

CSCI S-89C Deep Reinforcement Learning

Harvard Summer School

Dmitry Kurochkin

Summer 2020

Lecture 2

Contents

1 Finite Markov Decision Processes (MDP)

- Agent–Environment Interface
- Objective
- Policy
- Value Functions
 - State-value function $v_\pi(s)$ for policy $\pi(a|s)$
 - Action-value function $q_\pi(s, a)$ for policy $\pi(a|s)$
- Optimal Value Functions
 - Optimal state-value function $v_*(s)$
 - Existence of optimal policy $\pi_*(a|s)$
 - Optimal action-value function $q_*(s, a)$
- Bellman Equation
 - Bellman equation for $v_\pi(s)$
 - Bellman optimality equation for $v_*(s)$
 - Bellman optimality equation for $q_*(s, a)$

2 Dynamic Programming (DP)

- Iterative Policy Evaluation

Contents

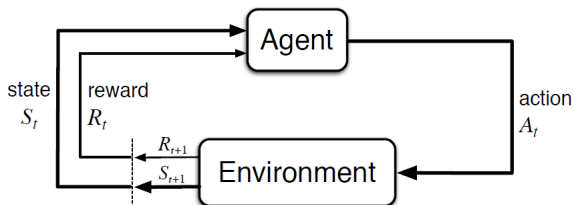
1 Finite Markov Decision Processes (MDP)

- Agent–Environment Interface
- Objective
- Policy
- Value Functions
 - State-value function $v_{\pi}(s)$ for policy $\pi(a|s)$
 - Action-value function $q_{\pi}(s, a)$ for policy $\pi(a|s)$
- Optimal Value Functions
 - Optimal state-value function $v_{*}(s)$
 - Existence of optimal policy $\pi_{*}(a|s)$
 - Optimal action-value function $q_{*}(s, a)$
- Bellman Equation
 - Bellman equation for $v_{\pi}(s)$
 - Bellman optimality equation for $v_{*}(s)$
 - Bellman optimality equation for $q_{*}(s, a)$

2 Dynamic Programming (DP)

- Iterative Policy Evaluation

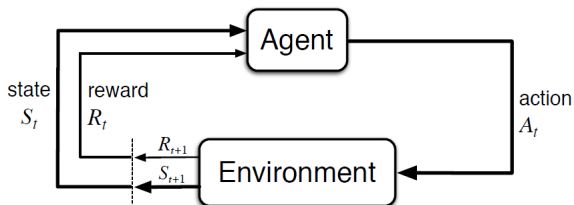
Agent–Environment Interactions



\mathcal{S} = set of all possible states S_t of the environment

$\mathcal{A}(s)$ = set of all admissible actions A_t the agent can take, given $S_t = s$

Agent–Environment Interactions



The *trajectory* will be

$\underline{S_0}, \underline{A_0}, \underline{R_1}, \underline{S_1}, \underline{A_1}, \underline{R_2}, \underline{S_2}, \underline{A_2}, \underline{R_3}, \underline{S_3}, \underline{A_3}, \dots, \underline{R_{t-1}}, \underline{S_{t-1}}, \underline{A_{t-1}}, \dots$

Assume **Markov property**, i.e. for each $s \in \mathcal{S}$ and $a \in \mathcal{A}(s)$:

$$P\{S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a, R_{t-1} = r_{t-1}, S_{t-2} = s_{t-2}, \dots\} = P\{S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a\} \doteq p(s', r | s, a),$$

for any history r_{t-1}, s_{t-2}, \dots

Source: *Reinforcement Learning: An Introduction* by R. Sutton and A. Barto

Agent-Environment Interactions

$$p(s', r|s, a) \doteq P\{S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a\}$$

completely defines the dynamics of the Markov decision processes

We notice that

$$\sum_{s', r} p(s', r|s, a) = 1$$

$$\sum_r p(s', r|s, a) = P\{S_t = s' | S_{t-1} = s, A_{t-1} = a\} \doteq p(s'|s, a)$$

Agent–Environment Interactions

$$p(s', r | s, a) \doteq P\{S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a\}$$

completely defines the dynamics of the Markov decision processes

We notice that

$$\sum_{s', r} p(s', r | s, a) = 1$$

$$\sum_r p(s', r | s, a) = P\{S_t = s' | S_{t-1} = s, A_{t-1} = a\} \doteq p(s' | s, a)$$

The expected reward depends on the current state s and action a only:

$$E[R_t | S_{t-1} = s, A_{t-1} = a, R_{t-1} = r_{t-1}, S_{t-2} = s_{t-2}, \dots] =$$

$$E[R_t | S_{t-1} = s, A_{t-1} = a] = \sum_r r \sum_{s'} p(s', r | s, a) \doteq r(s, a)$$

Contents

1 Finite Markov Decision Processes (MDP)

- Agent–Environment Interface
- Objective
- Policy
- Value Functions
 - State-value function $v_{\pi}(s)$ for policy $\pi(a|s)$
 - Action-value function $q_{\pi}(s, a)$ for policy $\pi(a|s)$
- Optimal Value Functions
 - Optimal state-value function $v_{*}(s)$
 - Existence of optimal policy $\pi_{*}(a|s)$
 - Optimal action-value function $q_{*}(s, a)$
- Bellman Equation
 - Bellman equation for $v_{\pi}(s)$
 - Bellman optimality equation for $v_{*}(s)$
 - Bellman optimality equation for $q_{*}(s, a)$

