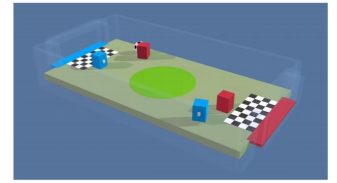




Cooperative-Competitive Multi-Agent Learning in Soccer Environments

Christopher Covert, Cameron McMillan, Patipan Pipatpinyong
[cwcovert, cmac12, pipat001]@stanford.edu



ABSTRACT

Inspired by multi-agent cooperation, we are investigating the use of multi-agent deep reinforcement learning networks to play simple cooperative games. This project utilizes a simulated soccer environment developed by Unity to test a suite of deep reinforcement learning algorithms against each other in a simple game benchmark. In the 2 vs 2 player format, where each team is equipped with a striker (offensive agent) and goalie (defensive agent), we explore how agents can learn to be simultaneously collaborative and competitive. To avoid the influence of consecutive examples' strong correlations and interdependence, a randomized experience replay scheme was used to train the network. We explored the training result between Deep Q-Networks (DQN) and Proximal Policy Optimization (PPO), each with small variations during training.

RESULTS

Table 1 shows the results on 500-run tests.

Model	Win (W) and Not Lose (NL) % Against:						
	Random	DQN _{base}	DQN ₁	DQN ₁ ^{mod}	DQN ₂	DQN ₂ ^{mod}	PPO _{base}
DQN _{base}	W: 13% NL: 81%	W: 0% NL: 100%	W: 18% NL: 32%	W: 14% NL: 41%	W: 5% NL: 38%	W: 10% NL: 59%	W: 21% NL: 55%
DQN ₁	W: 66% NL: 70%	W: 68% NL: 82%	W: 31% NL: 62%	W: 23% NL: 49%	W: 28% NL: 48%	W: 37% NL: 57%	W: 63% NL: 68%
DQN ₁ ^{mod}	W: 79% NL: 85%	W: 59% NL: 86%	W: 51% NL: 77%	W: 40% NL: 63%	W: 41% NL: 63%	W: 45% NL: 75%	W: 77% NL: 81%
DQN ₂	W: 77% NL: 81%	W: 62% NL: 95%	W: 52% NL: 72%	W: 37% NL: 59%	W: 41% NL: 62%	W: 48% NL: 76%	W: 70% NL: 77%
DQN ₂ ^{mod}	W: 53% NL: 74%	W: 41% NL: 90%	W: 43% NL: 63%	W: 25% NL: 55%	W: 24% NL: 52%	W: 39% NL: 59%	W: 49% NL: 71%
PPO _{base}	W: 53% NL: 70%	W: 45% NL: 79%	W: 32% NL: 37%	W: 19% NL: 23%	W: 23% NL: 30%	W: 29% NL: 51%	W: 46% NL: 58%

Table 1: Evaluation Results for the 2x2 Soccer Environment

DATASET

The dataset for training this model utilizes an experience replay method, meaning that the input set is generated at the time of training and added to a queue. For our application, the replay is a state-action set that is randomly selected from the queue to decouple the influence of consecutive runs.

FEATURES

The striker has 6 possible actions: forward, backward, left, right, clockwise rotation, and counterclockwise rotation. The keeper has only four possible actions: forward, backward, left, and right. As for the state space, each agent obtains a 336 element observation vector, that corresponds to 7 rays within the 180 degree view from the front of the agent, with 6 sets of rays stacked in the vertical direction. Each ray consists of 8 values: binary classification of seven possible object types and the object's distance from the agent.

MODELS

A Proximal Policy Optimization (PPO) model was used to provide a baseline for a series of reward-modified Deep Q-Network (DQN) tests. Both models utilize the basic reward structure of Table 2.

Proximal Policy Optimization (PPO)

The PPO baseline was implemented with the following loss function:

$$L_{\pi}(\theta) = r_{\pi}(\theta) \hat{A}_{\pi} - \beta \text{KL}[\pi_{\theta, \text{old}}, \pi_{\theta}]$$

Deep Q-Network (DQN)

The DQN models were implemented with the following target assignment:

$$y_j = \begin{cases} r_j & \text{if episode terminates} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta^-) & \end{cases}$$

DQN₁, the striker receives the original reward outlined earlier where scoring a goal rewards +1 and -0.1 penalty, corresponding to a "neutral" striker that focuses on scoring and may play some defense.

