



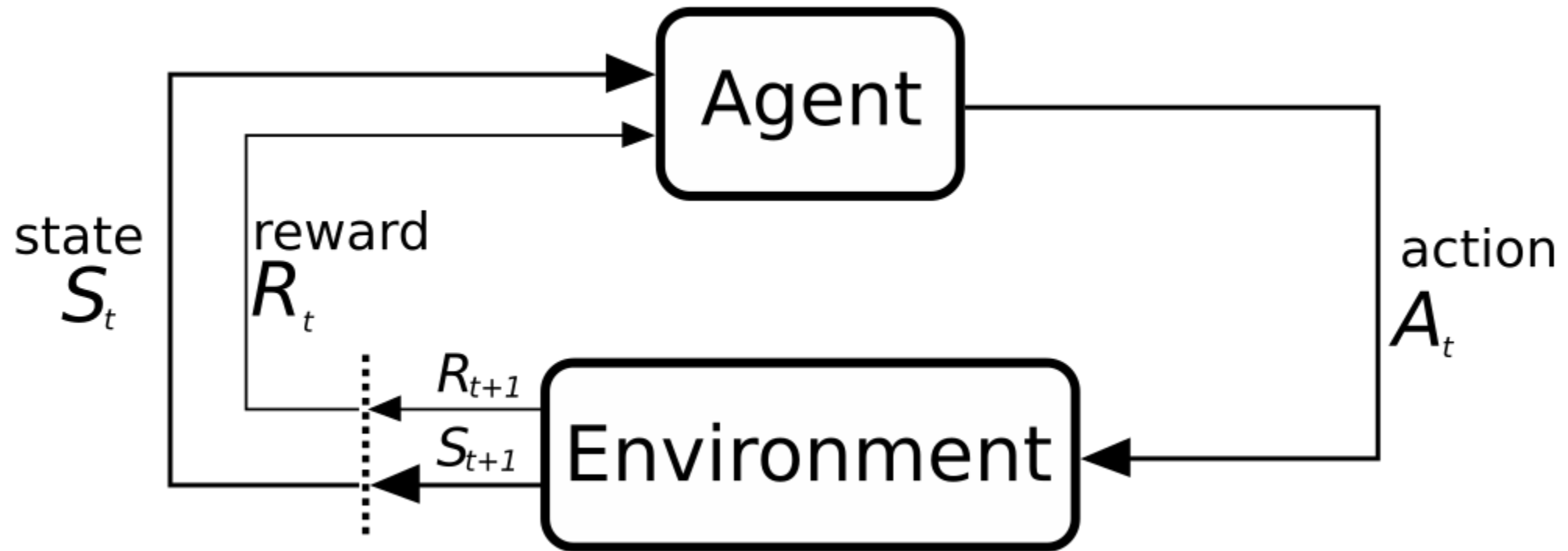
Exploration & Exploitation

Blink Sakulkueakulsuk

Recap

- A way to learn – a natural one
 - Learning through interactions
 - Active learning
 - Sequential interaction
 - Goal-directed
 - Learn even without the optimal sequence of actions

Recap



Rewards

It means, at time t , how good is the agent action

The ultimate goal of the agent is to maximize cumulative reward

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

This cumulative reward is called **return**

Values

We define the expected cumulative reward of a state s , the **value**.

$$\begin{aligned} v(s) &= \mathbb{E}[G_t \mid S_t = s] \\ &= \mathbb{E}[R_{t+1} + R_{t+2} + R_{t+3} + \dots \mid S_t = s] \end{aligned}$$

Action values

We can map value to both action and state

$$\begin{aligned} q(s, a) &= \mathbb{E}[G_t \mid S_t = s, A_t = a] \\ &= \mathbb{E}[R_{t+1} + R_{t+2} + R_{t+3} + \dots \mid S_t = s, A_t = a] \end{aligned}$$

Environment State

The internal state of the environment

Usually, it contains a lot of information. Those information may or may not be useful for the agent.

Sometimes, it is not even visible to the agent.

