

# Nash Q-Learning

2009.02.18 lab meeting



**GIST DCAS Lab, Han-Eol Kim**



GIST, Department of Mechatronics



## Index

### 1. Background of Reinforcement Learning

- Markov Decision Process
- Q-learning
- Stochastic Games

### 2. Nash Equilibrium

- Nash Q-function
- Nash Q-function Algorithm

### 3. Prove the Convergence of Nash Q-Learning

Keyword : Nash Q-Learning



GIST, Department of Mechatronics



## Markov Decision Process

Markov Decision Process game is a tuple  $\langle S, A, r, p \rangle$

In a Markov decision process, the objective of the agent is to find a strategy (policy  $\pi$ ) so as to maximize the expected sum of discounted rewards,

If it is optimal, solution of MDP, bellman equation can be described as

$$v(s, \pi^*) = \max_a \{r(s, a) + \beta \sum_{s'} p(s'|s, a)v(s', \pi^*)\}$$

$S_0$  is the initial state

$r_t$  is the reward at time  $t$

$0 \leq \beta < 1$  is the discount factor

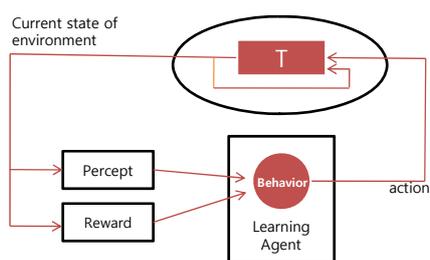
$\pi^*$  is the optimal strategy determined action by policy  $\pi$



GIST, Department of Mechatronics



## Q-Learning



At each time  $t$ , the agent chooses an action and observes its reward  
The agent then updates its Q-values based on the following equation:  
( $0 \leq \alpha_t < 1$  is the learning rate)

$$Q_t(s_t, a_t) = (1 - \alpha_t)Q_t(s_t, a_t) + \alpha_t[r_t + \beta \max_b Q_t(s_{t+1}, a)]$$



GIST, Department of Mechatronics



## Stochastic Games

In zero-sum game :  $r^1(s, a^1, a^2) + r^2(s, a^1, a^2) = 0$

In general-sum game :  $r^1(s, a^1, a^2) + r^2(s, a^1, a^2)$  is not 0 or any constant



GIST, Department of Mechatronics



Distributed Control and  
Autonomous System Laboratory

## Nash Equilibrium

### Definition of Nash Equilibrium

A Nash equilibrium is a joint strategy where each agent's is a best response to the others'. For a stochastic game, each agent's strategy is defined over the entire time horizon of the game.

In stochastic game, A Nash equilibrium point is a pair of strategies  $(\pi_s^1, \pi_s^2)$  such that for all  $s \in S$

$$\begin{aligned} r^1(\pi_s^1, \pi_s^2) &\geq r^1(\pi^1, \pi_s^2) \quad \text{for all } \pi^1 \in A^1 \\ r^2(\pi_s^1, \pi_s^2) &\geq r^2(\pi_s^1, \pi^2) \quad \text{for all } \pi^2 \in A^2 \end{aligned}$$

|        |         | agent1  |         |
|--------|---------|---------|---------|
|        |         | $a_1^1$ | $a_1^2$ |
| agent2 | $a_2^1$ | (3, 3)  | (0, 2)  |
|        | $a_2^2$ | (2, 0)  | (1, 1)  |

Example) Nash equilibrium point is  $(a_1^1, a_1^2)$



GIST, Department of Mechatronics



Distributed Control and  
Autonomous System Laboratory

## Nash Q-function

### Q-Learning Equation

$$Q_{t+1}(s_t, a_t) = (1 - \alpha_t)Q_t(s_t, a_t) + \alpha_t[r_t + \beta \max_b Q_t(s_{t+1}, a)]$$

### Nash Q-Function

$$\text{Nash}Q_t^i(s') = \pi^1(s') \times \dots \times \pi^n(s') \times Q_t^i(s')$$

### Nash Q-Learning

$$Q_{t+1}^i(s, a^1, \dots, a^n) = (1 - \alpha_t)Q_t^i(s, a^1, \dots, a^n) + \alpha_t[r_t^i + \beta \text{Nash}Q_t^i(s')]$$



GIST, Department of Mechatronics



Distributed Control and  
Autonomous System Laboratory

## Nash Q-Learning Algorithm

### Initialize:

Let  $t = 0$ , get the initial state  $s_0$ .

Let the learning agent be indexed by  $i$ .

For all  $s \in S$  and  $a^j \in A^j$ ,  $j = 1, \dots, n$ , let  $Q_t^j(s, a^1, \dots, a^n) = 0$ .

### Loop

Choose action  $a_t^i$ .

Observe  $r_t^1, \dots, r_t^n; a_t^1, \dots, a_t^n$ , and  $s_{t+1} = s'$

Update  $Q_t^j$  for  $j = 1, \dots, n$

$$Q_{t+1}^j(s, a^1, \dots, a^n) = (1 - \alpha_t)Q_t^j(s, a^1, \dots, a^n) + \alpha_t[r_t^j + \beta \text{Nash}Q_t^j(s')]$$

where  $\alpha_t \in (0, 1)$  is the learning rate, and  $\text{Nash}Q_t^k(s')$

Let  $t := t + 1$ .

$$\text{Nash}Q_t^i(s') = \pi^1(s') \times \dots \times \pi^n(s') \times Q_t^i(s')$$



GIST, Department of Mechatronics



Distributed Control and  
Autonomous System Laboratory

## Assumption

### Assumption 1

Every state  $s \in \mathcal{S}$  and action  $a^k \in A^k$  for  $k = 1, \dots, n$  are visited infinitely often.

