

복습)

Q-learning의 핵심

Q에는 1) 내가 있는 상태 state, 2) 내가 하는 행동 action input과 출력 3)quality(reward)가 있다.

frozen lake게임에서 process

$Q(s1, left)=0$, $Q(s1, RIGHT)=1.5$ $Q(s1, UP)=0$, $Q(s1, Down)=0.3$ 이 제공된다면 max value를 선택한다. max의 argument RIGHT를 선택한다

Max Q = $\max_{a'} Q(s, a')$ Q가 가질 수 있는 최대값 a'

$\pi^*(s) = \operatorname{argmax}_a Q(s, a)$ Q가 최대값이 되게끔 하는 argument $a \rightarrow$ optimal Π

Q러닝의 가정

S'에서 Q가 있어야 된다.

s에서 a를 하고 s'으로 갈거다. a를 하면 r을 받는다

$Q(s,a)$ 이라면 $Q(s',a')$ 이 존재할거다

$Q = r + \max Q(s',a')$ 이게 되고 이 Q를 학습하는 거다.



현재 $Q(s,a)=?$

S에서 $Q = r + \max Q(s',a')$

Future reward

$s_0, a_0, r_1, s_1, a_1, r_2, \dots, s_{n-1}, a_{n-1}, r_n, s_n$
 ↑ state ↑ action ↑ reward ↑ Terminal state

$$R = r_1 + r_2 + r_3 + \dots + r_n$$

$$R_t = r_t + r_{t+1} + r_{t+2} + \dots + r_n$$

Reward의 총합은 $R = r_1 + r_2 + r_3 + \dots + r_n$

$R_{t+1} = r_{t+1} + r_{t+2} + r_{t+3} + \dots + r_n$

$R_t = r_t + R_{t+1}$ 가 된다.

$R^*_t = r_t + \max R_{t+1}$ 가 되고, 결국

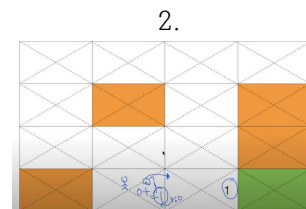
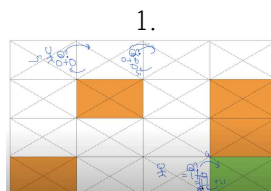
$Q(s,a) = r + \max Q(s',a')$ 가 된다.

Learning Q(s,a): Table
initial Q values are 0

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

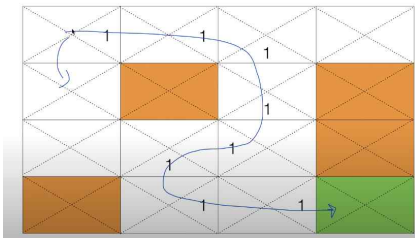
initial Q-table

0으로 초기화



1 -> 2 -> ... n번 진행하면서 각 상태의 Q를 update함

Exploit VS Exploration



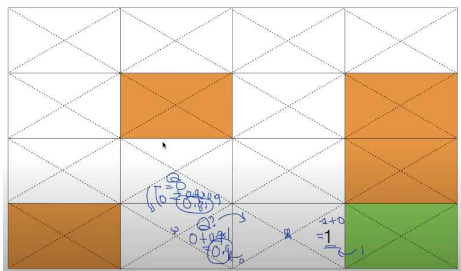
exploration을 통해 가본 길 말고 새로운 길을 탐색한다

"Select an action a and execute it"

