

The Prisoner's Dilemma

Can a Q-learning agent learn to play?

Irene Ferfaglia
UniTS - Reinforcement Learning exam
18.07.2022

The agenda

- Prisoner's Dilemma
 - Introduction
 - Formalisation
 - Player's strategies
- Q-Learning
 - TD learning
 - Policy control in TD(0)
 - The algorithm
 - Convergence
- Results

The Prisoner's Dilemma

The background of the slide is a solid dark blue. In the bottom right corner, there are two parallel diagonal stripes. The upper stripe is white, and the lower stripe is a lighter shade of blue. Both stripes extend from the bottom edge towards the top right corner.

Introduction

- Two prisoners with no means of speaking to each other
- Each can either cooperate with the other or give them up to the police
- They don't know each other's action
- No loyalty to each other: objective is to maximise their own reward



Multiple games

If the two players

- play more than once
- remember previous actions of their opponent
- change strategy accordingly

→ Iterated Prisoner's Dilemma

Formalisation

- Adversarial bandits problem
- Two choices of actions: Cooperate (C) or Defect (D)
- Four states: (A's previous action, B's previous action)
- Rewards: defined in a payoff matrix

Payoff matrix

| A \ B | Cooperate | Defect |
|-----------|-----------|--------|
| Cooperate | (R,R) | (S,T) |
| Defect | (T,S) | (P,P) |

- $T > R > P > S$, $2R > S + T$
- For each player it is beneficial to defect: Nash equilibrium
- **BUT** highest mutual payoff is with cooperation

Payoff matrix

| A \ B | Cooperate | Defect |
|-----------|-----------|--------|
| Cooperate | (3,3) | (0,5) |
| Defect | (5,0) | (1,1) |

- $5 > 3 > 1 > 0$, $6 > 5$
- For each player it is beneficial to defect: Nash equilibrium
- **BUT** highest mutual payoff is with cooperation

Strategies

- Always cooperate
- Always defect
- Random action
- Tit-for-Tat: copy opponent's previous action

Q-Learning

The background of the slide is a solid dark blue. In the bottom right corner, there are two parallel diagonal stripes. The top stripe is white, and the bottom stripe is a lighter shade of blue. Both stripes extend from the bottom right towards the top left.

Temporal-Difference (TD) learning

- Learn from experience, without a model of the environment, its rewards and next-state probability distributions
- Update estimates before final outcome is known (bootstrap)
- Simplest TD method, **TD(0)**:

$$V(S_t) \leftarrow V(S_t) + \alpha [R_{t+1} + \gamma V(R_{t+1}) - V(S_t)]$$

Policy control in TD(0)

| SARSA | Q-Learning |
|--|---|
| $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \underline{Q(S_{t+1}, A_{t+1})} - Q(S_t, A_t)]$ | $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \max_a \underline{Q(S_{t+1}, a)} - Q(S_t, A_t)]$ |
| On-policy | Off-policy |
| Learns near-optimal policy while exploring | Directly learns optimal policy |
| More conservative | More aggressive |

Q-Learning algorithm

Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$

Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+$, $a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

 Initialize S

 Loop for each step of episode:

 Choose A from S using policy derived from Q (e.g., ε -greedy)

 Take action A , observe R, S'

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

$S \leftarrow S'$

 until S is terminal

Requirements for convergence

- Every state-action pair continues to be visited
- Learning rate decreases over time
- Stationary and Markovian environment
- Problem! when vs Q-learner, environment is not stationary

ϵ -greedy policy

- Simple way to achieve balance of exploration and exploitation
- $0 \leq \epsilon \leq 1$

$$A_t = \begin{cases} \max_a Q(s, a) & \text{with probability } 1 - \epsilon \\ \text{random action} & \text{with probability } \epsilon \end{cases}$$

- Can be decaying by a factor

The choice of α

- TD algorithm convergence criterion:

$$\sum_t \alpha_t = \infty \qquad \sum_t \alpha_t^2 < \infty$$

So α is set to $\frac{1}{t+1}$.

