A simple learning strategy which realizes robust cooperation better than Pavlov in Iterated Prisoners' Dilemma

> Joe Yuichiro Wakano Norio Yamamura (Center for Ecological Research, Kyoto Univ.)

### Prisoner's Dilemma (PD)



C:cooperation D:defection

Defection is always the best strategy in a single round PD game.

When the game is repeated (Iterated Prisoner's Dilemma or IPD), adaptive strategy is not clear.

Pavlov is proposed as a leading strategy in IPD (Nowak & Sigmund 1993)

Pavlov := 
$$(1, 0, 0, 1)$$
  
Probability of  
playing C after CC CD DC DD

TFT (1,0,1,0) ends up very poor when noise is introduced.

Pavlov is invaded by All-D but Pavlov like strategy (1-e,0,0,1-e) is not. (Nowak & Sigmund 1993)

Is Pavlov like strategy an ESS in a strategy set of (p1,p2,p3,p4)?

## Performance of Classic Strategies

Payoff for (p1,p2,p3,p4) against (q1,q2,q3,q4) can be derived analytically (eigen vector in Markov process).

21<sup>4</sup> strategies are matched each other and payoff is calculated. By investigating the payoff matrix, we obtain ...

Pavlov like strategy (1-e,0,0,1-e) is invaded by (1,0,0,1). Only two ESSs found: All-D and GRIM (1,0,0,0) No strategy of type (p1,p2,p3,p4) can establish stable mutual cooperation. **Reinforcement learning:** 

"Repeat when rewarded, change when punished" or "Win-Stay, Lose-Shift"

Simple and intuitively understandable rule.

Pavlov can be considered as the most simple player adopting reinforcement learning. However, Pavlov considers the previous move only.

We will intoduce a learning player who considers past experiences (generalized Pavlov ?).

In order to adopt reinforcement learning, standard or aspiration level of evaluation must be defined.



To establish stable mutual cooperation, 1<s<3 is necessary.

# Model of Learning Players

Aspiration level *s* is genetically fixed. Internal state *h* is changed during reptition of PD. C is played when current internal state h>0, otherwise D.

$$\Delta h = a \cdot \text{sgn}(f - s) \cdot \text{sgn}(h)$$
$$\text{sgn}(x) \equiv \begin{cases} 1(x \ge 0) \\ -1(x < 0) \end{cases}$$

Examples:

If C is played and the resulting score f is larger than s, cooperation is affirmatively learned. That is, the player increases h to become more cooperative.

If C is played and the resulting score f is smaller than s, cooperation is negatively learned. That is, the player decreases h to become less cooperative.

If D is played and the resulting score f is larger than s, defection is affirmatively learned, leading decrease in h.

Each IPD consists of 10,000 PD games. Error or noise is introduced at probability *e*. When error happens, the opposite action expected from *h* value is played.

Internal state *h* is reset to a/2 every time the opponent changes, thus the learner always plays D at the first round.

→ The learner can learn to exploit All-C.

Internal state h can be considered as 'faith in cooperation'. Such faith is altered by past experience.

Payoffs of both players in IPD is determined only by aspiration levels.



$$\Delta h = 2 \cdot \text{sgn}(f - s) \cdot \text{sgn}(h)$$
  
$$\text{sgn}(x) \equiv \begin{cases} 1(x \ge 0) \\ -1(x < 0) \end{cases}$$

Only 5 independent variants;

DL(sf0), DL(0 < sf1), DL(1 < sf3), DL(3 < sf5), DL(s > 5)



## **Example of Digital Learner's behavior**

Table Comparison of digital learner  $(1 < s f_3, a = 2)$  and Pavlov

strategy	opponent	actual realization of game (h)	average
			payoff
DL(1 < sf3)	All-D	DD(-1), CD(1), DD(-1), CD(1), DD(-1),	0.5
	All-C	DC(-1), DC(-3), DC(-5), DC(-7), DC(-9),	5
	Pavlov	$DD(-1), CC(1), CC(3), CC(5), CC(7), \ldots$	3
	$\mathrm{TFT}$	DC(-1), DD(-3), DD(-1), CD(1), DC(-1),	1.75
	DL(1 <s<b>£3)</s<b>	DD(-1), $CC(1)$ , $CC(3)$ , $CC(5)$ , $CC(7)$ ,	3
Pavlov	All-D	DD, CD, DD, CD, DD,	0.5
	All-C	DC, DC, DC, DC,	5
	Pavlov	DD, CC, CC, CC, CC,	3
	$\mathrm{TFT}$	DC, DD, CD, DC, DD,	2
	DL(1 < sf3)	DD, CC, CC, CC, CC,	3

