

Adaptive Learning Rates for Multi-Agent Reinforcement Learning

Professor: Li-Der Chou

Presenter: Chih-Jou Tai

Conference on Autonomous Agents and Multiagent Systems
(AAMAS 2023), May 29 - June 2, 2023, London, UK

Outline



- Introduction
- Related Work
- Methodology
- Experiments
- Conclusion

Introduction



- MARL(Multi-Agent Reinforcement Learning)
- Adaptive Learning Rates: Introduction of adaptive learning rates to cooperative multi-agent reinforcement learning (MARL).
- AdaMa Method:
 - Dynamic Balancing Learning Rates
 - Second-Order Approximation
 - Performance Improvement & Cost Reduction

Outline



- Introduction
- Related Work
- Methodology
- Experiments
- Conclusion

Related Work

	AdaGrad	AdaDelta	RMSprop	WoLF	AdaMa
Learning Rate Adjustment	Accumulated gradients	Moving average of gradients	With decay factor introduced	Adjusts based on equilibrium strategy in games	Dynamically adjusts towards maximally improving Q-values
Training Process	Manual preset learning rate	Less manual presetting required	Less manual presetting required	Requires solving for equilibrium strategy	Adaptive adjustment reduces the need
Application Domain	General optimization problems	General optimization problems	General optimization problems	Stochastic games	Multi-agent reinforcement learning

