

Nash Q-Learning for General-Sum Stochastic Games

Hu, J. and Wellman, M.P.

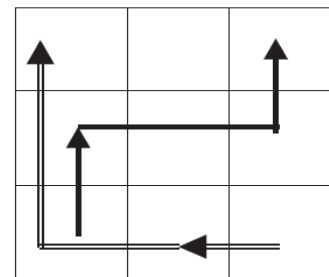
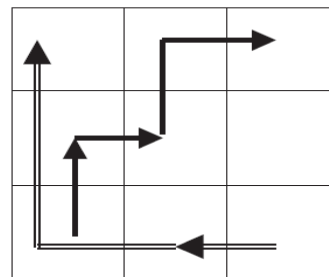
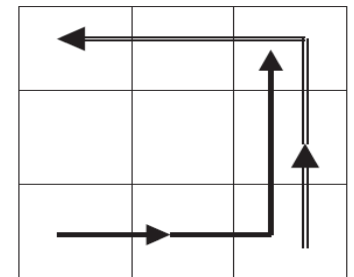
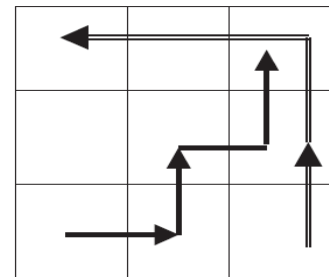
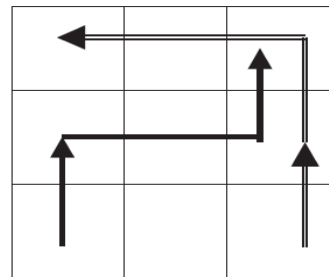
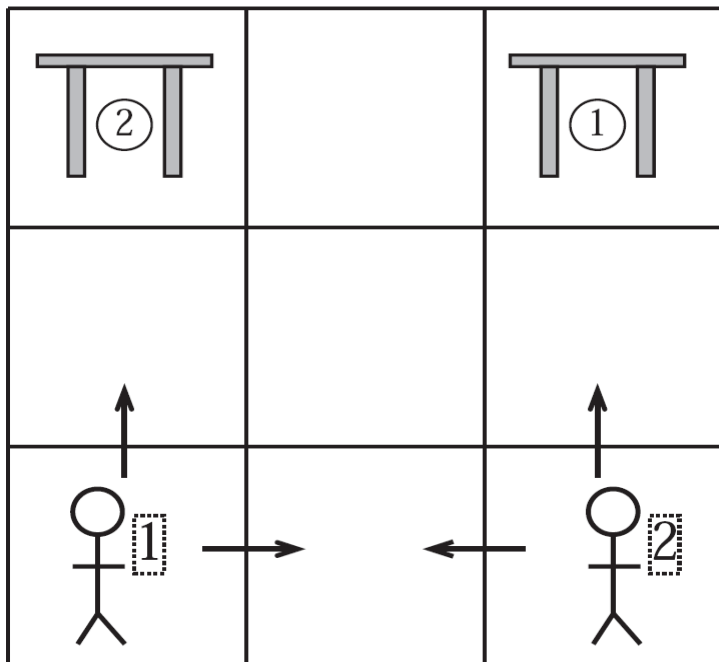
Lauri Lyly, llyly@cc.hut.fi

Stochastic general-sum games

- Stochasticity: Environment is in part formed by other agents
 - nondeterministic, noncooperative nature, no agreements
- Arbitrary relation between agents' rewards
 - Extends last time's topic zero-sum

Nash Equilibrium

- Best-response joint strategy
- Study limited to stationary strategies (policies)
- Rewards of others are perceived, strategies are not



Nash Q-Values

$$Q_*^i(s, a^1, \dots, a^n) = r^i(s, a^1, \dots, a^n) + \beta \sum_{s' \in S} p(s' | s, a^1, \dots, a^n) v^i(s', \pi_*^1, \dots, \pi_*^n)$$

	Multiagent	Single-Agent
Q-function	$Q(s, a^1, \dots, a^n)$	$Q(s, a)$
“Optimal” Q-value	Current reward + Future rewards when all agents play specified Nash equilibrium strategies from the next period onward	Current reward + Future rewards by playing the optimal strategy from the next period onward

Definitions of Q-values

Stage game

- One-period game as opposed to stochastic

Let σ^{-k} be the product of strategies of all agents other than k , $\sigma^{-k} \equiv \sigma^1 \dots \sigma^{k-1} \cdot \sigma^{k+1} \dots \sigma^n$.