# ESS among 5 variants of Digital Learner

able Results of matches among digital learners.

	player	opponent (wild type)				
	(mutant)	s f 0	0 <s<b>£1</s<b>	1 <s<b>£3</s<b>	$3 \le \mathbf{f5}$	s>5
a) 1% error	s f 0	1.030049	1.030053	2.965461	2.965461	2.985005
	0 <s£1< td=""><td>1.030048</td><td>1.030052</td><td>2.985288</td><td>2.965454</td><td>2.974166</td></s£1<>	1.030048	1.030052	2.985288	2.965454	2.974166
	$1 \le \mathbf{f}3$	0.535146	0.535152	2.984959	0.584596	1.519351
	$3 \le \mathbf{f}_5$	0.539876	0.53988	2.943325	1.738342	2.502725
	s>5	0.535053	0.545318	3.954128	1.258617	2.010029
o)10% error	s f 0	1.290936	1.290994	2.851567	2.709865	2.849876
	0 <s£1< td=""><td>1.290922</td><td>1.290977</td><td>2.851499</td><td>2.709804</td><td>2.788106</td></s£1<>	1.290922	1.290977	2.851499	2.709804	2.788106
	1 <s<b>£3</s<b>	0.848479	0.848549	2.865511	1.243602	1.848385
	$3 \le \mathbf{f}_5$	0.888843	0.88891	2.482726	1.910963	2.405803
	s>5	0.850024	0.98887	3.526731	1.601163	2.090827
· · · · · · · · · · · · · · · · · · ·						

cooperative variant

# **ESS among 5 variants of Digital Learner**

able Results of matches among digital learners.

	player			opponent (wild type)			
	(mutant)	$s \mathbf{f} 0$	0 <s<b>£1</s<b>	1 <s<b>£3</s<b>	$3 \le \mathbf{f}_5$	s>5	
a) 1% error	s f 0	1.030049	1.030053	2.965461	2.965461	2.985005	
	0 <s£1< td=""><td>1.030048</td><td>1.030052</td><td>2.985288</td><td>2.965454</td><td>2.974166</td></s£1<>	1.030048	1.030052	2.985288	2.965454	2.974166	
	$1 \le \mathbf{f}3$	0.535146	0.535152	2.984959	0.584596	1.519351	
	$3 \le \mathbf{f}_5$	0.539876	0.53988	2.943325	1.738342	2.502725	
	s>5	0.535053	0.545318	3.954128	1.258617	2.010029	
o)10% error	s f 0	1.290936	1.290994	2.851567	2.709865	2.849876	
	0 <s£1< td=""><td>1.290922</td><td>1.290977</td><td>2.851499</td><td>2.709804</td><td>2.788106</td></s£1<>	1.290922	1.290977	2.851499	2.709804	2.788106	
	1 <s<b>£3</s<b>	0.848479	0.848549	2.865511	1.243602	1.848385	
	$3 \le \mathbf{f}_{5}$	0.888843	0.88891	2.482726	1.910963	2.405803	
	s>5	0.850024	0.98887	3.526731	1.601163	2.090827	

ESS

cooperative variant

Cooperative variant DL(1<s<3) is not an ESS. DL(s<0) is almost same as All-D, thus defection spreads

# Analog Learner

The impact on learning process may depend on the payoff value (instead of a constant as in DL model). Analog evaluation model is developed.

$$\Delta h = 2 \cdot (f - s) \cdot \operatorname{sgn}(h)$$

Small difference in aspiration level *s* produces different learning process.

We used 56 different *s* values (s=-0.1, 0, 0.1, 0.2, ..., 5.4) in computer simulations.

### **Can Analog Learners invade classic strategies ?**



Analog Learner for s=1.6,...,3.0 invades All-C, All-D, TFT and Pavlov.