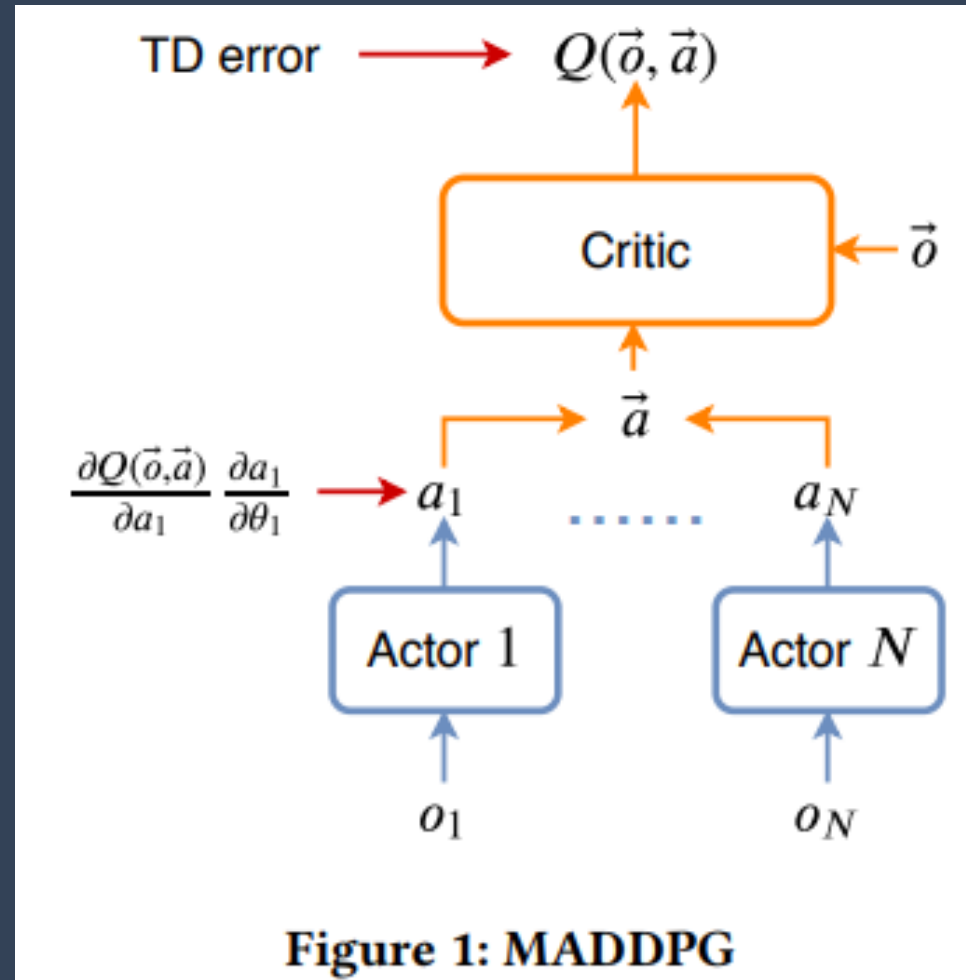
Outline



- Introduction
- Related Work
- Methodology
- Experiments
- Conclusion

Methodology

- Single-Critic MADDPG



Methodology

- Actor i learning rate: l_{a_i}
- Critic learning rate: l_c

- Adaptive \vec{l}_a Direction

$$\vec{l}_a = \alpha \vec{l}_a + (1 - \alpha) \eta \frac{\partial Q}{\partial \theta} \frac{\partial Q}{\partial \theta}^T / \left\| \frac{\partial Q}{\partial \theta} \frac{\partial Q}{\partial \theta}^T \right\| \quad (1)$$
$$\vec{l}_a = \vec{l}_a \frac{\eta}{\|\vec{l}_a\|},$$

- Adaptive l_c and $\|\vec{l}_a\|$

$$l_c = \alpha l_c + (1 - \alpha) l \cdot \text{clip}\left(\left| \frac{\partial \delta}{\partial \phi} \frac{\partial Q}{\partial \phi}^T \right| / m, \epsilon, 1 - \epsilon\right) \quad (2)$$
$$\eta = l - l_c.$$

- First-Order Approximation:
 - The actor i 's contribution to ΔQ is only related to the change of action i , without capturing the joint effect with other agents' updates.
 - When agents are strongly correlated, summing up individual updates from each actor doesn't adequately capture the increase in Q-value.
- Second-Order Approximation:
 - The second-order approximation allows the model to account for the interactions between agents' behaviors, enhancing prediction accuracy .

● Algorithm: AdaMa on MADDPG

Algorithm 1 AdaMa on MADDPG

- 1: Initialize critic network ϕ , actor networks θ_i , target networks, and the replay buffer \mathcal{D} .
- 2: Initialize the learning rates l_c and \vec{l}_a .
- 3: **for** episode = 1, ..., \mathcal{M} **do**
- 4: **for** $t = 1, \dots, \mathcal{T}$ **do**
- 5: Select action $a_t^i = \pi_i(o_t^i) + \mathcal{N}_t^i$ for each agent i
- 6: Execute action a_t^i , obtain reward r_t , and get new observation o_{t+1}^i for each agent i
- 7: Store transition $(\vec{o}_t, \vec{a}_t, r_t, \vec{o}_{t+1})$ in \mathcal{D}
- 8: **end for**
- 9: Sample a random minibatch of transitions from \mathcal{D}
- 10: Adjust l_c and $\|\vec{l}_a\|$ by (2).
- 11: Adjust \vec{l}_a by (1) (first order) or (3) (second order).
- 12: Update the critic ϕ by $\phi = \phi - l_c \frac{\partial \delta}{\partial \phi}$.
- 13: Update the actor θ_i by $\theta_i = \theta_i + l_{a_i} \frac{\partial Q(\vec{o}, \vec{a})}{\partial a_i} \frac{\partial a_i}{\partial \theta_i}$ for each agent.
- 14: Update the target networks.
- 15: **end for**

