

Exploration vs. Exploitation

Objectives

- ☐ Understand exploration and exploitation tradeoff
- ☐ Understand the method for balance the exploration and exploitation

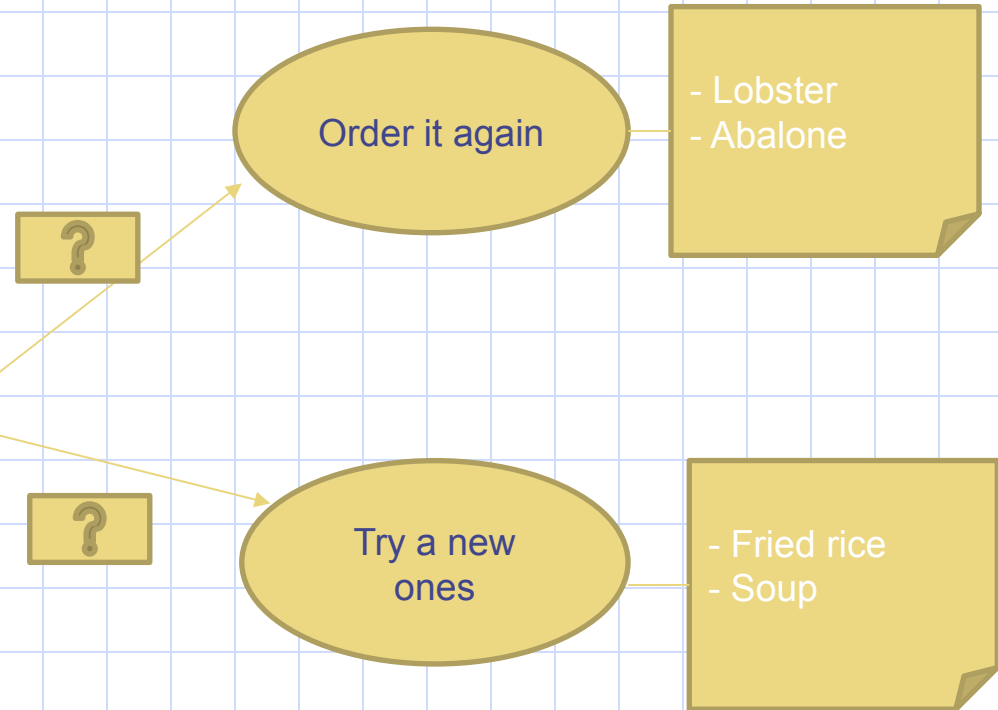
Exploration and Exploitation

- Discuss example:

- Mr. Manh is going out to eat with his lecturer Mr HoaDNT. When he get to the restaurant, he will have to decide what to order. He has been there a few times before and he always **order the same thing**. So he know that you will **be quite happy if he order it again**. Many of the **other items though look really tasty**.
- How does he decide when to order the same good meal again, or try something new?

Exploration and Exploitation

□ Discuss example:



□ What is Exploration and Exploitation? → discuss

Exploration and Exploitation

☐ Exploration:

- ☐ Exploration involves trying out new or less familiar options to gain more information about the environment or task.
- ☐ This is necessary to discover potentially better strategies or actions that may lead to higher rewards in the long run.
- ☐ Exploration is crucial for learning and discovering optimal solutions.
- ☐ Improve knowledge for long-term benefit

Exploration and Exploitation

- Example of exploratory behavior
 - Each of these plates represents a meal at your favorite restaurant, and you're trying to choose which meal to order.
 - Q of a is the estimated value for picking that meal.
 - N of a is the number of times you have picked that meal,

