

# Contracting Actions in Self-Interested Multi-Agent Reinforcement Learning

Paper #...

## ABSTRACT

Whereas agents in cooperative environments share the same reward function in self-interested settings agents maximize their individual returns. As a consequence agents learned policies tend to be overly greedy especially when shared resources between agents are scarce. Cooperative game theory provides a theoretical framework for self-interested settings where binding agreements between agents are possible. In this work we propose binding agreements for multi-agent reinforcement learning which allow agents to transfer reward in exchange for following a non greedy trajectory. Thus, agents are enabled to compensate each other for their occurred losses that arise from behaving non-greedy. We evaluate our proposed method in a smart-factory scenario where agents compete for available machines. The tasks either have low-priority or high-priority, giving small rewards or high-rewards respectively. We empirically demonstrate that RL agents stably learn agreements and thus originate higher returns compared with the pure self-interested case.

## KEYWORDS

multi-agent learning; cooperative game theory

## 5 EXPERIMENTS

### 5.1 Evaluation Environments

5.1.1 Smart Factory (SF).

### 5.2 Methods

### 5.3 Results

### 5.4 Discussion

## 6 CONCLUSION AND FUTURE WORK

## 1 INTRODUCTION

## 2 BACKGROUND

## 3 RELATED WORK

## 4 CONTRACTING

---

### Algorithm 1 Contracting

---

```
1: procedure CONTRACT( $Q^{GREEDY}, Q^{NON-GREEDY}, N_c$ )
2:   Observe current state  $s_1$ 
3:   Calculate compensation  $q_c \leftarrow \max_a(Q^{non-greedy}(s_1, a))$ 
4:   for  $t = 1, N_c$  do
5:      $a_{greedy} \leftarrow \operatorname{argmax}_a(Q^{greedy}(s, a))$ 
6:      $a_{non-greedy} \leftarrow \operatorname{argmin}_a(Q^{non-greedy}(s, a))$ 
7:     Execute  $a_t \leftarrow \langle a_{greedy}, a_{non-greedy} \rangle$ 
8:     Observe rewards  $r_{g,t}, r_{n,t}$  and new state  $s_{t+1}$ 
9:      $s \leftarrow s_{t+1}$ 
10:   Final compensation  $q_c \leftarrow q_c - \max_a(Q^{non-greedy}(s, a))$ 
11:   Final returns  $R_g \leftarrow -q_c + \sum_{t=1}^{N_c} r_{g,t}, R_n \leftarrow q_c + \sum_{t=1}^{N_c} r_{n,t}$ 
```

---