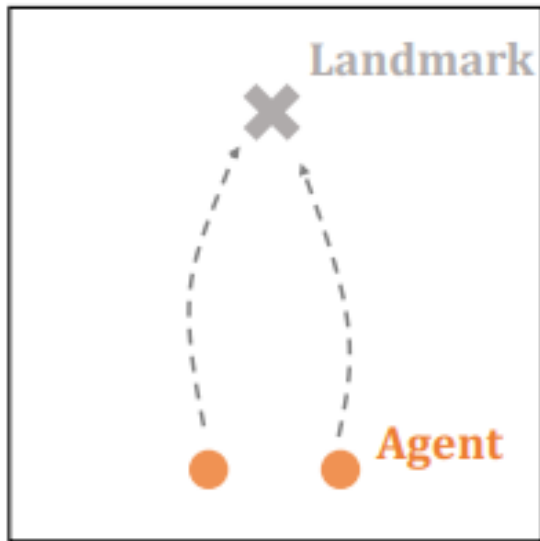
Outline



- Introduction
- Related Work
- Methodology
- Experiments
- Conclusion

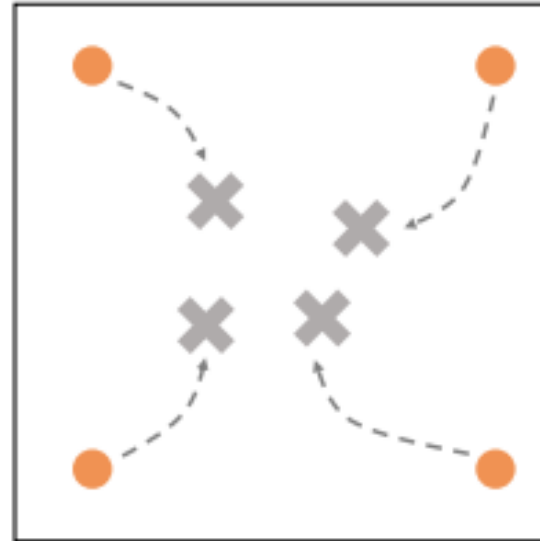
Experiments

- Four cooperative scenarios based on Multi-Agent Particle Environment (MPE)



