

1강 Q-learning

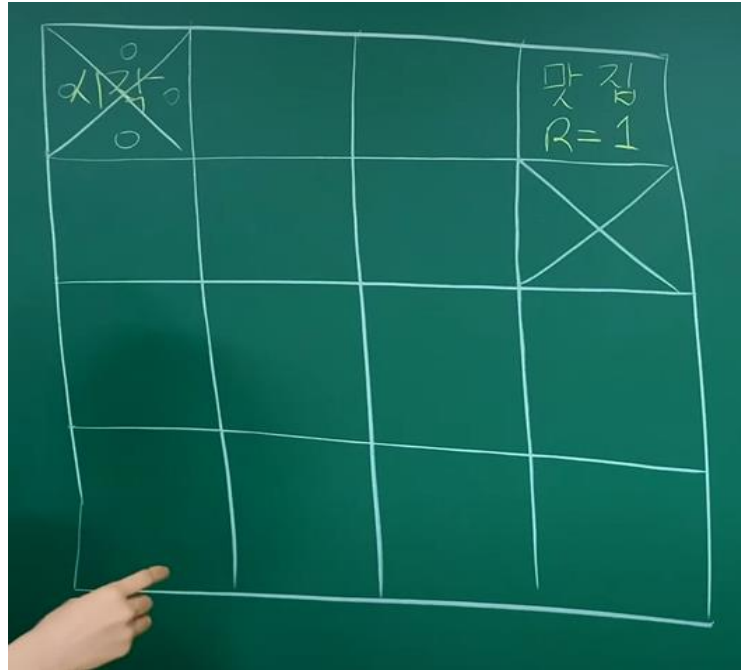
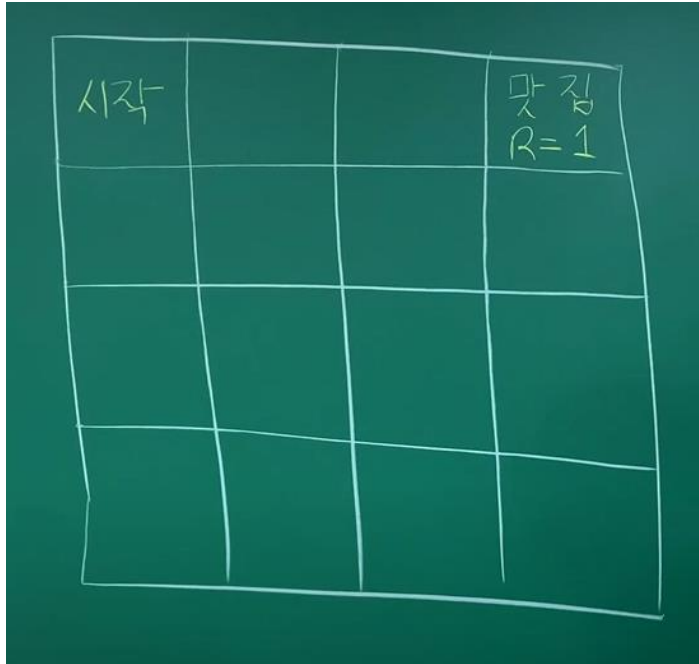
강화학습 연속된 액션을 찾아나가기 <- 보상 줌