- Mainly used in convergence proof
- Nash equilibrium for the stage game. M is a “payoff function”:

$$\sigma^k \sigma^{-k} M^k \geq \hat{\sigma}^k \sigma^{-k} M^k \quad \text{for all } \sigma^k \in \hat{\sigma}(A^k).$$

Update rule

- Same update rule for agent itself and its conjecture on other agent's Q-functions
- Q-functions can be initialized for example to 0
- Asynchronous updating: only entries pertaining to current state are updated

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

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

The 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 $\text{Nash} Q_t^k(s')$ is defined in (7)

Let $t := t + 1$.

Convergence proof requirements

- **Assumption 1:** Every state-action tuple is visited infinitely often
- **Assumption 2:** Learning rate $\alpha(t)$ satisfies:
 - Sum from goes towards infinity
 - Squared sum does not
 - $\alpha = 0$ if the element being updated doesn't correspond to current state-action tuple (asynchronous updating)

Proof basis and result

- Q-learning process updated by pseudo-contraction operator using the usual form:

$$Q = (1 - \alpha(t))Q(t) + \alpha(t)[P(t)Q(t)]$$

Contraction: Values approach optimal Q

- Link between stage games and stochastic games
- Goal: Show that NashQ is a pseudo-contraction operator
- Actually a real contraction operator in restricted conditions
 - Game with special types of Nash equilibrium points

Different Nash equilibria

- Global optimal point using stage game notation

$$\sigma M^k \geq \hat{\sigma} M^k \quad \text{for all } \hat{\sigma} \in \sigma(A).$$

- Saddle point

$$\sigma^k \sigma^{-k} M^k \geq \hat{\sigma}^k \sigma^{-k} M^k \quad \text{for all } \hat{\sigma}^k \in \sigma(A^k),$$

$$\sigma^k \sigma^{-k} M^k \leq \sigma^k \hat{\sigma}^{-k} M^k \quad \text{for all } \hat{\sigma}^{-k} \in \sigma(A^{-k})$$

- All equilibria chosen for update must be same type

Stage games compared

Global optimal point (*Up, Left*)

(Q_t^1, Q_t^2)	<i>Left</i>	<i>Right</i>
<i>Up</i>	10, 9	0, 3
<i>Down</i>	3, 0	-1, 2

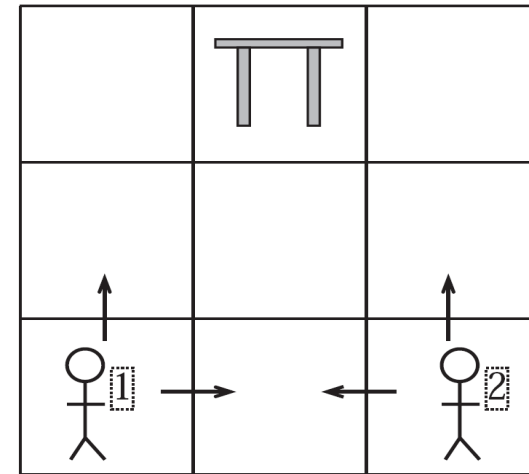
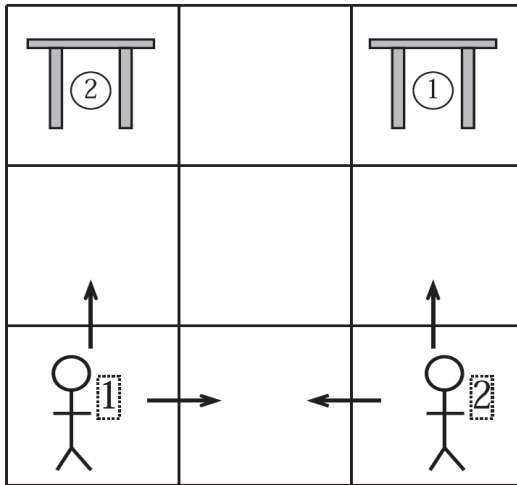
Saddle point (*Down, Right*)