### Assumption 2

The learning rate  $\alpha_t$  satisfies the following conditions for all  $s, t, a^1, \dots, a^n$  :

1.  $0 \leq \alpha_t(s, t, a^1, \dots, a^n) < 1$ ,  $\sum_{t=0}^{\infty} \alpha_t(s, t, a^1, \dots, a^n) = \infty$ ,  $\sum_{t=0}^{\infty} [\alpha_t(s, t, a^1, \dots, a^n)]^2 < \infty$
2.  $\alpha_t(s, t, a^1, \dots, a^n) = 0$  if  $(s, t, a^1, \dots, a^n) \neq (s_t, a_t^1, \dots, a_t^n)$



GIST, Department of Mechatronics



Distributed Control and  
Autonomous System Laboratory

## Assumption

### Assumption 3

A Nash equilibrium  $(\pi^1(s), \pi^2(s))$  for any 2-player stochastic game  $(Q^1(s), Q^2(s))$  satisfies one of the following properties.

1. Nash equilibrium is global optimal

$$\begin{array}{l} \forall \pi^1 \in A^1 \\ \forall \pi^2 \in A^2 \end{array} \quad \pi^1_*(s) Q^k(s) \pi^2_*(s) \geq \pi^1(s) Q^k \pi^2(s)$$

2. If the Nash equilibrium is not a global optimal, then an agent receives a higher payoff when the other agent deviates from the Nash equilibrium strategy.

$$\begin{array}{l} \forall \pi^1 \in A^1 \\ \forall \pi^2 \in A^2 \end{array} \quad \begin{array}{l} \pi^1_*(s) Q^k(s) \pi^2(s) \geq \pi^1(s) Q^k \pi^2(s) \\ \pi^1(s) Q^k(s) \pi^2_*(s) \geq \pi^1(s) Q^k \pi^2(s) \end{array}$$



GIST, Department of Mechatronics



Distributed Control and  
Autonomous System Laboratory

## Convergence

### Nash Q-Function

$$\text{Nash}Q_t^i(s') = \pi^1(s') \times \dots \times \pi^n(s') \times Q_t^i(s')$$

### Nash Q-Learning

$$\begin{aligned} Q_{t+1}^i(s, a^1, \dots, a^n) &= (1 - \alpha_t)Q_t^i(s, a^1, \dots, a^n) + \alpha_t[\gamma_t^i + \beta \text{Nash}Q_t^i(s')] \\ &= (1 - \alpha_t)Q_t^i(s, a^1, \dots, a^n) + \alpha_t[\gamma_t^i + \beta \pi^1(s') \times \dots \times \pi^n(s') \times Q_t^i(s')] \end{aligned}$$

### Simplifying Nash-Q

$$\begin{aligned} Q_{t+1}^i(s) &= (1 - \alpha_t)Q_t^i(s) + \alpha_t[P_t^i Q_t^i(s')] \\ P_t^i Q_t^i(s') &= \gamma_t^i + \beta \pi^1(s') \times \dots \times \pi^n(s') \times Q_t^i(s') \end{aligned}$$

### Convergence Condition

$$\|P_t Q - P_t Q_*\| \leq \beta \|Q - Q_*\|$$



## Convergence

### Prove

$$\begin{aligned} \|P_t^i Q^i - P_t^i Q_*^i\| &= \gamma_t^i + \beta \pi^1 \times \dots \times \pi^n \times Q^i - (\gamma_t^i + \beta \pi_*^1 \times \dots \times \pi_*^n \times Q_*^i) \\ &= \beta (\pi^1 \times \dots \times \pi^n \times Q^i - \pi_*^1 \times \dots \times \pi_*^n \times Q_*^i) \end{aligned}$$

New notation  $\sigma^i = \pi^i$ ,  $\sigma^{-i} \equiv \sigma^1 \times \dots \times \sigma^{i-1} \times \sigma^{i+1} \times \dots \times \sigma^n$

$$\begin{aligned} \sigma^i \sigma^{-i} Q^i - \sigma_*^i \sigma_*^{-i} Q_*^i &\leq \sigma^i \sigma^{-i} Q^i - \sigma^i \sigma^{-i} Q_*^i \\ &= \sum_{s'=1}^m \sigma^1 \times \dots \times \sigma^n (Q^i - Q_*^i) \\ \sum_{s'=1}^m p(s'|s, a^1, a^2) &= 1 \leftarrow \leq \sum_{s'=1}^m \sigma^1 \times \dots \times \sigma^n \|Q^i - Q_*^i\| = \|Q^i - Q_*^i\| \end{aligned}$$

Finally,  $\|P_t^i Q^i - P_t^i Q_*^i\| \leq \beta \|Q^i - Q_*^i\|$



# THE END

Thanks for Listening



GIST, Department of Mechatronics



Distributed Control and  
Autonomous System Laboratory