# Asynchronous Methods for Deep Reinforcement Learning

Volodymyr Mnih[1]                                          VMNIH@GOOGLE.COM
Adrià Puigdomènech Badia[1]                               ADRIAP@GOOGLE.COM
Mehdi Mirza[1,2]                                    MIRZAMOM@IRO.UMONTREAL.CA
Alex Graves[1]                                          GRAVESA@GOOGLE.COM
Tim Harley[1]                                          THARLEY@GOOGLE.COM
Timothy P. Lillicrap[1]                              COUNTZERO@GOOGLE.COM
David Silver[1]                                    DAVIDSILVER@GOOGLE.COM
Koray Kavukcuoglu [1]                                   KORAYK@GOOGLE.COM

[1] Google DeepMind
[2] Montreal Institute for Learning Algorithms (MILA), University of Montreal

## 【强化学习算法 5】A3C

**张楚珩** ✔
清华大学 交叉信息院博士在读

**原文传送门：**

Mnih, Volodymyr, et al. "Asynchronous methods for deep reinforcement learning." International conference on machine learning. 2016.

**特色**：发现异步并行地执行多个agent，让它们on-policy地去面对不同的状态，不仅加速了算法，而且有一种使得算法更加稳定的效果。

**分类**：Model-free、Policy-based（Actor-critic）、On-policy、Continuous State Space、Continuous Action Space、Support High-dim Input

**过程**：普通的Actor-critic方法加上baseline $V(s)$ ，只不过在每个节点上计算用到的 $\pi_\theta(a|s)$ 和 $V_{\theta_v}(s)$ 都用从master上面同步下来的权值，各自探索积累自己的梯度 $d\theta$ 和 $d\theta_v$ ，每过一阵子就推到master上去更新。

**算法**：

> **Algorithm S2** Asynchronous advantage actor-critic - pseudocode for each actor-learner thread.
> // Assume global shared parameter vectors $\theta$ and $\theta_v$ and global shared counter $T = 0$
> // Assume thread-specific parameter vectors $\theta'$ and $\theta'_v$
> Initialize thread step counter $t \leftarrow 1$
> **repeat**
>     Reset gradients: $d\theta \leftarrow 0$ and $d\theta_v \leftarrow 0$.
>     Synchronize thread-specific parameters $\theta' = \theta$ and $\theta'_v = \theta_v$
>     $t_{start} = t$
>     Get state $s_t$
>     **repeat**
>         Perform $a_t$ according to policy $\pi(a_t|s_t; \theta')$
>         Receive reward $r_t$ and new state $s_{t+1}$
>         $t \leftarrow t + 1$
>         $T \leftarrow T + 1$
>     **until** terminal $s_t$ **or** $t - t_{start} == t_{max}$
>     $R = \begin{cases} 0 & \text{for terminal } s_t \\ V(s_t, \theta'_v) & \text{for non-terminal } s_t \text{// Bootstrap from last state} \end{cases}$
>     **for** $i \in \{t-1, \ldots, t_{start}\}$ **do**
>         $R \leftarrow r_i + \gamma R$
>         Accumulate gradients wrt $\theta'$: $d\theta \leftarrow d\theta + \nabla_{\theta'} \log \pi(a_i|s_i; \theta')(R - V(s_i; \theta'_v))$
>         Accumulate gradients wrt $\theta'_v$: $d\theta_v \leftarrow d\theta_v + \partial (R - V(s_i; \theta'_v))^2 / \partial \theta'_v$
>     **end for**
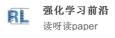>     Perform asynchronous update of $\theta$ using $d\theta$ and of $\theta_v$ using $d\theta_v$.
> **until** $T > T_{max}$

知乎 @张楚珩

编辑于 2018-09-19

强化学习 (Reinforcement Learning)　　算法（书籍）　　算法

▲ 赞同 1　▼　　💬 添加评论　　✈ 分享　　❤ 喜欢　　★ 收藏　　···

**文章被以下专栏收录**