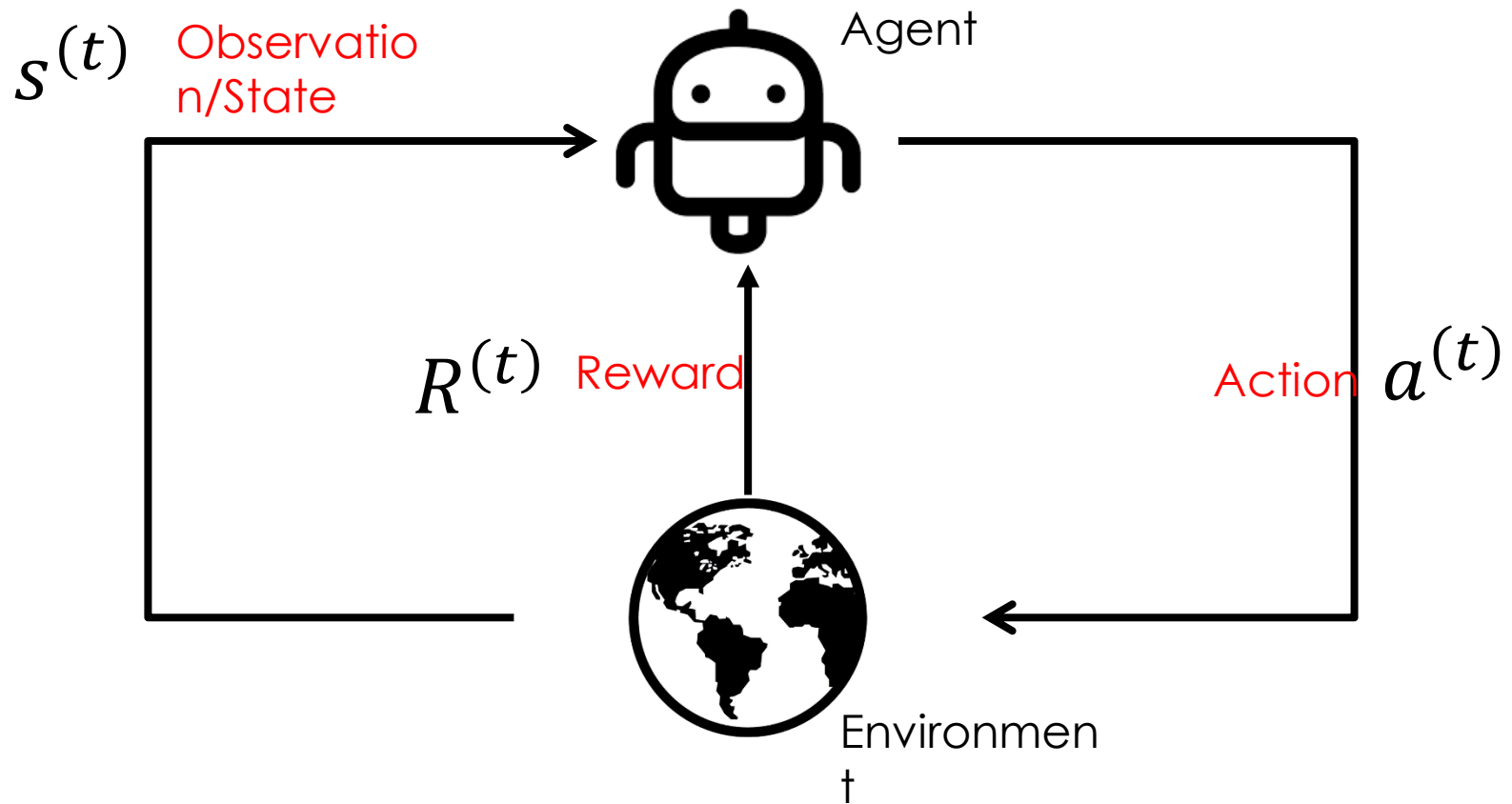


智慧化企業整合

Intelligent Integration of Enterprise

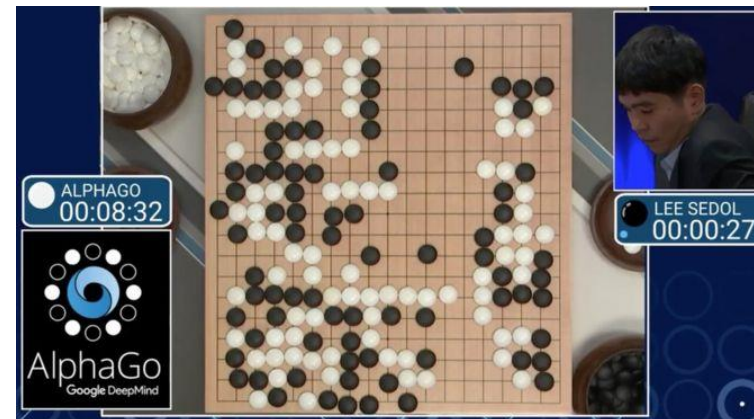
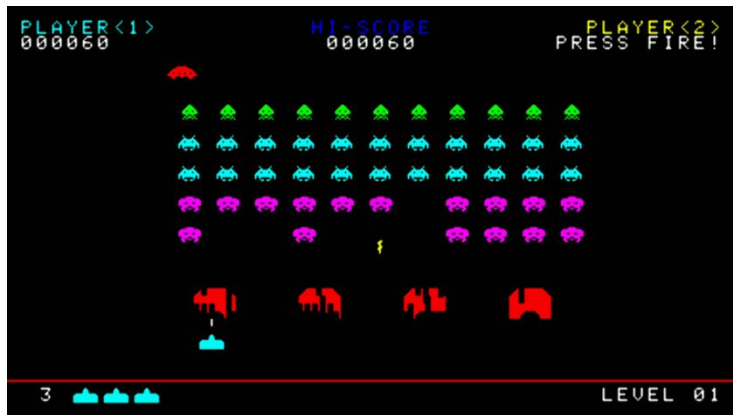
Reinforcement Learning

Reinforcement Learning



Reinforcement Learning