action → action → action → • • •

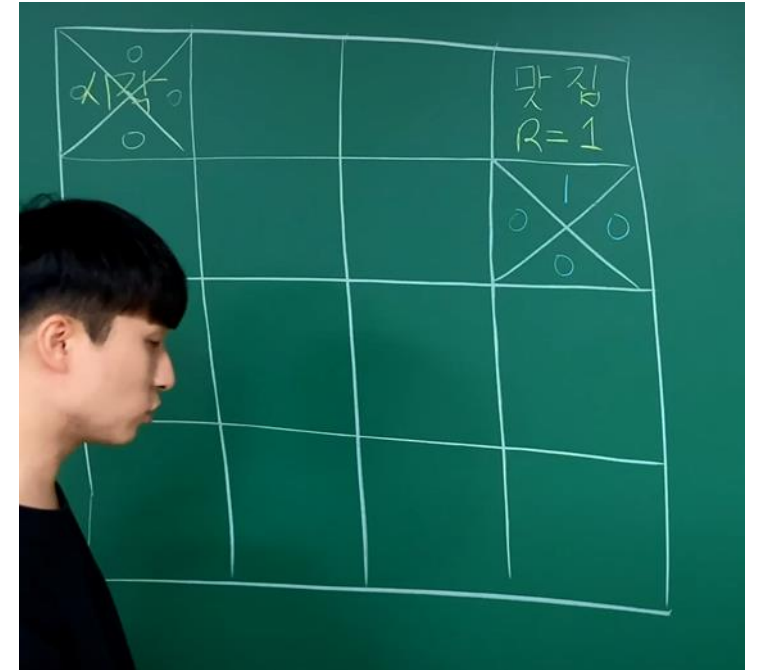
Goal = maximize Reward

Greedy Action

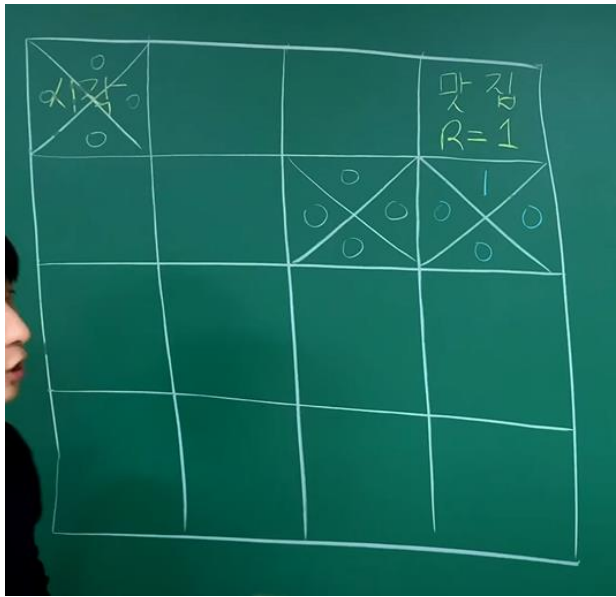
맛집 찾는 방법



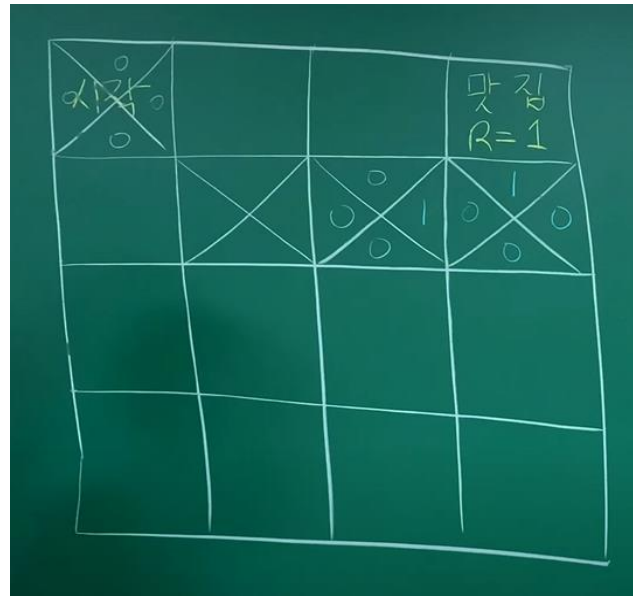
랜덤으로 움직임. 맛집 아래까지 왔다 가정



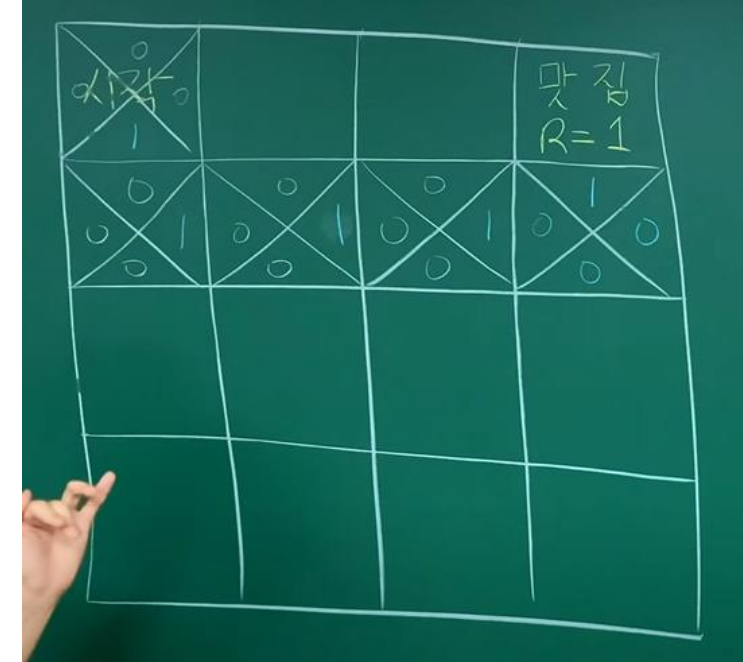
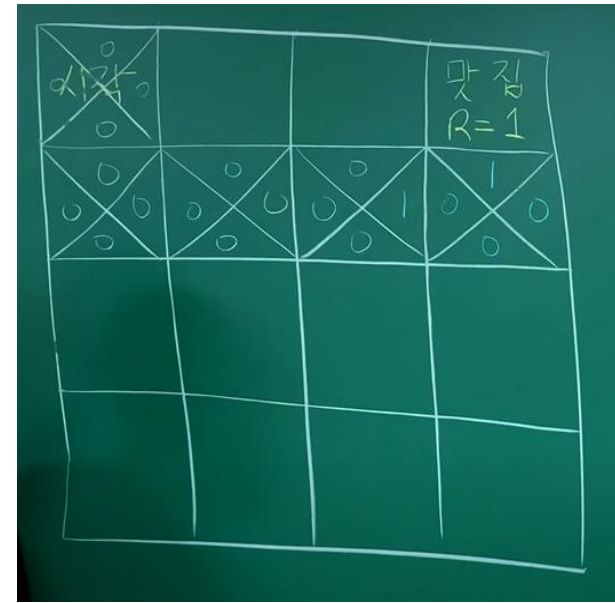
아래서 위로 갔을 때 맛집 도착
따라서 위쪽 그리드에서 1점이
부여됨



다시 랜덤으로 움직이다 점수
있는 칸 옆으로 감 아직 해당
칸은 다 0점



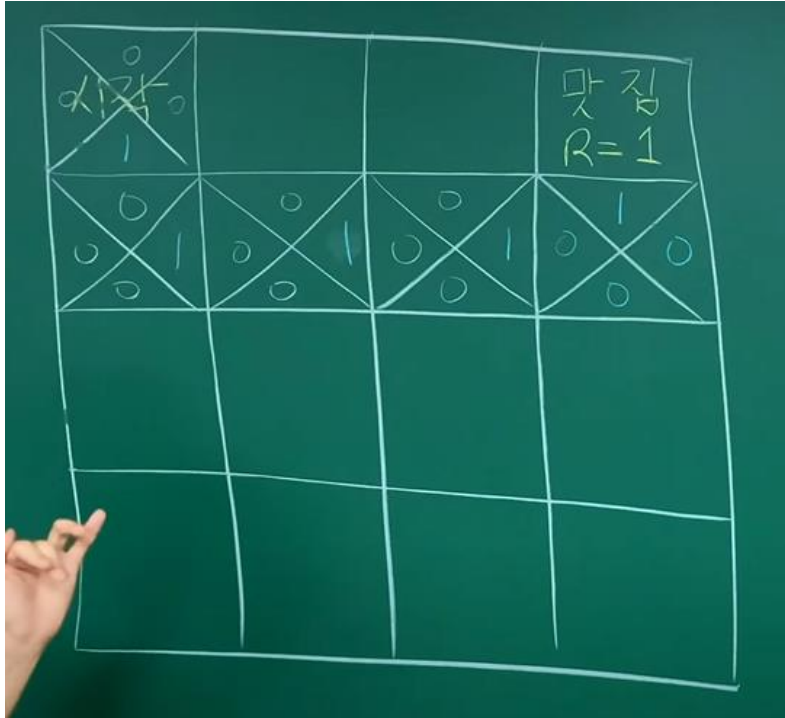
Q-learning 움직이면서 업데이
트 오른쪽 칸에 점수가 있으므
로 오른쪽 칸으로 움직이게 됨
그리고 오른쪽으로 가면 점수가
있으므로 현재 위치 내 오른쪽
그리드가 1점



