

Correlated  $Q$ -Learning  
&  
No-Regret  $Q$ -Learning

Amy Greenwald

with

David Gondek

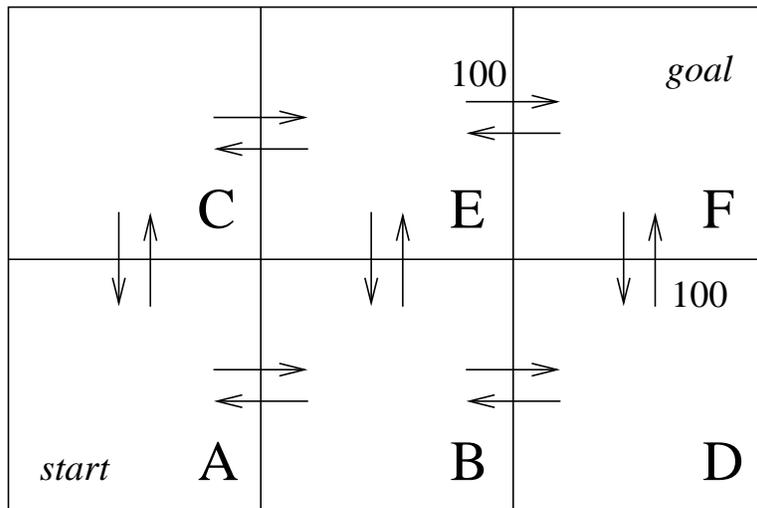
Keith Hall

Brown University

Economic Theory Workshop

March 4, 2002

# Reinforcement Learning



## Part I

### Multiagent $Q$ -Learning

- Correlated- $Q$  Learning
  - rational, empirically converges
- Nash- $Q$  [Hu and Wellman, 1998]
  - rational, does not converge
- Minimax- $Q$  [Littman, 1994]
  - converges, not rational

## Part II

### Approximate $Q$ -Learning

- No-regret  $Q$ -learning
  - No-external-regret
    - \* converge to minimax strategies in constant-sum games
  - No-internal-regret
    - \* converge to correlated equilibrium in general-sum games

# Markov Decision Processes (MDPs)

## Decision Process

- $S$  is a set of states ( $s \in S$ )
- $A$  is a set of actions ( $a \in A$ )
- $R : S \times A \rightarrow \mathbb{R}$  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:

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

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

## Bellman's Equations

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

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

## Value Iteration

VALUE\_ITERATION(MDP,  $\gamma$ )

Inputs discount factor  $\gamma$

Output optimal state-value function  $V^*$   
optimal action-value function  $Q^*$

Initialize  $V = Q = 0$

REPEAT

  for all  $s \in S$

    for all  $a \in A$

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

$$V(s) = \max_a Q(s, a)$$

FOREVER

## Q-Learning

Q\_LEARNING(MDP,  $\gamma, \alpha, \epsilon$ )

Inputs      discount factor  $\gamma$   
              rate of averaging  $\alpha$   
              rate of exploration  $\epsilon$

Output      optimal state-value function  $V^*$   
              optimal action-value function  $Q^*$

Initialize    $V = Q = 0$

REPEAT

  initialize  $s, a$

  WHILE  $s$  is nonterminal DO

    simulate action  $a$  in state  $s$

    observe reward  $R$  and next state  $s'$

    compute  $V(s') = \max_{a \in A(s)} Q(s, a)$

    update  $Q(s, a) = (1 - \alpha)Q(s, a) + \alpha[r + \gamma V(s')]$

    choose action  $a'$

$s = s', a = a'$

    decay  $\alpha$

FOREVER

**Theorem** [Watkins, 1989]

Q-learning converges to  $V^*$  and  $Q^*$

# Markov Games

## Stochastic Game

- $I$  is a set of  $n$  players ( $i \in I$ )
- $S$  is a set of states ( $s \in S$ )
- $A_i(s)$  is the  $i$ th player's set of actions at state  $s$   
let  $A(s) = A_1(s) \times \dots \times A_n(s)$  ( $\vec{a} \in A(s)$ )
- $P[s_{t+1}|s_t, \vec{a}_t, \dots, s_0, \vec{a}_0]$  is a probabilistic transition function that describes transitions between states, conditioned on past states and actions
- $R_i(s, \vec{a})$  is the  $i$ th player's reward at state  $s$  for action vector  $\vec{a}$

Markov Game = Stochastic Game + Markov Property:

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

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

## Bellman's Analogue

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

Foe-Q

$$V_1^*(s) = \max_{\sigma_1 \in \Sigma_1(s)} \min_{a_2 \in A_2(s)} Q_1^*(s, \sigma_1, a_2) = -V_2^*(s)$$

Friend-Q

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

Nash-Q

$$V_i^*(s) \in \text{Nash}_i(Q_1^*(s), \dots, Q_n^*(s))$$

CE-Q

$$V_i^*(s) \in \text{CE}_i(Q_1^*(s), \dots, Q_n^*(s))$$

## Multiagent $Q$ -Learning

MULTIQ(MGame,  $\gamma$ ,  $\alpha$ ,  $\epsilon$ )

REPEAT

  initialize  $s, a_1, \dots, a_n$

  WHILE  $s$  is nonterminal DO

    simulate actions  $a_1, \dots, a_n$  in state  $s$

    observe rewards  $R_1, \dots, R_n$  and next state  $s'$

    for all  $i \in I$

      compute  $V_i(s')$

      update  $Q_i(s, a_1, \dots, a_n)$

    (simultaneously) choose actions  $a'_1, \dots, a'_n$

$s = s', a_1 = a'_1, \dots, a_n = a'_n$

    decay  $\alpha$

FOREVER

FF- $Q$  converges [Littman 2001]