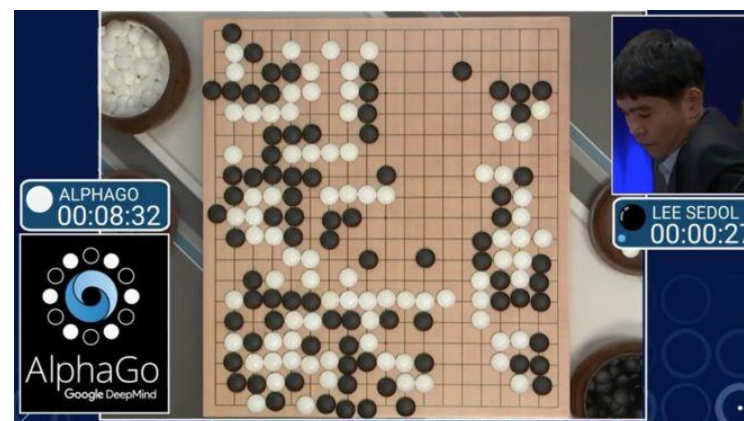
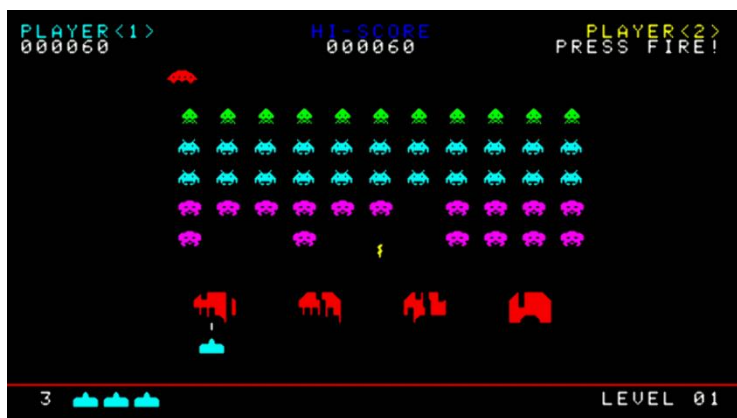
- An agent sees state, takes an action and receives a reward
- The state may change due to the action
- The state may also change without action
- Ex: space invader, AlphaGo...