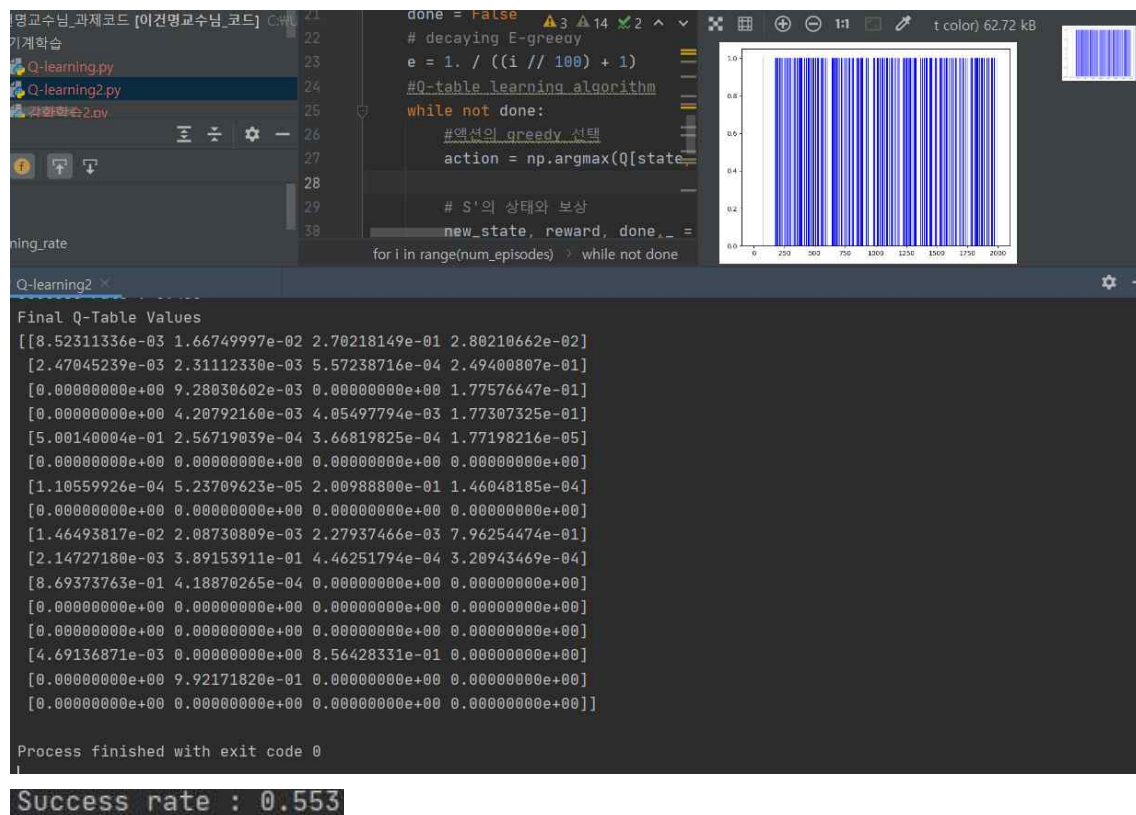
a를 선택할 때 exploit이나 exploration을 선택한다!

Discounted future reward

$$R_t = r_t + \Gamma r_{t+1} + \Gamma^2 r_{t+2} + \Gamma^3 r_{t+3} + \Gamma^4 r_{t+4} + \dots$$



코드 실행 결과



deterministic일땐 $r + \Gamma \max_{a'} Q(s', a')$ 이 Q 이고

스톱케스틱(미끄러운 lake) 일땐 $(1 - \Gamma)Q(s, a) + \Gamma(r + \Gamma \max_{a'} Q(s', a'))$

100*100에서 Q-table을 사용할 수 있을까?

$100 \times 100 \times 4(\text{action}) = 40000$ 개

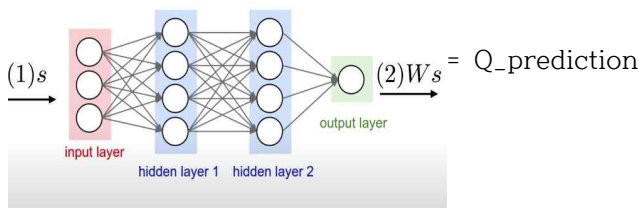
실제 80*80pixel + 2color에서는?

$2^{(80 \times 80)}$ 이다

Q-table로는 현실 세계의 문제를 못 푼다

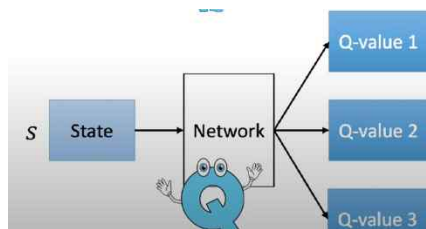
Q-network

상태 하나에 모든 가능한 액션값의 Q-value을 찾는다



$$\text{cost}(W) = (Ws - y)^2$$

$$Q^* = y = r + \gamma \max Q(s')$$



$$\text{cost}(W) = (Q_{\text{pred}} - Q^*)^2$$

cost를 minimize하도록 학습하는 것이 목표이다

$Q^{\wedge} = Q \text{ hat} =$ 예측값

$Q^{\wedge}(s,a|\theta)$ θ =weight 즉 weight에 따라서 s,a가 달라진다.

$$\rightarrow W_s = Q_{\text{prediction}} = Q^{\wedge}(s,a|\theta) \sim Q^*(s,a)$$

수학적 기호 표기

$$\min_{\theta} \sum_{t=0}^T [\hat{Q}(s_t, a_t | \theta) - (r_t + \gamma \max_{a'} \hat{Q}(s_{t+1}, a' | \theta))]^2$$

$\text{cost}(W) = (Q_{\text{pred}} - Q^*)^2$ 이다.

