

09

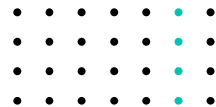
# PPO 알고리즘

## 1. 기본개념



# PPO 기본개념

## 중요도 샘플링

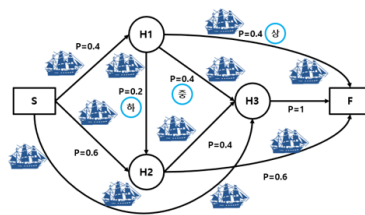


### 중요도 샘플링(Importance Sampling)

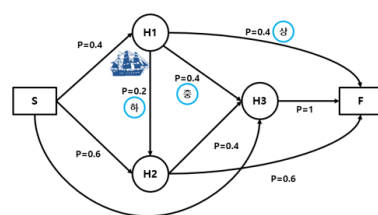
- $f(x)$ 의 기댓값을 계산하고자 하는 확률분포  $p(x)$ 를 알고 있지만  $p$ 에서 샘플을 생성하기 어려울 때, 비교적 샘플을 구하기 쉬운 확률분포  $q(x)$ 에서 샘플을 생성하여 확률분포  $p(x)$ 에서의  $f(x)$ 의 기댓값을 생성하는 것
- 중요도 샘플링을 활용해서 오프 폴리시에서 정책 간 연관성 문제를 해결한다.

$$\sum P(X)f(X) = \sum Q(X) \left[ \frac{P(X)}{Q(X)} f(X) \right]$$

<b>P(X)</b>	어떤 환경에서 변수 X의 확률분포 P
<b>Q(X)</b>	다른 환경에서 변수 X의 확률분포 Q
<b>f(X)</b>	X의 함수 어떤 함수도 가능(sin, cos, 2x+1 등)
<b><math>\sum P(X)f(X)</math></b>	변수X의 함수 f(X)에 대한 확률분포 P의 기댓값



기존 행로 데이터

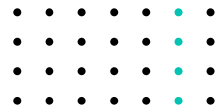


신규 행로



# PPO 기본개념

## 중요도 샘플링



### 활용방법

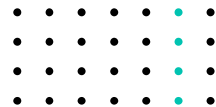
- 강화학습에서는  $\pi_{\theta_{old}}$  정책을 가지고 수집한 데이터  $X$ 를 가지고  $\pi_{\theta}$  정책을 가진 환경에서 기대값을 구하기 위해서는  $\pi_{\theta}$  정책의 기대값을 구하는 공식에다가 각각의 정책 비율을 곱하면 된다

$$\begin{aligned} E_{X \sim \pi_{\theta}}[f(X)] &= \sum \pi_{\theta} f(X) \\ &= \sum \pi_{\theta_{old}} \left[ \frac{\pi_{\theta}}{\pi_{\theta_{old}}} f(X) \right] \\ &= E_{X \sim \pi_{\theta_{old}}} \left[ \frac{\pi_{\theta}}{\pi_{\theta_{old}}} f(X) \right] \end{aligned}$$



# PPO 기본개념

A2C 오프 폴리시로 변경



Policy Gradient

$$\nabla_{\theta} J(\theta) = E_{\pi_{\theta}}[\nabla_{\theta} \log \pi_{\theta}(a|s) R_s^a]$$

$$\approx E_{\pi_{\theta}}[\nabla_{\theta} \log \pi_{\theta}(a|s) A_s^a]$$

$$= E_{\pi_{\theta}}\left[\frac{\nabla_{\theta} \pi_{\theta}(a|s)}{\pi_{\theta}(a|s)} A_s^a\right]$$

$$= E_{\pi_{\theta_{old}}}\left[\frac{\pi_{\theta}(a|s)}{\pi_{\theta_{old}}(a|s)} \frac{\nabla_{\theta} \pi_{\theta}(a|s)}{\pi_{\theta}(a|s)} A_s^a\right]$$

$$= E_{\pi_{\theta_{old}}}\left[\frac{\nabla_{\theta} \pi_{\theta}(a|s)}{\pi_{\theta_{old}}(a|s)} A_s^a\right]$$

