



Multi-Agent Reinforcement Learning

Theory and Examples

Quantitative Life Sciences - ICTP 2024

Cristopher Erazo

July 1, 2024



The Abdus Salam

**International Centre
for Theoretical Physics**



Table of Contents

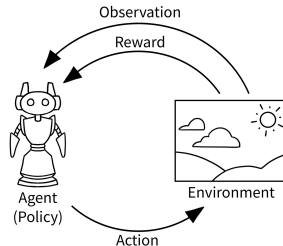
List

- ▶ Introduction to Single Agent
- ▶ Generalization MARL
- ▶ Particularities in MARL
- ▶ Matrix Games
- ▶ Numerical Examples
- ▶ Conclusions

Ingredients of Single Agent Reinforcement Learning

Markov Decision Process

- \mathcal{S} finite set of environment states $s \in \mathcal{S}$.
- \mathcal{A} finite set of actions $a \in \mathcal{A}$.
- $p(s'|s, a)$ state transition probability.
- $r(s, a)$ reward function.
- $\pi(a|s)$ policy.



At time t , the environment is at s_t . The agent observes the state and takes an action $a_t \sim \pi(\cdot|s_t)$. As a result: $s_{t+1} \sim p(\cdot|s_t, a_t)$ and the reward received is $r_{t+1} = r(s_t, a_t)$.

Goal of Reinforcement Learning

Value functions

The reward only evaluates the immediate effect of the action taken by the agent.

- Value Function: $V^\pi(s) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(S_t, A_t) | S_0 = s, \pi \right]$
- State-Action Value Function: $Q^\pi(s, a) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(S_t, A_t) | S_0 = s, A_0 = a, \pi \right]$
- Objective Function: $G(\pi) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(S_t, A_t) | \pi, \rho_0 \right]$ for $S_0 \sim \rho_0$

Goal

$$\pi^* = \arg \max_{\pi} G(\pi) \quad \Leftrightarrow \quad \pi^*(a|s) = \arg \max_a Q^*(s, a)$$



Table of Contents

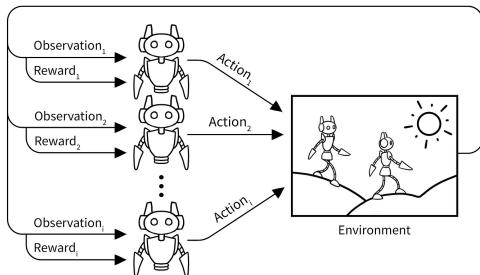
List

- ▶ Introduction to Single Agent
- ▶ Generalization MARL
- ▶ Particularities in MARL
- ▶ Matrix Games
- ▶ Numerical Examples
- ▶ Conclusions

Multi Agent Reinforcement Learning

Stochastic Games

- N agents.
- \mathcal{S} finite set of environment states.
- \mathcal{A}_n finite set of actions for agent n .
- $p(s'|s, \mathbf{a})$ state transition probability with $\mathbf{a} \in \prod_n \mathcal{A}_n$ being the joint action.
- $r_n(s, \mathbf{a})$ reward function for agent n that depends on the joint action taken.



The n -th policy is $\pi^n(a_n|s)$ with $a_n \in \mathcal{A}_n$ and the joint policy is $\pi(\mathbf{a}|s) = \prod_n \pi^n(a_n|s)$.

Value Functions in MARL

Direct Generalization

For agent n we define:

- Value Function: $V_n^\pi(s) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(S_t, \mathbf{A}_t) | S_0 = s, \pi \right]$
- State-Action Value Function: $Q_n^\pi(s, \mathbf{a}) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_n(S_t, \mathbf{A}_t) | S_0 = s, \mathbf{A}_0 = \mathbf{a}, \pi \right]$
- Objective Function: $G_n(\pi) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_n(S_t, \mathbf{A}_t) | \pi, \rho_0 \right]$

The value functions and the objective are computed with the joint policy and therefore now the returns are correlated among all the agents.

