Stars | 59k    Follow @labmlai

View code on Github

# Proximal Policy Optimization - PPO

This is a PyTorch implementation of Proximal Policy Optimization - PPO.

PPO is a policy gradient method for reinforcement learning. Simple policy gradient methods do a single gradient update per sample (or a set of samples). Doing multiple gradient steps for a single sample causes problems because the policy deviates too much, producing a bad policy. PPO lets us do multiple gradient updates per sample by trying to keep the policy close to the policy that was used to sample data. It does so by clipping gradient flow if the updated policy is not close to the policy used to sample the data.

You can find an experiment that uses it here. The experiment uses Generalized Advantage Estimation.

CO Open in Colab

```
28  import torch
29
30  from labml_helpers.module import Module
31  from labml_nn.rl.ppo.gae import GAE
```

## PPO Loss

Here's how the PPO update rule is derived.

We want to maximize policy reward

$$\max_\theta J(\pi_\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[ \sum_{t=0}^{\infty} \gamma^t r_t \right]$$

where $r$ is the reward, $\pi$ is the policy, $\tau$ is a trajectory sampled from policy, and $\gamma$ is the discount factor between $[0, 1]$.

$$\mathbb{E}_{\tau \sim \pi_\theta} \left[ \sum_{t=0}^{\infty} \gamma^t A^{\pi_{OLD}}(s_t, a_t) \right] =$$

$$\mathbb{E}_{\tau \sim \pi_\theta} \left[ \sum_{t=0}^{\infty} \gamma^t \left( Q^{\pi_{OLD}}(s_t, a_t) - V^{\pi_{OLD}}(s_t) \right) \right] =$$

$$\mathbb{E}_{\tau \sim \pi_\theta} \left[ \sum_{t=0}^{\infty} \gamma^t \left( r_t + V^{\pi_{OLD}}(s_{t+1}) - V^{\pi_{OLD}}(s_t) \right) \right] =$$

$$\mathbb{E}_{\tau \sim \pi_\theta} \left[ \sum_{t=0}^{\infty} \gamma^t \left( r_t \right) \right] - \mathbb{E}_{\tau \sim \pi_\theta} \left[ V^{\pi_{OLD}}(s_0) \right] = J(\pi_\theta) - J(\pi_{\theta_{OLD}})$$

So,

$$\max_\theta J(\pi_\theta) = \max_\theta \mathbb{E}_{\tau \sim \pi_\theta} \left[ \sum_{t=0}^{\infty} \gamma^t A^{\pi_{OLD}}(s_t, a_t) \right]$$

Define discounted-future state distribution,

$$d^\pi(s) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t P(s_t = s | \pi)$$

Then,

$$J(\pi_\theta) - J(\pi_{\theta_{OLD}}) = \mathbb{E}_{\tau \sim \pi_\theta} \left[ \sum_{t=0}^{\infty} \gamma^t A^{\pi_{OLD}}(s_t, a_t) \right]$$

$$= \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d^{\pi_\theta}, a \sim \pi_\theta} \left[ A^{\pi_{OLD}}(s, a) \right]$$

Importance sampling $a$ from $\pi_{\theta_{OLD}}$,

$$J(\pi_\theta) - J(\pi_{\theta_{OLD}}) = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d^{\pi_\theta}, a \sim \pi_\theta} \left[ A^{\pi_{OLD}}(s, a) \right]$$

$$= \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d^{\pi_\theta}, a \sim \pi_{\theta_{OLD}}} \left[ \frac{\pi_\theta(a|s)}{\pi_{\theta_{OLD}}(a|s)} A^{\pi_{OLD}}(s, a) \right]$$

Then we assume $d^{\pi_\theta}(s)$ and $d^{\pi_{\theta_{OLD}}}(s)$ are similar. The error we introduce to $J(\pi_\theta) - J(\pi_{\theta_{OLD}})$ by this assumption is bound by the KL divergence between $\pi_\theta$ and $\pi_{\theta_{OLD}}$. Constrained Policy Optimization shows the proof of this. I haven't read it.

