



Understanding Actor Critic Methods and A2C

Important Concepts in Deep Reinforcement Learning



Chris Yoon · Follow

Published in Towards Data Science

6 min read · Feb 6, 2019

Listen

Share

More

Preliminaries

In [my previous post](#), we derived policy gradients and implemented the REINFORCE algorithm (also known as Monte Carlo policy gradients). There are, however, some glaring issues with vanilla policy gradients: noisy gradients and high variance.

But why does that happen?

Recall the policy gradient:

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) G_t \right]$$

As in the REINFORCE algorithm, we update the policy parameter through Monte Carlo updates (i.e. taking random samples). This introduces inherent high variability in log probabilities (log of the policy distribution) and cumulative reward values, because each trajectories during training can deviate from each other at great degrees.

Consequently, the high variability in log probabilities and cumulative reward values will make noisy gradients, and cause unstable learning and/or the policy distribution skewing to a non-optimal direction.

Besides high variance of gradients, another problem with policy gradients occurs trajectories have a cumulative reward of 0. The essence of policy gradient is increasing the probabilities for “good” actions and decreasing those of “bad” actions in the policy distribution; both “goods” and “bad” actions with will not be learned if the cumulative reward is 0.

Overall, these issues contribute to the instability and slow convergence of vanilla policy gradient methods.

Improving policy gradients: Reducing variance with a baseline

One way to reduce variance and increase stability is subtracting the cumulative reward by a baseline:

Introducing baseline $b(s)$:

$$\nabla_{\theta} J(\theta) = \mathbb{E}\left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)(G_t - b(s_t))\right]$$

Intuitively, making the cumulative reward smaller by subtracting it with a baseline will make smaller gradients, and thus smaller and more stable updates.

Here is a very illustrative explanation, taken from [Jerry Liu's post to “Why does the policy gradient method have high variance”](#):

For the sake of hypothetical example, let's say that $\nabla_\theta \log \pi_\theta(\tau)$ is [0.5, 0.2, 0.3] respectively for three trajectories, and $r(\tau)$ is [1000, 1001, 1002].

Then the variance of the product of the two terms for these three samples is $Var(0.5 \times 1000, 0.2 \times 1001, 0.3 \times 1002)$ which according to WolframAlpha is around 23286.8.

Now what if we reduce all values of $r(\tau)$ by a constant, say, 1001? Then the variance of the product becomes $Var(0.5 \times 1, 0.2 \times 0, 0.3 \times 1)$, which is 0.1633, a much smaller value.

[Taken from Jerry Liu's post "Why does the policy gradient method have high variance"]

Summary of common baseline functions

The baseline can take various values. The set of equations below illustrates the classic variants of actor critic methods (with respect to REINFORCE). In this post, we will take a look at Q Actor Critic and Advantage Actor Critic.

$$\begin{aligned} \nabla_\theta J(\theta) &= \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) G_t] && \text{REINFORCE} \\ &= \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) \mathbf{Q^w}(s, a)] && \text{Q Actor-Critic} \\ &= \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) \mathbf{A^w}(s, a)] && \text{Advantage Actor-Critic} \\ &= \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) \delta] && \text{TD Actor-Critic} \end{aligned}$$

Image taken from CMU CS10703 lecture slides

Actor Critic Methods

Before we return to baselines, let's first take a look at the vanilla policy gradient again to see how the Actor Critic architecture comes in (and what is really is). Recall that:

$$\nabla_\theta J(\theta) = \mathbb{E}_\tau \left[\sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t | s_t) G_t \right]$$

We can then decompose the expectation into:

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{s_0, a_0, \dots, s_t, a_t} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right] \mathbb{E}_{r_{t+1}, s_{t+1}, \dots, r_T, s_T} [G_t]$$

The second expectation term should be familiar; it is the Q value! (If you did not already know this, I would suggest that you read up on value iteration and Q learning).

$$\mathbb{E}_{r_{t+1}, s_{t+1}, \dots, r_T, s_T} [G_t] = Q(s_t, a_t)$$

Plugging that in, we can rewrite the update equation as such:

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \mathbb{E}_{s_0, a_0, \dots, s_t, a_t} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right] Q_w(s_t, a_t) \\ &= \mathbb{E}_{\tau} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) Q_w(s_t, a_t) \right] \end{aligned}$$

As we know, the Q value can be learned by parameterizing the Q function with a neural network (denoted by subscript w above).