2 Dynamic Programming (DP)

- Iterative Policy Evaluation

Define Goal

$G_t \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1},$
where $\gamma \in [0, 1]$ is *discount* rate

Contents

1 Finite Markov Decision Processes (MDP)

- Agent–Environment Interface
- Objective
- Policy
- Value Functions
 - State-value function $v_{\pi}(s)$ for policy $\pi(a|s)$
 - Action-value function $q_{\pi}(s, a)$ for policy $\pi(a|s)$
- Optimal Value Functions
 - Optimal state-value function $v_{*}(s)$
 - Existence of optimal policy $\pi_{*}(a|s)$
 - Optimal action-value function $q_{*}(s, a)$
- Bellman Equation
 - Bellman equation for $v_{\pi}(s)$
 - Bellman optimality equation for $v_{*}(s)$
 - Bellman optimality equation for $q_{*}(s, a)$

2 Dynamic Programming (DP)

- Iterative Policy Evaluation

Define Policy

$$\pi(a|s) \doteq P\{A_t = a | S_t = s\}$$

Contents

1 Finite Markov Decision Processes (MDP)

- Agent–Environment Interface
- Objective
- Policy
- Value Functions
 - State-value function $v_{\pi}(s)$ for policy $\pi(a|s)$
 - Action-value function $q_{\pi}(s, a)$ for policy $\pi(a|s)$
- Optimal Value Functions
 - Optimal state-value function $v_{*}(s)$
 - Existence of optimal policy $\pi_{*}(a|s)$
 - Optimal action-value function $q_{*}(s, a)$
- Bellman Equation
 - Bellman equation for $v_{\pi}(s)$
 - Bellman optimality equation for $v_{*}(s)$
 - Bellman optimality equation for $q_{*}(s, a)$

2 Dynamic Programming (DP)

- Iterative Policy Evaluation

State-value function $v_\pi(s)$ for policy $\pi(a|s)$

$$\begin{aligned} v_\pi(s) &\doteq E_\pi [G_t | S_t = s] \\ &= E_\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \middle| S_t = s \right] \end{aligned}$$

Action-value function $q_\pi(s, a)$ for policy $\pi(a|s)$

$$\begin{aligned} q_\pi(s, a) &\doteq E_\pi [G_t | S_t = s, A_t = a] \\ &= E_\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \middle| S_t = s, A_t = a \right] \end{aligned}$$

Contents

1 Finite Markov Decision Processes (MDP)

- Agent–Environment Interface
- Objective
- Policy
- Value Functions
 - State-value function $v_{\pi}(s)$ for policy $\pi(a|s)$
 - Action-value function $q_{\pi}(s, a)$ for policy $\pi(a|s)$
- Optimal Value Functions
 - Optimal state-value function $v_{*}(s)$
 - Existence of optimal policy $\pi_{*}(a|s)$
 - Optimal action-value function $q_{*}(s, a)$
- Bellman Equation
 - Bellman equation for $v_{\pi}(s)$
 - Bellman optimality equation for $v_{*}(s)$
 - Bellman optimality equation for $q_{*}(s, a)$

2 Dynamic Programming (DP)

- Iterative Policy Evaluation

Optimal state-value function $v_*(s)$

Recall that the state-value function for policy $\pi(a|s)$ is defined as:

$$v_\pi(s) \doteq E_\pi [G_t | S_t = s]$$

Optimal state-value function:

$$\begin{aligned} v_*(s) &\doteq \max_{\pi} v_\pi(s) \\ &= \max_{\pi} E_\pi [G_t | S_t = s] \end{aligned}$$

Existence of $\pi_*(a|s)$

Theorem

For any MDP

- 1 there is a policy, denoted by $\pi_*(a|s)$, that is at least as good as any other policies, i.e. it maximizes $v_{\pi_*}(s)$ for all $s \in \mathcal{S}$ simultaneously
- 2 optimal policy does not have to be unique, i.e. non-equal policies $\pi_{*,1}(a|s)$ and $\pi_{*,2}(a|s)$ may result in the same state-value:
 $v_{\pi_{*,1}}(s) = v_{\pi_{*,2}}(s)$ for all $s \in \mathcal{S}$
- 3 there exists a deterministic optimal policy

Optimal action-value function $q_*(s, a)$

Optimal action-value function:

$$\begin{aligned}
 q_*(s, a) &\doteq \max_{\pi} q_{\pi}(s, a) \\
 &= \max_{\pi} E_{\pi} [G_t | S_t = s, A_t = a] \\
 &= \max_{\pi} E_{\pi} [R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a] \\
 &= \max_{\pi} \left[E_{\pi} [R_{t+1} | S_t = s, A_t = a] + \gamma E_{\pi} [G_{t+1} | S_t = s, A_t = a] \right] \\
 &= \max_{\pi} \left[E [R_{t+1} | S_t = s, A_t = a] \right. \\
 &\quad \left. + \gamma \sum_{s'} p(s' | s, a) E_{\pi} [G_{t+1} | S_{t+1} = s', S_t = s, A_t = a] \right] \\
 &= E [R_{t+1} + \gamma v_*(S_{t+1}) | S_t = s, A_t = a]
 \end{aligned}$$

