

## 总结二: GAE

2022年2月3日 12:27

- 在PPO、TRPO等类似算法中都常使用 GAE 技术, 这是一种对 advantage function 更稳定的估计方法
- 论文: [\[1506.02438\] High-Dimensional Continuous Control Using Generalized Advantage Estimation \(arxiv.org\)](#)
- 参考: [【强化学习技术 28】GAE - 知乎 \(zhihu.com\)](#)
- GAE (Generalized Advantage Estimation)
  - 在 policy gradient 中, 通过 advantage function, 提高稳定性, 寻求 unbiased 和 various 之间的平衡
  - 常见的 policy loss 有 (总结):

$$g = \mathbb{E} \left[ \sum_{t=0}^{\infty} \Psi_t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right], \quad (1)$$

where  $\Psi_t$  may be one of the following:

- |  |   |
|--|---|
| 1. $\sum_{t=0}^{\infty} r_t$ : total reward of the trajectory.                     | 4. $Q^{\pi}(s_t, a_t)$ : state-action value function.     |
| 2. $\sum_{t'=t}^{\infty} r_{t'}$ : reward following action $a_t$ .                 | 5. $A^{\pi}(s_t, a_t)$ : advantage function.              |
| 3. $\sum_{t'=t}^{\infty} r_{t'} - b(s_t)$ : baselined version of previous formula. | 6. $r_t + V^{\pi}(s_{t+1}) - V^{\pi}(s_t)$ : TD residual. |

$$V^{\pi}(s_t) := \mathbb{E}_{a_t: \infty}^{s_{t+1: \infty}} \left[ \sum_{l=0}^{\infty} r_{t+l} \right] \quad Q^{\pi}(s_t, a_t) := \mathbb{E}_{a_{t+1: \infty}}^{s_{t+1: \infty}} \left[ \sum_{l=0}^{\infty} r_{t+l} \right] \quad (2)$$

$$A^{\pi}(s_t, a_t) := Q^{\pi}(s_t, a_t) - V^{\pi}(s_t), \quad (\text{Advantage function}). \quad (3)$$

- 之后人们引入了 discount rate:  $\gamma$ , 即:

$$V^{\pi, \gamma}(s_t) := \mathbb{E}_{a_t: \infty}^{s_{t+1: \infty}} \left[ \sum_{l=0}^{\infty} \gamma^l r_{t+l} \right] \quad Q^{\pi, \gamma}(s_t, a_t) := \mathbb{E}_{a_{t+1: \infty}}^{s_{t+1: \infty}} \left[ \sum_{l=0}^{\infty} \gamma^l r_{t+l} \right] \quad (4)$$

$$A^{\pi, \gamma}(s_t, a_t) := Q^{\pi, \gamma}(s_t, a_t) - V^{\pi, \gamma}(s_t). \quad (5)$$

$$g^{\gamma} := \mathbb{E}_{a_{0: \infty}}^{s_{0: \infty}} \left[ \sum_{t=0}^{\infty} A^{\pi, \gamma}(s_t, a_t) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]. \quad (6)$$

需要注意引入  $\gamma$  之后, policy loss ( $g$ ) 的估计就有偏了

- 使用一个统计量估计这里的  $A$ , 要求替换后的 loss 求出的梯度无偏, 即所谓  $\gamma$ -just (不是  $g$  本身无偏)

**Definition 1.** The estimator  $\hat{A}_t$  is  $\gamma$ -just if

$$\mathbb{E}_{a_{0: \infty}}^{s_{0: \infty}} \left[ \hat{A}_t(s_{0: \infty}, a_{0: \infty}) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right] = \mathbb{E}_{a_{0: \infty}}^{s_{0: \infty}} [A^{\pi, \gamma}(s_t, a_t) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)]. \quad (7)$$

那满足  $\gamma$ -just 的  $A$  的估计有

- |  |  |
|--|--|
| <ul style="list-style-type: none"><li><math>\sum_{l=0}^{\infty} \gamma^l r_{t+l}</math></li><li><math>Q^{\pi, \gamma}(s_t, a_t)</math></li></ul> | <ul style="list-style-type: none"><li><math>A^{\pi, \gamma}(s_t, a_t)</math></li><li><math>r_t + \gamma V^{\pi, \gamma}(s_{t+1}) - V^{\pi, \gamma}(s_t)</math></li></ul> |
|--|--|

- 主要基于这样的观点: 赋予距离当前更远的时刻更小的权值, 可以减少扰动 (如 TD0 比 MC 更稳定)

(数学证明?)

- 论文提出了 GAE 也是一种  $A$  的符合  $\gamma$ -just 的估计, 是从 1-step 到  $n$ -step 的 TD Advantage function ( $\delta$ ) 的加权平均:

$$\hat{A}^{GAE(\gamma, \lambda)} = \frac{(\hat{A}^{(1)} + \lambda \hat{A}^{(2)} + \lambda^2 \hat{A}^{(3)} \dots)}{(1 - \lambda)(1 + \lambda + \lambda^2 \dots)} = \frac{(\hat{A}^{(1)} + \lambda \hat{A}^{(2)} + \lambda^2 \hat{A}^{(3)} \dots)}{(1 - \lambda)} \approx (1 - \lambda)(\hat{A}^{(1)} + \lambda \hat{A}^{(2)} + \lambda^2 \hat{A}^{(3)} \dots)$$

$$\begin{aligned}
\hat{A}_t^{GAE(\gamma,\lambda)} &:= (1-\lambda) \left( \hat{A}_t^{(1)} + \lambda \hat{A}_t^{(2)} + \lambda^2 \hat{A}_t^{(3)} + \dots \right) \\
&= (1-\lambda) \left( \delta_t^V + \lambda \left( \delta_t^V + \gamma \delta_{t+1}^V \right) + \lambda^2 \left( \delta_t^V + \gamma \delta_{t+1}^V + \gamma^2 \delta_{t+2}^V \right) + \dots \right) \\
&= \sum_{l=0}^{\infty} (\gamma \lambda)^l \delta_{t+l}^V
\end{aligned}$$

实际上是为更近的 A 赋予更大的权重