Nash- $Q$  is rational [Hu and Wellman 1998]

CE- $Q$  is rational and convergent (empirically)

## Correlated Equilibrium

Chicken

	$L$	$R$
$T$	6,6	2,7
$B$	7,2	0,0

CE

	$L$	$R$
$T$	1/2	1/4
$B$	1/4	0

$$\max 12\pi_{TL} + 9\pi_{TR} + 9\pi_{BL} + 0\pi_{BR}$$

subject to **probability** constraints

$$\begin{aligned} \pi_{TL} + \pi_{TR} + \pi_{BL} + \pi_{BR} &= 1 \\ \pi_{TL}, \pi_{TR}, \pi_{BL}, \pi_{BR} &\geq 0 \end{aligned}$$

& individual **rationality** constraints

$$\begin{aligned} 6\pi_{L|T} + 2\pi_{R|T} &\geq 7\pi_{L|T} + 0\pi_{R|T} \\ 7\pi_{L|B} + 0\pi_{R|B} &\geq 6\pi_{L|B} + 2\pi_{R|B} \\ 6\pi_{T|L} + 2\pi_{B|L} &\geq 7\pi_{T|L} + 0\pi_{B|L} \\ 7\pi_{T|R} + 0\pi_{B|R} &\geq 6\pi_{T|R} + 2\pi_{B|R} \end{aligned}$$

$$CE_i(Q_1(s), \dots, Q_n(s)) = \left\{ \sum_{\vec{a} \in A} \sigma^*(\vec{a}) Q_i(s, \vec{a}) \mid \sigma^* \text{ satisfies Eq. 1, 2, 3, or 4} \right\}$$

- **Utilitarian** maximize the **sum** of rewards

$$\sigma^* \in \arg \max_{\sigma \in CE} \sum_{\vec{a} \in A} \sigma(\vec{a}) \left( \sum_{i \in I} Q_i(s, \vec{a}) \right) \quad (1)$$

- **Egalitarian** maximize the **minimum** reward

$$\sigma^* \in \arg \max_{\sigma \in CE} \sum_{\vec{a} \in A} \sigma(\vec{a}) \left\{ \min_{i \in I} Q_i(s, \vec{a}) \right\} \quad (2)$$

- **Republican** maximize the **maximum** reward

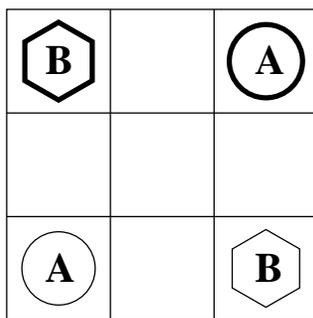
$$\sigma^* \in \arg \max_{\sigma \in CE} \sum_{\vec{a} \in A} \sigma(\vec{a}) \left\{ \max_{i \in I} Q_i(s, \vec{a}) \right\} \quad (3)$$

- **Libertarian**  $i$  maximizes only  $i$ 's rewards

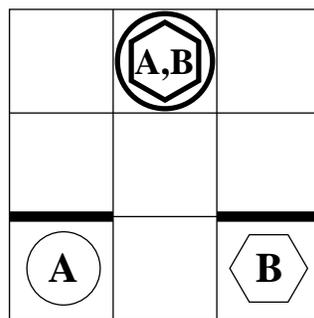
$$\sigma^i \in \arg \max_{\sigma \in CE} \sum_{\vec{a} \in A} \sigma(\vec{a}) Q_i(s, \vec{a}) \quad (4)$$

# Grid Games

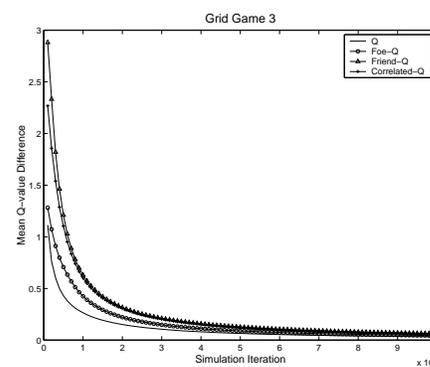
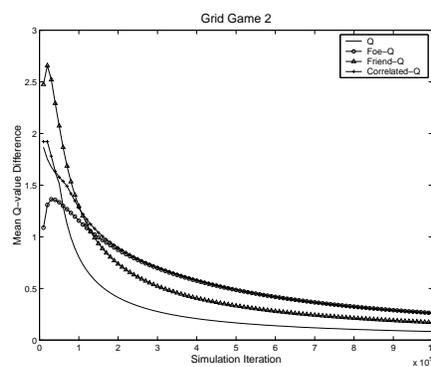
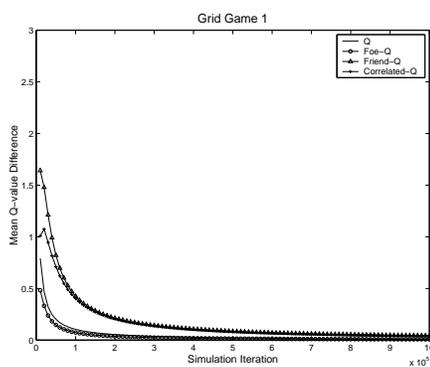
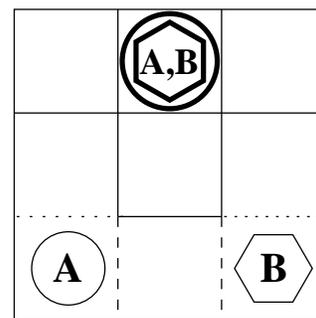
GG1



