# Adaptive Learning Rates for Multi-Agent Reinforcement Learning

Colin Acker

October 2023

## 1 Notes

For MADDPG, Deep Deterministic Policy Gradient extended to Multi-Agent systems, the agents share a single reward. Each agent learns an actor $\pi_i$ and a single shared critic. The critic takes $\vec{a} = (a_1, a_2...a_n)$ and $\vec{o}$, the observation vector, as input. The output of the critic is the Q-value. The critic, parameterized by $\phi$, is trained by minimizing the TD-error $\delta$

$$\mathbb{E}_{(\vec{o}, \vec{a}, r, \vec{o'})}[(Q(\vec{o}, \vec{a}) - y)^2]$$

$$\text{Where } y = r + \gamma Q^-(\vec{o'}, \pi_i^-(o_i'))$$

$Q^-$ is the target critic, $\pi_i^-$ is the target actor, and each $(\vec{o}, \vec{a}, r, \vec{o'})$ is from the replay buffer. The gradient of $Q$ with respect to $\theta_i$ is

$$\frac{\partial Q(\vec{o}, \vec{a})}{\partial a_i} \frac{\partial a_i}{\partial \theta_i}$$

The learning rates of each actor critic and the critic will be $l_{a_i}$ and $l_c$. If we leave the critic alone, we can approximate the gain in the Q-value contributed solely by the actors, using Taylor approximation.

$$\Delta Q = Q(\vec{o}, \vec{a} + \Delta\vec{a}) - Q(\vec{o}, \vec{a})$$

After some math, we get,

$$\Delta Q \approx \vec{l_a} \cdot \frac{\partial \vec{Q}}{\partial \theta} \frac{\partial Q^T}{\partial \theta}$$

So, the largest $\Delta Q$ is obtained when the direction of $\vec{l_a}$ is the same direction as $\frac{\partial \vec{Q}}{\partial \theta} \frac{\partial Q^T}{\partial \theta}$. This is basically gradient ascent on $\vec{l_a}$ to increase the Q-value. This only works when the critic is static, but in MADDPG, the critic and actor are trained simultaneously. So now,

$$\Delta Q = Q(\phi + \Delta\phi, \vec{o}, \vec{a} + \Delta\vec{a}) - Q(\phi, \vec{o}, \vec{a})$$

After some math, we get,

$$\Delta Q \approx \vec{l}_a \cdot \frac{\partial \vec{Q}}{\partial \theta} \frac{\partial Q^T}{\partial \theta} - l_c \frac{\partial \delta}{\partial \phi} \frac{\partial Q^T}{\partial \phi}$$

When the critic's update causes a large change in the Q-value, the critic is leading the learning. This is because the actor's learning is determined by the critic. When the critic is leading the learning, we should assign it a high learning rate. The actors struggle with a quickly changing critic, so we should assign the actors a low learning rate. On the other hand, when the critic has hit a plateau, we should increase the actor's learning rate. The learning rate for the critic, $l_c$ is found by updating $l_c$ using a clip function and many different hyperparameters.

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

$$\mu = l - l_c$$

I believe this is the most important result of the paper. The paper goes into details on how to tune these hyperparameters. In the paper, they called the algorithm AdaMa, which updates $l_c$ and $l_{a_i}$ using the equations above. They tested it on 4 different scenarios, all which I think would be good to test on. The key find during testing was adaptively updating $l_c$ and $l_{a_i}$, works faster than a fixed LR.

## 1.1 Algorithms

---
**Algorithm 1** AdaMa on MADDPG

---
1: Initialize critic network $\phi$, actor networks $\theta_i$, target networks, and the replay buffer $\mathcal{D}$.
2: Initialize the learning rates $l_c$ and $\vec{l_a}$.
3: **for** episode $= 1, \ldots, M$ **do**
4:     **for** $t = 1, \ldots, T$ **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 $\phi$ by $\phi = \phi - l_c \frac{\partial \delta}{\partial \phi}$.
13:     Update the actor $\theta_i$ 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**

---