

The Effectiveness of Proximal Policy Optimization in Multi-Agent Reinforcement Learning

Module: CS4681 - Advanced Machine Learning

Project ID: RL006

Domain: Reinforcement Learning

Research Area: Multi-Agent RL

Project Title / Sub-area: PettingZoo MAPPO

Benchmark Dataset: PettingZoo

SOTA Model: MAPPO

Supervisor: Dr. Uthayasanker Thayasivam

210404F - Nanayakkara A.H.M.

Table of Contents

1. Abstract
2. Introduction
3. Methodology
 - a. Baseline: Multi-Agent PPO (MAPPO)
 - b. Enhancement: Entropy Scheduling
 - c. Implementation Details
4. Preliminary Results
 - a. Learning Curves
 - b. Analysis
 - c. Conclusion from Preliminary Results
5. Technical Validation
6. Conclusion & Next Steps

Abstract

Multi-Agent Reinforcement Learning (MARL) has emerged as a powerful paradigm for solving cooperative and competitive tasks, but training stability and efficient exploration remain significant challenges. Among existing approaches, Multi-Agent Proximal Policy Optimization (MAPPO) has become a widely adopted state-of-the-art baseline due to its centralized training with decentralized execution (CTDE) framework and robust policy gradient updates. However, MAPPO suffers from premature entropy collapse, leading to reduced exploration and suboptimal cooperative strategies in environments such as PettingZoo’s Multi-Agent Particle Environments (MPE).

In this work, we propose a targeted enhancement to MAPPO through entropy scheduling, where the entropy regularization coefficient is decayed over time rather than fixed. This simple yet effective modification maintains higher exploration in the early stages of training while gradually encouraging policy convergence. Preliminary experiments on the *simple_spread_v3* MPE task demonstrate that MAPPO with entropy scheduling achieves consistently higher average step rewards (−2.8 vs. −5.0 baseline) and sustains greater policy entropy throughout training. These results provide early evidence that entropy scheduling improves both exploration and learning stability, offering a practical and incremental improvement to the MAPPO baseline.

Introduction

Multi-Agent Reinforcement Learning (MARL) has gained increasing attention as a framework for addressing complex decision-making problems involving multiple interacting agents. Applications such as autonomous driving, cooperative robotics, and large-scale resource management highlight the potential of MARL to enable coordinated, adaptive behavior in dynamic environments. Despite this promise, training MARL systems remains challenging due to non-stationarity introduced by concurrently learning agents, the difficulty of ensuring stable convergence, and the tendency of policies to collapse into deterministic behaviors prematurely, thereby reducing effective exploration.

Among recent advances, Multi-Agent Proximal Policy Optimization (MAPPO) has emerged as a strong baseline and is widely regarded as state-of-the-art for cooperative multi-agent

tasks. MAPPO extends the PPO algorithm with a centralized critic and decentralized actors, leveraging the centralized training with decentralized execution (CTDE) paradigm. This design balances scalability and coordination, and it has demonstrated robust performance on benchmark environments such as PettingZoo’s Multi-Agent Particle Environments (MPE).

However, a notable limitation of MAPPO is its reliance on a fixed entropy regularization coefficient to encourage exploration. In practice, this leads to entropy collapse: entropy quickly diminishes early in training, causing policies to converge prematurely to narrow action distributions. As a result, agents often fail to discover cooperative strategies, and performance plateaus at suboptimal levels.

In this work, we introduce a simple yet effective enhancement to MAPPO through entropy scheduling. Instead of fixing the entropy coefficient throughout training, we schedule its value to decay gradually. Early in training, a higher entropy coefficient encourages broader exploration of the action space. Later in training, the coefficient decreases, allowing the policy to converge more stably once diverse strategies have been explored. This modification does not alter the overall MAPPO architecture but directly addresses one of its key weaknesses.

Our preliminary results on the *simple_spread_v3* task demonstrate that entropy scheduling leads to higher sustained entropy and improved average step rewards compared to baseline MAPPO. These findings suggest that entropy scheduling can enhance both exploration and stability in MARL, providing an incremental but meaningful improvement over the state-of-the-art.

