

Notes: Asynchronous Methods for Deep Reinforcement Learning[1]

esgl Hu

March 2, 2018

1 Review

The simple online RL algorithms with deep neural networks was fundamentally unstable.

The sequence of observed data encountered by an online RL agent is non-stationary, and online RL updates are strongly correlated.

The experience replay has several drawbacks: it uses more memory and computation per real interaction; and it requires off-policy learning algorithms that can update from data generated by an older policy.

The loss function of one-step Q-learning is

$$L_i(\theta_i) = \mathbb{E}(r + \gamma \cdot \max_{a'} Q(s', a'; \theta_{i-1}) - Q(s, a; \theta_i))^2 \quad (1)$$

where s' is the state encountered after state s . One drawback of using one-step method is that obtaining a reward r only directly affects the value of the state action pair s, a that led to the reward. The values of other state action pairs are affected only indirectly through the updated value $Q(s, a)$. This can make the learning process slow since many updates are required to propagate a reward to the relevant preceding states and actions. One way of propagating rewards faster is by using n -step returns. In n -step Q-learning, $Q(s, a)$ is updated toward the n -step return defined as $r_t + \gamma r_{t+1} + \dots + \gamma^{n-1} r_{t+n-1} + \max_a \gamma^n Q(s_{t+n}, a)$.

Policy-based model-free methods directly parameterize the policy $\pi(a|s; \theta)$ and update the parameters θ by performing, typically approximate, gradient ascent on $\mathbb{E}[R_t]$. For example, standard **REINFORCE** updates the policy parameters θ in the direction $\nabla_{\theta} \log \pi(a_t|s_t; \theta) R_t$, which is an unbiased estimate of $\nabla_{\theta} \mathbb{E}[R_t]$. It is possible to reduce the variance of this estimate while keeping it unbiased by subtracting a learned function of the state $b_t(s_t)$, known as a baseline from the return. The resulting gradient is $\nabla_{\theta} \log \pi(a_t|s_t; \theta) (R_t - b_t(s_t))$.

A learned estimate of the value function is commonly used as the baseline $b_t(s_t) \approx V^{\pi}(s_t)$ leading to a much lower variance estimate of the policy gradient. When an approximate value function is used as the baseline, the quantity $R_t - b_t$ used to scale the policy gradient can be seen as an estimate of the *advantage* of the action a_t in state s_t , or $A(a_t, s_t) = Q(a_t, s_t) - V(s_t)$

2 Asynchronous one-step Q-learning

Algorithm 1 Asynchronous one-step Q-learning - pseudocode for each actor-learner thread.

```
1: // Assume global shared  $\theta, \theta^-$ , and counter  $T = 0$ .
2: Initialize thread step counter  $t \leftarrow 0$ 
3: Initialize target network weights  $\theta^- \leftarrow \theta$ 
4: Initialize network gradients  $d\theta \leftarrow 0$ 
5: Get initial state  $s$ 
6: repeat
7:   Take action  $a$  with  $\epsilon$ -greedy policy based on  $Q(s, a; \theta)$ 
8:   Receive new state  $s'$  and reward  $r$ 
9:    $y = \begin{cases} r & \text{for terminal } s' \\ r + \gamma \max_{a'} Q(s', a'; \theta^-) & \text{for non-terminal } s' \end{cases}$ 
10:  Accumulate gradients wrt :  $d\theta \leftarrow d\theta + \frac{\partial(y - Q(s, a; \theta))^2}{\partial \theta}$ 
11:   $s = s'$ 
12:   $T \leftarrow T + 1$  and  $t \leftarrow t + 1$ 
13:  if  $T \bmod I_{target} == 0$  then
14:    Update the target network  $\theta^- \leftarrow \theta$ 
15:  end if
16:  if  $t \bmod I_{asyncUpdate} == 0$  or  $s$  is terminal then
17:    Perform asynchronous update of  $\theta$  using  $d\theta$ .
18:    Clear gradients  $d\theta \leftarrow 0$ .
19:  end if
20: until  $T > T_{max}$ 
```

3 Asynchronous one-step Sarsa

The asynchronous one-step Sarsa algorithm is the same as asynchronous one-step Q-learning as given in algorithm 1 except that it uses a different target value for $Q(s, a)$. The target value used by one-step Sarsa is $r + \gamma Q(s', a'; \theta^-)$ where a' is the action taken in state s' .

