

Inducing Cooperation through Reward Reshaping based on Peer Evaluations in Deep Multi-Agent Reinforcement Learning

AAMAS 2020

Ming, Will, Yang

Intelligent Transportation Group
State Key Laboratory of Internet of Things for Smart City
Faculty of Science and Technology
University of Macau

2020/11/14



Contents

- 1 Motivation
- 2 Problem Setting
- 3 Solution
- 4 Experiment
- 5 Thoughts

Contents

- 1 Motivation
- 2 Problem Setting
- 3 Solution
- 4 Experiment
- 5 Thoughts

Motivation

semi-cooperative

agents may each have its own separate reward function (the joint reward scenario is a special case), but is willing to cooperate if an incentive for cooperation is appropriately provided.

social welfare

the sum of the rewards of each agent across the entire episode.

Contents

- 1 Motivation
- 2 Problem Setting
- 3 Solution
- 4 Experiment
- 5 Thoughts

Problem Setting

Stochastic game

a multi-agent task with each agent having its own individual reward by modeling it with a stochastic game

Semi-cooperative task

semi-cooperative tasks as the set of tasks where each agent may have a separate reward function but may benefit from cooperative strategies such as the prisoner's dilemma and the **stag hunt game**.

Goal: Maximizing the social welfare

the sum of the rewards of each agent across the entire episode.

$$\pi^* = (\pi_1^*(u_1 \mid o_1), \pi_2^*(u_2 \mid o_2), \dots, \pi_n^*(u_n \mid o_n))$$

$$\pi^* = \arg \max_{\pi} \sum_{a \in \mathcal{A}} \mathbb{E}_{s \sim \rho^{\pi}, u \sim \pi} [r_a(s, u)]$$

Contents

- 1 Motivation
- 2 Problem Setting
- 3 Solution**
- 4 Experiment
- 5 Thoughts

- ① Change of Games via Reward Reshaping
- ② Reward Update with Peer Evaluation
 - ① Peer evaluation signal
 - ② Reshaping reward from peer evaluation
- ③ Peer-evaluation based Dual DQN (PED-DQN)

Change of Games via Reward Reshaping

reward function

in the game G^t :

$$\hat{r}^t = (\hat{r}_a^t : a \in \mathcal{A})$$

at time step $t = 0, 1, \dots$

- 1 Compute the optimal π^t policy from G^t
- 2 Evaluate how well-coordinated G^t is by evaluating π^t
- 3 Update from G^t to G^{t+1} by updating from \hat{r}^t to \hat{r}^{t+1} , using the 'evaluation feedback' from (2)
- 4 Increment t and go to (1)

Framework

Policy Update: $\pi^{t+1} = F(\pi^t, \hat{r}^t)$

Reward Update: $\hat{r}^{t+1} = H(\hat{r}^t, \pi^t)$

Reward Update with Peer Evaluation

Peer evaluation signal:

counterfactual evaluation signal (CES) z_k^t

for agent k :

$$z_k^t(o_k^t, o_k^{t+1}, u_k^t, r_k^t) := r_k^t + \gamma Q_k^{\pi^t}(o_k^{t+1}, \pi_k^t(o_k^{t+1})) - Q_k^{\pi^t}(o_k^t, u_k^t)$$

Reshaping reward from peer evaluation:

agent a and agent k are peers.

For agent a :

$$Z_a^t[o_a^t, u_a^t] = \frac{1}{|K_a|} \sum_{k \in K_a} z_k^t$$

sample many times:

$$\hat{Z}_a^t [o_a^t, u_a^t] \approx \frac{1}{|K_a|} \sum_{k \in K_a} \times \mathbb{E}_{(o, o', u, r) \sim \pi^t: o_a = o_a^t, u_a = u_a^t} [z_k (o_k, o'_k, u_k, r_k)]$$

moving average:

$$\hat{Z}_a^{t+1} [o_a^t, u_a^t] \leftarrow (1 - \alpha) \hat{Z}_a^t [o_a^t, u_a^t] + \alpha Z_a^t [o_a^t, u_a^t]$$