반복하다보면 옆 그림처럼 됨
이때 더 좋은 길이 있는데 끝나는
문제점 발생

ϵ - Greedy

Exploration



<- 이 상황에서 오른쪽으로만 이동하는 길이 더 빠름

따라서 더 좋은 길을 찾기 위해 탐험을 감

ϵ 만큼은 Random하게, $1-\epsilon$ 만큼은 Greedy하게 움직인다 (ϵ 은 0~1의 사이의 값)

Exploration(랜덤하게 움직임 = 탐험) vs **Exploitation**(Greedy하게 움직임)

이를 통해서

1. 새로운 path를 찾을 수 있다.
2. 새로운 맛집을 찾을 수 있다.
3. ϵ 이 크면 계속 탐험 / 작으면 계속 Greedy

(decaying) ε - Greedy

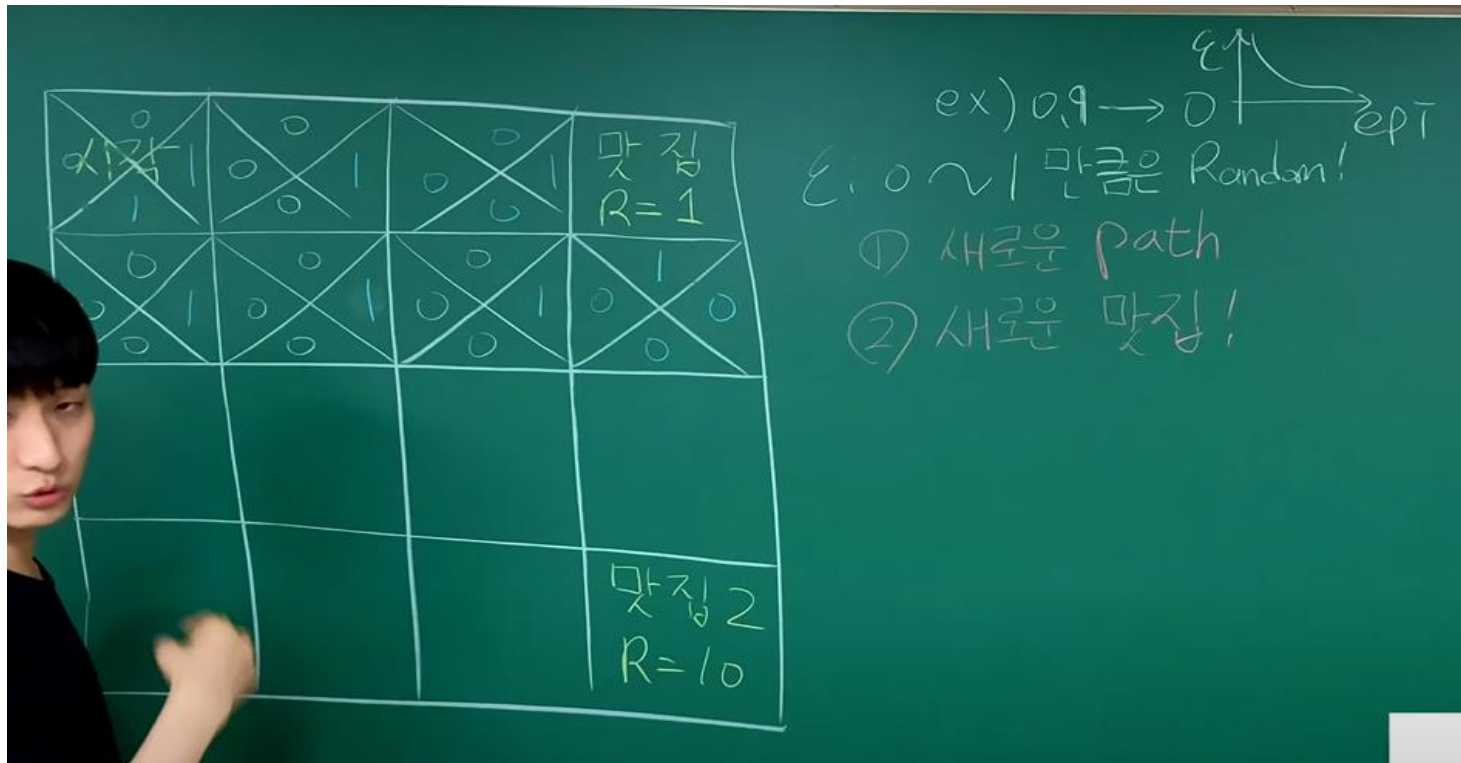
Exploration(랜덤하게 움직임 = 탐험) vs **Exploitation**(Greedy하게 움직임)

