

MARKOV DECISION PROCESSES

Consider an agent a , in an environment \mathcal{E} , restricted to the state-space \mathcal{S} , with an available action-space \mathcal{A} , and securing rewards from the reward-space \mathcal{R} .

Mathematically, we have:

$$\mathcal{S}, \mathcal{A} \in \text{Ob}(\text{Set}), \mathcal{E} = \mathcal{S} \times \mathcal{A}, \mathcal{R} \in \mathbb{R}^{\dim(\mathcal{E})}$$

Let the transition-probabilities be written as:

$$P_{x,y}[a] := P[Y_n = y \mid X_n = x, A_n = a]$$

Defⁿ: A **Markov Decision Process** \mathcal{M} , is described by the 4-tuple $(\mathcal{S}, \mathcal{A}, P, \mathcal{R})$.

- At $t=n$, let the agent register state $x_n \in \mathcal{X}$, choose an action $a \in \mathcal{A}$, and receive a probabilistic reward $r_n \in \mathbb{R}$, s.t.:

$$\mathbb{E}[r_n](x_n, a_n) =: p_{x_n}(a_n)$$

Let the discounting factor be $\gamma \in]0, 1[$.

- If $(R_t)_{t \in \mathbb{N} \cup \{0\}}$ describes the reward process and $(G_t)_{t \in \mathbb{N} \cup \{0\}}$ " " " " cumulative " " , then,

$$\begin{aligned} G_t &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \infty \\ &= R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \dots \infty) \\ &= R_{t+1} + \gamma G_{t+1} \end{aligned}$$

Defⁿ: let $x \in \mathcal{X}$. A **policy** function $\pi: \mathcal{X} \rightarrow \mathcal{A}$ is defined as the solution for M .

Defⁿ:

(i) For a given policy π , the **value** of state x , under π , written $V^\pi: \mathcal{S} \rightarrow \mathbb{R}$, is defined as:

$$\begin{aligned} V^\pi(x) &= \mathbb{E}[G_t | S_t = x] \\ &= \mathbb{E}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = x\right] \end{aligned}$$

• $V^\pi(x)$ is the same as the expected reward at x , ~~plus~~ plus the discounted V^π at next step:

$$V^\pi(x) = \int_x (\pi(x)) + \gamma \sum_{y \in \mathcal{S}(x)} P_{x,y}[\pi(x)] V^\pi(y)$$

(ii) For a given π , **Q-value** of $(x, a) \in \mathcal{E}$, is defined as:

$$Q^\pi(x, a) = \int_x(a) + \gamma \sum_{y \in \mathcal{S}(x)} P_{x,y}[\pi(x)] V^\pi(y)$$

TASK (Optimization problem):

Theorem [Bellman & Dreyfus, '62]: For a M.D.P. \mathcal{M} , described with, as above, $\exists \pi^*$ (**optimal policy**), s.t.:

$$V^* := V^{\pi^*}(x) = \max_{a \in A(x)} \left\{ r(x, a) + \gamma \sum_{y \in \mathcal{X}(x)} P_{x,y}[a] V^{\pi^*}(y) \right\} \quad \dots (3.5)$$

Objective: Determine Q -values for an optimal policy:

$$Q^*(x, a) := Q^{\pi^*}(x, a), \quad \forall x \in \mathcal{X}, \forall a \in A$$

Then, we have: $V^*(x) = \max_{a \in A(x)} \{Q^*(x, a)\},$

corresponding to an optimal policy $\pi^*(x)$.