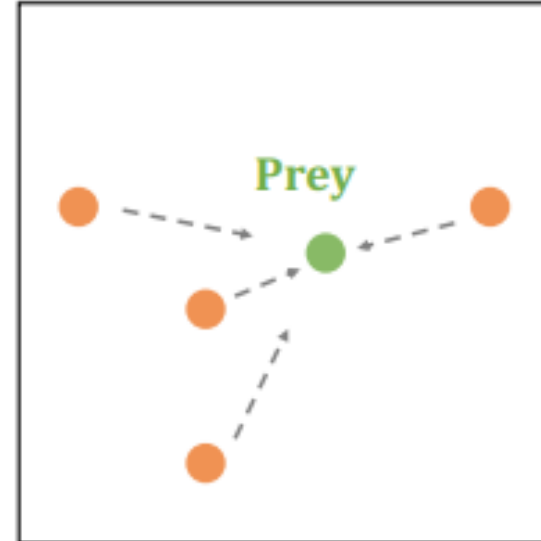
(a) going together

$$-0.5(d_i + d_j) - d_{ij}$$



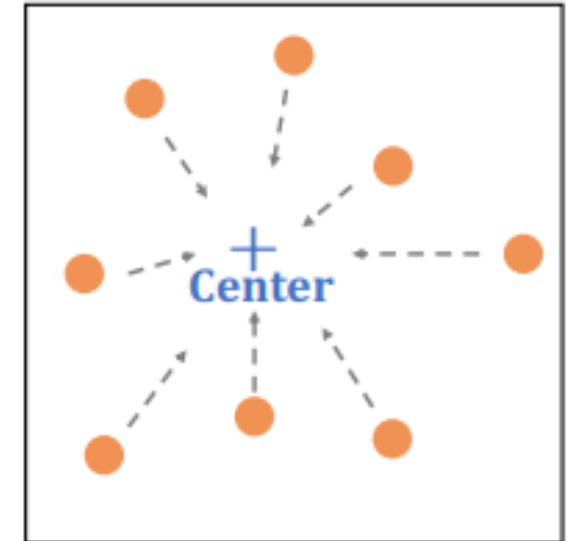
(b) cooperative navigation

$$-\max_i(d_i)$$



(c) predator-prey

$$+1 \text{ when touched}$$

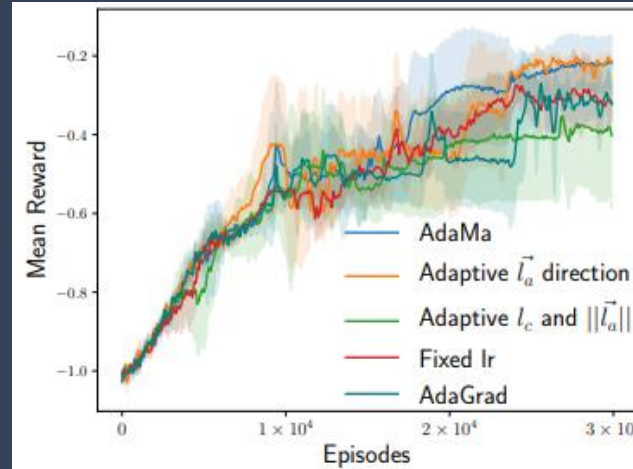


(d) clustering

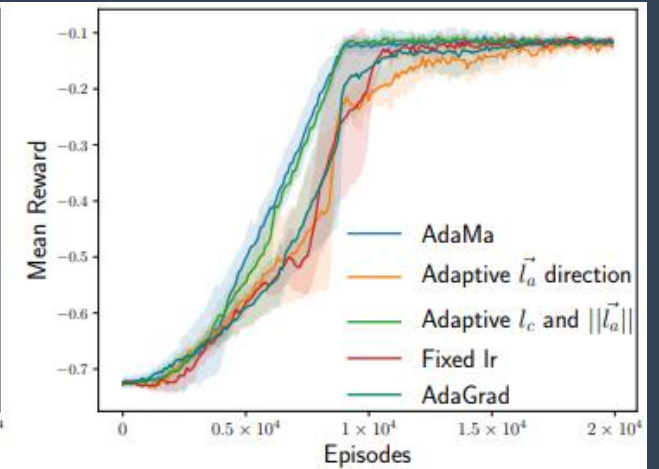
$$-\sum d_i$$

Experiments

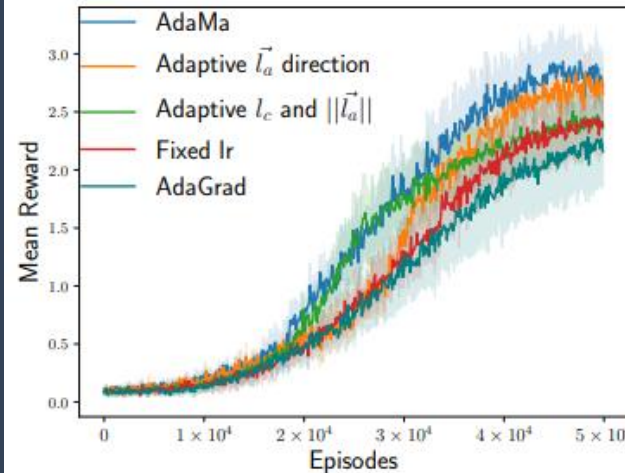
- Learning curves in the four scenarios.