4 Asynchronous n-step Q-learning

5 Asynchronous advantage actor-critic

References

- [1] Volodymyr Mnih, Adrià Puigdomènech Badia, Mehdi Mirza, Alex Graves, Timothy P. Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. CoRR, abs/1602.01783, 2016.

Algorithm 2 Asynchronous one-step Q-learning - pseudocode for each actor-learner thread.

```

1: // Assume global shared parameter vector  $\theta$ 
2: // Assume global shared target parameter vector  $\theta^-$ 
3: // Assume global shared counter  $T = 0$ 
4: Initialize thread step count  $t \leftarrow 1$ 
5: Initialize target network parameters  $\theta^- \leftarrow \theta$ 
6: Initialize thread-specific parameters  $\theta^- \leftarrow \theta$ 
7: Initialize network gradients  $d\theta \leftarrow 0$ 
8: repeat
9:   Clear gradient  $d\theta \leftarrow 0$ 
10:  Synchronous thread-specific parameters  $\theta' = \theta$ 
11:   $t_{target} = t$ 
12:  Get state  $s_t$ 
13:  repeat
14:    Take action  $a_t$  according to the  $\epsilon$ -greedy policy based on  $Q(s_t, a_t; \theta')$ 
15:    Receive reward  $r_t$  and new state  $s_{t+1}$ 
16:     $t \leftarrow t + 1$ 
17:     $T \leftarrow T + 1$ 
18:  until terminal  $s_t$  or  $t - t_{target} == t_{max}$ 
19:   $R = \begin{cases} 0 & \text{for terminal } s_t \\ \max_a Q(s^t, a; \theta^-) & \text{for non-terminal } s_t \end{cases}$ 
20:  for  $i \in t - 1, \dots, t_{start}$  do
21:     $R \leftarrow r_i + \gamma R$ 
22:    Accumulate gradients wrt  $\theta' : d\theta + \frac{\partial(R - Q(s_i, a_i; \theta'))^2}{\partial \theta'}$ 
23:  end for
24:  Performance asynchronous update of  $\theta$  using  $d\theta$ .
25:  if then  $T \bmod I_{target} == 0$ 
26:     $\theta^- \leftarrow \theta$ 
27:  end if
28: until  $T > T_{max}$ 

```

Algorithm 3 Asynchronous advantage actor-critic - pseudocode for each actor-learner thread.

```

1: // Assume global shared parameter vector  $\theta$  and  $\theta_v$  and global shared counter  $T = 0$ 
2: // Assume thread-specific parameters  $\theta'$  and  $\theta'_v$ 
3: Initialize thread step counter  $t \leftarrow 1$ 
4: repeat
5:   Reset gradients:  $d\theta \leftarrow 0$  and  $d\theta_v \leftarrow 0$ .
6:   Synchronous thread-specific parameters  $\theta' = \theta$  and  $\theta'_v = \theta_v$ 
7:    $t_{target} = t$ 
8:   Get state  $s_t$ 
9:   repeat
10:    Perform  $a_t$  according to policy  $\pi(a_t|s_t; \theta')$ 
11:    Receive reward  $r_t$  and new state  $s_{t+1}$ 
12:     $t \leftarrow t + 1$ 
13:     $T \leftarrow T + 1$ 
14:  until terminal  $s_t$  or  $t - t_{target} == t_{max}$ 
15:   $R = \begin{cases} 0 & \text{for terminal } s_t \\ V(s_t, \theta'_v) & \text{for non-terminal } s_t // \text{Bootstrap from last state} \end{cases}$ 
16:  for  $i \in t - 1, \dots, t_{start}$  do
17:     $R \leftarrow r_i + \gamma R$ 
18:    Accumulate gradients wrt  $\theta' : d\theta + \Delta_{\theta'} \log \pi(a_i|s_i; \theta')(R - V(s_i; \theta'_v))$ 
19:    Accumulate gradients wrt  $\theta'_v : d\theta_v + \frac{\partial(R - V(s_i; \theta'_v))^2}{\partial \theta'_v}$ 
20:  end for
21:  Perform asynchronous update of  $\theta$  using  $d\theta$  and of  $\theta_v$  using  $d\theta_v$ 
22: until  $T > T_{max}$ 

```