This leads us to *Actor Critic Methods*, where:

1. The “Critic” estimates the value function. This could be the action-value (the *Q value*) or state-value (the *V value*).
2. The “Actor” updates the policy distribution in the direction suggested by the Critic (such as with policy gradients).

and both the Critic and Actor functions are parameterized with neural networks. In the derivation above, the Critic neural network parameterizes the Q value — so, it is called *Q Actor Critic*.

Below is the pseudocode for Q-Actor-Critic:

Algorithm 1 Q Actor Critic

Initialize parameters s, θ, w and learning rates α_θ, α_w ; sample $a \sim \pi_\theta(a|s)$.

for $t = 1 \dots T$: **do**

- Sample reward $r_t \sim R(s, a)$ and next state $s' \sim P(s'|s, a)$
- Then sample the next action $a' \sim \pi_\theta(a'|s')$
- Update the policy parameters: $\theta \leftarrow \theta + \alpha_\theta Q_w(s, a) \nabla_\theta \log \pi_\theta(a|s)$; Compute the correction (TD error) for action-value at time t:
$$\delta_t = r_t + \gamma Q_w(s', a') - Q_w(s, a)$$
- and use it to update the parameters of Q function:
$$w \leftarrow w + \alpha_w \delta_t \nabla_w Q_w(s, a)$$
- Move to $a \leftarrow a'$ and $s \leftarrow s'$

end for

Adapted from Lilian Weng's post "Policy Gradient algorithms"

As illustrated, we update both the Critic network and the Value network at each update step.

Back to Baselines

To make an argument from authority (as I was not able to find the reason why), the state-value function makes an optimal baseline function. This is stated in the Carnegie Mellon CS10703 and Berekely CS294 lecture slides, but with no reason provided.

So, using the V function as the baseline function, we subtract the Q value term with the V value. Intuitively, this means *how much better it is to take a specific action compared to the average, general action at the given state*. We will call this value the *advantage value*:

$$A(s_t, a_t) = Q_w(s_t, a_t) - V_v(s_t)$$

Does that mean we have to construct two neural networks for both the Q value and the V value (in addition to the policy network)? No. That would be very inefficient. Instead, we can use the relationship between the Q and the V from the Bellman optimality equation:

$$Q(s_t, a_t) = \mathbb{E}[r_{t+1} + \gamma V(s_{t+1})]$$

So, we can rewrite the advantage as:

$$A(s_t, a_t) = r_{t+1} + \gamma V_v(s_{t+1}) - V_v(s_t)$$

Then, we only have to use one neural network for the V function (parameterized by v above). So we can rewrite the update equation as:

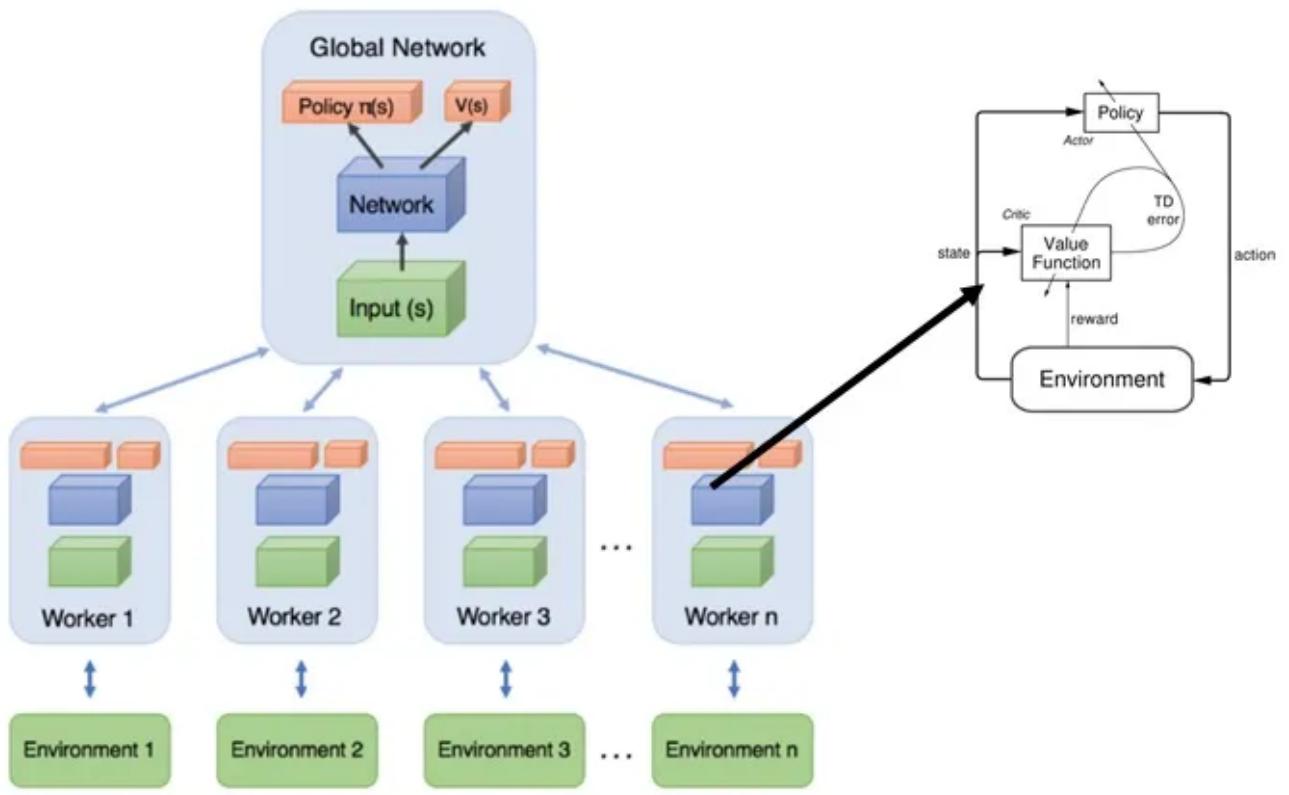
$$\begin{aligned}\nabla_\theta J(\theta) &\sim \sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t | s_t) (r_{t+1} + \gamma V_v(s_{t+1}) - V_v(s_t)) \\ &= \sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t | s_t) A(s_t, a_t)\end{aligned}$$

This is the *Advantage Actor Critic*.

Advantage Actor Critic (A2C) v.s. Asynchronous Advantage Actor Critic (A3C)

The *Advantage Actor Critic* has two main variants: the **Asynchronous Advantage Actor Critic (A3C)** and the **Advantage Actor Critic (A2C)**.

A3C was introduced in [Deepmind's paper "Asynchronous Methods for Deep Reinforcement Learning"](#) ([Mnih et al, 2016](#)). In essence, A3C implements *parallel training* where multiple workers in *parallel environments* *independently* update a *global value function*—hence “asynchronous.” One key benefit of having asynchronous actors is effective and efficient exploration of the state space.



High Level Architecture of A3C (image taken from [GroundAI blog post](#))

A2C is like A3C but without the asynchronous part; this means a single-worker variant of the A3C. It was empirically found that A2C produces comparable performance to A3C while being more efficient. According to [this](#) OpenAI blog post, researchers aren't completely sure if or how the asynchrony benefits learning:

After reading the paper, AI researchers wondered whether the asynchrony led to improved performance (e.g. “perhaps the added noise would provide some regularization or exploration?”), or if it was just an implementation detail that allowed for faster training with a CPU-based implementation ...

Our synchronous A2C implementation performs better than our asynchronous implementations — we have not seen any evidence that the noise introduced by asynchrony provides any performance benefit. This A2C implementation is more cost-effective than A3C when using single-GPU machines, and is faster than a CPU-only A3C implementation when using larger policies.

Anyhow, we will implement the A2C in this post as it is more simple in implementation. (This can easily be extended to A3C)

Implementing A2C

So, recall the new update equation, replacing the discounted cumulative award from vanilla policy gradients with the Advantage function:

$$\nabla_{\theta} J(\theta) \sim \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) A(s_t, a_t)$$

On each learning step, we update both the Actor parameter (with policy gradients and advantage value), and the Critic parameter (with minimizing the mean squared error with the Bellman update equation). Let's see how this looks like in code:

Below are the includes and hyper-parameters:

```
1 import sys
2 import torch
3 import gym
4 import numpy as np
5 import torch.nn as nn
6 import torch.optim as optim
7 import torch.nn.functional as F
8 from torch.autograd import Variable
9 import matplotlib.pyplot as plt
10 import pandas as pd
11
12 # hyperparameters
13 hidden_size = 256
14 learning_rate = 3e-4
15
16 # Constants
17 GAMMA = 0.99
18 num_steps = 300
19 max_episodes = 3000
```

a3c_params.py hosted with ❤ by GitHub

[view raw](#)