➔ 둘 중 하나로 치우치면 안 좋음 그래서 사용하는 방법이 (decaying) ϵ - Greedy

(decaying) ϵ - Greedy

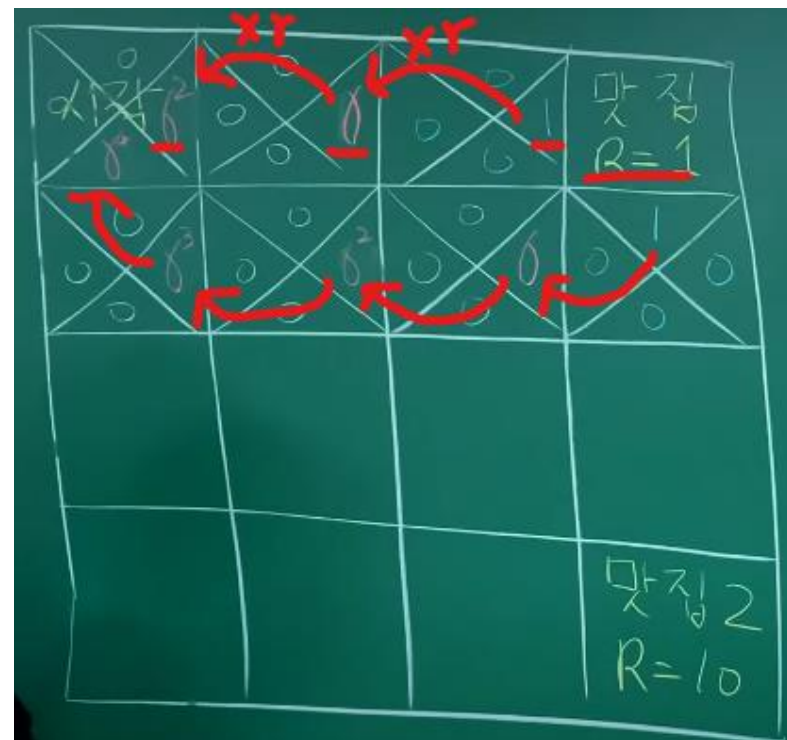
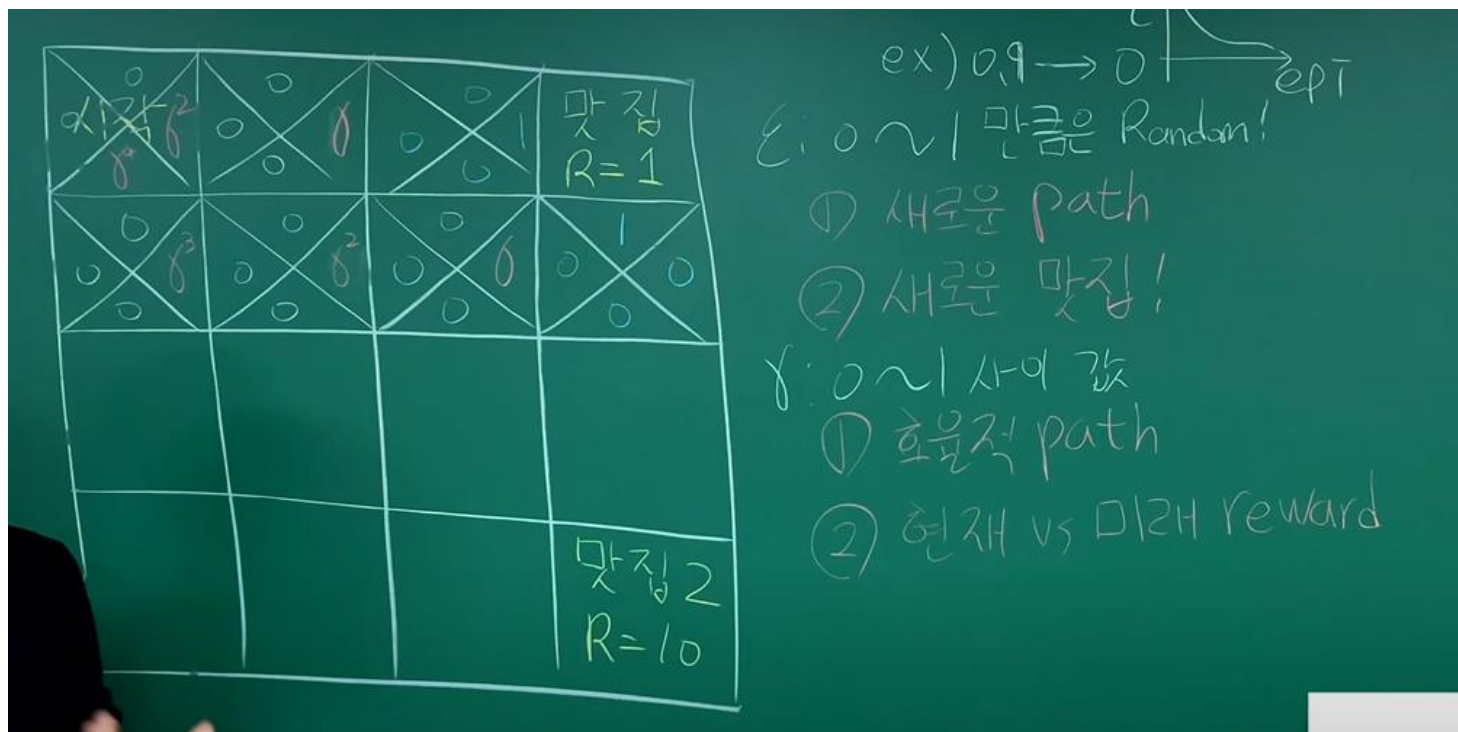
점점 ϵ 값을 줄여가며(탐험을 줄여가며) 학습

