

PPD

- 기존에 샘플 데이터를 조금씩 버리는 방식에서 여러번 사용하게 만들었다.
- On-policy 알고리즘이지만 데이터를 한번만 쓰고 버리지 않는다
- PPO는 TRPO를 발전시킨 모델이다.

TRPO

- ① = constant로 추정
- 기존의 Policy Gradient는 $\hat{\mathbb{E}}_t[\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \hat{A}_t]$ 값으로 정책을 업데이트 했다.
- 기존의 Loss 또한 $\hat{\mathbb{E}}_t[\log \pi_{\theta}(a_t | s_t) \hat{A}_t]$ 값으로 Loss를 만들었다.
- 만약 이전 θ 값과 old θ 의 차이가 심하게 많다면 Loss를 $\hat{\mathbb{E}}_t[\frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{old}}(a_t | s_t)} \hat{A}_t]$ 로 만들수 있다. (Importance sampling 방법이다.)

Chain rule를 이용해서 미분식을 만들고 미분을 뺀다.

$$\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) |_{\theta_{old}} = \frac{\nabla_{\theta} \pi_{\theta}(a_t | s_t) |_{\theta_{old}}}{\pi_{\theta}(a_t | s_t)} = \nabla_{\theta} \left(\frac{f(\theta)}{f(\theta_{old})} \right) |_{\theta_{old}}$$

- 위의 조건을 만족하도록 식을 작성하면 아래와 같다.

$$\begin{aligned} & \underset{\theta}{\text{maximize}} \quad \hat{\mathbb{E}}_t \left[\frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{old}}(a_t | s_t)} \hat{A}_t \right] \\ & \text{subject to} \quad \hat{\mathbb{E}}_t [\text{KL}[\pi_{\theta_{old}}(\cdot | s_t), \pi_{\theta}(\cdot | s_t)]] \leq \delta. \end{aligned}$$

식을 하나로 만들면 아래와 같다.

$$\underset{\theta}{\text{maximize}} \quad \hat{\mathbb{E}}_t \left[\frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{old}}(a_t | s_t)} \hat{A}_t \right] - \beta \hat{\mathbb{E}}_t [\text{KL}[\pi_{\theta_{old}}(\cdot | s_t), \pi_{\theta}(\cdot | s_t)]]$$

PPO

- TRPO에서 변형된 형태이다.
- 2가지 방식중에 clip 방식이 성능이 좋다.
- clip 방식의 Loss는 다음과 같다.
- 기존에 original Loss에 clip 함수를 추가하여 최대 최소 값을 정한다.
- 마지막에 2개중에 작은 값을 사용한다.