$$= E_{\pi_{\theta_{old}}}[\nabla_{\theta}(\frac{\pi_{\theta}(a|s)}{\pi_{\theta_{old}}(a|s)} A_s^a)]$$

$$= - \frac{\pi_{\theta}(a|s)}{\pi_{\theta_{old}}(a|s)} A_s^a$$

①

②

③

④

⑤

⑥

⑦

- 경험을 재사용하고 하나의 정책으로 많은 경험을 쌓아 인공 신경망을 업데이트 하기 위해 오프 폴리시로 변경

미분의 연쇄 법칙  
(Chain rule)

$$y = \log f(x)$$

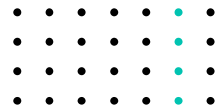
$$y' = \frac{f(x)'}{f(x)}$$

Cost Function



# PP0 기본개념

## 클리핑



### 빠른 학습속도의 문제점



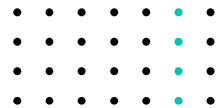
이미지 출처: <http://www.iautocar.co.kr>

- 변수의 변화 속도를 제어하면서 학습 효율을 높일 수 있다.
- PP0에서는 학습 속도를 제어하기 위해 클리핑(Clipping) 기법을 사용한다



# PPO 기본개념

## 클리핑



### 클리핑 개념

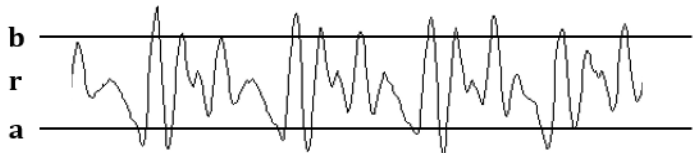
Clipping

`clip( r, a, b )`

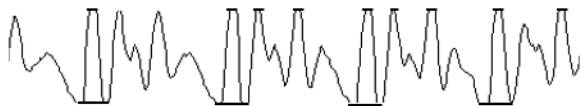
$r < a$  이면  $a$

$r > b$  이면  $b$

$a \leq r \leq b$  이면  $r$



`clip( r, a, b )`

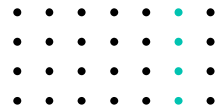


- 클리핑은 데이터의 하한(Lower)과 상한(Upper)을 정해 놓고 입력되는 데이터를 일정 범위 안으로 들어오도록 만드는 기법이다



# PP0 기본개념

## 클리핑



### 비용함수

- $\epsilon$  을 통해 클리핑의 상한과 하한을 지정할 수 있다. 강화학습 과정에서 효율적인  $\epsilon$  값을 선택하느냐는 효율적인 학습을 위한 아주 중요한 문제이다

$$r(\theta) = \frac{\pi_{\theta}(a|s)}{\pi_{\theta_{old}}(a|s)} \quad ①$$

### Cost Function

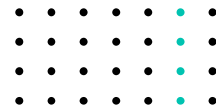
$$\min( \underbrace{r_t(\theta) A_t}_{②-1}, \underbrace{clip(r_t(\theta), 1 - \epsilon, 1 + \epsilon) A_t}_{②-2} ) \quad ②$$

Original Loss

Clipped Loss



# PPO 기본개념 GAE



Advantage

$$R_t + \gamma q_{t+1} - q_t$$

GAE

Generalized Advantage Estimation

delta

+

$\gamma$

GAE

감가율로 할인된 누적 Advantage

$\gamma$

$$R_t + \gamma q_{t+1} - q_t$$

- Generalized Advantage Estimation
- GAE : 감가율로 할인된 누적 어드벤처지
- 에이전트가 경험을 쌓을 때 얻는 보상(r)과 큐함수(q)를 모두 시간순으로 기록해 두었다가 어드벤처지를 계산할 때 사용

시간흐름

$r_1$	$q_1$
$r_2$	$q_2$
:	:
$r_n$	$q_n$

GAE  
계산순서