Agent State

A state of the agent can be reached by its **history**

$$H_t = O_0, A_0, R_1, O_1, \dots, O_{t-1}, A_{t-1}, R_t, O_t$$

* O = Observation

Once the agent reaches this state, the actions can be determined

Markov decision processes

$$p(r, s \mid S_t, A_t) = p(r, s \mid H_t, A_t)$$

The equation above essentially means “the future state is independent of the past state given the present state”

MDP provides baseline for reinforcement learning.

Agent State

A state of the agent in this case is **a function of history**

$$S_{t+1} = u(S_t, A_t, R_{t+1}, O_{t+1})$$

The agent state is typically much smaller than the environment state

Policy

Given a state, a policy $\pi(S)$ defines an agent's action

- Mapping from agent state to action

Deterministic policy: $\pi(S) = A$

Stochastic policy: $\pi(A|S) = p(A|S)$

Value function

state s อัตราริ่จจจจจจจจจจ π

$$\begin{aligned} v_{\pi}(s) &= \mathbb{E}[G_t \mid S_t = s, A_t \sim \pi(s)] \\ &= \mathbb{E}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) \mid S_t = s, A_t \sim \pi(s)] \end{aligned}$$

- $A_t \sim \pi(s)$ means that we choose action A from policy π
- $\gamma \in [0,1]$ is a discount factor

This is called a **Bellman equation**

If we want to find the optimal value, we can modify the equation to

$$v_*(s) = \max_a \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a]$$

Value function approximations

In many cases, we cannot find, or it is not feasible to find, the exact value from the value function.

- The states are infinitely many.

We need to approximate the function.

↳ bc state it's too big
use neural network

Model

A **model** predicts what will happen next in the environment

For example, a model P predicts the next state

$$P(s, a, s') \approx p(S_{t+1} = s' \mid S_t = s, A_t = a)$$

For example, a model R predicts the next reward

$$R(s, a) = \mathbb{E}[R_{t+1} \mid S_t = s, A_t = a]$$

Exploration & Exploitation

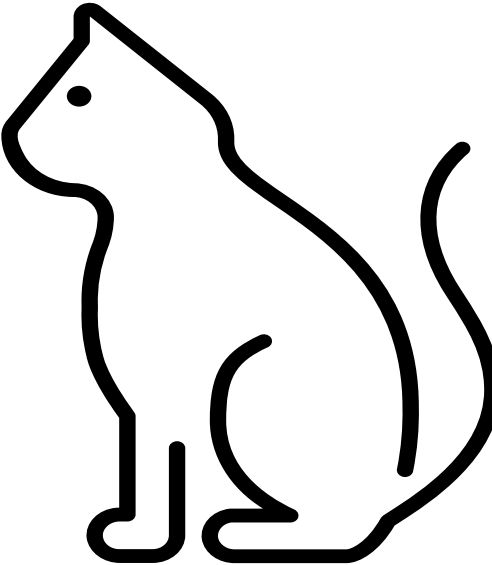




concept : ใ้ ใ้ ใ้ agent ใ้ ใ้ ใ้

Key Question

How can we **learn** without compromising the **reward**?

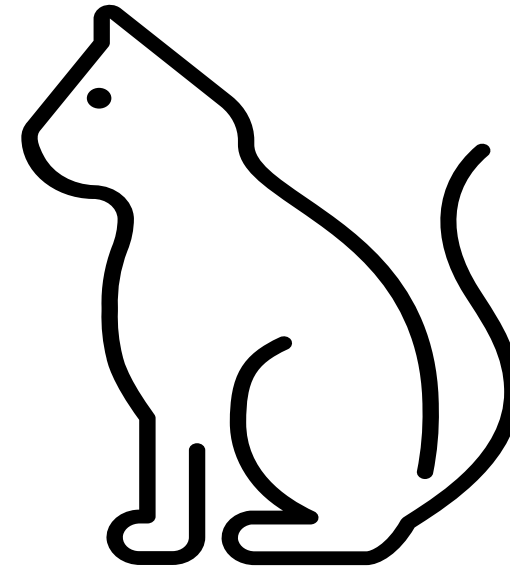
↳ reward ↓ take ใจทิ้ง: A รัง

Example

Day	Action	Reward	
1			
2			
3	?		