Reinforcement Learning

- Goal: to learn the best policy that maximize the total reward

$$\pi^*(s^{(t)}) = a^{(t)} \text{ that maximizes } \sum_t R^{(t)}$$



Markov process

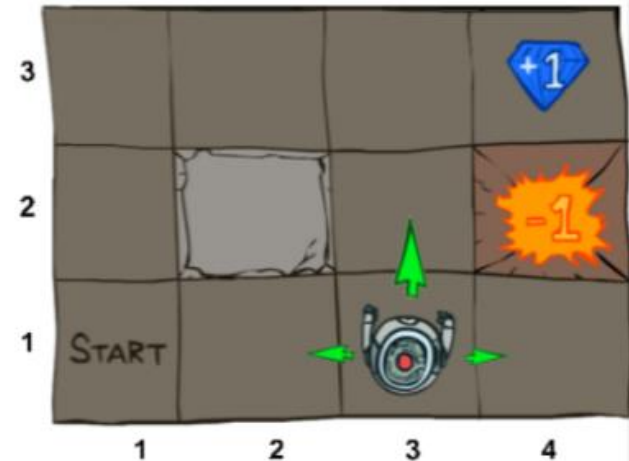
- A random process is a collection of time-indexed random variables
- A random process is called a Markov process if it satisfies the Markov property:

$$\mathbf{P}(\mathbf{s}^{(t+1)} | \mathbf{s}^{(t)}, \mathbf{s}^{(t-1)}, \dots) = \mathbf{P}(\mathbf{s}^{(t+1)} | \mathbf{s}^{(t)})$$

- The future is independent of the past given the present
- The state captures all relevant information from history

Define Markov Decision Process (MDP)

- S the state space, A the action space
- Start state s^0
- $P(s'|s; a)$ the transition distribution, fixed over t
- $R(S, a, s')$ the deterministic reward function
- $\gamma \in [0,1]$ the discount factor



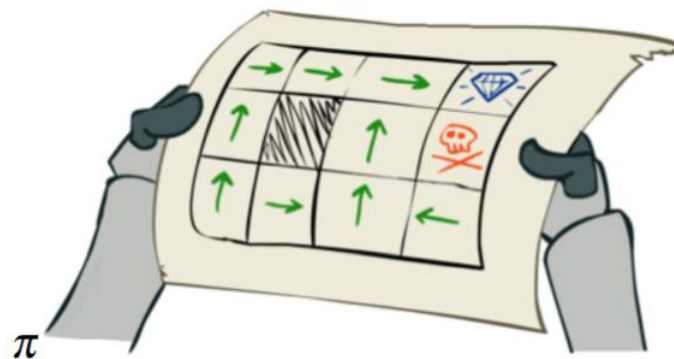
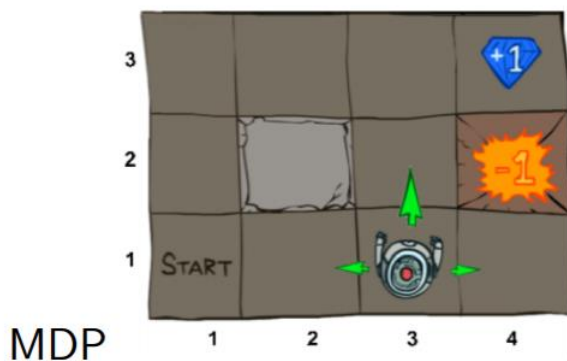
MDP

Given a policy $\pi(s) = \mathbf{a}$, an MDP proceeds as follows:

$$\mathbf{s}^{(0)} \xrightarrow{\mathbf{a}^{(0)}} \mathbf{s}^{(1)} \xrightarrow{\mathbf{a}^{(1)}} \dots \xrightarrow{\mathbf{a}^{(H-1)}} \mathbf{s}^{(H)},$$

with the accumulative reward

$$R(\mathbf{s}^{(0)}, \mathbf{a}^{(0)}, \mathbf{s}^{(1)}) + \gamma R(\mathbf{s}^{(1)}, \mathbf{a}^{(1)}, \mathbf{s}^{(2)}) + \dots + \gamma^{H-1} R(\mathbf{s}^{(H-1)}, \mathbf{a}^{(H-1)}, \mathbf{s}^{(H)})$$



