\pi$ is full rank.

An intuitive way to see this is:

according to Gershgorin circle th. \leftarrow since P^π is a stochastic matrix, the spectral radius of P^π should be 1, the the spectral radius of γP^π is strictly smaller than 1, which implies $\det(I - \gamma P^\pi) \neq 0$. \square

Lemma 1.6.

$$[(1 - \gamma)(I - \gamma P^\pi)^{-1}]_{(s, a), (s', a')} = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t P^\pi(s_t = s', a_t = a' | s_0 = s, a_0 = a)$$

$$\begin{aligned} \text{Pf: } (I - \gamma P^\pi) [I + \gamma P^\pi + \gamma^2 (P^\pi)^2 + \dots + \gamma^t (P^\pi)^t] \\ = I - \gamma^{t+1} (P^\pi)^{t+1} \end{aligned}$$

since $\max (P^\pi)^{t+1} = \max p^\pi (P^\pi)^t \leq \max (P^\pi)^t \leq \dots \leq \max p^\pi$
and $\gamma \in [0, 1)$,

when $t \rightarrow \infty$, we have:

$$\textcircled{1} \quad \sum_{t=0}^{\infty} \gamma^t (P^\pi)^t \text{ converges}$$

$$\textcircled{2} \quad (I - \gamma P^\pi) \left(\sum_{t=0}^{\infty} \gamma^t (P^\pi)^t \right) = I$$

$$\textcircled{2} \text{ implies } (I - \gamma P^\pi)^{-1} = \sum_{t=0}^{\infty} \gamma^t (P^\pi)^t.$$

On the other hand.

$$\begin{aligned} (P^\pi)^t_{(s,a),(s',a')} &= \sum_{(s_1,a_1)} P^\pi_{(s,a),(s_1,a_1)} (P^\pi)^{t-1}_{(s_1,a_1),(s',a')} \\ &= \dots \\ &= \sum_{(s_{t-1},a_{t-1})} \dots \sum_{(s_1,a_1)} P^\pi_{(s,a),(s_1,a_1)} \dots P^\pi_{(s_{t-1},a_{t-1}),(s',a')} \\ &= |p^\pi(s_t=s', a_t=a' | s_0=s, a_0=a)| \end{aligned}$$

then we have

$$[(I - \gamma P^\pi)^{-1}]_{(s,a),(s',a')} = \sum_{t=0}^{\infty} \gamma^t |p^\pi(s_t=s', a_t=a' | s_0=s, a_0=a)| \quad \square$$

1.3 Bellman Optimally Equations.

Now we proof the assertion at the beginning.

Theorem 1-7 Let Π be the set of all non-stationary and randomized policies. Define:

$$V^*(s) := \sup_{\pi \in \Pi} V^\pi(s)$$

$\Rightarrow V^*, Q^*$ are bounded by $\frac{1}{1-\gamma}$.

$$Q^*(s,a) := \sup_{\pi \in \Pi} Q^\pi(s,a)$$

There exists a stationary and deterministic policy π s.t. for

all $s \in S$ and $a \in A$,

$$V^\pi(s) = V^*(s),$$

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

Pf: let $(S_0, A_0, R_0, S_1, A_1, R_1, \dots, S_{t-1}, A_{t-1}, R_{t-1}, S_t)$

denote a random T_t . where S_i, A_j, R_k are r.v..

① First we show that

$$\sup_{\pi \in \Pi} E \left[\sum_{t=0}^{\infty} \gamma^t r(S_t, a_t) \mid \pi, (S_0, A_0, R_0, S_1) = (s, a, r, s') \right] = \gamma V^*(s'). \quad (\star)$$

For $\forall \pi \in \Pi$, define an "offset" policy $\pi_{(s, a, r)}$:

$$\pi_{(s, a, r)}(A_t = a \mid S_0 = s_0, A_0 = a_0, R_0 = r_0, \dots, S_t = s_t)$$

$$:= \pi(A_{t+1} = a \mid S_0 = s, A_0 = a, R_0 = r, S_1 = s_0, A_1 = a, R_1 = r_0, \dots, S_{t+1} = s_t)$$

