

강화학습 PPO

Wan Ju Kang

- Trust Region PO
- Maximum a Posteriori PO
- Proximal PO
- 모두 actor-critic 구조이지만 일부 자료에서는 policy gradient 계열로 소개하기도 함
 - 알고리즘의 독창성이 모두 actor 쪽의 구현이라서
- 기본적으로 학습을 진행하면서 정책이 나빠지지 않는 말게 하자는 공통된 목표를 위해 구현됨
 - Monotonic improvement guarantee
- 수렴 관련 성질의 우수함을 인정 받아 다양한 파생 연구가 진행됨

- Proximal; approximate; proximity
- = 주변
- $PPO = PG + Idea\{1, 2, 3\}$
- Idea 1: Actor – 현재 정책의 주변을 탐사해서 좀 더 좋다면 옮겨가자
- Idea 2: Actor – 주변을 최대한 면밀히 탐사하자
- Idea 3: Critic – 그냥 하던 거 하자

- RL slide 69
- Q값 자체 대신 Advantage라는 개념의 차용 → Advantage Actor-Critic (A2C)
 - $A = Q - V$

Policy Objective Function

$$L^{PG}(\theta) = E_t[\log \pi_{\theta}(a_t|s_t) * \underline{A_t}]$$

log probability of
taking that action at
that state

Advantage if $A > 0$, this action is
better than the other action
possible at that state

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}}[\nabla_{\theta} \log \pi_{\theta}(s, a) * Q_w(s, a)]$$

Idea 1: 주변만 탐사하기

- Clipped objective function
- 너무 급진적으로 다른 정책 말고 주변 가까운 정책부터 시도해보자
- TRPO는 clipping 대신 두 정책의 KL divergence를 계산
 - → 비쌈

PPO's Clipped surrogate objective function

$$L^{CLIP}(\theta) = \hat{\mathbb{E}}_t \left[\min(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t) \right]$$

The ratio function

$$r_t(\theta) = \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)}$$

Idea 2: 주변은 다 탐사하기

- 일반적으로 엔트로피 항 s 를 추가하면 보다 다양한 값을 바꿔 학습해보도록 진행됨

Final PPO's Actor Critic Objective Function

$$L_t^{CLIP+VF+S}(\theta) = \hat{\mathbb{E}}_t [L_t^{CLIP}(\theta) - c_1 L_t^{VF}(\theta) + c_2 S[\pi_\theta](s_t)]$$

c_1 and c_2 are coefficients.

Squared-error value loss: $(V_\theta(s_t) - V_t^{\text{targ}})^2$

Add an entropy bonus to ensure sufficient exploration

Idea 3: 비평가는 원래 하던 일

- 최종 손실 함수를 최대화하도록 PPO Agent 트레이닝
- L^{CLIP} : 새 정책의 advantage가 더 크되, epsilon 안에 있도록
- S : 엔트로피 크게 해서 epsilon 안은 최대한 모두 보도록
- L^{VF} : TD error

Final PPO's Actor Critic Objective Function

$$L_t^{CLIP+VF+S}(\theta) = \hat{\mathbb{E}}_t [L_t^{CLIP}(\theta) - c_1 L_t^{VF}(\theta) + c_2 S[\pi_\theta](s_t)]$$

c_1 and c_2 are coefficients.

Squared-error value loss: $(V_\theta(s_t) - V_t^{\text{targ}})^2$

Add an entropy bonus to ensure sufficient exploration

Thank you!

Any Questions?