MDP

Given a policy $\pi(s) = \mathbf{a}$, an MDP proceeds as follows:

$$\mathbf{s}^{(0)} \xrightarrow{\mathbf{a}^{(0)}} \mathbf{s}^{(1)} \xrightarrow{\mathbf{a}^{(1)}} \dots \xrightarrow{\mathbf{a}^{(H-1)}} \mathbf{s}^{(H)},$$

with the accumulative reward

$$R(\mathbf{s}^{(0)}, \mathbf{a}^{(0)}, \mathbf{s}^{(1)}) + \gamma R(\mathbf{s}^{(1)}, \mathbf{a}^{(1)}, \mathbf{s}^{(2)}) + \dots + \gamma^{H-1} R(\mathbf{s}^{(H-1)}, \mathbf{a}^{(H-1)}, \mathbf{s}^{(H)})$$

- To acquire rewards ASAP
- Different accumulative rewards in different trials

Goal

- Given a policy π , the expected accumulative rewards collected by taking actions can be expressed as:

$$V_{\pi} = \mathbb{E}_{\mathbf{s}^{(0)}, \dots, \mathbf{s}^{(H)}} \left(\sum_{t=0}^H \gamma^t R(\mathbf{s}^{(t)}, \pi(\mathbf{s}^{(t)}), \mathbf{s}^{(t+1)}); \pi \right)$$

- Our goal is to find the optimal policy
$$\pi^* = \arg \max_{\pi} V_{\pi}$$

Value Iteration

Optimal Value Function

$$\pi^* = \arg \max_{\pi} E_{\mathbf{s}^{(0)}, \dots, \mathbf{s}^{(H)}} \left(\sum_{t=0}^H \gamma^t R(\mathbf{s}^{(t)}, \pi(\mathbf{s}^{(t)}), \mathbf{s}^{(t+1)}); \pi \right)$$

- Optimal value function:
Maximum expected accumulative rewards when starting from state s and acting optimally for h steps

每一步都是最佳的policy

$$V^{*(h)}(\mathbf{s}) = \max_{\pi} E_{\mathbf{s}^{(1)}, \dots, \mathbf{s}^{(h)}} \left(\sum_{t=0}^h \gamma^t R(\mathbf{s}^{(t)}, \pi(\mathbf{s}^{(t)}), \mathbf{s}^{(t+1)}) \mid \mathbf{s}^{(0)} = \mathbf{s}; \pi \right)$$

Optimal Value Function

$$V^{*(h)}(s) = \max_{\pi} E_{s^{(1)}, \dots, s^{(h)}} \left(\sum_{t=0}^h \gamma^t R(s^{(t)}, \pi(s^{(t)}), s^{(t+1)}) \mid s^{(0)} = s; \pi \right)$$

- Having $V^{*(H-1)}(s)$ for each s , we can solve π^* by

由前一步可以找出現在這一步的最佳policy

$$\pi^*(s) = \arg \max_a \sum_{s'} P(s' | s; a) [R(s, a, s') + \gamma V^{*(H-1)}(s')], \forall s$$

直接根據 accumulative reward最高的的policy行動

Dynamic Programming

$$V^{*(h)}(\mathbf{s}) = \max_{\pi} E_{\mathbf{s}^{(1)}, \dots, \mathbf{s}^{(H)}} \left(\sum_{t=0}^h \gamma^t R(\mathbf{s}^{(t)}, \pi(\mathbf{s}^{(t)}), \mathbf{s}^{(t+1)}) \mid \mathbf{s}^{(0)} = \mathbf{s}; \pi \right)$$

- $h = H - 1$:

$$V^{*(H-1)}(\mathbf{s}) = \max_{\mathbf{a}} \sum_{\mathbf{s}'} P(\mathbf{s}' | \mathbf{s}; \mathbf{a}) [R(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma V^{*(H-2)}(\mathbf{s}')], \forall \mathbf{s}$$

- $h = H - 2$:

$$V^{*(H-2)}(\mathbf{s}) = \max_{\mathbf{a}} \sum_{\mathbf{s}'} P(\mathbf{s}' | \mathbf{s}; \mathbf{a}) [R(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma V^{*(H-3)}(\mathbf{s}')], \forall \mathbf{s}$$