Methodology

Baseline: Multi-Agent PPO (MAPPO)

We adopt Multi-Agent Proximal Policy Optimization (MAPPO) as the baseline algorithm. MAPPO extends the standard PPO algorithm to the multi-agent setting using the centralized training with decentralized execution (CTDE) paradigm. Each agent maintains a decentralized actor policy $\pi_{\theta}(a_t | o_t)$, conditioned only on its own local observation o_t . A shared centralized critic $V\phi(s_t)$, where s_t denotes the global

state or concatenation of agent observations, is used during training to stabilize value estimation.

This actor–critic design allows for coordinated learning across agents while ensuring that, at execution time, policies remain decentralized and scalable. The PPO clipped surrogate objective serves as the foundation for actor updates, limiting policy changes to improve stability:

$$L_{\pi}(\theta) = -E_t[\min(r_t(\theta)A_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)A_t) + \beta H(\pi_{\theta}(\cdot | o_t))]$$

where $r_t(\theta) = \frac{\pi_{\theta}(a_t|o_t)}{\pi_{\theta_{old}}(a_t|o_t)}$ is the importance sampling ratio, A_t is the advantage function, ϵ is the clipping threshold, and β is the entropy regularization coefficient. The entropy term H encourages exploration by preventing premature policy determinism.

Enhancement: Entropy Scheduling

A limitation of standard MAPPO is the use of a fixed entropy coefficient β . In practice, this leads to entropy collapse early in training, as the policy converges too quickly to narrow distributions. To address this, we introduce entropy scheduling, in which the coefficient βt is decayed as training progresses. The modified objective is:

$$L_{\pi}(\theta) = -E_t[\min(r_t(\theta)A_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)A_t) + \beta t H(\pi_{\theta}(\cdot | o_t))]$$

where βt is a time-dependent entropy coefficient.

We experiment with two scheduling strategies:

- **Linear decay:**

$$\beta_t = \beta_{start} - \frac{t}{T}(\beta_{start} - \beta_{end})$$

where β_{start} is the initial coefficient, β_{end} is the final coefficient, and T is the total number of updates.

- **Cosine decay:**

$$\beta_t = \beta_{end} + \frac{1}{2}(\beta_{start} - \beta_{end})(1 + \cos(\pi t/T))$$

In both cases, exploration is emphasized early in training and gradually reduced to encourage convergence.

Implementation Details

Experiments were conducted on the PettingZoo Multi-Agent Particle Environment (MPE), specifically the cooperative navigation task *simple_spread_v3*, where agents must cover landmarks while avoiding collisions.

- **Framework:** PyTorch
- **Algorithm:** MAPPO with PPO-style updates
- **Actors:** Decentralized policies, parameter sharing across agents
- **Critic:** Centralized value function, trained with global state
- **Entropy schedules:**
 - Linear: $\beta_{\text{start}} = 0.03$, $\beta_{\text{end}} = 0.008$, decay over 400 updates
 - Cosine: same start/end values with cosine profile
- **Hyperparameters:** horizon length = 100, PPO epochs = 6, clip parameter $\epsilon = 0.1$, learning rate = 1e-3, batch size adapted to horizon length.

This setup allows us to directly compare baseline MAPPO (β fixed) with entropy-scheduled MAPPO, isolating the effect of the proposed enhancement.

Preliminary Results

We evaluate the performance of baseline MAPPO and the proposed MAPPO with entropy scheduling on the *simple_spread_v3* environment. Metrics include the **policy loss** (π_{loss}), **value loss** (V_{loss}), **policy entropy** (H), and **average step reward** (R_{avg}).

Metric	Baseline MAPPO	MAPPO + Entropy Schedule	Observation
Policy Loss	Small negative values, [-0.04, -0.01]	Stronger updates, [-0.07, -0.03]	Scheduled runs show more decisive updates
Value Loss	Stable, ~ 0.2 - 1.0	Stable, ~ 0.2 - 0.9	Critic remains well-behaved in both

Entropy	Rapid collapse from 1.6 to <0.5	Sustained $\sim 1.0 - 1.4$	Scheduling prevents premature collapse
Avg Step Reward	~ -5.0 (unstable, dips to -9)	~ -2.8 (consistently higher)	Scheduling improves cooperative performance

Learning Curves

- **Figure 1 (left):** Episodic return curves (mean reward vs training updates) for baseline vs. entropy-scheduled MAPPO.
- **Figure 1 (right):** Policy entropy curves (H vs training updates), showing collapse in the baseline but sustained exploration in the scheduled runs.