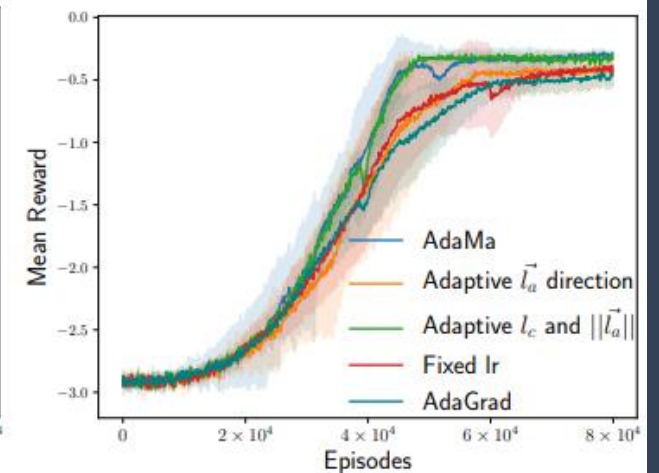
(a) going together



(b) cooperative navigation



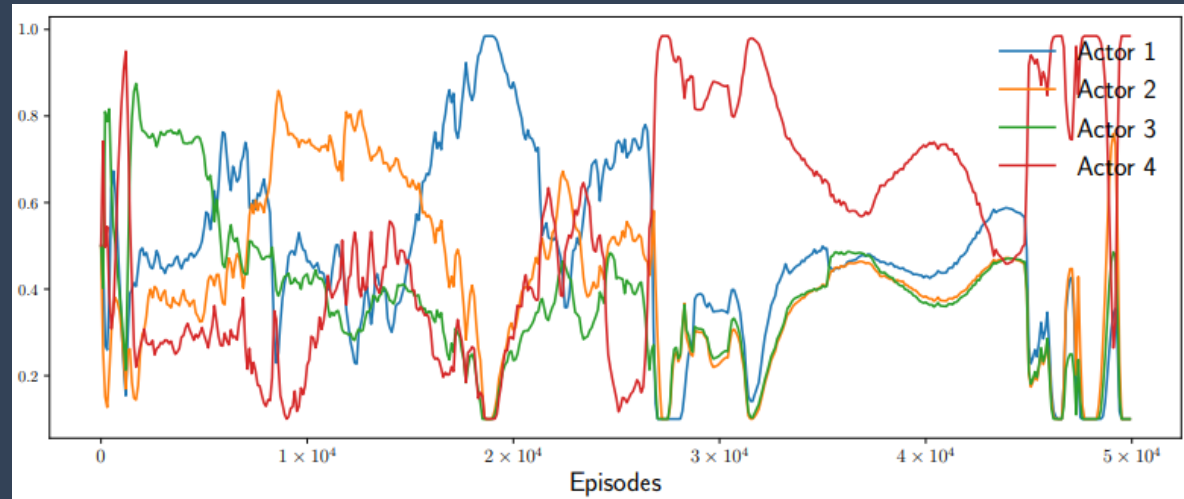
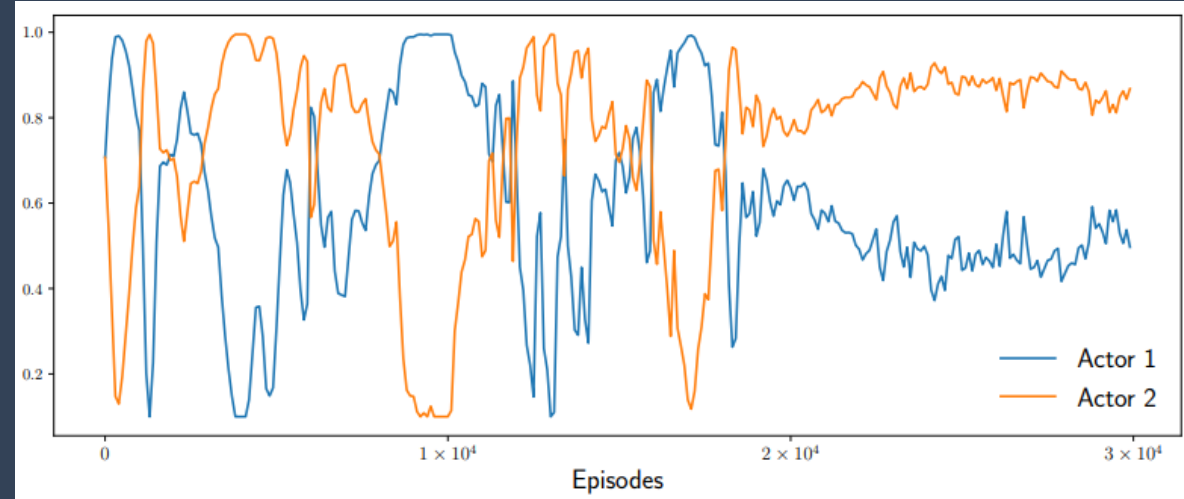
(c) predator-prey



