



Department of Computer Engineering

# Artificial Intelligence

## Assignment 7 Part 2

Dr. Rohban

Parsa Mohammadian — 98102284

January 7, 2022

## Contents

<b>1</b>		<b>1</b>
1.1	.....	1
1.2	.....	1
1.3	.....	1

# 1

## 1.1

V-value is state-value function. It means that  $V^\pi(s)$  is the value of starting from state  $s$  and obeying policy  $\pi$ . At the other hand, Q-value is action-value function. So  $Q^\pi(s, a)$  is the value of taking action  $a$  from state  $s$  and obeying policy  $\pi$ . Since  $R$  and  $P$  are not available, we can only use Q-value for this kind of problems (model-free).

## 1.2

In  $\epsilon$ -greedy method, we choose currently best action with probability  $1 - \epsilon$  and randomly choose other actions with probability  $\epsilon$ . For the normal policy iteration ( $\epsilon = 1$  or always choose the current best action), we only exploit our current knowledge and never explore. In this fashion we might stuck in local optimum. But when  $\epsilon$  is not 1, we explore with probability  $1 - \epsilon$  and seek for better local optimum even if we have already found one, and finally we might get to the global optimum.

## 1.3

$$\begin{aligned}
 \forall s : \mathbb{E}_{a \sim \pi'}[Q^\pi(s, a)] &\geq \mathbb{E}_{a \sim \pi}[Q^\pi(s, a)] \\
 \mathbb{E}[X] &= \sum_{x \in X} p(x)x \\
 \Rightarrow \sum_{a \in A} P(a)Q^{\pi'}(s, a) &\geq \sum_{a \in A} P(a)Q^\pi(s, a) \\
 P(a) &\propto \pi(a|s) \\
 \Rightarrow \sum_{a \in A} \pi(a|s)Q^{\pi'}(s, a) &\geq \sum_{a \in A} \pi(a|s)Q^\pi(s, a) \\
 V^\pi(s) &= \sum_{a \in A(s)} \pi(a|s)Q^\pi(s, a) \\
 \Rightarrow V^{\pi'}(s) &\geq V^\pi(s)
 \end{aligned}$$