First, let's start with implementing the Actor Critic Network with the following configurations:

```
1  class ActorCritic(nn.Module):
2      def __init__(self, num_inputs, num_actions, hidden_size, learning_rate=3e-4):
3          super(ActorCritic, self).__init__()
4
5          self.num_actions = num_actions
6          self.critic_linear1 = nn.Linear(num_inputs, hidden_size)
7          self.critic_linear2 = nn.Linear(hidden_size, 1)
8
9          self.actor_linear1 = nn.Linear(num_inputs, hidden_size)
10         self.actor_linear2 = nn.Linear(hidden_size, num_actions)
11
12     def forward(self, state):
13         state = Variable(torch.from_numpy(state).float().unsqueeze(0))
14         value = F.relu(self.critic_linear1(state))
15         value = self.critic_linear2(value)
16
17         policy_dist = F.relu(self.actor_linear1(state))
18         policy_dist = F.softmax(self.actor_linear2(policy_dist), dim=1)
19
20         return value, policy_dist
```

a2c_model.py hosted with ❤ by GitHub

[view raw](#)

The main loop and update loop, as outlined above:

```

1  def a2c(env):
2      num_inputs = env.observation_space.shape[0]
3      num_outputs = env.action_space.n
4
5      actor_critic = ActorCritic(num_inputs, num_outputs, hidden_size)
6      ac_optimizer = optim.Adam(actor_critic.parameters(), lr=learning_rate)
7
8      all_lengths = []
9      average_lengths = []
10     all_rewards = []
11     entropy_term = 0
12
13     for episode in range(max_episodes):
14         log_probs = []
15         values = []
16         rewards = []
17
18         state = env.reset()
19         for steps in range(num_steps):
20             value, policy_dist = actor_critic.forward(state)
21             value = value.detach().numpy()[0,0]
22             dist = policy_dist.detach().numpy()
23
24             action = np.random.choice(num_outputs, p=np.squeeze(dist))
25             log_prob = torch.log(policy_dist.squeeze(0)[action])
26             entropy = -np.sum(np.mean(dist) * np.log(dist))
27             new_state, reward, done, _ = env.step(action)
28
29             rewards.append(reward)
30             values.append(value)
31             log_probs.append(log_prob)
32             entropy_term += entropy
33             state = new_state
34
35             if done or steps == num_steps-1:
36                 Qval, _ = actor_critic.forward(new_state)
37                 Qval = Qval.detach().numpy()[0,0]
38                 all_rewards.append(np.sum(rewards))
39                 all_lengths.append(steps)
40                 average_lengths.append(np.mean(all_lengths[-10:]))
41                 if episode % 10 == 0:
42                     sys.stdout.write("episode: {}, reward: {}, total length: {}, average length: {} \n".format(episode, np.sum(rewards), steps, np.mean(all_lengths[-10:])))
43                     break
44
45             # compute Q values
46             Qvals = np.zeros_like(values)
47             for t in reversed(range(len(rewards))):
48                 Qval = rewards[t] + GAMMA * Qval

```

```

46     Qval = rewards[1] + gamma * Qval
47     Qvals[t] = Qval
48
49
50         #update actor critic
51         values = torch.FloatTensor(values)
52         Qvals = torch.FloatTensor(Qvals)
53         log_probs = torch.stack(log_probs)
54
55
56         advantage = Qvals - values
57         actor_loss = (-log_probs * advantage).mean()
58         critic_loss = 0.5 * advantage.pow(2).mean()
59         ac_loss = actor_loss + critic_loss + 0.001 * entropy_term
60
61         ac_optimizer.zero_grad()
62         ac_loss.backward()
63         ac_optimizer.step()
64
65
66
67     # Plot results
68     smoothed_rewards = pd.Series.rolling(pd.Series(all_rewards), 10).mean()
69     smoothed_rewards = [elem for elem in smoothed_rewards]
70     plt.plot(all_rewards)
71     plt.plot(smoothed_rewards)
72     plt.plot()
73     plt.xlabel('Episode')
74     plt.ylabel('Reward')
75     plt.show()
76
77     plt.plot(all_lengths)
78     plt.plot(average_lengths)
79     plt.xlabel('Episode')
80     plt.ylabel('Episode length')
81     plt.show()