Q-learning ALGORITHM:

Defⁿ: Consider a M.D.P. described as above. We define the Q-learning, as the sequence $(Q_n(\cdot, \cdot): \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R})_{n \in \mathbb{N} \cup \{0\}}$:

$$Q_n(x, a) := \begin{cases} (1 - \alpha_n) Q_{n-1}(x, a) + \alpha_n (r_n + \gamma V_{n-1}(y_n)), & \text{if } (x, a) = (x_n, a_n) \\ Q_{n-1}(x, a) & , \text{ otherwise} \end{cases} \quad \text{--- (1)}$$

Here, (α_n) , (r_n) are, resp., the learning-rate sequence and the reward-sequence, and $V_{n-1}(y)$ is:

$$V_{n-1}(y) = \max_{b \in \mathcal{A}} \{Q_{n-1}(y, b)\} \quad \text{--- (2)}$$

Notⁿ: Define $n^i(x, a)$: index of i^{th} time, a is tried in x .
Define $I := \{n^i : i \in \mathbb{N}\}$

Theorem (C) [convergence theorem]:

Given bdd rewards, $|R_n| \leq R$, and learning rates $0 \leq \alpha_n < 1$, with $(\alpha_n)_{n \in \mathbb{I}} \notin l^1(I)$, and $(\alpha_n)_{n \in \mathbb{I}} \in l^2(I)$.

$$\left[\sum_{i=1}^{\infty} \alpha_{n_i(x,a)} = +\infty \right]$$

$$\left[\sum_{i=1}^{\infty} \alpha_{n_i(x,a)}^2 < +\infty \right]$$

Then, we have:

$$Q_n(x,a) \xrightarrow{n \rightarrow +\infty} Q^*(x,a), \quad \text{w.p. 1.}$$

• Strategy of Proof - \mathcal{C} :

Step 1: Action-Replay-Process, A , description

Step 2: • lemma A: Q_n optimal for A

• lemma B1-B3: { preparatory lemmas, and,
 $A \rightarrow$ real process

Step 3: use lemmas A-B3, to prove \mathcal{C} .

Action-Replay-Process [A]:

Consider the state-space $\{ \langle x, n \rangle \mid x \in \mathcal{X} \}$, action-space \mathcal{A} .
Let the discount factor $\gamma \in]0, 1[$ for A, be the same as well.
Recall n^i :

$$n^i := \underline{n^i(x, a)} : \begin{cases} \text{index of the } i^{\text{th}}\text{-time, action } a \\ \text{was tried at space } x \end{cases}$$

Also, define i_* as:

$$i_* := \begin{cases} \operatorname{argmax}_i \{ n^i < n \}, & \text{if } (x, a) \text{ expected before 'n'} \\ 0, & \text{otherwise} \end{cases}$$

$\Rightarrow n^{i_*}$ is the last time before episode 'n', that (x, a) was expected

Now, if $i_* = 0$, then the reward is set as $Q_0(\pi, a)$
 $\Rightarrow A$ absorbs

otherwise,

let i_e be the index of the episode replayed as :

$$i_e = \begin{cases} i_* & , \quad \text{with prob. } \alpha_{n i_*} \\ i_* - 1 & , \quad \text{with prob. } (1 - \alpha_{n i_*}) \alpha_{n i_* - 1} \\ i_* - 2 & , \quad \text{with prob. } (1 - \alpha_{n i_*}) (1 - \alpha_{n i_* - 1}) \alpha_{n i_* - 2} \\ \vdots & \\ 0 & , \quad \text{with prob. } \prod_{i=1}^{i_*} (1 - \alpha_{n i}) \end{cases}$$

Lemma A: [Q_n optimal for A]

$Q_n(x, a)$ are the optimal action-values for A , with state $\langle x, n \rangle$ and action a :

$$Q_n(x, a) = Q_A^*(\langle x, n \rangle, a), \quad \forall (a, x) \in \mathcal{E}, \forall n \geq 0$$

Proof. From the construction of A : $Q_0(x, a) = Q_A^*(\langle x, 0 \rangle, a)$

Hence, the theorem holds for $n=0$.