이제 cost를 최소화 하는 θ 를 구하는 것이 문제이다.

deep Q learning Algorithm

Q 네트워크의 w를 초기화 한다

1번째 상태를 가져온다. 그리고 전처리 한다 (ϕ) [단순히 S라고 생각하면 된다]

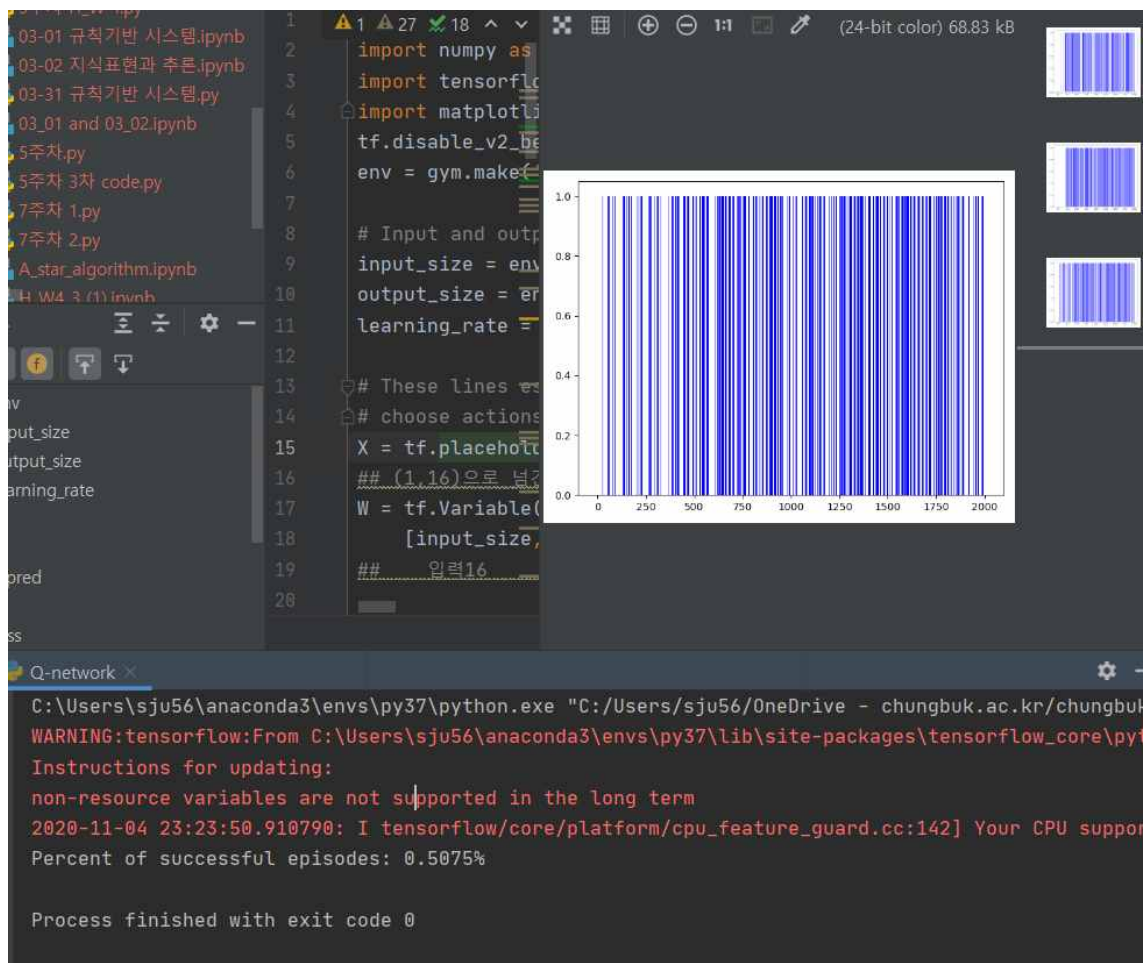
for

에그리디로 랜덤하게 action을 취하거나 현재 Q에서 max값을 고른다

그리고 a을 한다. 그리고 S_{t+1} 와 r_t 를 받는다

$$y_j = \begin{cases} r_j & \text{마지막 일 때} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{마지막이 아닐 때} \end{cases}$$

그리고 y_i 를 구하고 cost값을 최소화 하는 θ 를 찾으면 된다 $(y_j - Q(\phi_j, a_j; \theta))^2$



Q-table보다 성능이 안 좋다. 0.5%정도 나온다.