A₃C

Asynchronous Advantage Actor-Critic (A3C) is a powerful reinforcement learning algorithm that combines the strengths of actor-critic methods with asynchronous parallel training. It allows for faster and more efficient learning by training multiple agents in parallel on different environments.

Algorithm Steps:

1. Initialization:

- Initialize the global actor and critic networks.
- Create multiple worker agents.

2. Worker Agent Loop:

- For each worker agent:
 - Initialize the local actor and critic networks as copies of the global network.
 - While training:
 - Collect a sequence of experiences by interacting with the environment.
 - Calculate the advantage function for each experience:

where

A_t is the advantage function, R_t is the cumulative discounted reward from time step t onwards, y is the discount factor, and V(s) is the value function.

$$A_t = R_t + \gamma V(s_{t+1}) - V(s_t)$$

A3C 1

Update the local actor and critic networks using gradient descent:

where

ullet are the actor's parameters, \underline{w} are the critic's parameters, and $\underline{\alpha}$ is the learning rate.

$$egin{aligned}
abla heta J(heta) &=
abla heta log\pi(a_t|s_t, heta)A_t \ & heta \leftarrow heta + lpha
abla heta J(heta) \ \\
abla w J(w) &=
abla w V(s_t)(r_t + \gamma V(s_{t+1}) - V(s_t)) \ & heta \leftarrow w + lpha
abla w J(w) \end{aligned}$$

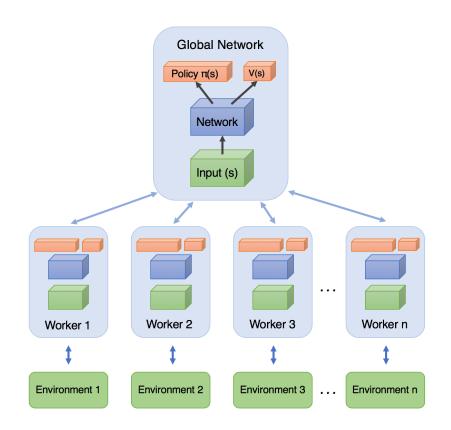
Synchronize the local networks with the global network.

Advantages of A3C:

- **Faster Training:** A3C can train significantly faster than single-threaded actor-critic methods by leveraging parallel computation.
- **Improved Stability:** Asynchronous updates can help to reduce the variance of the gradients, leading to more stable training.
- **Scalability:** A3C can be easily scaled to handle large-scale problems by increasing the number of worker agents.

Challenges of A3C:

- **Synchronization Overhead:** Synchronizing the global network with the worker agents can introduce overhead, especially for large-scale problems.
- **Exploration-Exploitation Trade-off:** Balancing exploration (trying new actions) with exploitation (repeating actions that have worked well in the past) is still a challenge in A3C.



A3C 3