

Advanced Machine Learning: Assignment 1

Seung-Hoon Na

October 23, 2017

1 Policy improvement theorem

Let $v_\pi(s)$ be the *state-value function for policy π* , which is the expected return when starting in s and following π thereafter. $v_\pi(s)$ is formally defined as follows:

$$v_\pi(s) = \mathbf{E}_\pi [G_t | S_t = s] = \mathbf{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s \right]$$

where $\mathbf{E}_\pi[\cdot]$ denotes the expected value of a random variable given that the agent follows policy π , and t is any time step.

Let $q_\pi(s, a)$ be the *action-value function for policy π* , which is the expected return starting from s , taking the action a , and thereafter following policy π . $q_\pi(s, a)$ is formally defined as follows:

$$q_\pi(s, a) = \mathbf{E}_\pi [G_t | S_t = s, A_t = a] = \mathbf{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a \right]$$

The *policy improvement theorem* is stated as:

Theorem 1 (Policy improvement theorem) *Let π and π' be any pair of deterministic policies such that for all $s \in S$, $q_\pi(s, \pi'(s)) \geq v_\pi(s)$. Then, the policy π' is better than π for all states $s \in S$, i.e., $v_{\pi'}(s) \geq v_\pi(s)$.*

Prove the policy improvement theorem.

2 Bellman equations

2.1 Bellman equation for state-value function

The Bellman equation for $v_\pi(s)$ is formulated as follows:

$$v_\pi(s) = \mathbf{E}_\pi [R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s]$$

Derive the Bellman equation for state-value function.

2.2 Bellman equation for action-value function

$$q_\pi(s, a) = \mathbf{E}_\pi [R_{t+1} + \gamma q_\pi(S_{t+1}, A_{t+1}) | S_t = s, A_t = a]$$

Derive the Bellman equation for action-value function.

3 Temporal Difference

3.1 Forward view TD(λ)

Describe the update equation of forward view TD(λ) with details.

3.2 Backward view TD(λ)

Describe the update equation of forward view TD(λ) with details

3.3 Equivalence between forward view and Backward view TD(λ)

Provide a full proof for the equivalence between the offline updates of forward view and backward view TD(λ), i.e., the sum of offline updates is identical for forward-view and backward-view TD(λ).

4 Sarsa and Q-learning

4.1 Cliff Walking task (implementation)

1. Read the cliff walking task of example 6.6 in the textbook and implement Sarsa and Q-learning methods with ϵ -greedy action selection for the cliff walking task (using python). (ϵ is fixed to 0.1).
2. Compare Sarsa and Q-learning methods on the cliff walking task by increasing the number of episodes and draw a comparison result between them like Figure 6.5 (using python).
3. Write additional codes for printing learned policies on console display (no GUI) (using python). It may also be necessary to store a learned policy (or action-value function) in a separate file.
4. Submit all the python codes with a short report.