

Outline

- A. Introduction
- B. Single Agent Learning
- C. Game Theory
- D. Multiagent Learning
 - Definition of the Problem: Stochastic Games
 - Equilibrium Learners
 - Best-Response Learners
 - Regret Minimizing Algorithms
 - Learning to Coordinate
- E. Future Issues and Open Problems

Equilibrium Learners

Extend Q table to values on joint actions, one for each learner.

Replace “max” with other operators.

Minimax-Q (Littman 94)

- Converges, zero-sum equilibrium (Littman & Szepesvári 96)

Nash-Q (Hu & Wellman 98)

- Applies to general-sum scenarios; works ok sometimes.

Friend-or-Foe-Q (Littman 01)

- Opponents set as friends (use max; Claus & Boutilier 98), foes (use minimax); converges, equilibria if saddlepoint/global optima.

CE-Q (Hall & Greenwald 02)

- Use correlated equilibria.

SG Analogies to MDPs

In the zero-sum case, results analogous to MDPs:

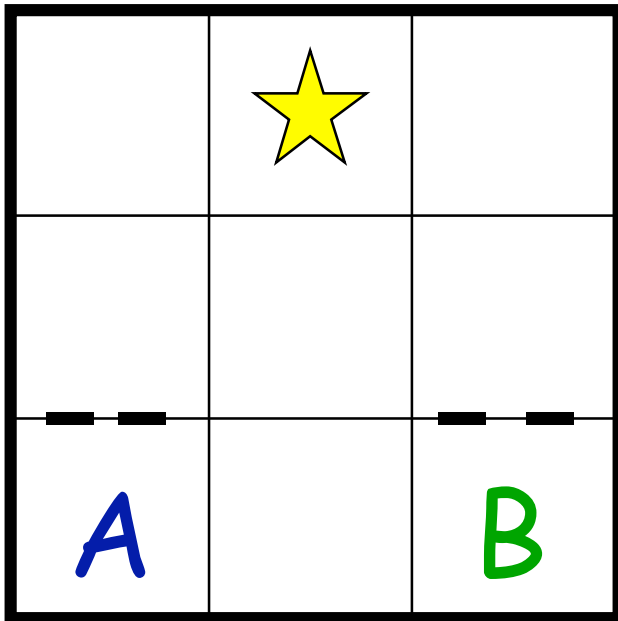
- optimal value function, policy, Q function
- can be found via simulation, search, DP (not LP!)
- can define Q-learning like algorithm

Failed analogies for general-sum games:

- optimal value function need not be unique
- Q-learning like algorithm doesn't converge
- no efficient algorithm known

Active area of research. What's the right thing to do?

Grid Game 3 (Hu & Wellman 01)



U, D, R, L, X

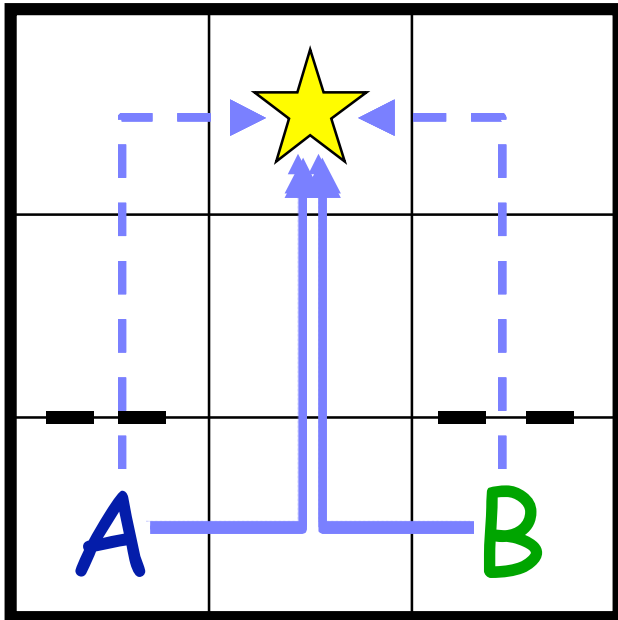
No move on collision

Semiwalls (50%)

-1 for step, -10 for collision, +100 for goal

Both can get goal.

Nash in Grid Game



Average total:

- (97, 48)
- (48, 97)
- (- , -) (not Nash)
- (64, 64) (not Nash)
- (75, 75)?

Collaborative Solution

A	★	
A	A	
---	---	---
A	A	B

Average total:

– (96, 96) (not Nash)

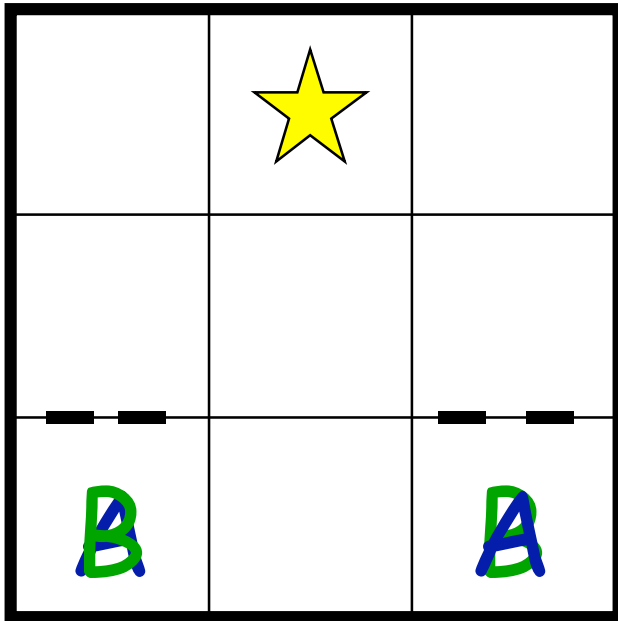
A won't wait.

B changes incentives.

Can we formalize collaboration like this?

Simpler setting: matrix games

Symmetric Markov Game



Episodic

Roles chosen randomly

Algorithm:

- Maximize sum (MDP)
- Security-level (0-sum)
- Choose max if better

Converges to Nash.

Regret Minimizing Algorithms

- Freund and Schapire
- Hart and Mas-Colell
- No internal regret
- No external regret
- Connection to minimax and correlated equilibria