By the Markov property:

$$\text{LHS of } (\star) = \sup_{\pi \in \Pi} \gamma E \left[\sum_{t=0}^{\infty} \gamma^t r(S_t, a_t) \mid \pi_{(s, a, r)}, S_0 = s' \right]$$

$$= \sup_{\pi \in \Pi} \gamma V^{\pi_{(s, a, r)}}(s')$$

we have for all (s, a, r) , $\{\pi_{(s, a, r)} \mid \pi \in \Pi\} = \Pi$.

$\forall \pi \in \Pi$, let $\pi'_{(s, a, r)} = \pi$

and $\pi'(s)$ is one-point distribution.

then $\pi' \in \Pi \Rightarrow \{\pi_{(s, a, r)} \mid \pi \in \Pi\} \supset \Pi$

$$\text{We have LHS of } (\star) = \gamma \cdot \sup_{\pi \in \Pi} V^\pi(s')$$

$$= \gamma \cdot V^*(s')$$

$$= \text{RHS of } (\star)$$

② We now show the deterministic and stationary policy $\hat{\pi}$ is optimal.

(2) we now show the deterministic and stationary policy $\hat{\pi}$ is optimal:

$$\hat{\pi}(s) = \arg \sup_{a \in A} E[r(s, a) + \gamma V^*(s_1) | (s_0, a_0) = (s, a)]$$

For this, we have

$$\begin{aligned} V^*(s_0) &= \sup_{\pi \in \Pi} E[r(s_0, a_0) + \sum_{t=1}^{\infty} \gamma^t r(s_t, a_t) | \pi, s_0 = s_0] \\ &= \sup_{\pi \in \Pi} E[r(s_0, a_0) + E[\sum_{t=1}^{\infty} \gamma^t r(s_t, a_t) | \pi, (s_0, a_0, r_0, s_1) = (s_0, a_0, r_0, s_1)]] \\ &\leq \sup_{\pi \in \Pi} E[r(s_0, a_0) + \sup_{\pi' \in \Pi} E[\sum_{t=1}^{\infty} \gamma^t r(s_t, a_t) | \pi', (s_0, a_0, r_0, s_1) = (s_0, a_0, r_0, s_1)]] \\ &= \sup_{\pi \in \Pi} E[r(s_0, a_0) + \gamma V^*(s_1)] \\ &= \sup_{a_0 \in A} E[r(s_0, a_0) + \gamma V^*(s_1)] \\ &= E[r(s_0, a_0) + \gamma V^*(s_1) | \hat{\pi}] \end{aligned}$$

By applying the argument recursively:

$$\begin{aligned} V^*(s_0) &\leq E[r(s_0, a_0) + \gamma V^*(s_1) | \hat{\pi}] \\ &\leq E[r(s_0, a_0) + \gamma E[r(s_1, a_1) + \gamma V^*(s_2) | \hat{\pi}] | \hat{\pi}] \\ &= E[r(s_0, a_0) + \gamma r(s_1, a_1) + \gamma^2 V^*(s_2) | \hat{\pi}] \\ &\leq \dots \\ &\leq V^{\hat{\pi}}(s_0) \end{aligned}$$

since $V^{\hat{\pi}}(s_0) \leq V^*(s_0)$ by definition,

we have $V^{\hat{\pi}} = V^*$ \square

Notation ①: $\pi_Q := \arg \max_{a \in A} Q(s, a)$

②: $V_Q(s) := \max_{a \in A} Q(s, a)$

Theorem 1.8 [Bellman optimality equations , BOE]

A vector $Q \in \mathbb{R}^{|S||A|}$ satisfies the BOE if

$$Q(s,a) = r(s,a) + \gamma \mathbb{E}_{s' \sim p(\cdot|s,a)} \left[\max_{a' \in A} Q(s',a') \right].$$

Then for any $Q \in \mathbb{R}^{|S||A|}$, $Q = Q^* \Leftrightarrow Q$ satisfies BOE.