(Q_*^1, Q_*^2)	<i>Left</i>	<i>Right</i>
<i>Up</i>	5, 5	0, 6
<i>Down</i>	6, 0	2, 2

Figure 2: Two stage games with different types of Nash equilibria

Experimentation framework

- Two grid-world games



- Motivation for grid games: state-specific actions, qualitative transitions, immediate and long-term rewards

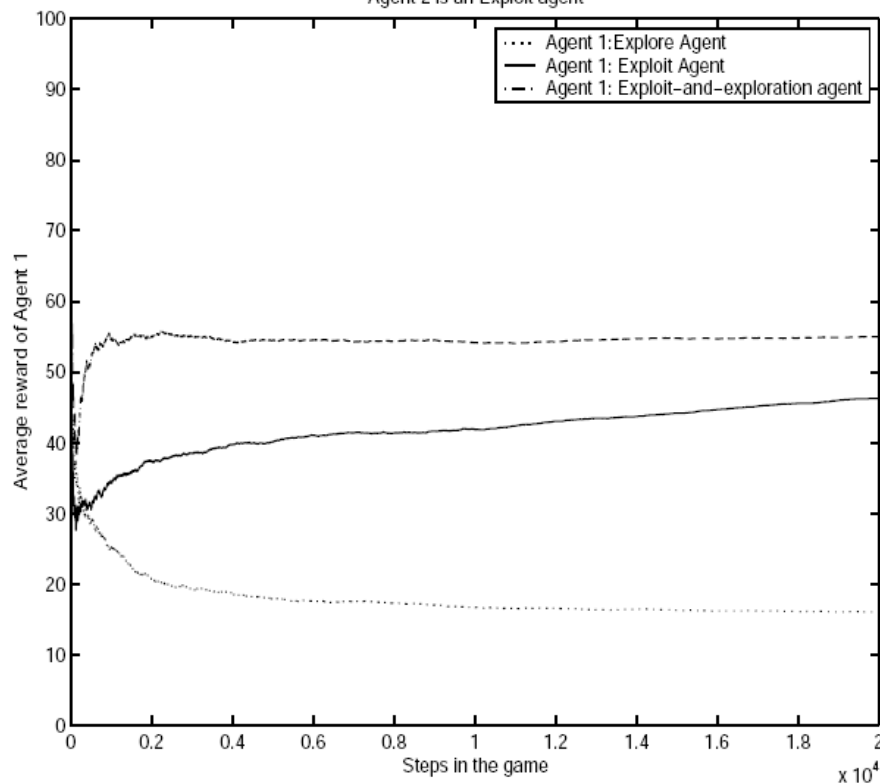
Learning process

- Violates assumption 3 of monotonic selection of global optima or saddle points.
- Still converges in most cases regardless of selection
- Offline and online learning rates separately