GG2



GG3



## Equilibrium Policies

Grid Games	GG1		GG2		GG3	
Algorithm	Score	Games	Score	Games	Score	Games
$Q$	100,100	2500	49,100	3333	100,125	3333
Foe- $Q$	0,0	0	67,68	3003	120,120	3333
Friend- $Q$	$-\infty, -\infty$	0	$-\infty, -\infty$	0	$-\infty, -\infty$	0
$u$ CE- $Q$	100,100	2500	50,100	3333	117,117	3333
$e$ CE- $Q$	100,100	2500	51,100	3333	118,118	3333
$r$ CE- $Q$	100,100	2500	100,49	3333	125,100	3333
$l$ CE- $Q$	100,100	2500	100,51	3333	$-\infty, -\infty$	0

## Part I

### Correlated- $Q$ Learning

- rational
- convergent (empirically)

### Conjecture

Correlated- $Q$  learning converges to  $Q$ -values that support equilibrium policies in Markov games.

## Part II

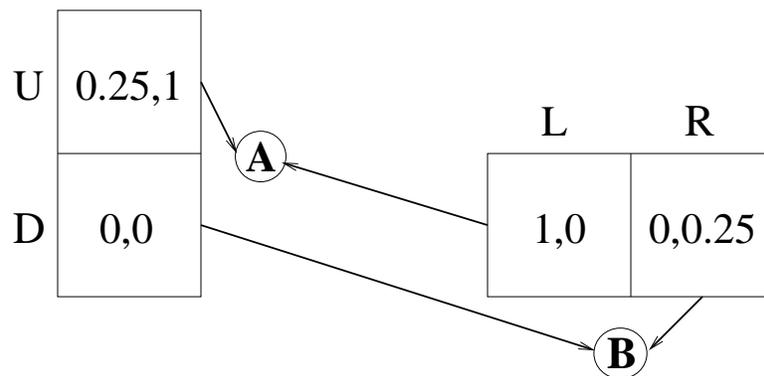
### No-regret $Q$ -Learning

- No-external-regret
  - converge to minimax strategies in constant-sum games
- No-internal-regret
  - converge to correlated equilibrium in general-sum games

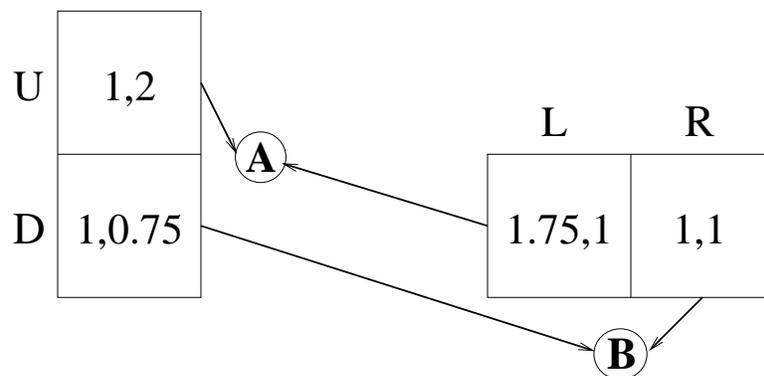
# Existence Game

Marty Zinkevich

## Payoffs



## Q-Values



Unique Mixed Strategy Equilibrium

$$\pi_r(U) = 7/15 \text{ and } \pi_c(L) = 4/9$$

## Repeated Games

A **game** is a tuple  $\Gamma = (I, (A_i, R_i)_{i \in I})$  where

- $I$  is a set of **players** ( $i \in I$ )
- $A_i$  is a set of **pure actions** ( $a_i \in A_i$ )
- $R_i : A \rightarrow \mathbb{R}$  is a **reward function** ( $a \in A = \prod_i A_i$ )

A **repeated game** is a sequence of tuples  $\Gamma^T$  or  $\Gamma^\infty$

## No-Regret Definitions

**Regret** is the difference in rewards for playing action  $a'_i$  rather than  $a_i$  at time  $t$ :

$$\rho_i^t(a_i, a'_i) = \pi_i^t(a'_i) [R_i(a_i, a_{-i}^t) - R_i(a'_i, a_{-i}^t)]$$

A learning algorithm exhibits **no-external-regret** iff it generates weights  $\{\pi_i^t\}$  s.t. for all opposing policies, there exists  $T$  s.t. for all  $T > T_0$ ,

$$\max_{a_i \in A_i} \frac{1}{T} \sum_{t=1}^T \sum_{a'_i \in A_i} \pi_i^t(a'_i) \rho_i^t(a_i, a'_i) \leq \text{ERR}(T)$$

where  $\text{ERR}(T) \rightarrow 0$  as  $T \rightarrow \infty$ .

A learning algorithm exhibits **no-internal-regret** iff it generates weights  $\{\pi_i^t\}$  s.t. for all opposing policies, there exists  $T$  s.t. for all  $T > T_0$ ,

$$\max_{a_i \in A_i} \frac{1}{T} \sum_{a'_i \in A_i} \left( \sum_{t=1}^T \pi_i^t(a'_i) \rho_i^t(a_i, a'_i) \right)^+ \leq \text{ERR}(T)$$

where  $\text{ERR}(T) \rightarrow 0$  as  $T \rightarrow \infty$  and  $X^+ = \max\{X, 0\}$ .