$q \leq$



$$\begin{aligned} q(a) &= 0 \\ N(a) &= 0 \\ q_*(a) &= 3 \end{aligned}$$



$$\begin{aligned} q(a) &= 0 \\ N(a) &= 0 \\ q_*(a) &= 4 \end{aligned}$$

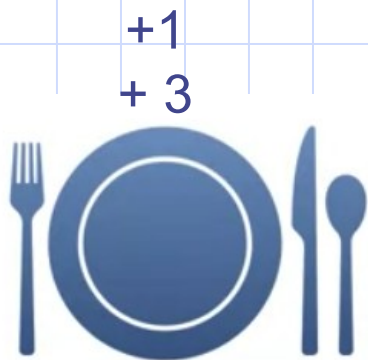


$$\begin{aligned} q(a) &= 0 \\ N(a) &= 0 \\ q_*(a) &= 2 \end{aligned}$$

exploration vs. Exploitation
Tradeoff

Exploration and Exploitation

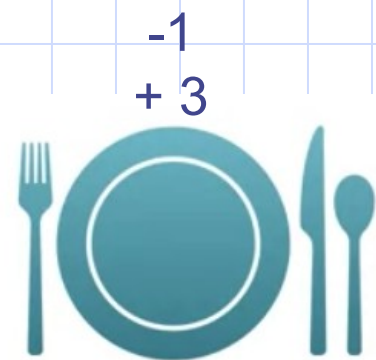
- Example of exploratory behavior
 - Each time you visit this restaurant, you follow a strict regimen and choose each meal in a Round Robin fashion.



$$\begin{aligned} q(a) &= 2 \\ N(a) &= 2 \\ q_*(a) &= 3 \end{aligned}$$



$$\begin{aligned} q(a) &= 4 \\ N(a) &= 2 \\ q_*(a) &= 4 \end{aligned}$$



$$\begin{aligned} q(a) &= 1 \\ N(a) &= 2 \\ q_*(a) &= 2 \end{aligned}$$

Exploration and Exploitation

- Example of exploratory behavior
 - After some time, you'll find the best meal to order



$$q(a) = 2$$

$$N(a) = 2$$

$$q_*(a) = 3$$



$$q(a) = 4$$

$$N(a) = 2$$

$$q_*(a) = 4$$



$$q(a) = 1$$

$$N(a) = 2$$

$$q_*(a) = 2$$

Exploration and Exploitation

☐ Exploitation:

- ☐ Exploitation involves maximizing immediate rewards by choosing actions that have been found to be effective based on past experiences or knowledge.
- ☐ Exploitation exploits the current knowledge to make decisions that are likely to yield the highest immediate rewards.
- ☐ Exploitation may lead to suboptimal decisions if there are better but undiscovered options.
- ☐ Exploit knowledge for short- term benefit

Exploration and Exploitation

- Example of exploitation behavior



$$\begin{aligned} q(a) &= 0 \\ N(a) &= 0 \\ q_*(a) &= 3 \end{aligned}$$



$$\begin{aligned} q(a) &= 0 \\ N(a) &= 0 \\ q_*(a) &= 4 \end{aligned}$$

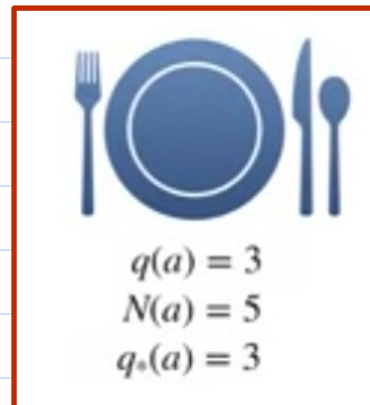


$$\begin{aligned} q(a) &= 0 \\ N(a) &= 0 \\ q_*(a) &= 2 \end{aligned}$$

Exploration and Exploitation

- Example of exploitation behavior
 - The agent received a positive reward making the value for that meal higher.
 - The estimated values for the other actions are zero.
 - The greedy action is always the same, to pick the first meal

+1
+3
+5
+1
+5



Exploration and Exploitation

- Example of exploitation behavior
 - The agent never saw any samples for the other meals.
 - The estimated values for the other two actions remain far from the true values, which means the agent never

