

REINFORCE 알고리즘

- 알고리즘 이름이진짜 REINFORCE 알고리즘이다.

Policy Gradient 관계

- Policy Gradient를 바탕으로 한다.

- 기댓값에서 큰수의 법칙으로 'N'이 충분하면 아래처럼 표시 가능하다.

$$\int x \cdot p(x) dx \approx \frac{1}{N} \sum_{i=1}^N x_i$$

- 그렇다면 Policy Gradient도 적용시킬수있다.

$$\theta \leftarrow \theta + \alpha \left[\int_{\tau=0}^{\infty} \left(\nabla_{\theta} \ln P_{\theta}(a_t | s_t) \right) \times G_t \cdot P_{\theta}(\tau) d\tau \right]$$

• 만약 충분한(N)이 없다면 큰수의 법칙 적용

$$\rightarrow \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^{\infty} \nabla_{\theta} \ln P_{\theta}(a_t | s_t) G_t$$

- 위의 처럼하면 가능한 하지만 1개의 에피소드 만해도 수만번(5)④ 때문에 계산이 힘들다.

그래서(④)는 힘들지만 (N=1)을 사용해서 미분정행한다.

$$\bullet \theta \leftarrow \theta + \alpha \sum_{t=0}^{\infty} \nabla_{\theta} \ln P_{\theta}(a_t | s_t) G_t$$

- 알고리즘 순서

① 샘플 τ 생성 (즉 에피소드 생성) from $P_{\theta}(a_t | s_t)$

$$\textcircled{2} \theta \leftarrow \theta + \alpha \sum_{t=0}^{\infty} \nabla_{\theta} \ln P_{\theta}(a_t | s_t) G_t$$

문제점

- (n=1) 이다보니 결과값이 biased, High Variance 때문에 찾아가기 힘들다.

- Actor-Critic 으로 해결했다.