# No-Regret Algorithms

Average Regret Matrix

$$\mathcal{R}_i^t(a'_i, a_i) = \frac{1}{t} \sum_{x=1}^t \rho_i^x(a'_i, a_i)$$

Hart and Mas-Colell (HMC)

$$\pi_i^{t+1}(a_i) = \frac{[\sum_{a'_i \in A_i} \mathcal{R}_i^t(a'_i, a_i)]^+}{\sum_{a_i \in A_i} [\sum_{a'_i \in A_i} \mathcal{R}_i^t(a'_i, a_i)]^+}$$

NER learning approximates minimax equilibria

[Freund and Schapire, 1996]

PEACE Probably Empirically Approximating CE

$$\pi_i^{t+1}(a_i) = \frac{\sum_{a'_i \in A_i} [\mathcal{R}_i^t(a'_i, a_i)]^+}{\sum_{a_i \in A_i} \sum_{a'_i \in A_i} [\mathcal{R}_i^t(a'_i, a_i)]^+}$$

NIR learning approximates correlated equilibria

[Foster and Vohra, 1997]

# Normal Form Games

## Matching Pennies

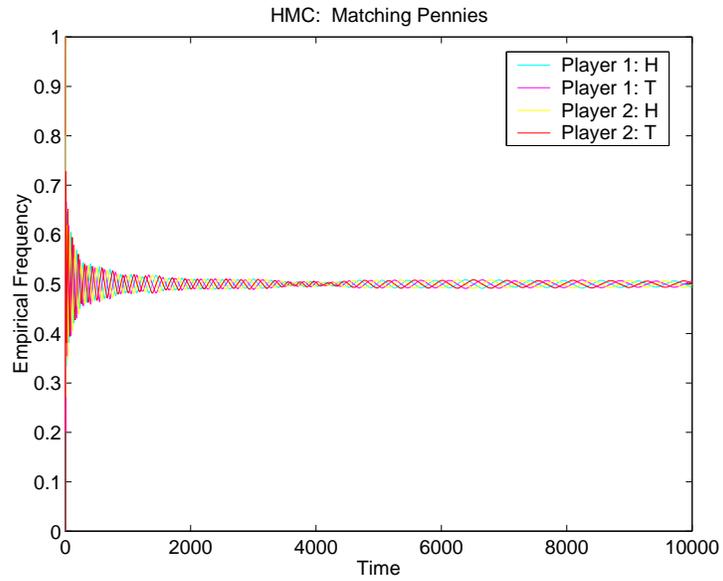
$1 \backslash 2$	$H$	$T$
$H$	1,0	0,1
$T$	0,1	1,0

## Shapley Game

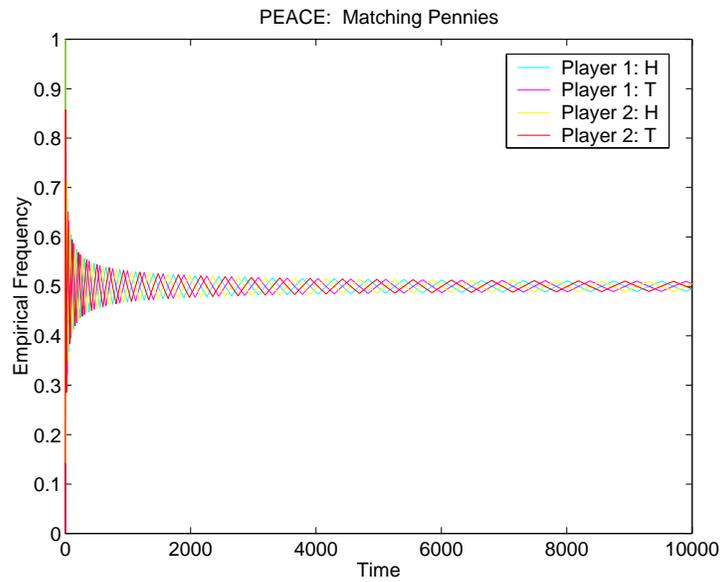
$1 \backslash 2$	$L$	$C$	$R$
$T$	1,0	0,1	0,0
$M$	0,0	1,0	0,1
$B$	0,1	0,0	1,0