+3
+5
+1
+5



$q(a) = 3$
 $N(a) = 5$
 $q_*(a) = 3$



$q(a) = 0$
 $N(a) = 0$
 $q_*(a) = 4$



$q(a) = 0$
 $N(a) = 0$
 $q_*(a) = 2$

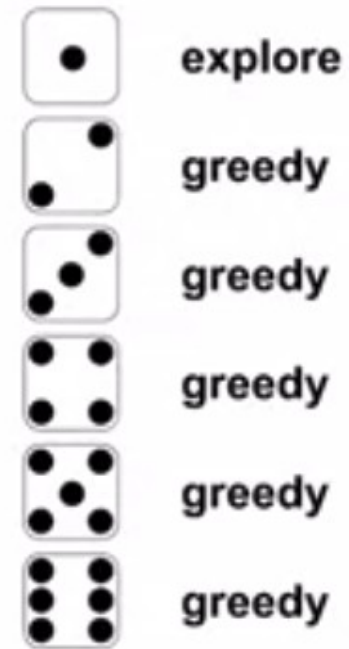
Exploration and Exploitation

- ☐ How do we choose explore or exploit?
 - ☐ When we explore:
 - ☐ Improve knowledge for long-term benefit
 - ☐ We get more accurate estimates of our values.
 - ☐ When we exploit,
 - ☐ Exploit knowledge for short-term benefit
 - ☐ We might get more reward.
 - ☐ We cannot however choose to do both simultaneously.

Exploration and Exploitation

□ Choose randomly- Epsilon Greedy method: We could choose to exploit most of the time with a small chance of exploring

- We could roll a dice. If it lands on one, then we'll explore.
- Otherwise, we'll choose the greedy action.
- Where epsilon refers to the probability of choosing to explore. In this case, epsilon will be equal to one over six.



Exploration and Exploitation

□ Epsilon Greedy method:

□ The action that we select on time-step t is:

- The greedy action with probability one minus epsilon
- or a random action with probability epsilon.

$$A_t \leftarrow \begin{cases} \operatorname{argmax}_a Q_t(a) & \text{with probability } 1 - \epsilon \\ a \sim \operatorname{Uniform}(\{a_1 \dots a_k\}) & \text{with probability } \epsilon \end{cases}$$

Exploration and Exploitation

- Example:

- Python code implementing epsilon-greedy exploration and exploitation. This code creates an agent that interacts with a multi-armed bandit environment and uses epsilon-greedy strategy to balance exploration and exploitation.

Exploration and Exploitation

□ Example

```

1 # 1.3 Example for Exploration and Exploitation- HoaDNT@fe.edu.vn
2 import numpy as np
3
4 class EpsilonGreedyAgent:
5     def __init__(self, num_actions, epsilon=0.1):
6         self.num_actions = num_actions
7         self.epsilon = epsilon
8         self.action_values = np.zeros(num_actions)
9         self.action_counts = np.zeros(num_actions)
10
11     def select_action(self):
12         if np.random.rand() < self.epsilon:
13             # Randomly choose an action for exploration
14             action = np.random.randint(self.num_actions)
15         else:
16             # Choose the greedy action for exploitation
17             action = np.argmax(self.action_values)
18         return action
19
20     def update_value(self, action, reward):
21         self.action_counts[action] += 1
22         # Update action-value estimate using incremental update rule
23         self.action_values[action] += (1 / self.action_counts[action]) * (reward - self.action_values[action])
24
25 # Create a simple multi-armed bandit environment
26 class MultiArmedBandit:
27     def __init__(self, num_arms):
28         self.num_arms = num_arms
29         self.true_action_values = np.random.normal(0, 1, num_arms)
30
31     def get_reward(self, action):
32         # Reward is sampled from a normal distribution with mean true action value and unit variance
33         return np.random.normal(self.true_action_values[action], 1)
34

```

Exploration vs. Exploitation
Tradeoff

Exploration and Exploitation

□ Example

```

34
35 # Initialize the environment and agent
36 num_arms = 10
37 num_steps = 1000
38 agent = EpsilonGreedyAgent(num_arms)
39
40 # Interaction loop
41 bandit = MultiArmedBandit(num_arms)
42 total_rewards = 0
43 for step in range(num_steps):
44     action = agent.select_action()
45     reward = bandit.get_reward(action)
46     agent.update_value(action, reward)
47     total_rewards += reward
48
49 print("Total rewards obtained:", total_rewards)
50 print("Estimated action values:", agent.action_values)

```

```

Total rewards obtained: 1133.394282414202
Estimated action values: [ 0.66969484  0.75865086  0.77213978  0.76244285 -1.14527773  0.97416503
-2.01990784  0.86064529 -0.50302905  1.25977875]

```

Summary

- ☐ Understand exploration and exploitation tradeoff
- ☐ Understand the method for balance the exploration and exploitation

Q & A

Exploration vs. Exploitation
Tradeoff