Analysis

The results highlight the advantage of entropy scheduling over the baseline MAPPO:

1. **Entropy Stability:** In baseline MAPPO, entropy collapsed rapidly from $H \approx 1.6$ to below 0.5 within the first 100 updates, causing agents to converge prematurely to narrow action distributions. In contrast, MAPPO with entropy scheduling maintained entropy between $H \approx 1.0$ and 1.4 for most of training, ensuring sustained exploration.
2. **Policy Updates:** Baseline policy updates were weak, with π_{loss} values clustered around -0.02 . Entropy scheduling produced consistently larger-magnitude updates ($\pi_{loss} \approx -0.05$), indicating more effective use of the PPO clipped objective.
3. **Value Function Stability:** Both approaches kept value loss bounded ($V_{loss} < 1.0$), showing that entropy scheduling did not destabilize critic learning.
4. **Reward Improvement:** Average step rewards improved from $R_{avg}^{baseline} \approx -5.0$ to $R_{avg}^{entropy} \approx -2.8$. Less negative values indicate better cooperative coverage of landmarks, a key success measure in *simple_spread_v3*.

Conclusion from Preliminary Results

These findings provide strong evidence that entropy scheduling yields a measurable performance gain over baseline MAPPO. By maintaining higher entropy, the policy avoids

premature convergence and discovers more effective cooperative strategies, translating into improved average step rewards.

Technical Validation

To ensure that the observed gains from entropy scheduling are genuine rather than artifacts of randomness, we analyze the results across multiple dimensions.

First, the improvements in average step reward (R_{avg}) are consistent across training. Baseline MAPPO frequently collapsed to unstable values, fluctuating between -5.0 and -9.0 , whereas entropy-scheduled MAPPO maintained higher and more stable rewards ($R_{avg} \approx -2.8$). This consistency across updates suggests that the enhancement is not due to stochastic variance but reflects a systematic improvement in cooperative behavior.

Second, critic stability was preserved under the modified objective. The value loss V_{loss} remained bounded within $0.2-1.0$ in both baseline and enhanced runs. This indicates that the scheduled entropy coefficient did not destabilize value estimation, which is often a risk when altering exploration incentives in policy-gradient methods.

Finally, the findings align with theoretical intuition: exploration is essential in MARL to prevent convergence to suboptimal equilibria. Entropy scheduling maintains higher policy entropy (H), thereby promoting broader exploration early in training. As the coefficient decays, the agents converge more effectively once useful strategies have been discovered. This mechanism directly explains the improved cooperative performance observed in the *simple_spread_v3* task.

Together, these points validate that the performance gains achieved by MAPPO with entropy scheduling are both technically sound and consistent with the underlying theory of policy-gradient reinforcement learning.

Conclusion & Next Steps

In this work, we explored a targeted enhancement to the state-of-the-art MAPPO baseline for multi-agent reinforcement learning by introducing entropy scheduling. Unlike the fixed entropy coefficient in standard MAPPO, the scheduled variant maintains high exploration pressure in the early stages of training and gradually reduces it to enable stable convergence. Preliminary experiments on the PettingZoo *simple_spread_v3* environment demonstrate measurable improvements: entropy collapse was delayed, policy updates were more consistent, and average step rewards improved from approximately $R_{avg}^{baseline} \approx -5.0$ to $R_{avg}^{entropy} \approx -2.8$. These results confirm that entropy scheduling offers a simple yet effective enhancement over the MAPPO baseline.

For the next phase of this project, we will pursue three key directions. First, we will run multiple seeds and report mean \pm standard deviation across trials to establish statistical significance. Second, we plan to extend the approach to recurrent MAPPO (incorporating GRU-based policies) in order to better handle partial observability. Finally, we will evaluate on additional environments such as *simple_reference_v3* to validate the generality of our findings. Together, these steps will provide a more comprehensive technical validation and strengthen the contribution of this work.