# Matching Pennies

## Frequencies: HMC

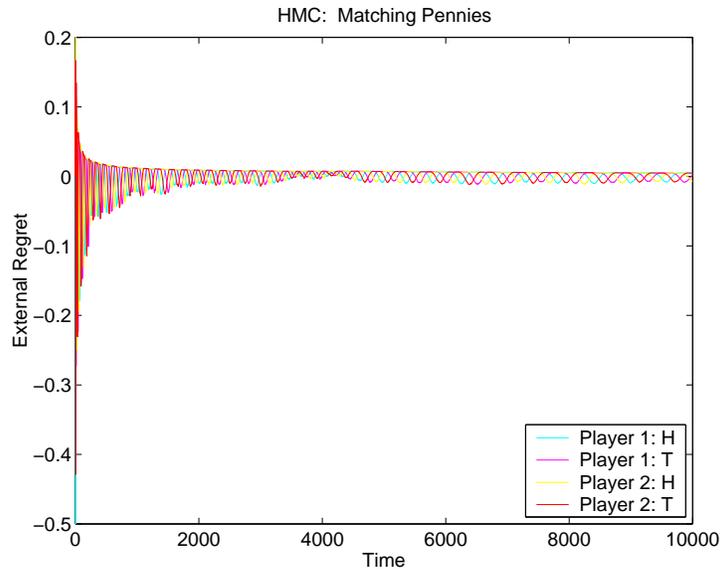


## Frequencies: PEACE

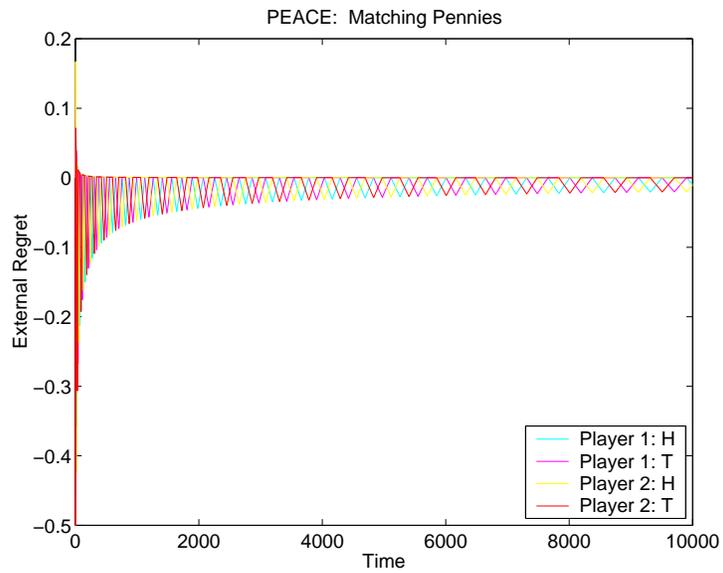


# Matching Pennies

## External Regret: HMC

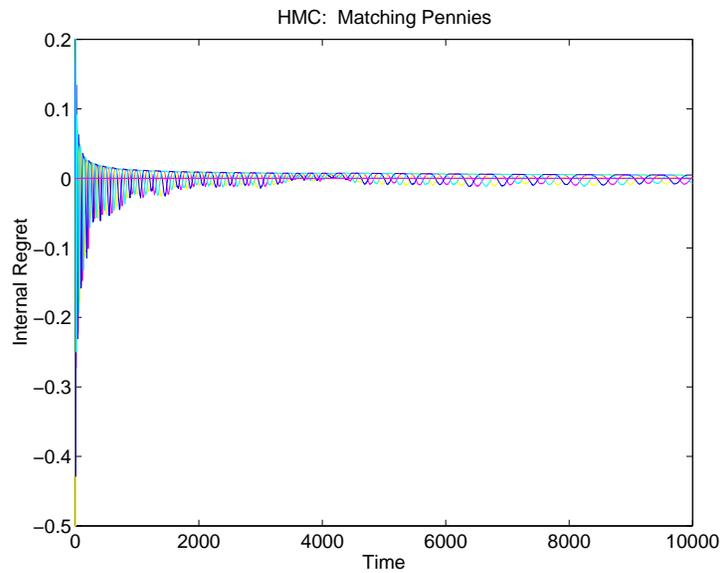


## External Regret: PEACE

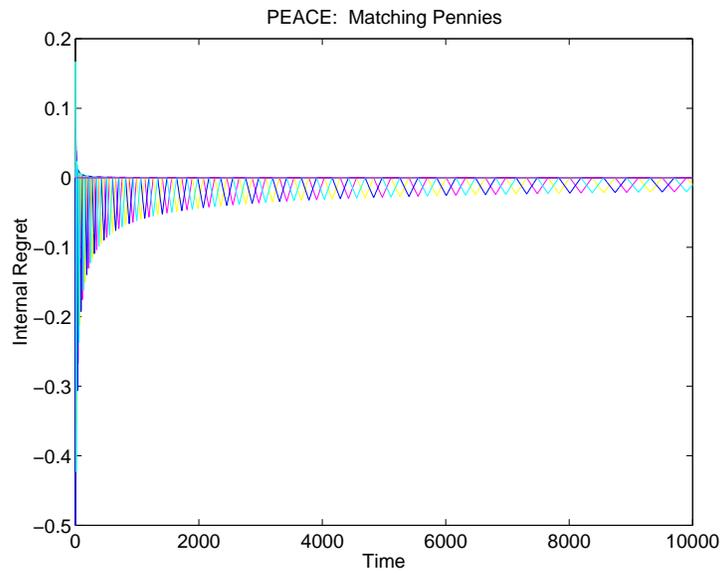


# Matching Pennies

## Internal Regret: HMC

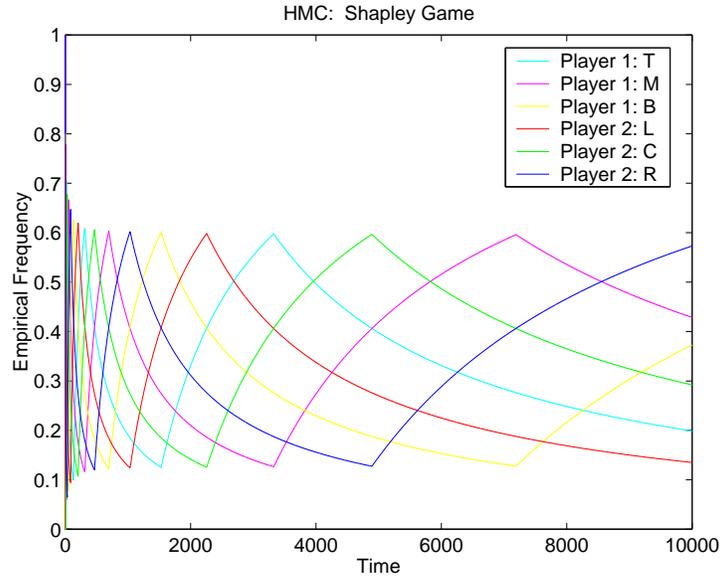


## Internal Regret: PEACE

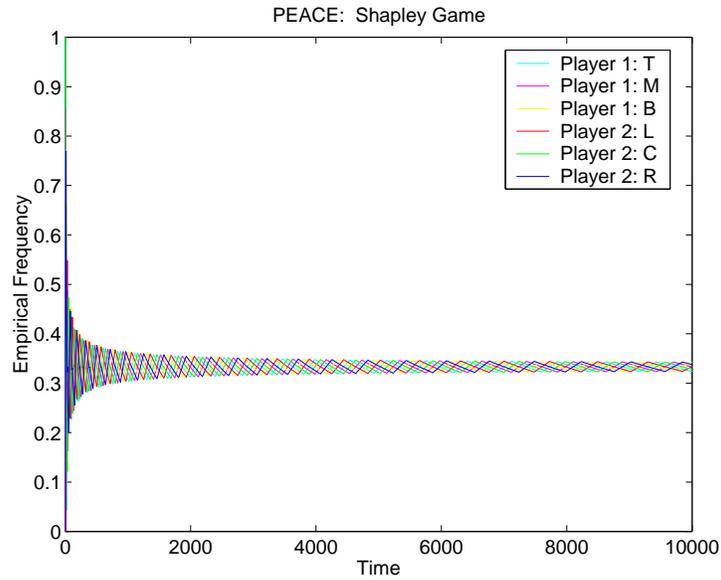


# Shapley Game

## Frequencies: HMC

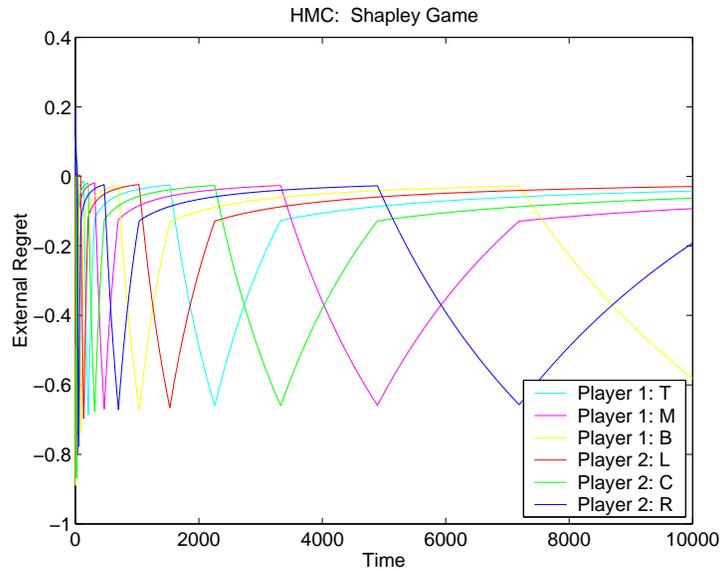


## Frequencies: PEACE

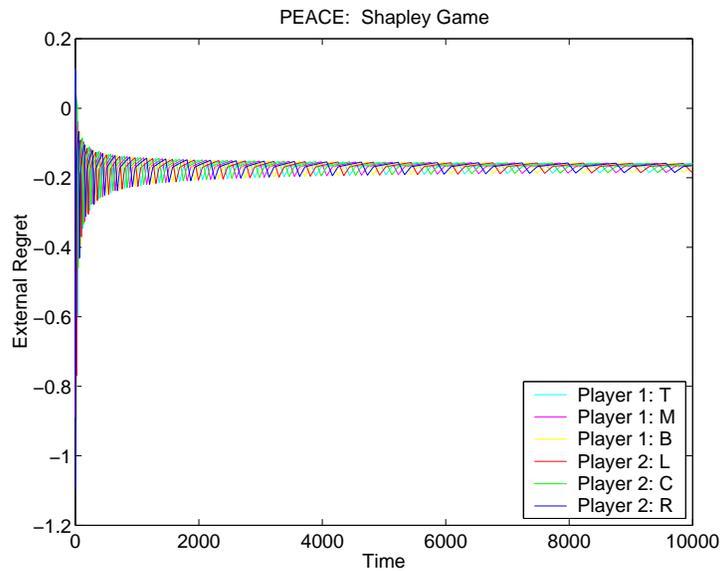


# Shapley Game

## External Regret: HMC

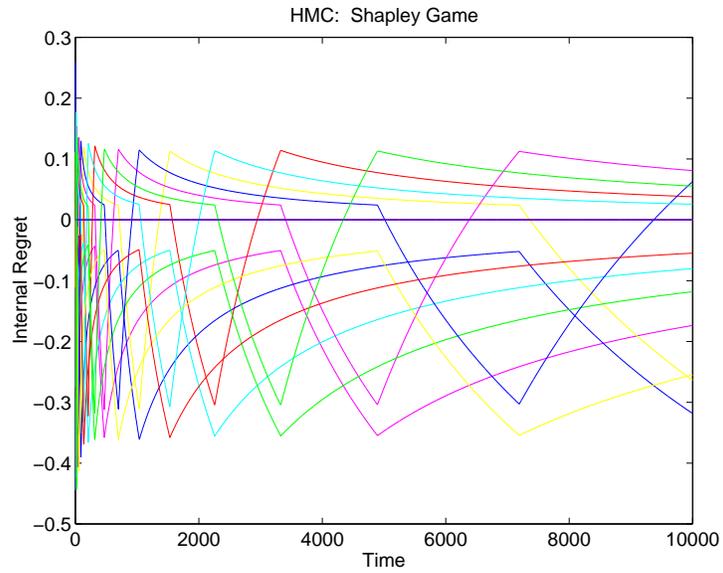


## External Regret: PEACE

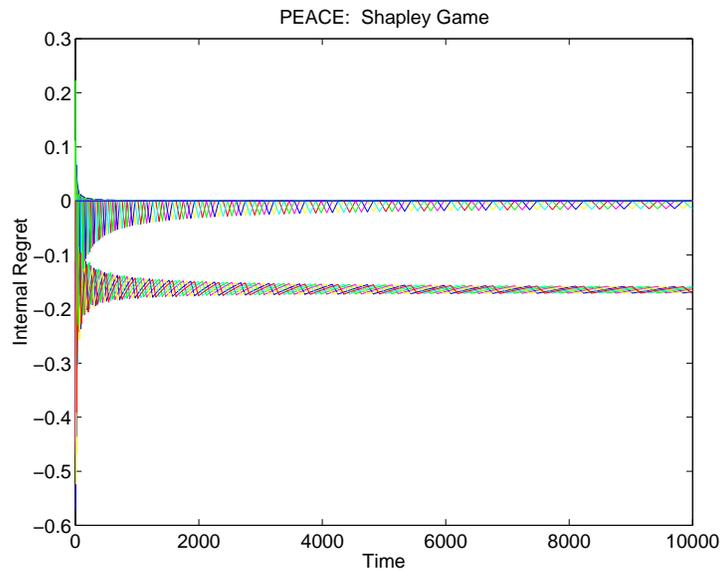


# Shapley Game

## Internal Regret: HMC



## Internal Regret: PEACE



## Naive No-Regret Learning

$$\hat{R}_i(a_i, a_{-i}^t) = \begin{cases} \frac{R_i(a_i, a_{-i}^t)}{\hat{\pi}_i^t(a_i)} & \text{if } a_i^t = a_i \\ 0 & \text{otherwise} \end{cases}$$