```

Running the code, we can see how the performance improves on the OpenAI Gym CartPole-v0 environment:

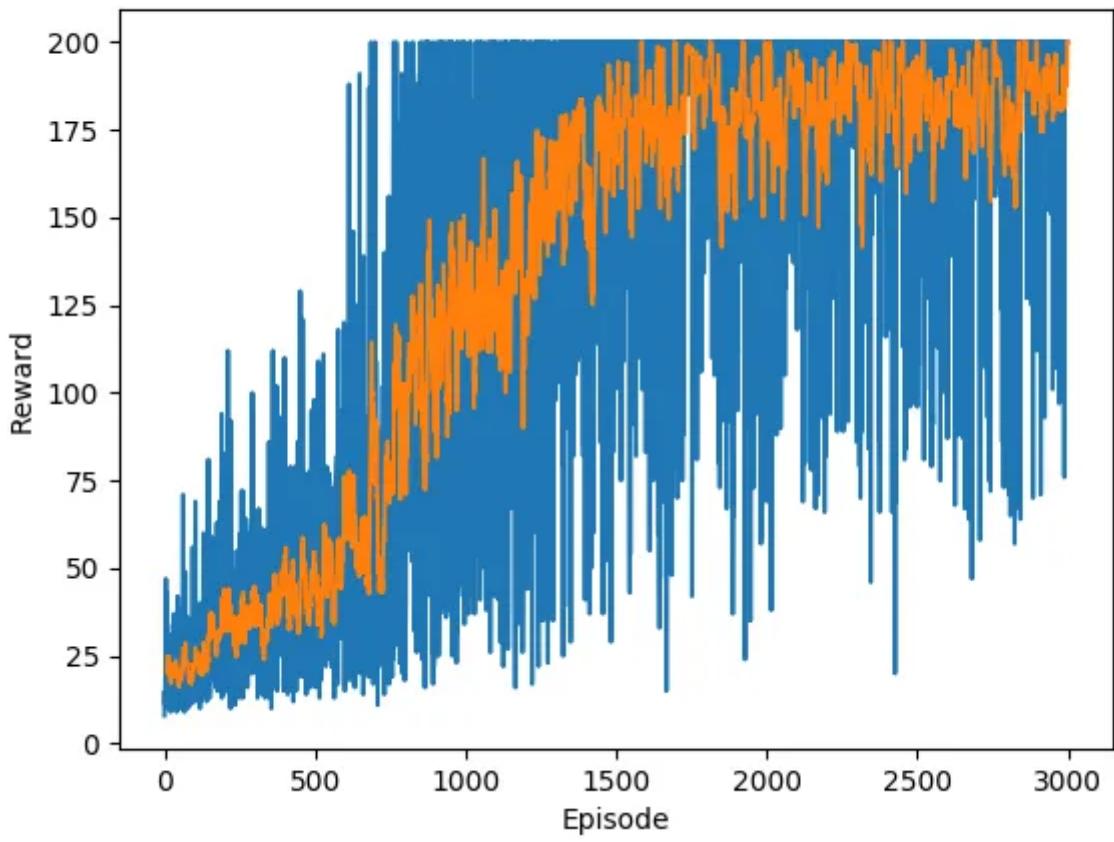
```

1 if __name__ == "__main__":
2     env = gym.make("CartPole-v0")
3     a2c(env)

```

[main.py](#) hosted with ❤ by GitHub

[view raw](#)



Blue: Raw rewards; Orange: Smoothed rewards

Find the full implementation here:

<https://github.com/thechrisyoon08/Reinforcement-Learning/>

References:

1. [UC Berkeley CS294 Lecture Slides](#)
2. [Carnegie Mellon University CS10703 Lecture Slides](#)
3. [Lilian Weng's Post on Policy Gradient Algorithms](#)
4. [Jerry Liu's answer on Quora Post "Why does the policy gradient method have a high variance"](#)
5. [Naver D2 RLCode Lecture Video](#)
6. [OpenAI blog post on A2C and ACKTR](#)
7. Diagram from [GroundAI's blog post](#) "Figure 4: Schematic representation of Asynchronous Advantage Actor Critic algorithm (A3C) algorithm."