- $h = 0$:

$$V^{*(0)}(\mathbf{s}) = \max_{\mathbf{a}} \sum_{\mathbf{s}'} P(\mathbf{s}' | \mathbf{s}; \mathbf{a}) [R(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma V^{*(-1)}(\mathbf{s}')], \forall \mathbf{s}$$

- $h = -1$:

$$V^{*(-1)}(\mathbf{s}) = 0, \forall \mathbf{s}$$

Policy Iteration

Policy Iteration

- Given a policy π , define value function:

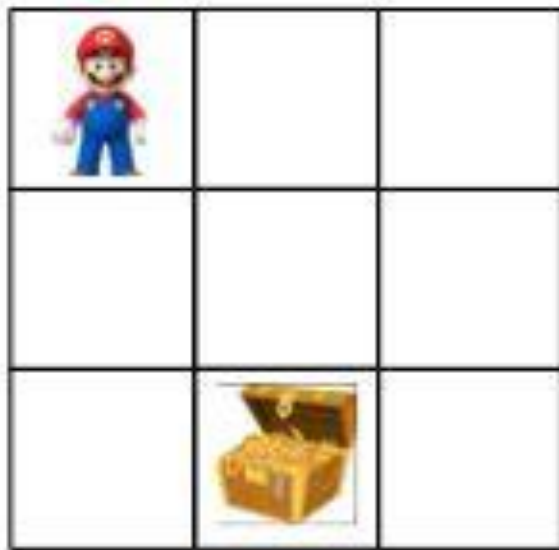
$$V_{\pi}(\mathbf{s}) = \mathbb{E}_{\mathbf{s}^{(1)}, \dots} \left(\sum_{t=0}^{\infty} \gamma^t R(\mathbf{s}^{(t)}, \pi(\mathbf{s}^{(t)}), \mathbf{s}^{(t+1)}) \mid \mathbf{s}^{(0)} = \mathbf{s}; \pi \right)$$

- Start from state s and act based on policy π

Policy Iteration

1. 隨機選擇一個策略當作初始值
2. 根據目前的策略計算Value function的值
3. 根據Value function的值計算目前狀態中最好的action
4. 修改策略

Policy Iteration



1. 假設有一個 3×3 的尋寶遊戲
2. 初始化: 不管在哪, 一律往下走
3. 策略評估:
 - 如果寶藏在正下方, 期望值較高
 - 如果寶藏不在正下方, 期望值較低
4. 策略改善:
 - 如果寶藏在正下方, 策略不變
 - 如果寶藏不在正下方, 最好的策略是向右一步
5. 策略更新: 現在的策略是往下或往右走

這邊的policy的樣子就會類似一張告訴Agent「在什麼位置, 要做什麼動作」的地圖

Policy Iteration

Algorithm: Policy Iteration

Input: MDP $(\mathcal{S}, \mathcal{A}, P, R, \gamma, H \rightarrow \infty)$

Output: $\pi(s)$'s for all s 's

For each state s , initialize $\pi(s)$ randomly;

repeat

 For each state s , initialize $V_\pi(s) \leftarrow 0$;

 repeat

 foreach s do

$V_\pi(s) \leftarrow \sum_{s'} P(s'|s; \pi(s)) [R(s, \pi(s), s') + \gamma V_\pi(s')]$;

 end

 until $V_\pi(s)$'s converge;

 foreach s do

$\pi(s) \leftarrow \arg \max_a \sum_{s'} P(s'|s; a) [R(s, a, s') + \gamma V_\pi(s')]$;

 end

until $\pi(s)$'s converge;

Algorithm 5: Policy iteration.

Randomly
pick

Policy Iteration

Algorithm: Policy Iteration

Input: MDP $(\mathcal{S}, \mathcal{A}, P, R, \gamma, H \rightarrow \infty)$

Output: $\pi(s)$'s for all s 's

For each state s , initialize $\pi(s)$ randomly;

repeat

 For each state s , initialize $V_\pi(s) \leftarrow 0$;

repeat

foreach s **do**

$V_\pi(s) \leftarrow \sum_{s'} P(s'|s; \pi(s)) [R(s, \pi(s), s') + \gamma V_\pi(s')]$;

end

until $V_\pi(s)$'s converge;

foreach s **do**

$\pi(s) \leftarrow \arg \max_a \sum_{s'} P(s'|s; a) [R(s, a, s') + \gamma V_\pi(s')]$;

end

until $\pi(s)$'s converge;

Evaluate
 $V_\pi(s)$

Algorithm 5: Policy iteration.