DQN₂ corresponds to an "aggressive" striker that receives reward of +1 for scoring and 0 penalty if the team gets scored on.

DQN₃ corresponds to a "defensive" striker that gets reward of +1 for scoring and penalty of -0.6 if it gets scored on.

The mod tag represents models that have no penalty for not observing the ball.

Agent	Event		
	Ball Enters Opponent's Goal	Ball Enters Own Team's Goal	Existing
Striker	+1	-0.1	-0.001
Goalie	+0.1	-1	+0.001

Table 2: Basic Agent Reward Structure

DISCUSSION

The PPO model initially performed better against random and trained opponents. Intermediate rewards were then introduced to help with training DQN models and made a significant improvement. One interesting result is that the penalty for the striker not seeing the ball doesn't seem to have much effect on the model's performance. We were able to see that the more aggressive strikers performed better in this game overall. We also found that introducing intermediate reward actually hurt the performance when training PPO models.

FUTURE

Future work includes further fine tuning of the reward structure. We would also like to expand the problem to include more agents, less environmental boundaries, more complex agent roles, and have different models trained against each other. Lastly, there are other deep reinforcement learning methods to be explored such as double DQN, dueling DQN, and imitation learning.

REFERENCES

[1] Google, Google, David Rosenberg, and Corina Truoc. "Self-supervised training controllers for autonomous control systems using deep reinforcement learning." arXiv preprint arXiv:1902.02821, 2019. <https://arxiv.org/abs/1902.02821>

[2] Denis Sjoqvist, Maxim Agreus, and Rafael Aichewitz. "Cooperative multi-agent control using deep reinforcement learning." International Conference on Autonomous Agents and Multiagent Systems, Springer, 2017. https://doi.org/10.1007/978-3-319-62222-2_10

[3] Iqbal, Ahsan. "Multi-agent deep reinforcement learning." 2019. <https://arxiv.org/abs/1902.02821>

[4] Hester, Timothy, et al. "Learning multi-agent deep reinforcement learning." Proceedings of the 17th International Conference on Autonomous Agents and Multiagent Systems, International Foundation for Autonomous Agents and Multiagent Systems, 2016. <https://arxiv.org/abs/1609.02984>

[5] Reinforcement learning systems laboratory [2015]. "Control repository." <https://github.com/rl-solutions>

[6] Mita, Y., Kawahara, K., Sato, O., Hagi, A., Iwama, S., Bekker, M., et al. and Petersen, S. (2015). "Human-level control through deep reinforcement learning." 2015: 1200-1205.

[7] Hester, T., Sutton, R., and Elmer, J. (2016). "Mastering deep reinforcement learning with double q-learning in StarCraft." Conference on Artificial Intelligence. <https://arxiv.org/abs/1603.04712>

[8] Hester, T., Sutton, R., and Elmer, J. (2016). "Mastering deep reinforcement learning with double q-learning in StarCraft." Conference on Artificial Intelligence. <https://arxiv.org/abs/1603.04712>

[9] Hester, T., Sutton, R., Elmer, J., and Elmer, J. (2015). "Proximal policy optimization algorithms." arXiv preprint arXiv:1509.02984. <https://arxiv.org/abs/1509.02984>

[10] Schaul, T., Leibo, J., Ghahramani, Z., and Silver, D. (2015). "Trust Region Policy Optimization in Atari." In: pp. 188-197. <https://arxiv.org/abs/1509.02984>

[11] Sutton, R., Barto, A., Silver, D., Zhou, Y., Powell, M., and Ghahramani, Z. (2018). "Using a recurrent framework for intelligent agents." arXiv preprint <https://arxiv.org/abs/1808.07312>

[12] Vucelja, M. (2016). "Reinforcement learning." <https://github.com/rl-solutions>