Other Posts:

Check out my other posts on Reinforcement Learning:

- [Deriving Policy Gradients and Implementing REINFORCE](#)
- [Understanding Actor Critic Methods](#)
- [Deep Deterministic Policy Gradients Explained](#)

Thanks for reading!

Machine Learning

Reinforcement Learning

Deep Learning

Data Science

Artificial Intelligence



tds

Follow



Written by Chris Yoon

651 Followers · Writer for Towards Data Science

Student in NYC. <https://www.linkedin.com/in/chris-yoon-75847418b/>

More from Chris Yoon and Towards Data Science

$$\vartheta) = \mathbb{E} \left[\sum_{t=0}^{T-1} r_t \right]$$

 Chris Yoon

Deriving Policy Gradients and Implementing REINFORCE

Policy gradient methods are ubiquitous in model free reinforcement learning algorithms—they appear frequently in reinforcement learning...

4 min read · Dec 30, 2018

 1.1K  19



 Bex T. in Towards Data Science

130 ML Tricks And Resources Curated Carefully From 3 Years (Plus Free eBook)

Each one is worth your time

★ · 48 min read · Aug 1

👏 3.1K

💬 10



...



 Maxime Labonne  in Towards Data Science

Fine-Tune Your Own Llama 2 Model in a Colab Notebook

A practical introduction to LLM fine-tuning

★ · 12 min read · Jul 25

👏 1.8K

💬 33



...

Parameters:

θ^Q : Q network

θ^μ : Deterministic policy function

$\theta^{Q'}$: target Q network

$\theta^{\mu'}$: target policy network

 Chris Yoon in Towards Data Science

Deep Deterministic Policy Gradients Explained

Reinforcement Learning in Continuous Action Spaces

6 min read · Mar 20, 2019

 1.1K

 9

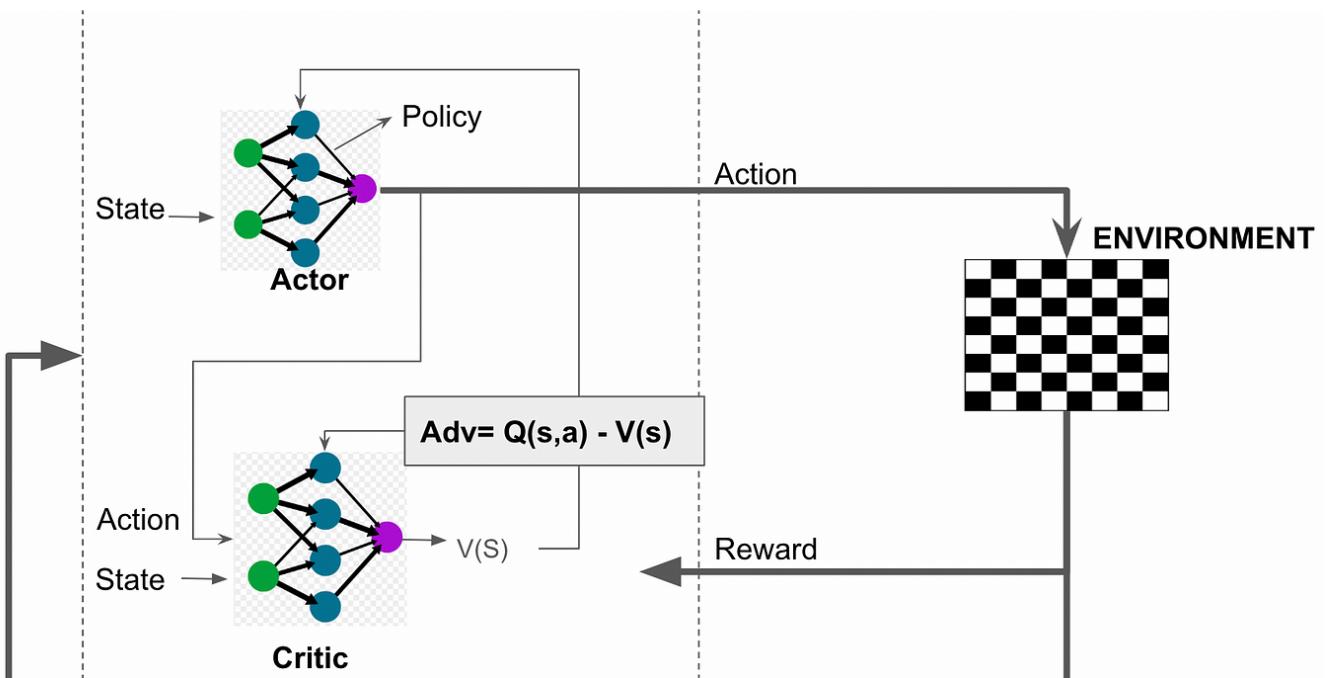


...