(d) clustering

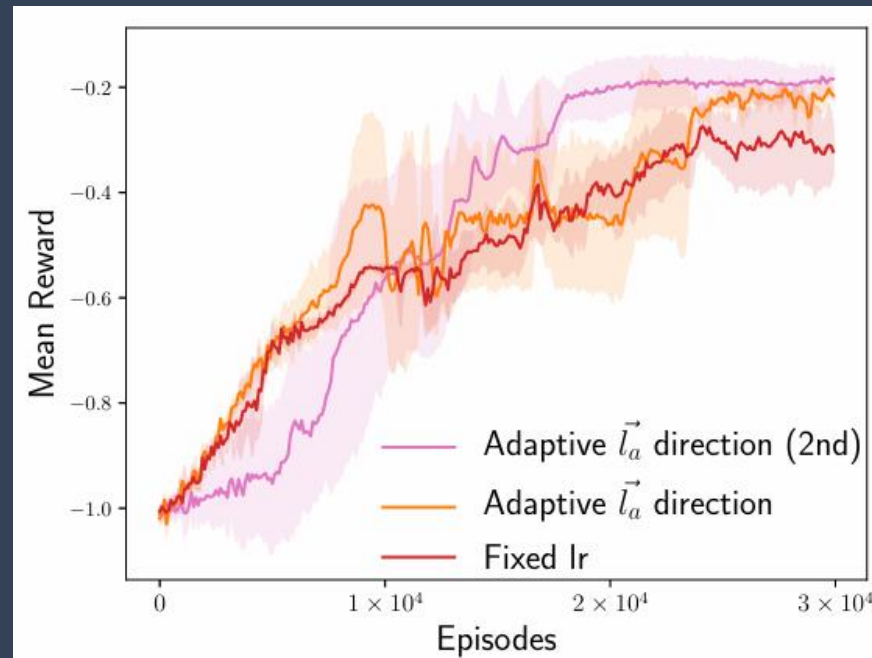
Experiments

- Normalized actors' learning rates during the training in going together.
- Normalized actors' learning rates during the training in predator-prey.

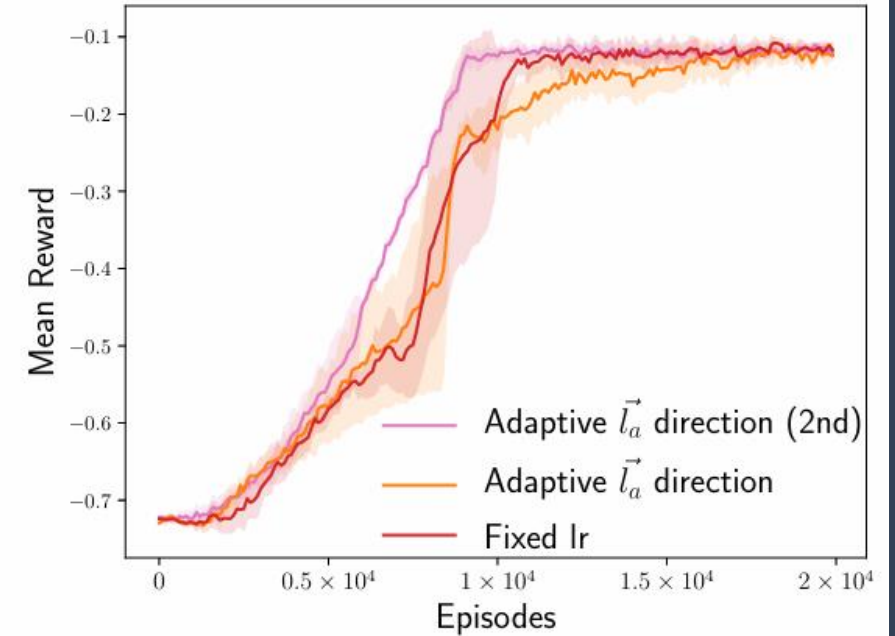


Experiments

- Learning curves with the second-order approximation.