Ex) 처음 0.9로 시작 탐험을 더 많이 했다면 반복하면서 0.8, 0.7..0.1 점점 줄여서 탐험 확률을 줄임



<- 첫번째 경로와 두번째 경로 점수 똑같음

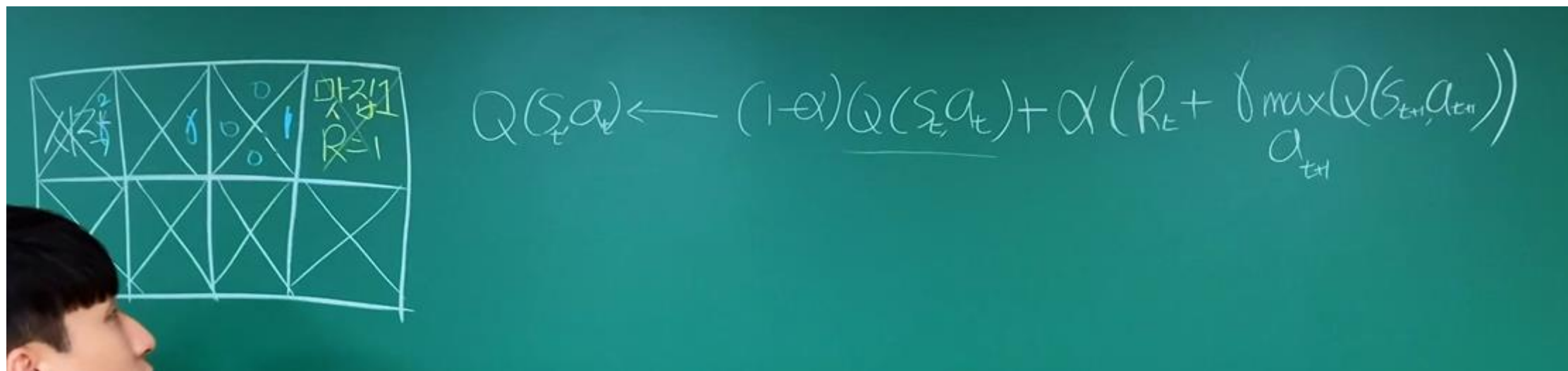
Discount factor



reward 업데이트 할 때 $\text{MAX}(\text{reward})$ 에서 γ 만큼 곱해서 가져온다(γ 는 0~1 사이의 값)
이를 통해서

1. 효율적 path (동일 reward값이라면 경로가 짧은 것)
2. 현재에 관심(γ 가 클 때) vs 미래에 관심(γ 가 작을 때) reward

Q - Update



$$Q(s_t, a_t) \leftarrow (1 - \alpha) \cdot \underbrace{Q(s_t, a_t)}_{\text{old value}} + \underbrace{\alpha}_{\text{learning rate}} \cdot \left(\underbrace{r_t}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \cdot \underbrace{\max_a Q(s_{t+1}, a)}_{\text{estimate of optimal future value}} \right)$$

learned value

Q - value 는 old value + future value 의 합
 α = 새로운 걸(future value) 얼마나 받아 들이냐?

2-1강

Markov Decision Process (MDP)

Markov **Decision** Process

action → **action** → **action** → • • •

Markov Decision Process -> 액션을 해나가는 프로세스

Markov Decision Process

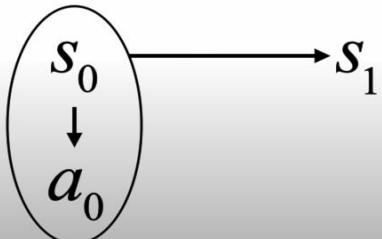
action → action → action → ...

s_0

s_0 : 시작 상태 ex (0,0) 이나 맛집 탐방에서 시작 칸

Markov Decision Process

action → action → action → ...



따라서 s_0 에서 a_0 행동 $\rightarrow s_1$ 이 된다.

Markov Decision Process

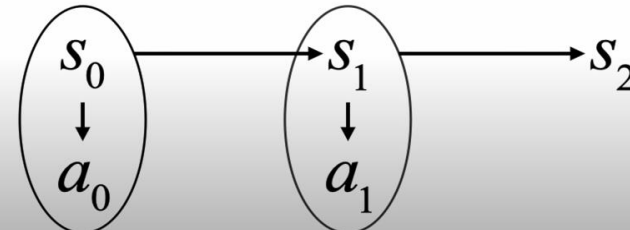
action → action → action → ...



s_0 에서 a_0 행동은 함 ex) 맛집탐방 일 때 액션을 한다 그럼 다른 칸으로 움직이게 됨 다른칸 s_1

Markov Decision Process

action → action → action → ...



$$1. P(a_1 | s_0, a_0, s_1)$$

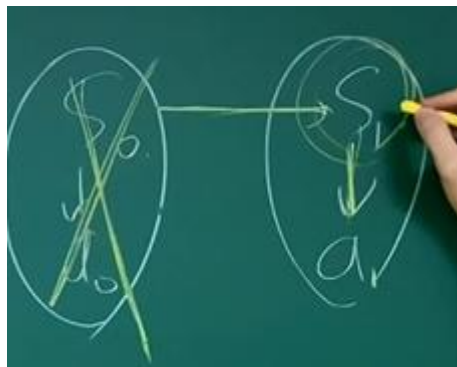
s_0, a_0, s_1 이 주어질 때 a_1 이 일어날 상황(확률)

$$1. P(a_1 | \cancel{s_0}, \cancel{a_0}, s_1)$$

s_0, a_0 이 필요 없어짐

→ s_1 는 s_0, a_0 의해서 넘어간 값

→ s_1 는 s_0, a_0 의해 넘어간 값이기 때문에 s_0, a_0 몰라도 됨



s_1 는 s_0, a_0 정보 흡수되어 있음

$$2. P(s_2 | s_0, a_0, s_1, a_1)$$

s_0 에서 a_0 행동 해서 s_1 로 넘어가게 되고 s_1 에서 a_1 행동했을때 s_2 가 될 상황(확률)



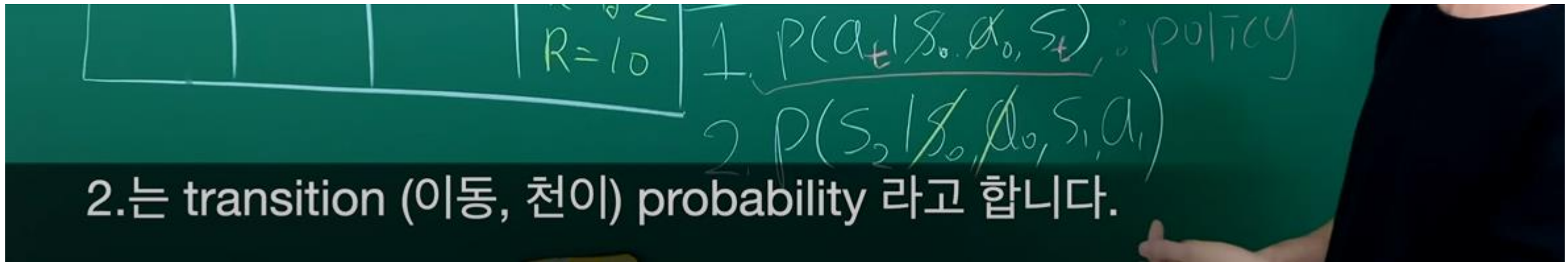
s_0 에서 a_0 행동 해야 s_1 로 넘어갈 수 있듯 둘은 세트 이 세트를 해야 다음 상태로 변함

