Cyclic Equilibria in Markov Games

Amy Greenwald Brown University

with Michael Littman and Martin Zinkevich

Games Research and Development February 8, 2006

Agenda

Theorem

Value iteration converges to a stationary optimal policy in Markov decision processes.

Question

Does multiagent value iteration converge to a stationary equilibrium policy in Markov games?

Markov Decision Processes (MDPs)

Decision Process

- $\circ~S$ is a set of states
- \circ A is a set of actions
- $\circ \ R:S\times A\to \mathbb{R}$ is a reward function
- $P[s_{t+1} | s_t, a_t, \dots, s_0, a_0]$ is a probabilistic transition function that describes transitions between states, conditioned on past states and actions
- MDP = Decision Process + Markov Property:

$$P[s_{t+1} \mid s_t, a_t, \dots, s_0, a_0] = P[s_{t+1} \mid s_t, a_t]$$

 $\forall t, \forall s_0, \ldots, s_t \in S, \forall a_0, \ldots, a_t \in A$

Bellman's Equations

$$Q^{*}(s,a) = R(s,a) + \gamma \sum_{s'} P[s' \mid s,a] V^{*}(s')$$
(1)

$$V^*(s) = \max_{a \in A} Q^*(s, a)$$
 (2)

Value Iteration

 $\begin{array}{ll} \mathsf{VI}(\mathsf{MDP},\gamma) \\ \text{Inputs} & \text{discount factor } \gamma \\ \text{Output} & \text{optimal state-value function } V^* \\ & \text{optimal action-value function } Q^* \\ \text{Initialize} & V \text{ arbitrarily} \\ \end{array}$ $\begin{array}{l} \mathsf{REPEAT} \\ \text{for all } s \in S \\ & \text{for all } a \in A \\ & Q(s,a) = R(s,a) + \gamma \sum_{s'} P[s' \mid s,a] V(s') \\ & V(s) = \max_a Q(s,a) \\ \end{array}$ $\begin{array}{l} \mathsf{FOREVER} \end{array}$

Markov Games

Stochastic Game

- $\circ~N$ is a set of players
- $\circ~S$ is a set of states
- \circ A_i is the *i*th player's set of actions
- $R_i(s, \vec{a})$ is the *i*th player's reward at state s given action vector \vec{a}
- $P[s_{t+1} | s_t, \vec{a}_t, \dots, s_0, \vec{a}_0]$ is a probabilistic transition function that describes transitions between states, conditioned on past states and actions

Markov Game = Stochastic Game + Markov Property:

$$P[s_{t+1} \mid s_t, \vec{a}_t, \dots, s_0, \vec{a}_0] = P[s_{t+1} \mid s_t, \vec{a}_t]$$

$$\forall t, \forall s_0, \dots, s_t \in S, \forall \vec{a}_0, \dots, \vec{a}_t \in A$$

Bellman's Analogue

$$Q_i^*(s, \vec{a}) = R_i(s, \vec{a}) + \gamma \sum_{s'} P[s' \mid s, \vec{a}] V_i^*(s')$$
(3)

$$V_i^*(s) = \sum_{\vec{a} \in A} \pi^*(s, \vec{a}) Q_i^*(s, \vec{a})$$
(4)

Foe-VI $\pi^*(s) = (\sigma_1^*, \sigma_2^*)$, a minimax equilibrium policy
[Shapley 1953, Littman 1994]Friend-VI $\pi^*(s) = e_{\vec{a}^*}$ where $\vec{a}^* \in \arg \max_{\vec{a} \in A} Q_i^*(s, \vec{a})$
[Littman 2001]Nash-VI $\pi^*(s) \in \operatorname{Nash}(Q_1^*(s), \dots, Q_n^*(s))$
[Hu and Wellman 1998]CE-VI $\pi^*(s) \in \operatorname{CE}(Q_1^*(s), \dots, Q_n^*(s))$
[G and Hall 2003]

Multiagent Value Iteration

MULTI-VI(Inputs	(MGame, γ , f) discount factor γ selection mechanism f		
Output	equilibrium state-value function V^* equilibrium action-value function Q^*		
Initialize	Varbitrarily		
REPEAT			
for all $s \in S$			
for all $\vec{a} \in A$			
for all $i \in N$			
$Q_i(s, \vec{a}) = R_i(s, \vec{a}) + \gamma \sum_{s'} P[s' \mid s, \vec{a}] V_i(s')$			
$\pi(s) \in f(Q_1(s), \dots, Q_n(s))$			
for all $i \in N$			
$V_i(s) = \sum_{\vec{a} \in A} \pi(s, \vec{a}) Q_i(s, \vec{a})$			
FOREVER			

Friend-or-Foe-VI always converges [Littman 2001] Nash-VI and CE-VI converge to stationary equilibrium policies in zero-sum & common-interest Markov games [GZ and Hall 2005]

An Example



Observation

This game has no stationary deterministic equilibrium policy when $\gamma > \frac{1}{2}$.

Proof

 $(A \text{ quits}, B \text{ quits}) \Rightarrow A \text{ prefers send to quit } (2\gamma > 1)$ $(A \text{ sends}, B \text{ quits}) \Rightarrow B \text{ prefers send to quit } (0 > -1)$ $(A \text{ sends}, B \text{ sends}) \Rightarrow A \text{ prefers quit to send } (1 > 0)$ $(A \text{ quits}, B \text{ sends}) \Rightarrow B \text{ prefers quit to send } (-1 > -2)$



Observation

This game has a deterministic cyclic equilibrium policy when $\gamma = \frac{2}{3}$.

Example

PolicyV(A)V(B)1(A quits, B sends)(1, -2) $(\frac{8}{9}, -\frac{4}{9})$ 2(A sends, B sends) $(\frac{4}{3}, -\frac{2}{3})$ $(\frac{8}{9}, -\frac{4}{9})$ 3(A sends, B quits) $(\frac{4}{3}, -\frac{2}{3})$ (2, -1)4(A quits, B quits)(1, -2)(2, -1)

More Cyclic Policies



γ	(A sends, B sends)	Total
$\frac{2}{3}$	1	4
0.9	5	8
0.999999	693146	693149

Random Markov Games

$$\begin{split} |N| &= 2 \\ |A| \in \{2,3\} \\ |S| \in \{1,\ldots,10\} \\ \text{Random Rewards} \in [0,99] \\ \text{Random Deterministic Transitions} \\ \gamma &= \frac{3}{4} \end{split}$$



Multiagent Q-Learning

Minimax-Q Learning [Littman 1994]

 provably converges to stationary minimax equilibrium policies in zero-sum Markov games

Nash-Q Learning [Hu and Wellman 1998] Correlated-Q Learning [G and Hall 2003]

 converge empiricially to stationary equilibrium policies on a testbed of general-sum Markov games