$\pi(s) \in \arg \max_{a \in A} Q^*(s,a)$ is an optimal policy.

By the notation: the optimal policy $\pi^* = \pi_{Q^*}$

$$T_M : \mathbb{R}^{|S||A|} \rightarrow \mathbb{R}^{|S||A|}, \quad TQ := r + \gamma P V_Q$$

Bellman optimality operator

we can write BOE: $Q = TQ$. so the theorem states that

$Q = Q^* \Leftrightarrow Q$ is a fixed point of T .

Pf: ① First we want to show $V^*(s) = \max_a Q^*(s,a)$ (A)

Let π^* be an optimal stationary and deterministic policy

$V^*(s) = \sup_{\pi \in \Pi} V^\pi(s)$, consider such a policy π' : take action a

then follow π^* . $\pi' \in \Pi$

$$\Rightarrow V^*(s) \geq V^{\pi'}(s) = Q^{\pi'}(s,a) \stackrel{\text{Th 1.7}}{=} Q^*(s,a) \quad (\forall a \in A)$$

$$\Rightarrow V^*(s) \geq \max_a Q^*(s,a)$$

On the other hand, by Th 1.7 & Lem 1.4 (since π^* is stationary)

$$\begin{aligned} V^*(s) &= V^{\pi^*}(s) = Q^{\pi^*}(s, \pi^*(s)) \\ &\leq \max_a Q^{\pi^*}(s,a) \\ &= \max_a Q^*(s,a) \end{aligned}$$

we prove (A) holds.

② sufficiency: Q^* satisfies $Q^* = TQ^*$

for all actions $a \in A$:

$$\begin{aligned} Q^*(s, a) &= Q^{\pi^*}(s, a) \\ &= r(s, a) + \gamma E_{s' \sim p(s, a)} [V^{\pi^*}(s')] \\ &= r(s, a) + \gamma E_{s' \sim p(s, a)} [V^*(s')] \\ &= r(s, a) + \gamma E_{s' \sim p(s, a)} [\max_{a'} Q^*(s', a')] \\ &= \mathcal{T} Q^* . \end{aligned}$$

③ necessity : assume $Q = \mathcal{T} Q$, we show that $Q = Q^*$

$$\text{Let } \pi = \pi_Q \quad . \quad Q = \mathcal{T} Q = r + \gamma P V_Q$$

since $\pi = \pi_Q = \arg \max_{\alpha} Q(s, a)$ is stationary and deterministic

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

$$\Rightarrow Q = r + \gamma P^{\pi} Q$$

$$\Rightarrow Q = (I - \gamma P^{\pi})^{-1} r \stackrel{\text{Cor 1.5}}{=} Q^{\pi}$$

i.e. Q is action value of π_Q .

\forall stationary and deterministic π' , we have

$$[(P^{\pi} - P^{\pi'}) Q^{\pi}]_{s, a} = E_{s' \sim p(\cdot | s, a)} [Q^{\pi}(s', \pi(s)) - Q^{\pi}(s', \pi'(s))] \geq 0$$

$$\begin{aligned} \Rightarrow Q - Q^{\pi'} &= Q^{\pi} - Q^{\pi'} \\ &= Q^{\pi} - (I - \gamma P^{\pi'})^{-1} r \\ &= (I - \gamma P^{\pi'})^{-1} [(I - \gamma P^{\pi'}) - (I - \gamma P^{\pi})] Q^{\pi} \\ &= \gamma (I - \gamma P^{\pi'})^{-1} (P^{\pi} - P^{\pi'}) Q^{\pi} \\ &\geq 0 \end{aligned}$$

By Lem 1.6 $\gamma (I - \gamma P^{\pi'})^{-1} \geq 0$.

Specifically, $Q^{\pi} = Q \geq Q^{\pi^*} = Q^* \geq Q^{\pi}$

$$\Rightarrow Q = Q^{\pi} = Q^* \quad \square$$