$$2. P(s_2 | \cancel{s_0}, \cancel{a_0}, s_1, a_1)$$

s_1, a_1 둘 다 남아있음

정리

어떤 액션을 할지 정하는 것이 policy



$$1. P(a_t | s_{t-1}, a_{t-1}, s^t) = P(a_t | s_t) = \text{Policy}$$

어떤 상태에서 어떤 액션을 취할까? = Policy

s_0, a_0 와 상관 없다. 왜냐하면 s_0, a_0 는 s_1 에 모두 담겨 있기 때문

$$2. P(s_{t+1} | s_{t-1}, a_{t-1}, s_t, a_t) = P(s_{t+1} | s_t, a_t)$$

=transition probability

Return

Markov **Decision** Process

action → action → action → . . .

Goal = maximize **Reward**

Markov **Decision** Process

action → action → action → . . .

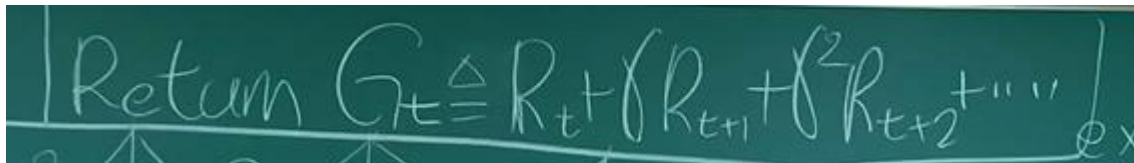
Goal = maximize **Return**



Markov **Decision** Process

action → action → action → . . .

Goal = maximize **Expected Return**


$$\text{Return } G_t \triangleq R_t + \gamma(R_{t+1} + \gamma^2 R_{t+2} + \dots)$$

리턴 = 리워드의 합

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

Return G_t 는 시간 t 이후부터 얻을 수 있는 reward의 합을 의미하며, discount factor γ 를 통해 위 식과 같이 정의된다.

결국 평균 return을 maximize하는 policy를 찾는게 강화학습의 목표

2-2강

상태 가치 함수 V & 행동 가치
함수 Q & Optimal policy

강화학습은 expected return을 최대화 하는 것이 목표이고, 그 목표를 달성하는 action을 뽑아주는 policy들을 찾아냄

expected return을 잘 표현하는 2가지 방법인 state value function, action value function

State value function

지금부터 기대되는 Return

따라서 지금 state에 대한 가치 평가를 내려줌

Action value function

지금 행동으로부터 기대되는 Return

지금 이 state에서의 어떤 행동으로부터 기대되는 Return

Optimal policy

State Value Function 즉 지금부터 기대되는 Return를 최대화 하는것

state value function

$$\text{Return } G_t = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots \quad E[f(x)] = \int f(x) p(x) dx$$

$$\textcircled{1} V(s_t) \triangleq \int_{a_t: a_\infty} G_t p(a_t, s_{t+1}, a_{t+1}, \dots | s_t) da_t: a_\infty$$

행동 하나 하나는 이산적인 과정이지만 그 경우의 수는 셀 수 없을 만큼 많기 때문에 기대값 E 를 계산할 때에 적분을 한다.

지금 상태에서의 a_t 행동의 경우의 수 중에 하나를 뽑고, 그 다음 s_{t+1} 상태에서 가능한 a_{t+1} 경우의 수 중에서 하나를 뽑고... 이것을 무한대로 이어나가는 경우의 수들 각각의 확률을 구해서 리턴과 곱한 뒤 적분해주면 svf 가 나온다.

action value function

$$② Q(s_t, a_t) \triangleq \int_{s_{t+1}, a_{t+1}}^{a_t, a_{\infty}} G_t p(s_{t+1}, a_{t+1}, s_{t+2}, a_{t+2}, \dots | s_t, a_t) ds_{t+1} da_{t+1}$$

현재 state에서 어떤 행동을 했을때부터 쪽 진행한 기대되는 return

Optimal Policy

Return $G_t = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots$ $E[f(x)] = \int f(x) p(x) dx$

$$\textcircled{1} V(s_t) \triangleq \int_{a_t: a_\infty} G_t p(a_t, s_{t+1}, a_{t+1}, \dots | s_t) da_t: a_\infty \stackrel{?}{=} \text{maximize over } \left(\begin{matrix} p(a_t | s_t) \\ p(a_{t+1} | s_{t+1}) \\ \vdots \\ p(a_\infty | s_\infty) \end{matrix} \right) \text{ \textit{Optimal Policy}}$$

$$\textcircled{2} Q(s_t, a_t) \triangleq \int_{s_{t+1}: a_\infty} G_t p(s_{t+1}, a_{t+1}, s_{t+2}, a_{t+2}, \dots | s_t, a_t) ds_{t+1}: a_\infty$$

$V(s_t)$ 를 maximize하는 action들의 모임

2-3강. 벨만 방정식 (Bellman equation)

Bellman equation

Return $G_t = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots$ $E[f(x)] = \int f(x) p(x) dx$

① $V(s_t) \triangleq \int_{a_t, a_{\infty}} G_t p(a_t, s_{t+1}, a_{t+1}, \dots | s_t) da_t : a_{\infty} \stackrel{\text{maximize}}{\Rightarrow} \left(\begin{matrix} p(a_t | s_t) \\ p(a_{t+1} | s_{t+1}) \\ \vdots \\ p(a_{\infty} | s_{\infty}) \end{matrix} \right)$ Optimal Policy