Suppose Q_{n-1} -values (by Q-learning) are optimal for A :

$$Q_{n-1}(x, a) = Q_A^*(\langle x, n-1 \rangle, a), \quad \forall (a, x) \in \mathcal{A} \times \mathcal{S}$$

[Induction hypothesis]

$$\begin{aligned} \Rightarrow V^*(\langle x, n-1 \rangle) &= V_{n-1}(x) \\ &= \max_{a \in \mathcal{A}(x)} \{Q_{n-1}(x, a)\} \end{aligned}$$

Case (i): $(x, a) \neq (x_n, a_n)$

$$\begin{aligned}\Rightarrow Q_n(x, a) &= Q_{n-1}(x, a) \\ &= Q_{\mathbb{A}}^*(\langle x, n-1 \rangle, a) \\ &= Q_{\mathbb{A}}^*(\langle x, n \rangle, a)\end{aligned}$$

Case ii: $(x, a) = (x_n, a_n)$

$$\begin{aligned}Q_{\mathbb{A}}^*(\langle x_n, n \rangle, a_n) &= (1 - \alpha_n) Q_{\mathbb{A}}^*(\langle x_n, n-1 \rangle, a_n) + \alpha_n (r_n + \gamma V^*(\langle y_n, n-1 \rangle)) \\ &= (1 - \alpha_n) Q_{n-1}(x_n, a_n) + \alpha_n (r_n + \gamma V_{n-1}(y_n)) \\ &= Q_n(x_n, a_n)\end{aligned}$$

~~Hence~~, from the induction hypothesis and Q_n -iteration formula (1).

Hence, $Q_n(x, a) = Q_{\mathbb{A}}^*(\langle x, n \rangle, a)$, $\forall (a, x) \in \mathbb{A} \times \mathbb{X}$



Lemma B.1: [Discounting infinite sequence]

consider a finite Markov Process, with discounting factor γ , bounded-rewards ($|r_n| < R$), transition-probabilities $P_{x,y}[a]$.
Let $X_s = (x_0, x_1, \dots, x_s)$ be the fixed s -steps. Then,

$$\forall \pi: \quad |V^\pi(X_s) - V^\pi(x_0)| \xrightarrow{s \rightarrow +\infty} 0$$

Proof.

Ignoring the value of $(s+1)^{\text{th}}$ -state, incurs the penalty:

$$\delta := \gamma^s \sum_{y_{s+1} \in \mathcal{S}(y_s)} P_{y_s, y_{s+1}}[a_s] V^\pi(y_{s+1})$$

Since $|x_n| < R$, $\forall n \in \mathbb{N} \cup \{0\}$:

$$\begin{aligned} V^\pi(y_{s+1}) &< R + \gamma R + \gamma^2 R + \dots \infty \\ &= R \left(\frac{1}{1-\gamma} \right) = \frac{R}{1-\gamma} \end{aligned}$$

$$\Rightarrow V^\pi(y_{s+1}) < \frac{R}{1-\gamma}$$

$$\begin{aligned} \therefore |\delta| &< \gamma^s \sum_{y_{s+1} \in \mathcal{Y}(y_s)} P_{y_s, y_{s+1}}[a_s] \left(\frac{R}{1-\gamma} \right) \\ &= \gamma^s \left(\frac{R}{1-\gamma} \right) (1) = \frac{\gamma^s R}{1-\gamma} \end{aligned}$$

$$\Rightarrow |\delta| < \frac{\gamma^s R}{1-\gamma} \xrightarrow{s \rightarrow +\infty} 0$$



Lemma B.2 : [Rewards & transⁿ probabilities converge]

w. P. 1,
$$P_{A(x,y)}^{(n)}[a] \xrightarrow{n \rightarrow +\infty} P_{(x,y)}^{(n)}[a]$$

and
$$\int_{A_x}^{(n)}[a] \xrightarrow{n \rightarrow +\infty} \int_x^{(n)}[a], \quad \forall a \in \mathcal{A}$$

Proof.

