

Value-based, Policy-based

Value-based

- 가치함수를 Maximization 하는 action을 고르는 것

- (예) DQN

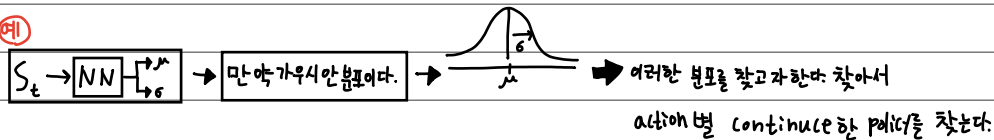
Greedy \Rightarrow $\arg \max_{a_t} Q(s_t, a_t) \Rightarrow$ 가치값이 큰 action을 고른다.

ϵ -Greedy \Rightarrow Greedy 방식에서 ϵ 를 추가한 것

Policy-based

- 강화학습 자체가 좋은 policy $[P(a_t|s_t)]$ 를 찾는 건데 Policy-based는 그 확률 분포를 찾는 것이다.

- (예)



- Policy-based를 많이 쓰는 이유

1) continuous action을 고를 때 좋다.

\Rightarrow Value-based 기반의 DQN을 보면 CNN을 거쳐서 $\begin{bmatrix} Q(s, a) \\ Q(s, \text{왼}) \end{bmatrix}$ 이런 식으로 discrete 하게 뽑아야 한다.

그러다보니 continuous한 action (자동차 핸들)을 control 하기 힘들다.

2) stochastic policy를 찾을 수 있다.

\Rightarrow (예)

s_0	s_1	s_2	s_1	s_3
Die	①	Goal	②	Die

Value-based 입장

• s_0, s_2, s_3 에서는 문제가 생기지 않는다.

• s_1 상태에서는 ① 위치에서 컴퓨터 입장에서 보았을 때 우측으로 가는 것이 Value 값이 크다.

$$\Rightarrow Q(s_1, \text{오}) < Q(s_1, \text{우})$$

하지만 ② 위치에서는 좌측으로 가는 것이 Value 값이 크다.

$$\Rightarrow Q(s_1, \text{오}) > Q(s_1, \text{우})$$

\rightarrow 이런 상황이 계속 반복된다.

Policy-based 입장

• 오른쪽 왼쪽 확률을 $\frac{1}{2}, \frac{1}{2}$ 로 만들면 충분히 가능하다.