② $Q(s_t, a_t) \triangleq \int_{s_{t+1}, a_{\infty}} G_t p(s_{t+1}, a_{t+1}, s_{t+2}, a_{t+2}, \dots | s_t, a_t) ds_{t+1} : a_{\infty}$

v라는 것을 Q라고도 표현할 수 있고, 이를 바꿔서도 가능하고
다음 t+1로도 표현할 수 있게 해주는 방정식

함수 뒤에 뺄 수 있음 Bayes rule 적용하기

$$p(x, y) = p(x|y) p(y)$$

$$p(x, y|z) = p(x|y, z) p(y|z)$$


위 방법으로 a_t 를 뺀

$$\textcircled{1} V(s_t) \triangleq \int_{a_t: a_{\infty}} G_t \frac{p(a_t, s_{t+1}, a_{t+1}, \dots | s_t)}{\textcircled{2} p(s_{t+1}, a_{t+1}, \dots | s_t, a_t) p(a_t | s_t)} da_t : a_{\infty}$$

$$\textcircled{1} V(s_t) \triangleq \int_{a_t: a_{\infty}} G_t \frac{p(a_t, s_{t+1}, a_{t+1}, \dots | s_t)}{\textcircled{2} p(s_{t+1}, a_{t+1}, \dots | s_t, a_t) p(a_t | s_t)} da_t : a_{\infty} \stackrel{?}{=} \max$$

$$\textcircled{2} Q(s_t, a_t) \triangleq \int_{s_{t+1}: a_{\infty}} G_t p(s_{t+1}, a_{t+1}, s_{t+2}, a_{t+2}, \dots | s_t, a_t) ds_{t+1}: a_{\infty}$$

$$\textcircled{1} - \textcircled{1} = \int_{a_t} \underbrace{\int_{s_{t+1}: a_{\infty}} G_t p(s_{t+1}, a_{t+1}, \dots | s_t, a_t) ds_{t+1}: a_{\infty}}_{Q(s_t, a_t)} p(a_t | s_t) da_t$$

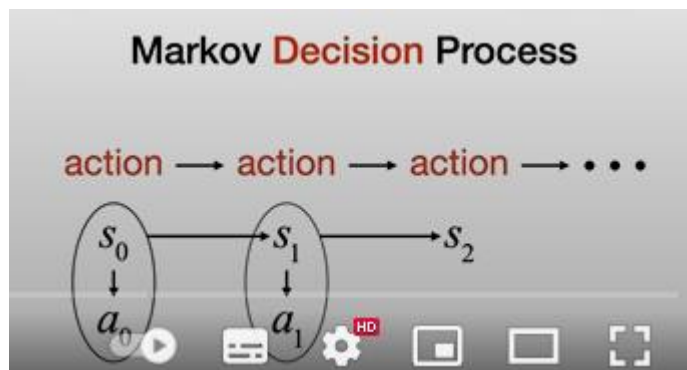
$$\begin{aligned}
 ① V(s_t) &\triangleq \int_{a_t: a_0} G_t \frac{p(a_t, s_{t+1}, a_{t+1}, \dots | s_t)}{p(a_t | s_t)} da_t : a_0 \stackrel{\text{maximize}}{=} \\
 ② Q(s_t, a_t) &\triangleq \int_{s_{t+1}: a_0} G_t p(s_{t+1}, a_{t+1}, s_{t+2}, a_{t+2}, \dots | s_t, a_t) ds_{t+1} : a_0 \\
 ① - ① &= \int_{a_t} \int_{s_{t+1}: a_0} G_t \underbrace{p(s_{t+1}, a_{t+1}, \dots | s_t, a_t)}_{Q(s_t, a_t)} p(s_{t+1}: a_0) p(a_t | s_t) da_t \\
 &= \int_{a_t} Q(s_t, a_t) p(a_t | s_t) da_t
 \end{aligned}$$


V to Q

State에서 기대되는 return 값은 아래처럼 이해 가능
state안의 action들의 Q값을 평균을 내는 것.

$$\textcircled{2} P(a_{t+1}, \dots | s_t, a_t, s_{t+1})$$

2번째 a_t 와 s_{t+1} 도 같이 뺌



Markov Decision process에 따라 s_{t+1} 은 s_t at모두의 값을 가지고 있으므로 이 둘을 생략이 가능

$$\textcircled{2} P(a_{t+1}, \dots | s_t, a_t, s_{t+1})$$

전체 식

$$\textcircled{2} P(a_{t+1}, \dots | s_t, a_t, s_{t+1}) P(a_t, s_{t+1} | s_t)$$

$$\textcircled{1}-\textcircled{2} = \int_{a_t, s_{t+1}} \int_{a_{t+1}, a_{\infty}} (R_t + \gamma G_{t+1}) p(a_{t+1}, \dots | s_{t+1}) da_{t+1}, a_{\infty} p(a_t, s_{t+1} | s_t) da_t, s_{t+1}$$


위 감마 G_{t+1} 어디서 $\rightarrow G_t$ 를 옆과 같이 다르게 표현가능

$$\text{Return } G_t = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots$$

① $V(s_t) \triangleq \int_{a_t, a_{\infty}} G_t p(a_t, s_{t+1}, a_{t+1}, \dots | s_t) da_t, a_{\infty} \equiv \text{maximize 하는 policy + optimal policy}$

② $Q(s_t, a_t) \triangleq \int_{s_{t+1}, a_{\infty}} G_t p(s_{t+1}, a_{t+1}, s_{t+2}, a_{t+2}, \dots | s_t, a_t) ds_{t+1}, a_{\infty}$

①-① = $\int_{a_t} \underbrace{\int_{s_{t+1}, a_{\infty}} G_t p(s_{t+1}, a_{t+1}, \dots | s_t, a_t) ds_{t+1}, a_{\infty}}_{Q(s_t, a_t)} p(a_t | s_t) da_t$

= $\int_{a_t} Q(s_t, a_t) p(a_t | s_t) da_t$ 

①-② = $\int_{a_t, s_{t+1}} \int_{a_{t+1}, a_{\infty}} (R_t + \gamma G_{t+1}) p(a_{t+1}, \dots | s_{t+1}) da_{t+1}, a_{\infty} p(a_t, s_{t+1} | s_t) da_t, s_{t+1}$