Optimal action-value function $q_*(s, a)$

Corollary

For any MDP

- any optimal policy, i.e. the policy that maximizes $v_\pi(s)$, achieves the optimal action-value $q_*(s, a)$

Contents

1 Finite Markov Decision Processes (MDP)

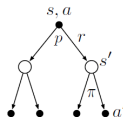
- Agent–Environment Interface
- Objective
- Policy
- Value Functions
 - State-value function $v_\pi(s)$ for policy $\pi(a|s)$
 - Action-value function $q_\pi(s, a)$ for policy $\pi(a|s)$
- Optimal Value Functions
 - Optimal state-value function $v_*(s)$
 - Existence of optimal policy $\pi_*(a|s)$
 - Optimal action-value function $q_*(s, a)$
- Bellman Equation
 - Bellman equation for $v_\pi(s)$
 - Bellman optimality equation for $v_*(s)$
 - Bellman optimality equation for $q_*(s, a)$

2 Dynamic Programming (DP)

- Iterative Policy Evaluation

Bellman equation for $v_\pi(s)$

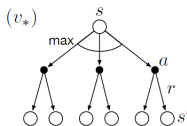
$$\begin{aligned}
 v_\pi(s) &\doteq E_\pi [G_t | S_t = s] \\
 &= E_\pi [R_{t+1} + \gamma G_{t+1} | S_t = s] \\
 &= \sum_a \pi(a|s) E_\pi [R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a] \\
 &= \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma E_\pi [G_{t+1} | S_{t+1} = s']] \\
 &= \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma v_\pi(s')]
 \end{aligned}$$



Source: *Reinforcement Learning: An Introduction* by R. Sutton and A. Barto

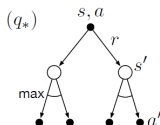
Bellman optimality equation for $v_*(s)$

$$\begin{aligned}
 v_*(s) &= \max_{a \in \mathcal{A}(s)} q_{\pi_*}(s, a) \\
 &= \max_a E_{\pi_*} [G_t | S_t = s, A_t = a] \\
 &= \max_a E_{\pi_*} [R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a] \\
 &= \max_a E_{\pi_*} [R_{t+1} + \gamma v_*(S_{t+1}) | S_t = s, A_t = a] \\
 &= \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma v_*(s')]
 \end{aligned}$$



Bellman optimality equation for $q_*(s, a)$

$$\begin{aligned} q_*(s, a) &= E \left[R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') \middle| S_t = s, A_t = a \right] \\ &= \sum_{s', r} p(s', r | s, a) \left[r + \gamma \max_{a'} q_*(s', a') \right] \end{aligned}$$



Contents

1 Finite Markov Decision Processes (MDP)

- Agent–Environment Interface
- Objective
- Policy
- Value Functions
 - State-value function $v_{\pi}(s)$ for policy $\pi(a|s)$
 - Action-value function $q_{\pi}(s, a)$ for policy $\pi(a|s)$
- Optimal Value Functions
 - Optimal state-value function $v_{*}(s)$
 - Existence of optimal policy $\pi_{*}(a|s)$
 - Optimal action-value function $q_{*}(s, a)$
- Bellman Equation
 - Bellman equation for $v_{\pi}(s)$
 - Bellman optimality equation for $v_{*}(s)$
 - Bellman optimality equation for $q_{*}(s, a)$

2 Dynamic Programming (DP)

- Iterative Policy Evaluation

Iterative Policy Evaluation

Recall the Bellman equation for $v_\pi(s)$:

$$v_\pi(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a) [r + \gamma v_\pi(s')]$$

Iterative Policy Evaluation

Recall the Bellman equation for $v_\pi(s)$:

$$v_\pi(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma v_\pi(s')]$$

Let $v_k(s)$, $k = 0, 1, 2, \dots$ denote an estimate of $v_\pi(s)$. The fixed-point iteration can be written as follows:

$$v_{k+1}(s) \doteq \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma v_k(s')]$$

Iterative Policy Evaluation

Let $v_k(s)$, $k = 0, 1, 2, \dots$ denote an estimate of $v_\pi(s)$. The fixed-point iteration can be written as follows:

$$v_{k+1}(s) \doteq \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a) [r + \gamma v_k(s')]$$

Iterative Policy Evaluation, for estimating $V \approx v_\pi$

Input π , the policy to be evaluated

Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation

Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(\text{terminal}) = 0$

Loop:

$\Delta \leftarrow 0$

Loop for each $s \in \mathcal{S}$:

$v \leftarrow V(s)$

$V(s) \leftarrow \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a) [r + \gamma V(s')]$

$\Delta \leftarrow \max(\Delta, |v - V(s)|)$

until $\Delta < \theta$