Policy Iteration

Algorithm: Policy Iteration

Input: MDP $(\mathcal{S}, \mathcal{A}, P, R, \gamma, H \rightarrow \infty)$

Output: $\pi(s)$'s for all s 's

For each state s , initialize $\pi(s)$ randomly;

repeat

 For each state s , initialize $V_\pi(s) \leftarrow 0$;

repeat

foreach s **do**

$V_\pi(s) \leftarrow \sum_{s'} P(s'|s; \pi(s)) [R(s, \pi(s), s') + \gamma V_\pi(s')]$;

end

until $V_\pi(s)$'s converge;

foreach s **do**

$\pi(s) \leftarrow \arg \max_a \sum_{s'} P(s'|s; a) [R(s, a, s') + \gamma V_\pi(s')]$;

end

until $\pi(s)$'s converge;

Improve $\pi(s)$

Algorithm 5: Policy iteration.

Q-Learning

Q Function

$$Q^*(s, a) = \max_{\pi} E_{s^{(1)}, \dots} \left(R(s, a, s^{(1)}) + \sum_{t=1}^{\infty} \gamma^t R(s^{(t)}, \pi(s^{(t)}), s^{(t+1)}); s, a, \pi \right)$$

$$Q^*(s, a) = \sum_{s'} P(s' | s; a) [R(s, a, s') + \gamma \max_{a'} Q^*(s', a')], \forall s$$

Policy iteration:

會有一張告訴Agent「在什麼位置，要做什麼動作」的地圖

Q learning:

會有一個由state跟action組成的表格以儲存Q值

Q-table

Q-Table	a1	a2
s1	$q(s1,a1)$	$q(s1,a2)$
s2	$q(s2,a1)$	$q(s2,a2)$
s3	$q(s3,a1)$	$q(s3,a2)$

Q table

兩小時後要meeting，我們得到的初始Q-table可能是這個樣子

	a1: 看paper	a2: 玩楓之谷
s1	5	-3
s2	4	-2

1. 假設我們在s1，決定看paper (因為Q高)
2. 到達s2，不採取動作，而是看 $Q(s2,a1)$ 與 $Q(s2,a2)$ 哪個大
3. $Q(s2,a1)$ 比較大，將它乘 γ ，再加上s1到s2拿到的reward(5)
→得到現實中 $Q(s1,a1)$ 的值，並和Q-table中原本的值(5)比較
4. 更新表格

Temporal Difference

- We want to estimate:

$$Q^*(s, a) = \sum_{s'} P(s'|s; a) [R(s, a, s') + \gamma \max_{a'} Q^*(s', a')], \forall s$$

Temporal difference (TD) estimation of $Q^*(s, a)$ for exploitation policy π :

- ① $\hat{Q}^*(s, a) \leftarrow$ random value, $\forall s, a$
- ② Repeat until converge **for each action $a^{(t)}$** :

$$\hat{Q}^*(s^{(t)}, a^{(t)}) \leftarrow \hat{Q}^*(s^{(t)}, a^{(t)}) + \eta \left[(R(s^{(t)}, a^{(t)}, s^{(t+1)}) + \gamma \max_a \hat{Q}^*(s^{(t+1)}, a)) - \hat{Q}^*(s^{(t)}, a^{(t)}) \right]$$

現實的情況下可能拿到的reward
(注意: 這裡只是想像未來)

Q-table中存的reward值

References

- [清大資工吳尚鴻老師的講義](#)
- [清大資工吳尚鴻老師的影片](#)
- [台大電機李宏毅老師的影片](#)
- [Policy iteration介紹](#)
- [Q-learning介紹](#)

Homework

- Please see the attached files (`run.ipynb/RL_brain.py/maze_env.py`) carefully and write comments to illustrate how these code work in `run.ipynb/RL_brain.py/maze_env.py`(optional).
- (Optional) You are also encouraged to modify the code under different situation (different maze size, episodes, learning rate, gamma or epsilon) and illustrate your observations.
- Turn in your work with the package (*.zip) that contains the three files above.
- Folder name: hw8_Your Chinese Name