= $\int_{a_t, s_{t+1}} (R_t + \gamma V(s_{t+1})) p(a_t, s_{t+1} | s_t) da_t, s_{t+1}$

$p(x, y) = p(x|y) p(y)$
 $p(x, y | z) = p(x|y, z) p(y|z)$

$$\begin{aligned}
 \textcircled{1} &= \int_{a_t} \int_{s_{t+1}: a_{\infty}} G_t \underbrace{p(s_{t+1}, a_{t+1}, \dots | s_t, a_t)}_{Q(s_t, a_t)} \underbrace{p(a_{t+1} | s_t)}_{\text{transition policy}} da_{t+1} ds_{t+1} \\
 &= \int_{a_t} Q(s_t, a_t) p(a_t | s_t) da_t \quad \boxed{V(s_{t+1})} \\
 \textcircled{1} - \textcircled{2} &= \int_{a_t, s_{t+1}} \left(R_t + \gamma \underbrace{G_{t+1}}_{\text{red}} \right) \underbrace{p(a_{t+1}, \dots | s_{t+1})}_{\text{red}} da_{t+1} a_{\infty} \underbrace{p(a_t, s_{t+1} | s_t)}_{\text{red}} da_t ds_{t+1} \\
 &= \int_{a_t, s_{t+1}} (R_t + \gamma V(s_{t+1})) \underbrace{p(a_t, s_{t+1} | s_t)}_{\text{red}} da_t ds_{t+1}
 \end{aligned}$$

$p(x, y) = p(x|y)p(y)$
 $p(x, y|z) = p(x|y, z)p(y|z)$


아래식처럼 최대화 하는 policy 구하는 것이 중요한데 위 적분식 다르게 표현한 식(빨강색)을 보면 policy 구할 수 있게 됨

$$\textcircled{1} V(s_t) \triangleq \int_{a_t: a_{\infty}} G_t \underbrace{p(a_t, s_{t+1}, a_{t+1}, \dots | s_t)}_{\text{red}} da_{t+1} a_{\infty} \triangleq \text{maximize over policy} \quad \text{optimal policy}$$

Q value Bellman equation

Q to Vst+1

Q를 next V로 표현이 가능, 숨어있는 V를 찾기



Return $G_t = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots$, $E[f(x)] = \int f(x)p(x)dx$

① $V(s_t) \triangleq \int G_t \underbrace{p(a_{t+1}, s_{t+1}, a_{t+2}, \dots | s_t)}_{\text{① } p(s_{t+1}, a_{t+1}, \dots | s_t, a_t) p(a_t | s_t)} \underbrace{d a_{t+1} : a_{\infty}}_{\text{② } p(a_{t+1}, \dots | s_t, a_t, s_{t+1}) p(a_t, s_{t+1} | s_t)} \triangleq \text{maximize 하는 policy} \quad \text{+ optimal policy}$

② $Q(s_t, a_t) \triangleq \int G_t \underbrace{p(s_{t+1}, a_{t+1}, s_{t+2}, a_{t+2}, \dots | s_t, a_t)}_{\text{① } p(a_{t+1}, s_{t+1}, a_{t+2}, \dots | s_t, a_t, s_{t+1}) p(s_{t+1} | s_t, a_t)} d s_{t+1} : a_{\infty}$

$p(x, y) = p(x|y) p(y)$
 $p(x, y|z) = p(x|y, z) p(y|z)$

Markov Decision process에 따라 s_{t+1} 은 s_t at 모두의 값을 가지고 있으므로 이 둘을 생략이 가능

Return $G_t = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots$, $E[F(x)] = \int f(x) p(x) dx$

① $V(s_t) \triangleq \int G_t p(a_{t+1}, s_{t+1}, a_{t+2}, \dots | s_t) da_{t+1} a_{t+2} \dots \stackrel{?}{=} \text{maximize over policy} \rightarrow \text{optimal policy}$

② $Q(s_t, a_t) \triangleq \int G_t p(s_{t+1}, a_{t+1}, s_{t+2}, a_{t+2}, \dots | s_t, a_t) ds_{t+1} a_{t+2} \dots$

$V(s_{t+1}) \leftarrow \int a_{t+1} a_{t+2} \dots \int p(a_{t+1}, s_{t+2}, a_{t+2}, \dots | s_t, a_t, s_{t+1}) p(s_{t+1} | s_t, a_t) p(x, y) = p(x|y) p(y)$

② - ① = $\int \int_{s_{t+1}} \int_{a_{t+1} a_{t+2} \dots} (R_t + \gamma G_{t+1}) p(a_{t+1} a_{t+2} \dots | s_{t+1}) da_{t+1} a_{t+2} \dots p(s_{t+1} | s_t, a_t) ds_{t+1} p(x, y | z) = p(x|y, z) p(y | z)$

= $\int_{s_{t+1}} (R_t + \gamma V(s_{t+1})) p(s_{t+1} | s_t, a_t) ds_{t+1}$

마지막 식 Q를 next V로 표현된걸 확인 가능

Qst+1로 표현

$$\begin{aligned} \textcircled{2} - \textcircled{2} &= \int \int (R_t + \beta Q(s_{t+1}, a_{t+1})) p(s_{t+1}, a_{t+1} | s_t, a_t) ds_{t+1}, a_{t+1} \\ &= \int (R_t + \beta Q(s_{t+1}, a_{t+1})) p(s_{t+1}, a_{t+1} | s_t, a_t) ds_{t+1}, a_{t+1} \end{aligned}$$

$$\int_{s_{t+1}, a_{t+1}} (R_{t+1} + \gamma Q(s_{t+1}, a_{t+1})) p(s_{t+1}, a_{t+1} | s_t, a_t) ds_{t+1}, da_{t+1}$$

풀어쓰기, 적분할 거리를 줄여준다.

$$\frac{p(a_{t+1} | s_t, a_t, s_{t+1})}{\times p(s_{t+1} | s_t, a_t)}$$

풀어쓰기, 적분할 거리를 줄여준다.
transition과 policy 나타낼 수 있음