See all from Chris Yoon

See all from Towards Data Science

Recommended from Medium



 Renu Khandelwal

Unlocking the Secrets of Actor-Critic Reinforcement Learning: A Beginner’s Guide

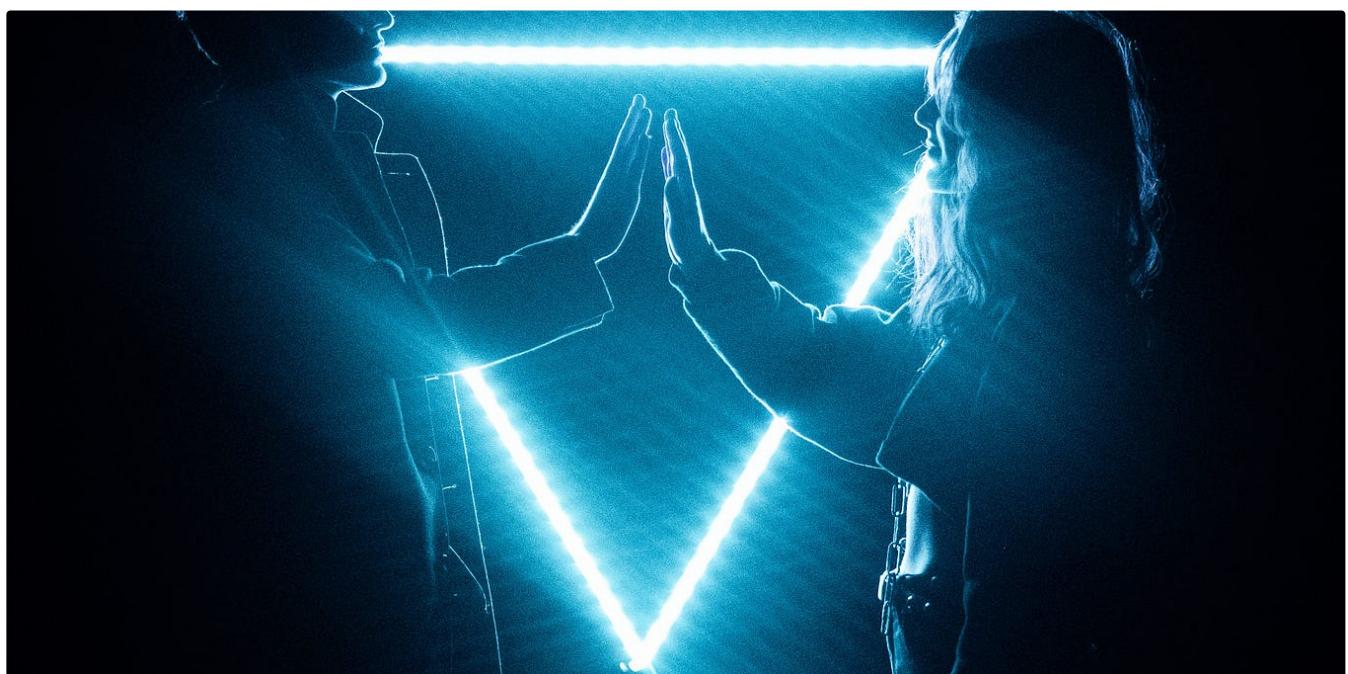
Understanding Actor-Critic Mechanisms, Different Flavors of Actor-Critic Algorithms, and a Simple Implementation in PyTorch

★ · 6 min read · Feb 21

 59  1



...



 Wouter van Heeswijk, PhD in Towards Data Science

Proximal Policy Optimization (PPO) Explained

The journey from REINFORCE to the go-to algorithm in continuous control

◆ · 13 min read · Nov 29, 2022

191

3

+

...