A static game has no states ($\mathcal{S} = \emptyset$). Is described only by the set of actions \mathcal{A}_n and the reward functions $r_n(\mathbf{a})$. The policy becomes a strategy $\pi^n(a_n)$ as in Game Theory.

Expected utility: $u_n(\pi) = \mathbb{E} [r_n(\mathbf{a}) | \pi]$ where $\pi(\mathbf{a}) = \prod_n \pi^n(a_n)$ is the joint strategy.

Can be separated:

$$u_n(\pi) = u_n(\pi^n, \pi^{-n}) = \sum_{a_n, \mathbf{a}_{-n}} r_n(a_n, \mathbf{a}_{-n}) \pi^n(a_n) \pi^{-n}(\mathbf{a}_{-n}) \quad (1)$$

where $(-n)$ refers to all the other agents appart from n .

Many MARL algorithms are designed for static games or work stage-wise.

Stage Game = Static Game that arises in a fixed state of the Stochastic Game.

Reward Functions = Q functions projected in action space at fixed state.

Strategy = Projection of the policy in action space at fixed state.

MARL learn strategies stage-wise and then aggregate them to obtain the overall policy.

Important Concepts in Game Theory

The **best response** for agent n is:

$$BR^n(\pi^{-n}) = \arg \max_{\pi^n} u_n(\pi^n, \pi^{-n})$$

Profile π^* is a **Nash Equilibrium** if:

$$\pi_*^n = BR^n(\pi_*^{-n}) \quad \forall n$$



Table of Contents

List

- ▶ Introduction to Single Agent
- ▶ Generalization MARL
- ▶ **Particularities in MARL**
- ▶ Matrix Games
- ▶ Numerical Examples
- ▶ Conclusions



Benefits and Challenges

In a Nutshell

Benefits:

- Experience sharing can help agents with similar tasks to learn faster.
- MARL can be implemented with parallel computation when the task is decentralized.
- MARL is more robust because when one agent fails other one can take over some of the tasks.

Challenges:

- Curse of dimensionality more severe because of the exponential growth in the state-action space with the number of agents.
- Specifying a good goal in MARL is difficult. The agent returns are correlated and cannot be maximized independently.
- Non-stationarity arises in MARL because all the agents learn simultaneously. Each agent faces a moving-target learning problem.

Specifying a good general MARL goal is **challenging** because the agent's returns are different and correlated. In the literature, the goals are generally tracked within two edges: **Stability** and **Adaptation**.

- **Stability:** Convergence to equilibrium requirement. **Nash equilibrium** are the most frequently used, but there is not a consensus in weather it is useful. A normal objection is that the relation between stage-wise convergence and convergence in the stochastic game is still unclear.
- **Adaptation:** One criteria is rationality: convergence to a best response equilibrium if the others remain stationary. In this cases Nash equilibrium arises naturally combining convergence and rationality. Another technique is **no-regret** which means that the agent achieves a return that is at least as good as any stationary strategy.



MARL Goal

In a Nutshell

Some other properties can be related to **stability** and **adaptation**.

Opponent-independent algorithms converges to a stable strategies regardless of what the other agents are doing

Opponent-aware algorithms learn models of the other agents and use them in some way of best response and is related to adaptation.

A good combination of both axis should be designed according to each problem. Stability is needed because stable agents are simpler to analyze and have more meaningful performance guarantees. Adaptation is needed because the behaviour of other agents is generally unpredictable.



MARL Algorithms

Classification

The classification of MARL algorithms encompasses several dimensions and in a simplified way can be separated by **task type**.

In this case ($r_1 = \dots = r_n$) and the goal is to maximize the common objective function. This is reduced to a MDP if a centralized controller is available $\pi(a_1, \dots, a_n | s)$. Complications might arise because, multiple joint actions may be optimal. In absence of coordination mechanisms, the agents might pick sub-optimal joint action. Several solutions are proposed:

- **Team Q learning:** avoids coordination assuming that the optimal joint actions are unique (not always true).
- **Distributed Q learning:** solves cooperation without assuming coordination, only works in deterministic problems with non-negative rewards.
- **Coordination based:** when global Q functions can be decomposed additive into local Q functions of subset of agents.