$$\hat{\pi}_i^t = (1 - \epsilon)\pi_i^t + \frac{\epsilon}{|A_i|}$$

### Theorem

If an informed learning algorithm  $\mathcal{A}_i$  exhibits no-regret, then the naive learning algorithm  $\hat{\mathcal{A}}_i$  exhibits  $\epsilon$ -no-regret.

# Multiagent Algorithms

NRQ(MGame,  $\gamma, \alpha, \epsilon$ )

Inputs      discount factor  $\gamma$   
              rate of averaging  $\alpha$   
              rate of exploration  $\epsilon$

Output      equilibrium state-value function  $V^*$   
              equilibrium action-value function  $Q^*$

Initialize    $V = Q = 0$

REPEAT

  initialize  $s, a_1, \dots, a_n$

  WHILE  $s$  is nonterminal DO

    simulate actions  $a_1, \dots, a_n$  in state  $s$

    observe rewards  $R_1, \dots, R_n$  and next state  $s'$

    for all  $i \in I$

      let  $\pi(s', \vec{a}) = \prod_i \pi_i(s', a_i)$

      compute  $V_i(s') = \sum_{\vec{a} \in A} \pi(s', \vec{a}) Q_i(s', a)$

      update  $Q_i(s, \vec{a}) = (1 - \alpha) Q_i(s, \vec{a}) + \alpha P_i(s, a_i)$

      where  $P_i(s, a_i) = R_i + \gamma V_i(s')$

      update policies

        informed

        naive

    (simultaneously) choose actions  $a'_1, \dots, a'_n$

$s = s', a_1 = a'_1, \dots, a_n = a'_n$

    decay  $\alpha$

FOREVER

# Grid Game 1

GG1

$u\text{CE-}Q$		$S$	$C$
	$S$	72.90,72.90	72.90,72.90
	$C$	72.90,72.90	64.61,64.61

$\text{CE-}Q$		$S$	$C$
	$S$	72.90,72.90	72.90,72.90
	$C$	72.90,72.90	64.61,64.61

<b>Informed NER-<math>Q</math></b>		$S$	$C$
	$S$	72.90,72.90	72.90,72.90
	$C$	72.90,72.90	64.61,64.61

<b>Naive NIR-<math>Q</math></b>		$S$	$C$
	$S$	71.41,71.40	71.77,71.76
	$C$	72.05,72.05	63.36,63.35

## Grid Game 2

GG2

	<i>S</i>	<i>C</i>
<i>S</i>	58.55,49.77	39.46,81.00
<i>C</i>	81.00,40.22	71.76,33.34

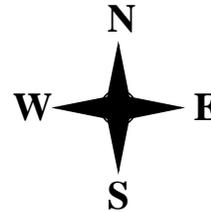
	<i>S</i>	<i>C</i>
<i>S</i>	49.27,59.64	40.19,81.00
<i>C</i>	38.37,42.63	34.47,71.00