Example

Day	Action	Reward	
1			
2			
3			อย่าทำแบบนั้น
4	?	ทำตาม pattern	



Exploration vs Exploitation

Exploration: gain knowledge

exploitation

Exploitation: gain maximum reward

ឆ្លើយ រក្សាទុក រកចំណេញ reward

But do we have enough knowledge to know the action to gain maximum reward?

- It is a trade-off. We need to learn, but not too much.

Introducing Multi-Armed Bandit ^{เครื่อง}

⁶ เครื่องเล่นสุ่ม so popular

AKA Slot Machine

At each step t the agent performs an action $A_t \in A$ ^{เลือก 1 action} set of all action

- The agent receives a reward R_t ^{เลือกมา 1 set ที่ timestep t ใน action ทั้งหมด}
- The agent doesn't know the reward rate for each action.
(each machine has its own reward distribution $p(r|a)$ ^{ค. เหน้รางวัลที่ได้รับ reward ใน 1 action})
- The agent wants the maximum return over the lifetime.

* Agent ไม่รู้

แต่ละเครื่องมี reward distribution ต่างกัน



Action values → บาท 1000 คือ ได้ avg reward เท่าไหร่

Since we have only 1 state, the action value is the following equation.

$$q(a) = \mathbb{E}[R_t | A_t = a]$$

t = 1000	return
S ₁ = 600	1,200 บาท → $q(A_1) = \frac{1200}{600} = 2$
S ₂ = 400	1,200 บาท → $q(A_2) = \frac{1200}{400} = 3$

on average $S_1 = 2$
 $S_2 = 3$
 value มากกว่า

$$Q_t(a) = \frac{\text{sum of reward}}{\text{numbers of action}}$$

Expectation can be found as the average reward of the action.

$$Q_t(a) = \frac{\text{sum of reward} \quad \text{ว. action a}}{\sum_{n=1}^t R_n I(A_n = a)} \quad \text{identity}$$

if $A_n = a, I = 1$
 $A_n \neq a, I = 0$
 ได้ว่า a คือ action ที่สนใจ

Given $I(\text{True}) = 1$ and $I(\text{False}) = 0$

Action values

* เก็บประวัติได้ history 100; เก็บ memory 8; เก็บจุดคุ้ม

Or it can be updated iteratively:

error

error ที่เก็บค่าควร: Adjust มาจาก

or different $\times \alpha_t$

$$Q_t(A_t) = \begin{cases} Q_{t-1}(A_t) + \alpha_t (R_t - Q_{t-1}(A_t)) & \text{for } a = A_t \\ Q_{t-1}(a) & \text{for } a \neq A_t \end{cases}$$

$$\left. \begin{array}{l} \text{เก็บค่า} \\ \bar{X} - \text{state average} \\ N - \text{จำนวนครั้ง} \end{array} \right\} \frac{\bar{X}N + \text{new state}}{N+1}$$

Given

$$\alpha_t = \frac{1}{N_t(A_t)}, \quad N_t(A_t) = N_{t-1}(A_t) + 1, \quad N_t(a) = 0, \forall a$$

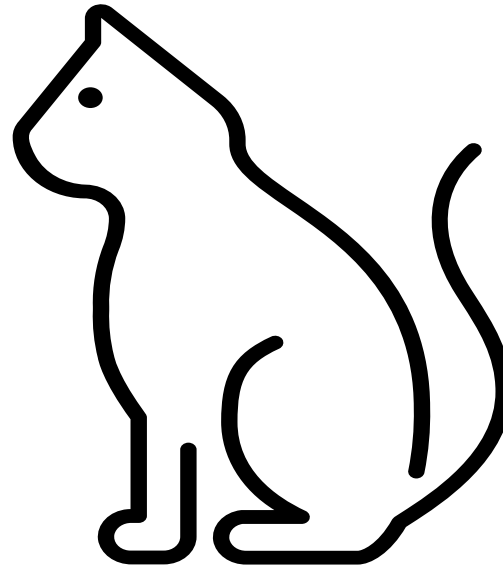
initial condition เริ่มต้นที่ 0

\downarrow
 α_t คือ จำนวนครั้งที่เลือกตัวที่ n ไว้ $\text{avg reward} = 2$ แต่เลือกได้ 3 , extra = 1 $\rightarrow \frac{1}{10}$