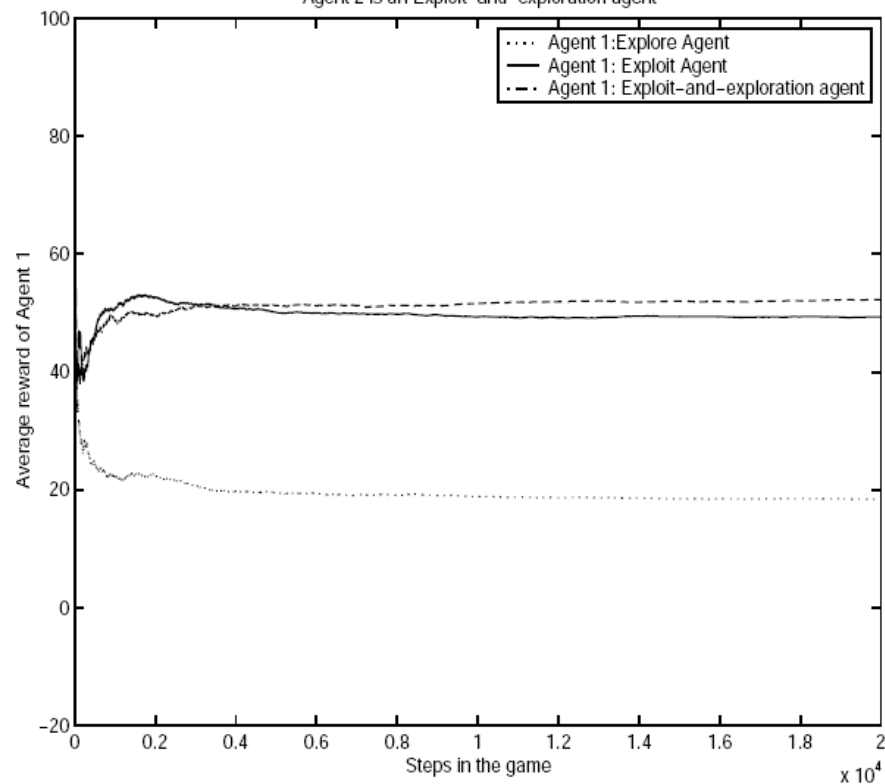
LEARNING STRATEGY		RESULTS OF LEARNING
AGENT 1	AGENT 2	PERCENT THAT REACH A NASH EQUILIBRIUM
SINGLE	SINGLE	20%
SINGLE	FIRST NASH	60%
	SECOND NASH	50%
	BEST EXPECTED NASH	76%
FIRST NASH	SECOND NASH	60%
	BEST EXPECTED NASH	76%
SECOND NASH	BEST EXPECTED NASH	84%
BEST EXPECTED NASH	BEST EXPECTED NASH	100%
FIRST NASH	FIRST NASH	100%
SECOND NASH	SECOND NASH	100%

Table 11: Learning performance in Grid Game 1

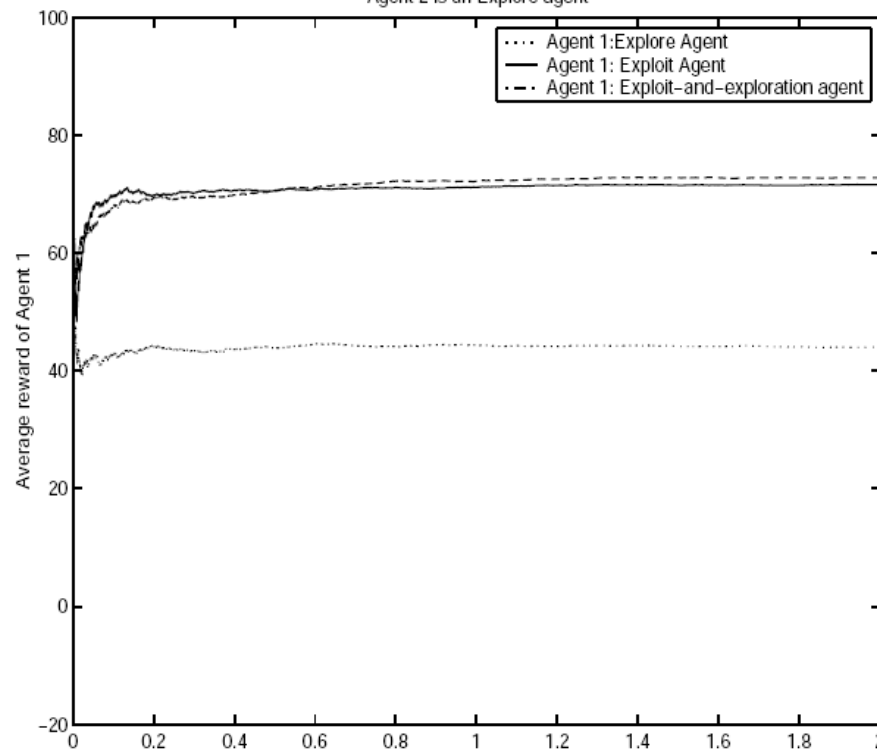
Agent 2 is an Exploit agent



Agent 2 is an Exploit-and-exploration agent



Agent 2 is an Explore agent



Conclusions

- No current method provides performance guarantees for general-sum stochastic games
- Works as a starting point
 - Other promising variants
 - Nash equilibrium itself can be refined

References

- [7] Hu, J. and Wellman, M.P. (2003). Nash Q-Learning for General-Sum Stochastic Games. *Journal of Machine Learning Research* 4:1039–1069.