In this case (for $N = 2$ is equivalent to $r_1 = -r_2$) the minimax principle can be applied: maximize your benefit under the worst case assumption that the opponent will always try to minimize it. This principle suggests a opponent-independent algorithm.

- **Minimax Q:** similar to Q learning introducing minimax operator.

$$\mathbf{m}_1(Q, s) = \max_{\pi_1(\cdot|s)} \min_{a_2} \sum_{a_1} \pi^1(a_1|s) Q_1(s, a_1, a_2) \quad (2)$$

The schematic update rule will be:

$$Q_1^{t+1} \leftarrow Q_1^t + \alpha_t (Q_1^t + \gamma \mathbf{m}_1(Q_1^t) - Q_1^t) \quad (3)$$

No constraints imposed on the reward functions of the agents. Self-interested, but not necessarily competing agents. Here the influence of Game Theory is the strongest.

- **Single - agent RL:** Algorithms for single agents like Q-learning can be applied directly. The non-stationarity of the problem invalidates most of the single-agent RL theoretical guarantees, nonetheless, it has a significant number of applications because of its simplicity.
- **Agent - independent methods:** Nash Q-learning modifies the Q-learning update rule to include a convergence to equilibrium (probably Nash) in each state.
- **Agent - tracking methods:** Estimate models of the other agents strategies or policies and act using some form of best response.
- **Agent - aware methods:** They target convergence as well as adaptation to other agents.



Table of Contents

List

- ▶ Introduction to Single Agent
- ▶ Generalization MARL
- ▶ Particularities in MARL
- ▶ **Matrix Games**
- ▶ Numerical Examples
- ▶ Conclusions

Application in a Matrix Game

Natural Policy Gradient

N independent agents with action sets \mathcal{A}_i . The reward functions are $r_n(\mathbf{a})$ and depends on the joint action \mathbf{a} . The policy is:

$$\pi_{\theta^n}^n(a) = \frac{e^{\theta_a^n}}{\sum_{a' \in \mathcal{A}_n} e^{\theta_{a'}^n}} \text{ where } \theta^n = (\theta_a^n)_{a \in \mathcal{A}_n} \text{ is the set of parameters} \quad (4)$$

The goal for each agent is to maximize its expected utility:

$$u_n(\theta) = u_n(\theta^n, \theta^{-n}) = u_n(\pi_{\theta^n}^n, \pi_{\theta^{-n}}^{-n}) = \sum_{a_n, \mathbf{a}_{-n}} r_n(a_n, \mathbf{a}_{-n}) \pi_{\theta^n}^n(a_n) \pi_{\theta^{-n}}^{-n}(\mathbf{a}_{-n}) \quad (5)$$

Application in a Matrix Game

Natural Policy Gradient

The update rule for the parameters will follow the natural gradient:

$\theta_a^n \leftarrow \theta_a^n + \eta \tilde{\nabla}_a u_n(\theta^n)$ where $\tilde{\nabla} = \mathbb{J}^{-1} \nabla$ is the natural gradient and \mathbb{J} is the Fisher matrix of the policy and η is the learning rate ($\tilde{\nabla}_a \equiv \tilde{\nabla}_{\theta_a^n}$)

If we apply the natural gradient to the expected utility we will obtain:

$$\tilde{\nabla}_a u_n(\theta^n) = \sum_{\mathbf{a}_n, \mathbf{a}_{-n}} r_n(\mathbf{a}_n, \mathbf{a}_{-n}) \left[\tilde{\nabla}_a \pi_{\theta^n}^n(a_n) \right] \pi_{\theta^{-n}}^{-n}(\mathbf{a}_{-n})$$