Results

Can a Q-learning agent learn to play:

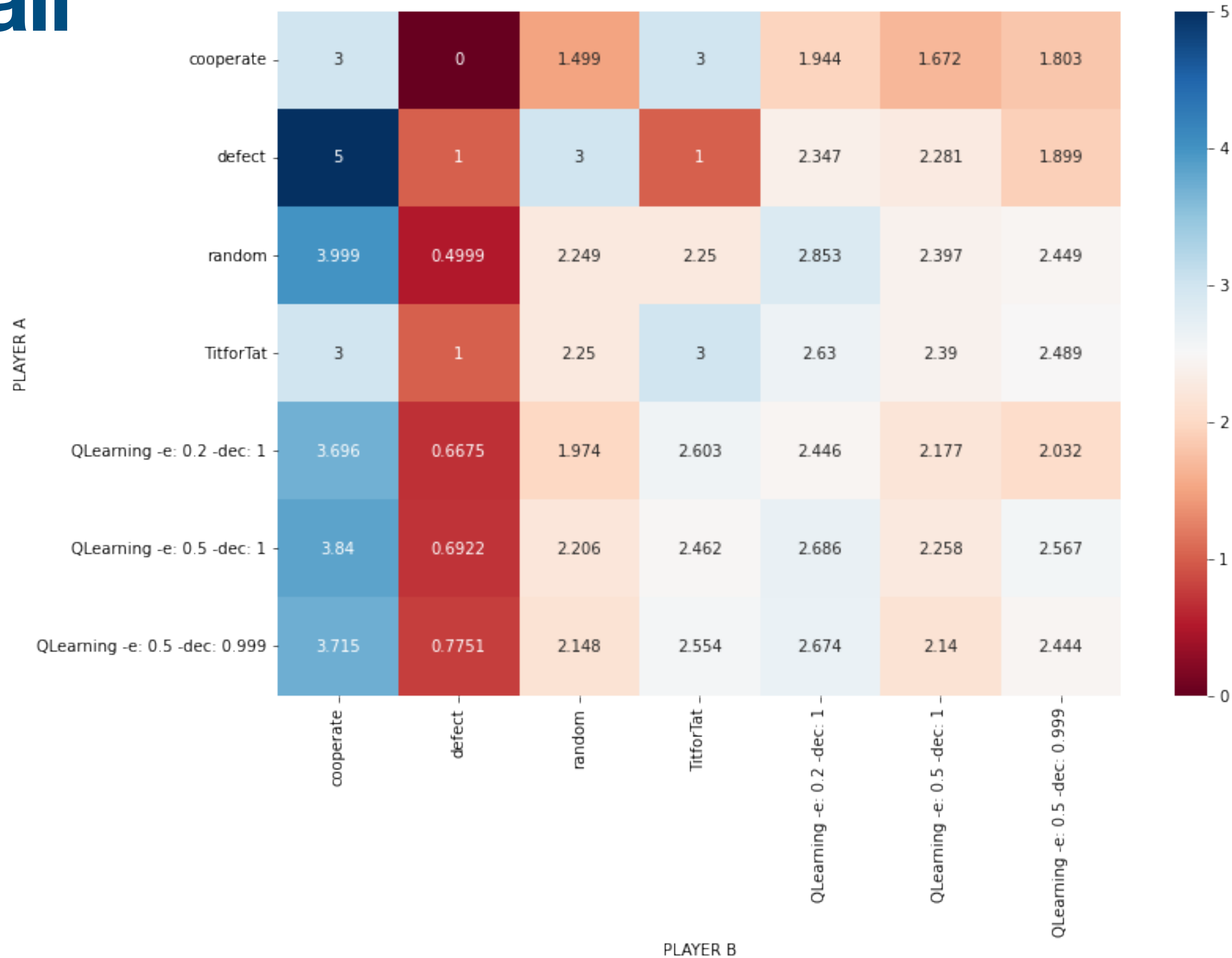
- vs a fixed strategy?
- vs another Q-learning agent?