### **Can Analog Learners invade classic strategies ?**



Analog Learner for s=1.6,...,3.0 invades All-C, All-D, TFT and Pavlov.

#### Can classic strategies invade Analog Learners?



Analog Learner for s=2.6,...,3.0 is stable against invasion of All-C, All-D, TFT and Pavlov.

#### Can classic strategies invade Analog Learners?



Analog Learner for s=2.7,...,2.8 is stable against invasion of All-C, All-D, TFT and Pavlov.

Summary of classic strategies vs. Analog Learner

Though Analog Learner cannot invade All-D population, cooperative variants of Analog Learner do very well against classic strategies.

1% error:

Analog Learner for s=2.6,...,3.0 invades and is stable against All-C, TFT, Pavlov (and All-D).

10% error:

Analog Learner for s=2.7,...,2.8 invades and is stable against All-C, TFT, Pavlov (and All-D).

## **Evolutionary dynamics of aspiration level**

F(s,s\*):= payoff value of AL(s=s) when matched against AL(s=s\*)

 $F(s,s^*)$  is obtained by compter simulation. By calculating many sets of  $(s,s^*)$ , we obtain the payoff matrix (56 x 56 matrix).

### **Evolutionary dynamics of aspiration level**



10.00		
	2.9000	-3.000
	2.8000	-2.900
	2.7000	-2.800
	2.6000	-2.700
	2.5000	-2.600
	2.4000	-2.500
	2.3000	-2.400
	2.2000	-2.300
	2.1000	-2.200
	2.0000	-2.100
	1.9000	-2.000
	<b>1</b> .8000	-1.900
	<b>1</b> .7000	-1.800
	1.6000	-1.700
	<b>1.5000</b>	-1.600
	<b>1</b> .4000	-1.500
	<b>1</b> .3000	-1.400
	1.2000	-1.300
	1.1000	-1.200
	1.0000	-1.100
	0.9000	-1.000
	0.8000	-0.900
	0.7000	-0.800
	0.6000	-0.700
	0.5000	-0.600
1.00		

# **Evolutionary dynamics of aspiration level**

From this payoff matrix, Pairwise Invasibility Plot (PIP) is drawn.

By surveying PIP, it becomes clear whether evolutionarily stable strategy (ESS) exists or not.



ESS among Analog Learners

Pairwise Invasibility Plot

s\*

ESS among Analog Learners



Pairwise Invasibility Plot

## Analog Learner can establish robust cooperation

- Under both 1% and 10% error rate, Analog Learner with aspiration level around 2.7 can invade All-C, TFT or Pavlov populations and establish very robust mutual cooperation.
- Such population is more cooperative than Pavlov population.

	Average scores of populations for infinitely					
$\mathbf{F}(2,7,2,7)$	iterated Prisoners' Dilemma					
F(2./,2./):	erro	All-C	Pavlov	TFT	All-I	
Average score of	r			(1,0,1,0)	(0,0,0	
AL(s=2.7) population	rate	(1,1,1,1)	(1,0,0,1)			
	0%	3.000	3.000	3.000	1.00	
1% 2.9891	1%	2.990	2.951	2.010	1.03	
	5%	2.948	2.777	2.048	1.14	
10% 2.8491	10%	2.890	2.602	2.090	1.29	
	20%	2.760	2.376	2.160	1.56	

Due to strong faith in cooperation, it allows its partner's defection played by error without punishment.

## Difference between Digital and Analog Learner

# When matched against All-D (or DL(s<1) or AL(s<1));

	-1	1	-1	1	1
DL(1 <s<3)< th=""><th>DD</th><th>CD</th><th>DD</th><th>CD</th><th>DD</th></s<3)<>	DD	CD	DD	CD	DD
AL(s=2)	DD	CD	DD	DD	CD
	- 1	1	-3	- 1	1

All-D strategy's payoff =  $\begin{cases} 3 \text{ (Digital Learner)} \\ 2.3 \text{ (Analog Learner)} \\ \land \\ \text{Average payoff of AL(s=2) population (=3)} \end{cases}$ 

Quantified evaluation ability is essential to avoid being exploited by always defecting strategies. Memory decay; more recent result might affect internal state more strongly

Analog Learner with non-linear evaluation