สมมติ $N_t(A_t)$ ค่าของ error จะทำให้อัตรา Swing ของ $N_t(A_t)$ ของค่าจะสูงขึ้น

Example

Day	Action	Reward
1		
2		
3		
4		?



reward
 Food: $R = 1$

No food: $R = -1$

ณ วันที่ 3

↓

$$Q_3(\text{white}) = 0 \rightarrow \text{Reward } 1 \text{ and } -1 ; \frac{1 + (-1)}{2}$$

$$Q_3(\text{black}) = -1 \rightarrow \text{Reward } -1 \text{ and } -1 ; \frac{-1 + (-1)}{2}$$

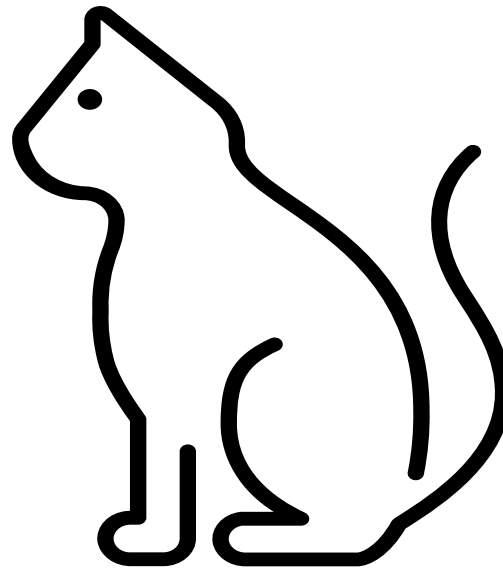
↓

$$\frac{\bar{x}_N + x}{N+1}$$

Example

Day Action Reward

1		
2		
3		
4		
5		
6		



exploration → ลองกด; ไปกดดำ

exploitation → ขาวเพราะเคยได้ ดูจากผลลัพธ์ที่เห็นไปแล้ว

Food: $R = 1$

No food: $R = -1$

$$Q_6(\text{white}) = -0.6 \rightarrow \frac{-0.6 + (-1)}{2} = -0.8$$

$$Q_6(\text{black}) = -1$$

What choice to choose next?

Introducing Regrets → for balancing exploration & exploitation

At the point, everyone knows that we need to balance between exploration and exploitation.

- How can we balance it?
- Can we set some assumptions and follow them?
- Can we even do it optimally?

Regrets

The optimal value

$$v_* = \max_{a \in A} q(a) = \max_a \mathbb{E}[R_t | A_t = a]$$

Regret is the loss occurred when taking sub-optimal action

$$v_* - q(A_t)$$

Noted that the agent doesn't know the regret, because it doesn't know the optimal value

Regrets

We can sum up all the regrets in the agent's lifetime

$$L_t = \sum_{i=1}^t (v_* - q(a_i))$$

By minimizing the cumulative regret, we maximize the cumulative reward

Regret Example

Day	Action	Reward	behavior : <u>ทำซ้ำเพราะกลัว</u> ↑ greedy
1			
2			
3			
4			
5			
6			

กลัวว่าได้ >>>> regret จะเพิ่มขึ้น
เพราะไม่เลือกจะเปลี่ยนใจ

Food: $R = 1$ / No food: $R = -1$

$$p(\text{food}|\text{white}) = 0.1$$

$$p(\text{food}|\text{black}) = 0.9$$

$$v_* = q(\text{black}) = \overset{\text{food}}{0.9} * 1 + \overset{\text{value ของการเลือกซ้ำๆ}}{0.1} * (-1) = 0.8$$

$$q(\text{white}) = \underset{\text{No food}}{0.9} * (-1) + \underset{\text{กลัวจึงเลือกซ้ำๆ}}{0.1} * 1 = -0.8$$

If the cat continues to choose white (because of the earlier episodes), it will get regrets of $1.6t$

Regrets

action a action b action c
 การกระทำ + การกระทำ + การกระทำ
 regret in this regret in this regret in this

$$L_t = \sum_{i=1}^t (v_* - q(a_i)) = \sum_{a \in A} N_t(a) (v_* - q(a)) = \sum_{a \in A} \underbrace{N_t(a)}_{\text{exploration}} \underbrace{\Delta_a}_{\text{exploitation}}$$

Now, minimizing the equation above shows the balance between exploration and exploitation.