final the reshaped reward $\hat{r}_a^t [o_a^t, u_a^t]$:

$$\hat{r}_a^t [o_a^t, u_a^t] = r_a^t + \beta \hat{Z}_a^t [o_a^t, u_a^t]$$



Example: Evaluation feedback exchange

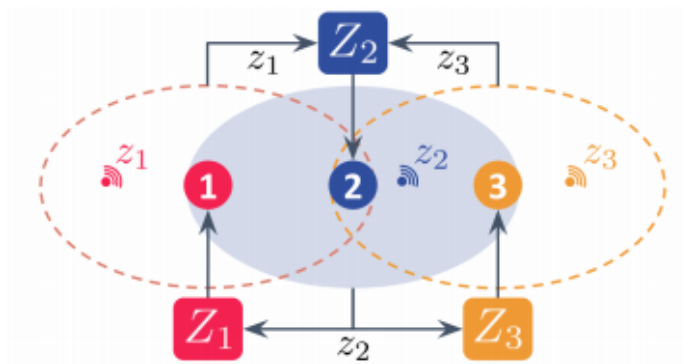
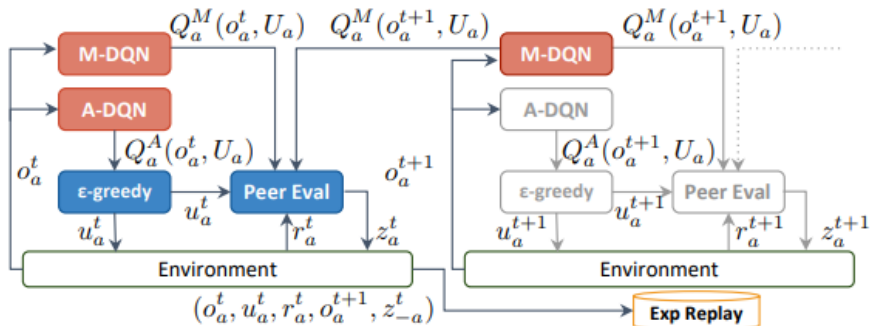


Figure: Example: Evaluation feedback exchange

Peer-evaluation based Dual DQN (PED-DQN)



(b) The architecture of PED-DQN

Figure: The architecture of PED-DQN

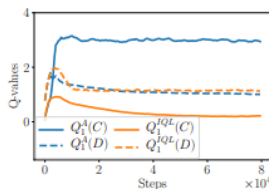
Contents

- 1 Motivation
- 2 Problem Setting
- 3 Solution
- 4 Experiment**
- 5 Thoughts

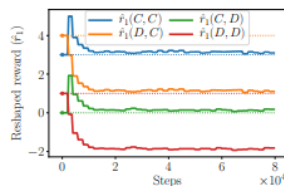
Case study: Prisoner's dilemma

Table 1: Prisoner's dilemma

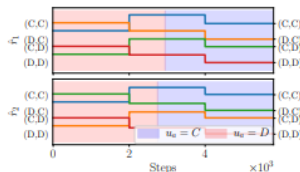
(a) original payoff			(b) reshaped payoff		
	C	D		C	D
C	3, 3	0, 4	C	3.11, 3.13	0.10, 1.13
D	4, 0	1, 1	D	1.11, 0.13	-1.90, -1.87



(a) IQL learned defection while ours learned cooperation



(b) Reshaped rewards $\hat{r}_1(D, C)$ and $\hat{r}_1(D, D)$ penalize defection



(c) Deep evaluations steered the agents' policy



Experiments

Algorithms:

- QMIX
- IDQ PED-DQN
- PE:single network
- Pro DQN:directly use reward

$$z_a = r_a$$

$$\hat{r}_a = r_a + \frac{\beta}{|K_a|} \sum_{k \in K_a} r_k$$

Environments:

- Resource share
- Partially cooperative pursuit (PCP)



Experiments

please read the paper to check the experiments details.

Contents

- 1 Motivation
- 2 Problem Setting
- 3 Solution
- 4 Experiment
- 5 Thoughts

communication & reward shaping