	<i>S</i>	<i>C</i>
<i>S</i>	57.62,48.32	41.32,39.96
<i>C</i>	81.00,41.39	71.90,36.57

	<i>S</i>	<i>C</i>
<i>S</i>	30.47,68.36	41.67,81.00
<i>C</i>	80.30,29.21	29.75,61.57

# Grid Soccer

A			B
A			B
A			B
A			B
A	B	(A)	B
A			B
A			B



## Policies

### MM- $Q$

	N	E	W	stick
$P1$	0.290	0	0	0.710
$P2$	0.716	0	0	0.284

### Informed NER- $Q$

	N	E	W	stick
$P1$	0.358	0.001	0.001	0.641
$P2$	0.657	0.001	0.001	0.342

### Naive NIR- $Q$

	N	E	W	stick
$P1$	0.739	0.013	0.044	0.204
$P2$	0.236	0.026	0.013	0.725

# Grid Soccer

## Q-Values

MM-Q

	N	E	W	stick
N	26.21, -26.21	54.59, -54.59	81.00, -81.00	32.54, -32.54
E	-100.00, 100.00	-100.00, 100.00	-100.00, 100.00	-100.00, 100.00
W	36.71, -36.71	-0.74, 0.74	95.26, -95.26	-29.52, 29.52
stick	28.75, -28.75	26.12, -26.12	31.10, -31.10	26.15, -26.15

Informed NER-Q

	N	E	W	stick
N	24.43, -26.50	56.15, -58.14	81.00, -81.00	28.85, -33.91
E	-100.00, 100.00	-100.00, 100.00	-100.00, 100.00	-100.00, 100.00
W	26.08, -29.02	-0.35, -2.12	95.23, -95.15	-27.07, 23.64
stick	27.25, -30.46	24.04, -26.41	28.11, -27.39	24.04, -26.41

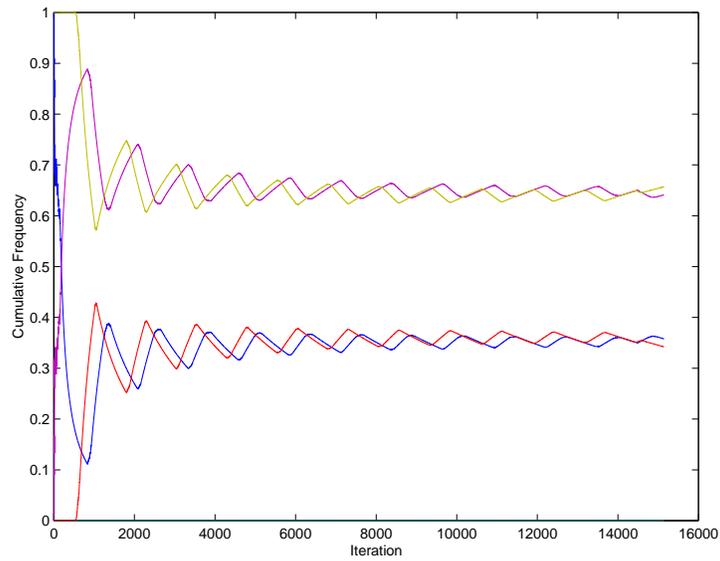
Naive NIR-Q

	N	E	W	stick
N	27.74, -27.73	50.70, -50.71	80.48, -80.48	31.27, -31.28
E	-100.00, 100.00	-100.00, 100.00	-100.00, 100.00	-100.00, 100.00
W	28.93, -28.94	18.44, -18.44	97.53, -97.53	-19.26, 19.25
stick	30.67, -30.66	26.98, -26.98	34.00, -34.01	26.36, -26.35

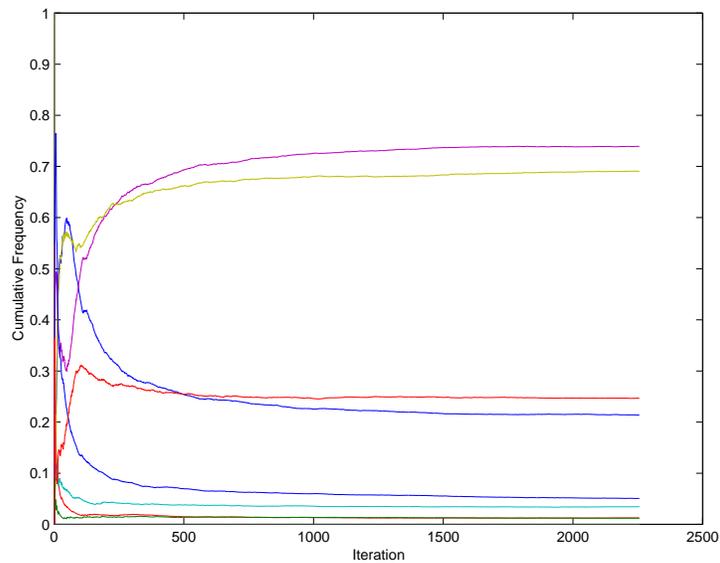
# Grid Soccer

## Empirical Frequencies

### Informed NER-Q



### Naive NIR-Q



## Conjectures

- **Correlated- $Q$**  learning converges to  $Q$ -values that support equilibrium policies in Markov games.
- **WAR** and **PEACE** exhibit no-internal regret.
- **NER  $Q$ -Learning** converges to minimax strategies in constant-sum Markov games.
- **NIR  $Q$ -Learning** converges to correlated equilibrium in general-sum Markov games.