The natural gradient of the policy as $\tilde{\nabla}_a \pi_{\theta^n}^n(a_n) = \delta_{a, a_n}$ and we obtain:

$$\tilde{\nabla}_a u_n(\theta^n) \equiv u^n(a, \pi^{-n}) = \sum_{\mathbf{a}_{-n}} r_n(a, \mathbf{a}_{-n}) \pi_{\theta^{-n}}^{-n}(\mathbf{a}_{-n})$$

The great advantage of the softmax parametrization along with the natural gradient in this case is that now is very easy to construct an estimator of the gradient based on samples: if in a given step, agent n plays action A_n when the other agents play the joint action \mathbf{A}_{-n} , the estimator of the gradient will be:

$\widehat{\tilde{\nabla}}_a u_n = r_n(A_n, \mathbf{A}_{-n})$ and the update rule becomes:

$$\theta_a^n \leftarrow \theta_a^n + \eta r_n(A_n, \mathbf{A}_{-n}) \delta(A_n = a)$$

Application in a Matrix Game

Replicator Dynamics

From a parallel point of view, we can check what is the update rule for the policy $\pi_{\theta^n}^n(a)$ instead of just its parameters.

$$\begin{aligned}
 \pi_{\theta^n}^n(a) &\leftarrow \frac{e^{\theta_a^n + \eta u^n(a, \pi^{-n})}}{\sum_{a' \in \mathcal{A}_n} e^{\theta_{a'}^n + \eta u^n(a', \pi^{-n})}} \approx \frac{e^{\theta_a^n} (1 + \eta u^n(a, \pi^{-n}))}{\sum_{a' \in \mathcal{A}_n} e^{\theta_{a'}^n} (1 + \eta u^n(a', \pi^{-n}))} = \\
 &= \frac{e^{\theta_a^n} (1 + \eta u^n(a, \pi^{-n}))}{\sum_{a' \in \mathcal{A}_n} e^{\theta_{a'}^n} + \eta \sum_{a' \in \mathcal{A}_n} e^{\theta_{a'}^n} u^n(a', \pi^{-n})} = \pi_{\theta^n}^n(a) \frac{1 + \eta u^n(a, \pi^{-n})}{1 + \eta u^n(\pi)} \approx \\
 &\pi_{\theta^n}^n(a) (1 + \eta u^n(a, \pi^{-n}) - \eta u^n(\pi))
 \end{aligned}$$

In summary, the update rule is:

$$\pi_{\theta^n}^n(a) \leftarrow \pi_{\theta^n}^n(a) + \eta \pi_{\theta^n}^n(a) [u^n(a, \pi^{-n}) - u^n(\pi)]$$

which in the continuous limit $\eta \rightarrow 0$ recovers the structure of the replicator dynamics:

$$\frac{d\pi^n(a)}{dt} = \pi^n(a) [u^n(a, \pi^{-n}) - u^n(\pi)]$$

that is a set of equations found in Evolutionary Game Theory that describes the evolution of a population of individuals with different fitnesses.



Table of Contents

List

- ▶ Introduction to Single Agent
- ▶ Generalization MARL
- ▶ Particularities in MARL
- ▶ Matrix Games
- ▶ Numerical Examples
- ▶ Conclusions



Table of Contents

List

- ▶ Introduction to Single Agent
- ▶ Generalization MARL
- ▶ Particularities in MARL
- ▶ Matrix Games
- ▶ Numerical Examples
- ▶ **Conclusions**

- MARL bring more challenges than the single agent case because the returns are correlated and cannot be optimized independently.
- Defining a good MARL goal is not trivial and depends on the task we want to solve.
- Game Theory gives us a powerful mathematical toolbox to study problems in MARL and many algorithms are based on Game Theory concepts.
- The continuous time limit of the Natural Policy Gradient algorithm recovers the Replicator Dynamics with some differences in particular cases due to the local observability of the agents.



Thank you for listening!
Any questions?