(a) going together

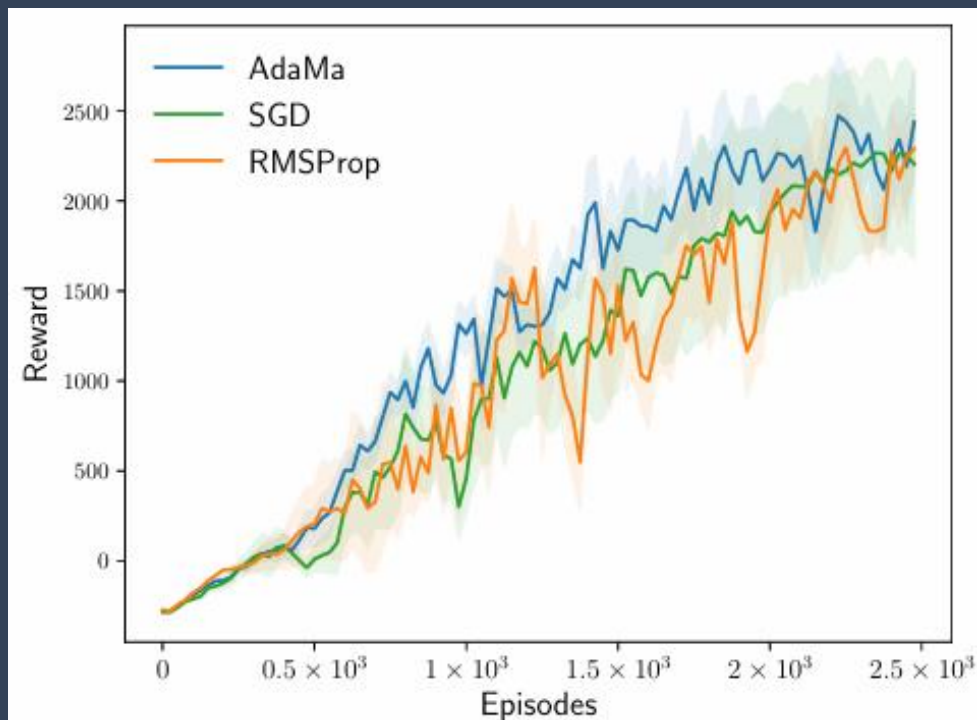


(b) cooperative navigation

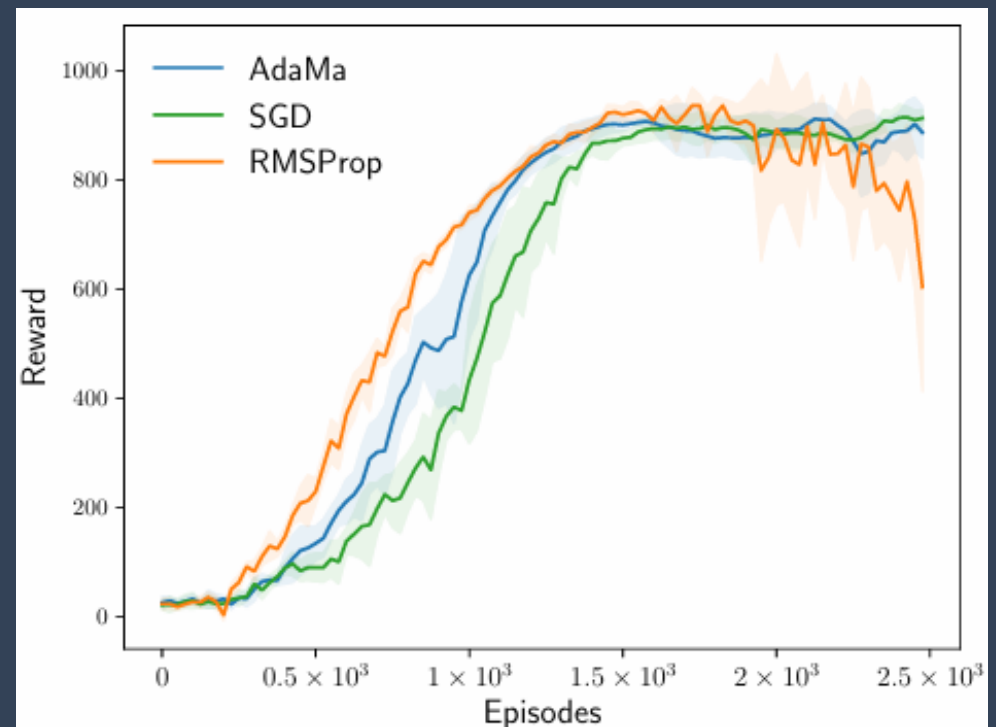
Experiments

- Learning curves of AdaMa on multi-agent mujoco.

HalfCheetah-MADDPG



Ant-MADDPG



Outline



- Introduction
- Related Work
- Methodology
- Experiments
- Conclusion

Conclusion



- Pros:

- AdaMa adaptively update the learning rate in multi-agent environments, effectively accelerating learning.
- AdaMa can be applied to various multi-agent scenarios with a single critic.

- Cons:

- In some environments, using AdaMa may result in less effective learning in the early stages compared to other methods with adaptive learning rates.

Conclusion



- Future work:
 - The current work on AdaMa does not cover environments with multiple critics in multi-agent systems, therefore, AdaMa can be modified in the future to suit environments with multiple critics in multi-agent systems.

Q & A