By the theorem of Kushner & Clark, 1978, " if $(X_n)_{n \in \mathbb{N} \cup \{0\}}$:

$$X_{n+1} = X_n + \beta_n (\xi_n - X_n),$$

with $0 \leq \beta_n < 1$, $(\beta_n) \notin l^1$, $(\beta_n) \in l^2$, ξ_n bdd, $\mathbb{E}[\xi_n] = \Xi$,
then,

$$X_n \xrightarrow{n \rightarrow +\infty} \Xi, \quad \text{w. P. 1.}$$

In our case:

$$f_{\langle x, n^{i+1} \rangle}(a) = f_{\langle x, n^i \rangle}(a) + \alpha_{n^{i+1}} (r_{n^{i+1}} - f_{\langle x, n^i \rangle}(a)),$$

and,

$$P_{A(x,y)}^{(n^{i+1})}[a] = P_{A(x,y)}^{(n^i)}[a] + \alpha_{n^{i+1}} (\mathbb{1}_{\{y_n=y\}} - P_{A(x,y)}^{(n^i)}[a]),$$

such that,

$$\mathbb{E}[r_n] = f_x(a), \quad \mathbb{E}[\mathbb{1}_{\{y_n=y\}}] = P_{x,y}[a]$$

hence, the theorem applies to both, and we have proved our lemma.



Lemma B.3: [close reward & probabilities \Rightarrow close values]

Consider an s -step Markov chain, formed according to the probabilities $(P_{(x_n, x_{n+1})}^{(n)}[a_n])_{n \in \{1, 2, \dots, s\}}$. Let $\eta > 0$ be given.

Let $\bar{Q}(x, a_1, \dots, a_s)$: expected reward for the real process
 $\bar{Q}'(x, a_1, \dots, a_s)$: " " " " Markov Chain

If we have: $|P^i[a] - P_{x,y}[a]| < \frac{\eta}{R}$, $\forall a, x, y, \forall i \in \{1, \dots, s\}$

and, $|\int_x^{(i)}(a) - \int_x(a)| < \eta$, $\forall a, x, \forall i \in \{1, \dots, s\}$

then, $|\bar{Q}'(x, a_1, \dots, a_s) - \bar{Q}(x, a_1, \dots, a_s)| < \frac{s(s+1)}{2} \eta$

Proof. We have: $\bar{Q}(x, a_1, a_2) = f_x(a_1) + \gamma \sum_{y \in \mathcal{S}(x)} P_{x,y} f_y(a_2)$

and, $\bar{Q}'(x, a_1, a_2) = f_x^{(1)}(a_1) + \gamma \sum_{y \in \mathcal{S}(x)} P_{x,y}^{(1)} f_y^{(2)}(a_2)$

$$\begin{aligned} \Rightarrow |\bar{Q}'(x, a_1, a_2) - \bar{Q}(x, a_1, a_2)| &\leq |f_x^{(1)}(a_1) - f_x(a_1)| \\ &\quad + \left| \gamma \sum_{y \in \mathcal{S}(x)} P_{x,y}^{(1)} (f_y^{(2)} - f_y)(a_2) \right| \\ &\quad + \left| \gamma \sum_{y \in \mathcal{S}(x)} f_y(a_2) (P_{x,y}^{(1)} - P_{x,y}) \right| \\ &\leq \eta + \gamma \cdot \eta \cdot 1 + \gamma \cdot R \cdot \frac{\eta}{R} \\ &\stackrel{(\gamma \leq 1)}{\leq} \eta + \eta + \eta = 3\eta \end{aligned}$$

Similarly, for the s -step chain, we get $\frac{s(s+1)}{2}$ terms and we get the result.



Proof of \mathcal{Q} :

Putting the above proved lemmas together, we get the following:

From lemma B.2:
$$P_{A(x,y)}^{(n)}[a] \xrightarrow{n \rightarrow +\infty} P_{x,y}^{(n)}[a]$$

$$P_{A_x}^{(n)}(a) \xrightarrow{n \rightarrow +\infty} P_x^{(n)}(a), \quad \forall a \in A$$

From lemma B.3:

$$Q_A^*(x, a_1, \dots, a_r) \longrightarrow Q^*(x, a_1, \dots, a_r)$$

$$\begin{aligned} &\Rightarrow Q_A^*(x, a) \longrightarrow Q^*(x, a) \\ &(\because \text{lemma B.1}) \end{aligned}$$

From lemma A: $Q_n(x, a) \longrightarrow Q_A^*(x, a)$

$$\therefore Q_n(x, a) \longrightarrow Q^*(x, a)$$