Lists



Predictive Modeling w/ Python

20 stories · 313 saves



Natural Language Processing

551 stories · 174 saves



Practical Guides to Machine Learning

10 stories · 344 saves



ChatGPT prompts

24 stories · 296 saves

A decorative graphic featuring a Tic-Tac-Toe board with yellow 'O's and pink 'X's, the words "tic tac toe" in yellow, blue, and pink, and a video game controller icon.

REINFORCEMENT LEARNING



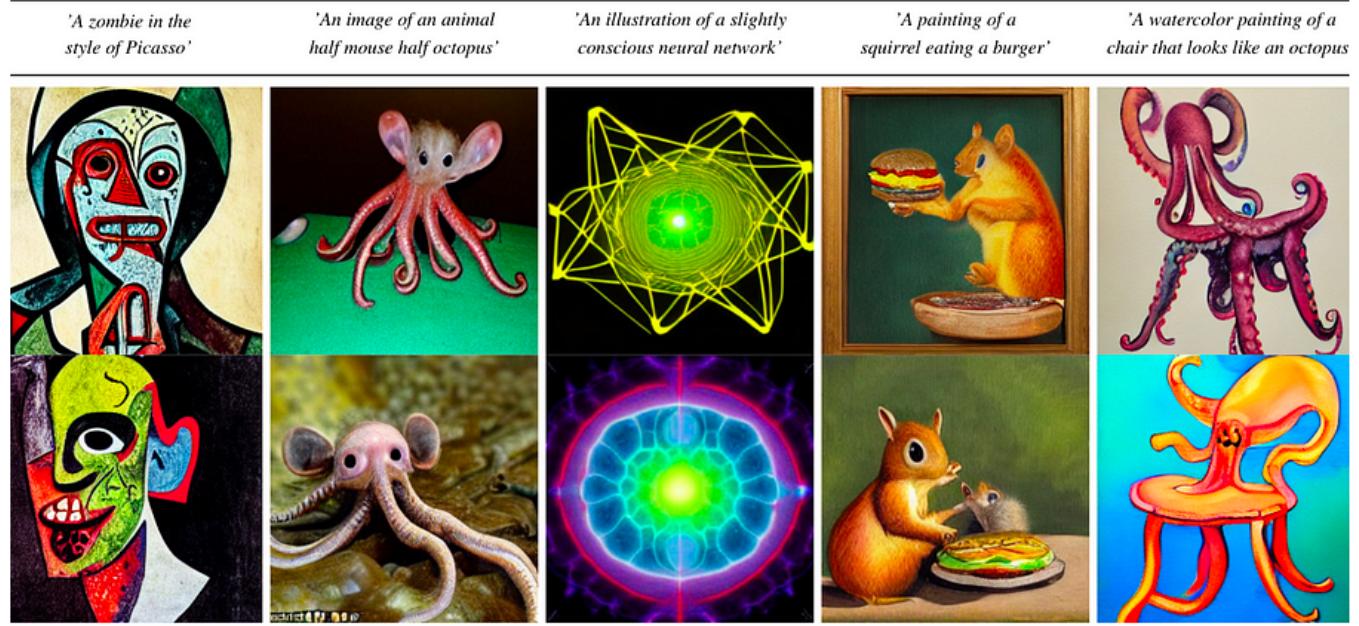
Waleed Mousa in Artificial Intelligence in Plain English

Building a Tic-Tac-Toe Game with Reinforcement Learning in Python: A Step-by-Step Tutorial

Welcome to this step-by-step tutorial on how to build a Tic-Tac-Toe game using reinforcement learning in Python. In this tutorial, we will...

9 min read · Mar 13

👏 39 ⚡ 1



👤 Onkar Mishra

Stable Diffusion Explained

How does Stable diffusion work? Explaining the tech behind text to image generation.

6 min read · Jun 8

👏 291 ⚡ 2



cal Programmatic Reinforcement ia Learning to Compose Program

¹ En-Pei Hu ^{*1} Pu-Jen Cheng ¹ Hung-Yi Lee ¹

 Ming-Hao Hsu

[RL] Hierarchical Programmatic Reinforcement Learning via Learning to Compose Programs (ICML23)

Paper Link: [Hierarchical Programmatic Reinforcement Learning via Learning to Compose Programs](#)

4 min read · Jul 21



Solving the CartPole Environment

A (brief) introduction to reinforcement learning, in plain English.

10 min read · Apr 10



...

[See more recommendations](#)