

## PP0 알고리즘

(Proximal Policy Optimization) 중심부 정책 최적화

- **TRPO** Trust Region Policy Optimization (신뢰 영역 정책 최적화)을 연산에 용이하도록 알고리즘 수정한 것
- Actor 와 Critic 을 통해서 최적화

※ Reinforcement Learning 에서의 목표 함수 Expected Reward

$$\eta(\pi) = \hat{E}_t [\sum_{t=0}^{\infty} \gamma^t \log \pi_{\theta}(a_t | s_t) \hat{A}_t] \rightarrow t \text{의 시점에서 추정되는 Advantage}$$

$\rightarrow$  Stochastic Policy [확률 정책]

$\rightarrow$  파라미터  $\theta$ 에 대해서 Gradient를 사용하여 최대화하는 방향으로 업데이트한다.

[illegible]

[illegible]

[illegible]