$$J(\pi_\theta) - J(\pi_{\theta_{OLD}}) = \frac{1}{1-\gamma} \mathop{\mathbb{E}}_{\substack{s \sim d^{\pi_\theta} \\ a \sim \pi_{\theta_{OLD}}}} \left[ \frac{\pi_\theta(a|s)}{\pi_{\theta_{OLD}}(a|s)} A^{\pi_{OLD}}(s,a) \right]$$

$$\approx \frac{1}{1-\gamma} \mathop{\mathbb{E}}_{\substack{s \sim d^{\pi_{\theta_{OLD}}} \\ a \sim \pi_{\theta_{OLD}}}} \left[ \frac{\pi_\theta(a|s)}{\pi_{\theta_{OLD}}(a|s)} A^{\pi_{OLD}}(s,a) \right]$$

$$= \frac{1}{1-\gamma} \mathcal{L}^{CPI}$$

```
34  class ClippedPPOLoss(Module):
```

```
136     def __init__(self):
137         super().__init__()
```

```
139     def forward(self, log_pi: torch.Tensor, sampled_log_pi: torch.Tensor,
140                 advantage: torch.Tensor, clip: float) -> torch.Tensor:
```

ratio $r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{OLD}}(a_t|s_t)}$; *this is different from rewards* $r_t$.

```
143         ratio = torch.exp(log_pi - sampled_log_pi)
```

## Cliping the policy ratio

$$\mathcal{L}^{CLIP}(\theta) = \mathbb{E}_{a_t, s_t \sim \pi_{\theta_{OLD}}} \left[ min\Big( r_t(\theta)\bar{A}_t, clip\big(r_t(\theta), 1-\epsilon, 1+\epsilon\big)\bar{A}_t \Big) \right]$$

The ratio is clipped to be close to 1. We take the minimum so that the gradient will only pull $\pi_\theta$ towards $\pi_{\theta_{OLD}}$ if the ratio is not between $1-\epsilon$ and $1+\epsilon$. This keeps the KL divergence between $\pi_\theta$ and $\pi_{\theta_{OLD}}$ constrained. Large deviation can cause performance collapse; where the policy performance drops and doesn't recover because we are sampling from a bad policy.

Using the normalized advantage $\bar{A}_t = \frac{\hat{A}_t - \mu(\hat{A}_t)}{\sigma(\hat{A}_t)}$ introduces a bias to the policy gradient estimator, but it reduces variance a lot.

```
172         clipped_ratio = ratio.clamp(min=1.0 - clip,
173                                     max=1.0 + clip)
174         policy_reward = torch.min(ratio * advantage,
175                                   clipped_ratio * advantage)
176
177         self.clip_fraction = (abs((ratio - 1.0)) > clip).to(torch.float).mean()
178
179         return -policy_reward.mean()
```

# Clipped Value Function Loss

Similarly we clip the value function update also.

$$V_{CLIP}^{\pi_\theta}(s_t) = clip\Big(V^{\pi_\theta}(s_t) - \hat{V}_t, -\epsilon, +\epsilon\Big)$$

$$\mathcal{L}^{VF}(\theta) = \frac{1}{2}\mathbb{E}\Big[max\Big(\big(V^{\pi_\theta}(s_t) - R_t\big)^2, \big(V_{CLIP}^{\pi_\theta}(s_t) - R_t\big)^2\Big)\Big]$$

Clipping makes sure the value function $V_\theta$ doesn't deviate significantly from $V_{\theta_{OLD}}$.

```python
182  class ClippedValueFunctionLoss(Module):
```

```python
204      def forward(self, value: torch.Tensor, sampled_value: torch.Tensor, sampled_return: torch.Tensor, clip: float):
205          clipped_value = sampled_value + (value - sampled_value).clamp(min=-clip, max=clip)
206          vf_loss = torch.max((value - sampled_return) ** 2, (clipped_value - sampled_return) ** 2)
207          return 0.5 * vf_loss.mean()
```

labml.ai