Set Up

- 100 independent meetings, to avoid random uncertainty
- 10'000 games each meeting
- Q learner with varying epsilons and decays
- Strategies are unknown to one another
- Objective: learn optimal strategy to maximise own avg reward

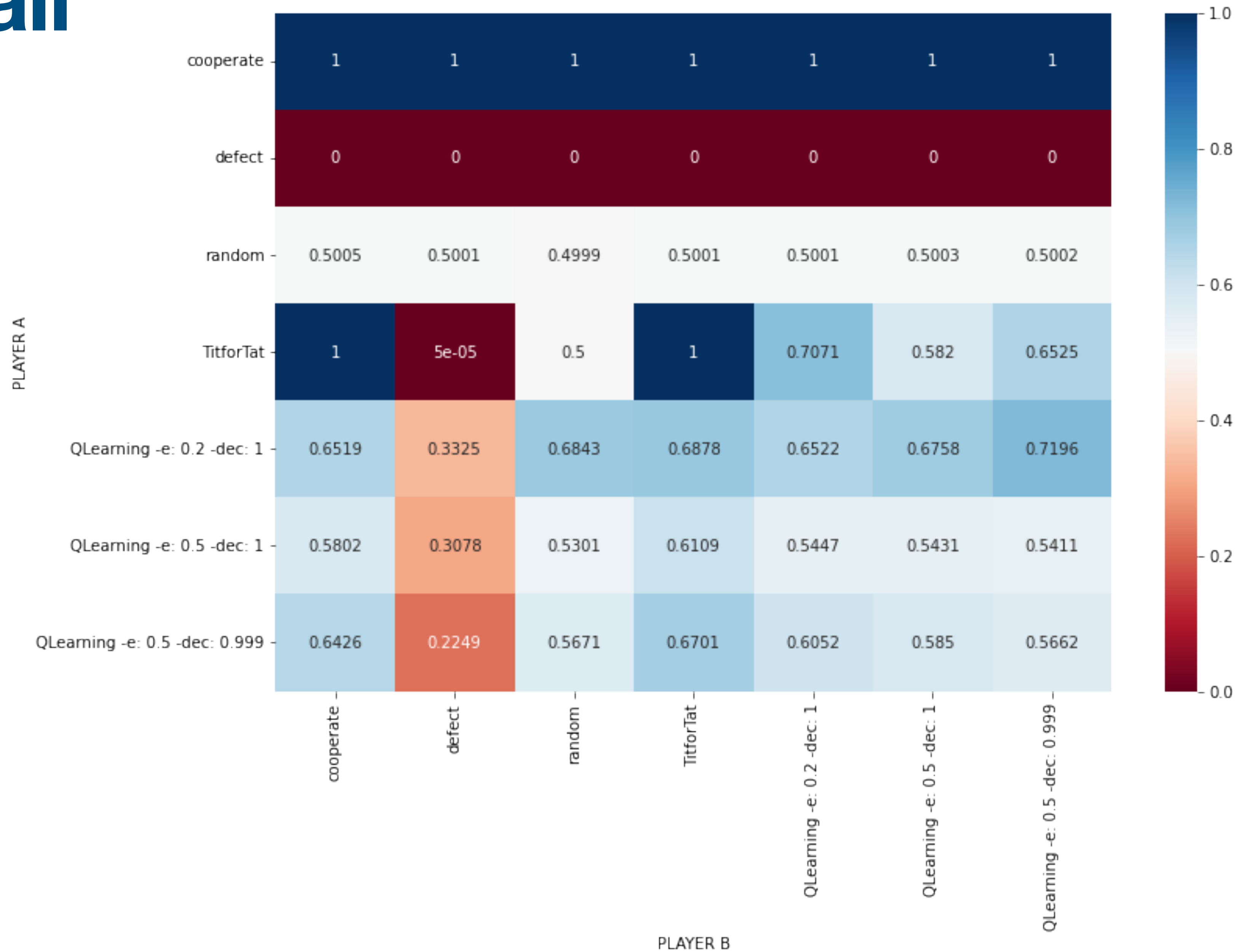
All vs all