If we want to take an action a many times, we need to reduce the action regret Δ_a .

Option 1: epsilon-greedy

exploration ผ่านค่า epsilon ไม่รู้ว่าตัวค่าไหนดี
แต่พอรู้ก็

A simple, yet powerful, algorithm

- $\epsilon \in [0,1]$ is a variable for exploration
- We select random action with probability ϵ
 - $a = \text{random}(A)$ random ผ่านค่า epsilon if $\epsilon > \text{threshold}$; random กับใน state action value
- We select the best action with probability $1 - \epsilon$
 - $a = \text{argmax}_{a \in A} Q_t(a)$ if $\epsilon < \text{threshold}$; เลือก state action value and select best action

เลือก epsilon อย่างไร ?

- ไม่มีค่าตายตัว
-

Option 2: Upper Confidence Bounds

idea มาจาก theory ใน มั่ว Probability

Add an upper confidence term $U_t(a)$ to the action value and adjust it through learning

Then we select an action that minimizing the adjusted action value

$$a_t = \underset{\text{max value ที่เกิดขึ้น}}{\underset{\text{action ในหน่วยนี้}}{\underset{\text{max}}{\underset{\text{แล้ว return ๑:1}}{\arg\max_{a \in A} Q_t(a) + U_t(a)}}}}$$

$a_t \leftarrow$ action ในหน่วยที่ ถ้า max

Option 2: Upper Confidence Bounds

$$a_t = \operatorname{argmax}_{a \in A} Q_t(a) + U_t(a)$$

bc เพื่อจะได้อ้าง argmax มา bc ไม่มั่นใจใน a

When we use action a a few times, $U_t(a)$ is large. (not sure if the value is correct) → จังจกนี้ Q_t จะ dominate U_t

When we use action a many times, $U_t(a)$ is small. (pretty sure that the value is correct)

* What is $U_t(a)$

Option 2: Upper Confidence Bounds

Assumption 1: we want to minimize this:

$$L_t = \sum_{i=1}^t (v_* - q(a_i)) = \sum_{a \in A} N_t(a) (v_* - q(a)) = \sum_{a \in A} N_t(a) \Delta_a$$

Assumption 2: we want to implement this:

When we use action a a few times, $U_t(a)$ is large. (not sure if the value is correct)

When we use action a many times, $U_t(a)$ is small. (pretty sure that the value is correct)

Option 2: Upper Confidence Bounds represent regret ไม่ค่อย no care

$$U_t(a) = \sqrt{\frac{\log t}{2N_t(a)}}$$

จน. ความน่าจะเป็น action a

Proved by using Hoeffding's Inequality

If we use the term above, we can implement an algorithm to both explore and exploit at the same time.

adjust แบบ dynamic to threshold like epsilon

$$a_t = \underset{\substack{\text{ถ้า state ใน} \\ \text{return max}}}{\operatorname{argmax}}_{a \in A} Q_t(a) + c \sqrt{\frac{\log t}{N_t(a)}}$$

เพื่อแบบ linear } ยิ่ง $N_t(a)$ มากเท่าไร ยิ่ง
พจน์นี้ จะน้อยลงเรื่อยๆ

↓
เป็นวิธี compensate a_t ว่าควร explore, exploit

Option 2: Upper Confidence Bounds

โง่กว่าตัว

```
n0 = 0
n1 = 0
sum = 0
sum_n0 = 0
sum_n1 = 0
```

```
for i in range (100000):
    x = np.random.uniform()
    if x < 0.9: epsilon
        sum_n0 += -1
        n0 += 1
    else:
        sum_n1 += 1
        n1 += 1
```

} state action value

```
print(sum_n0)
print(sum_n1)
print(n0)
print(n1)
```