Average reward per game of Player A vs Player B

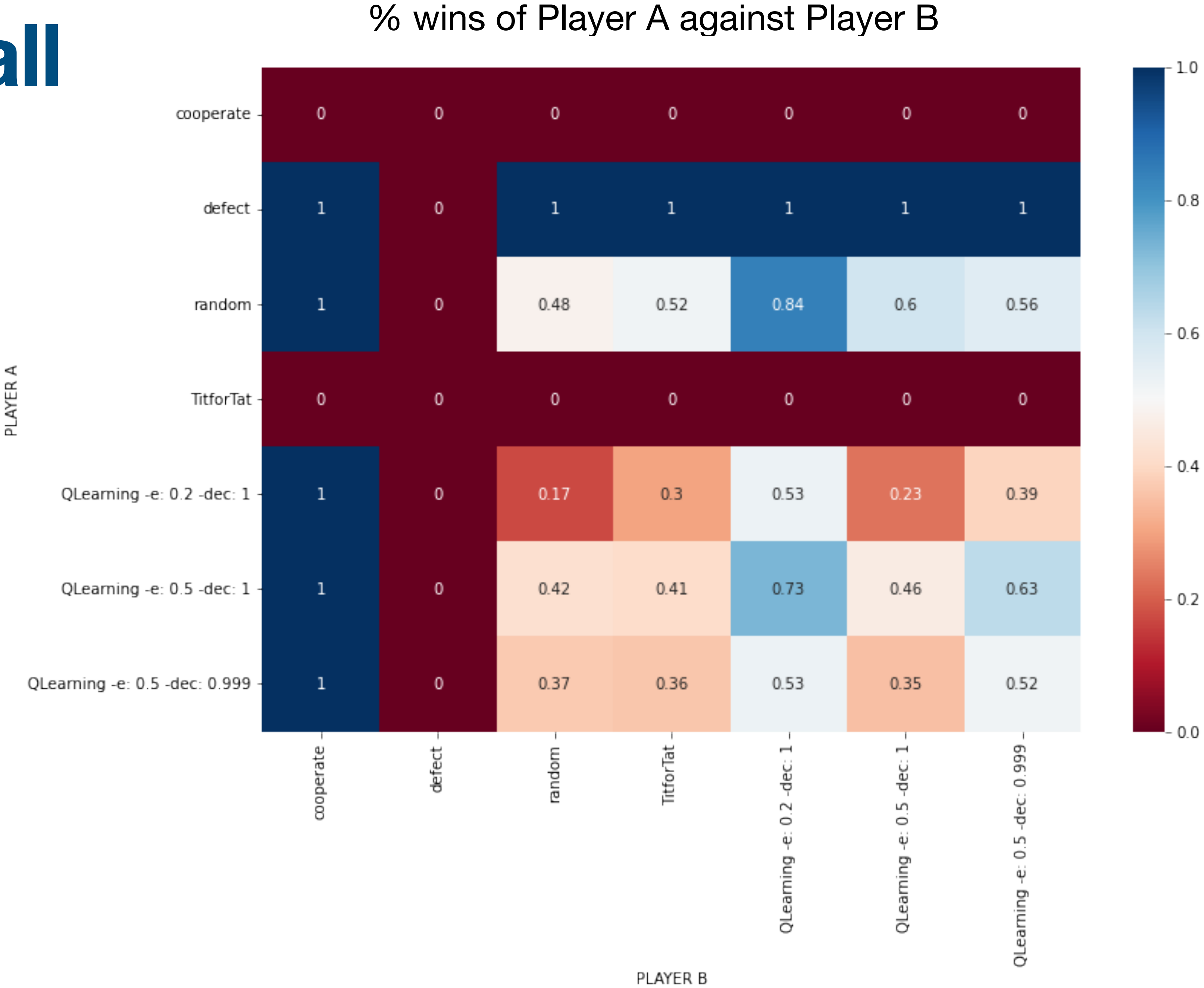


All vs all

Average % of cooperations per game of Player A vs Player B



All vs all



Conclusions

Can a Q-learning agent learn to play the Iterated Prisoner's Dilemma...

- Vs a fixed strategy?
- Vs another Q-learning agent?

Kind of

Thank you!