

① **k-armed bandit** : (making decisions under uncertainty) \Rightarrow you can't be sure what you would like, but make the best choice you can.

每次只能选 1 种治疗 - \Rightarrow treatment 有 3 种 (1, 2, 3) 但只能选 1 种.

- 但是能 轮流 选 treatment 3 种 (1, 2, 3). 有 3 种治疗 - \Rightarrow treatment 有 3 种 (1, 2, 3) 但是需要轮流使用.

(治疗 1, 2, 3 交替使用 3 种 treatment 的策略).

each treatment is an action.

- have an agent (an agent) to choose between k actions (the treatments 1, 2, 3) and receives a reward (rewards) from the treatment (1, 2, 3) based on its choose.

- agent decide which action is best. \Rightarrow find the value of taking each action. RT Action-Values. (Real Action-Values)

Defined the value of choosing an action as the expected reward: (goal: maximize the expected reward)

expected: $q_*(a) \stackrel{?}{=} E[R_t | A_t=a]$ $\forall a \in \{1, \dots, k\}$

referred to as selected action

for each possible action 1 to k.

If the agent (1, 2, 3) chose the action has the highest value, it achieves the goal. This procedure is called

Argmax $q_*(a)$.

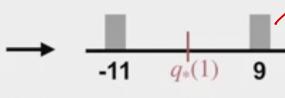
the mean of the distribution of each action.

$$= \sum_r p(r|a) r \Rightarrow \text{the sum of all possible rewards}$$

(multiplying possible reward with the probability of observing that reward).

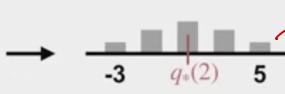
Calculating $q_*(a)$

What is an action to
in the distribution.



根据分布, 11.5 分。

$q_*(a) = .5 \times -11 + .5 \times 9 = -1$
the expected value of (分布) distribution



根据分布

$$q_*(a) = 1$$



$$q_*(a) = 3$$

② Sample-Average Method:

Sample-Average Method

$$Q_t(a) \stackrel{?}{=} \frac{\text{sum of rewards when } a \text{ taken prior to } t}{\text{number of times } a \text{ taken prior to } t}$$

$$= \frac{\sum_{i=1}^{t-1} R_i}{t-1} \Rightarrow \text{sum of rewards when taking action } a. \\ \Rightarrow \text{total # of times of action } a \text{ has been taken}$$

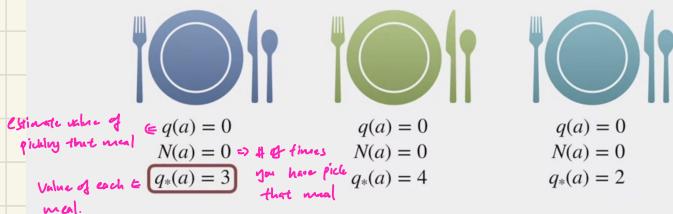
Action Selection



⑤

Exploration and Exploitation

- Exploration - improve knowledge for long-term benefit



grdually方法不是完全导致你的解!! 因为只有对 agent 会奖励 [奖励这个 meal 选饭] 才能知道这个 meal 没有发现其他 meal 的价值。

Exploration versus Exploitation

- .. 选择探索 (选择探索的策略). {
- Exploration - improve knowledge for long-term benefit \Rightarrow get more accurate estimate to our value
 - Exploitation - exploit knowledge for short-term benefit \Rightarrow might get more reward

是探索?

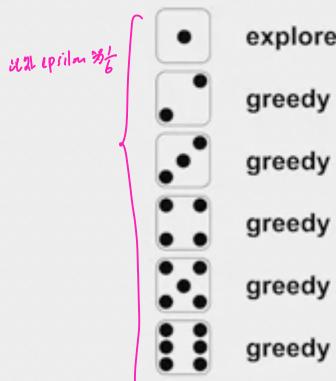
(选择一个策略)

- .. need to choose when to explore and when to exploit?

when to exploit?

Epislon: 选择探索的概率:

Epsilon-Greedy Action Selection



Epsilon-Greedy Action Selection

$A_t \leftarrow \{$

the action that we select at time step t

$\underbrace{\arg\max_a Q_t(a)}$ is the greedy action (maximize current value estimate). \hookrightarrow **explore** **with probability** $1 - \epsilon$

$\underbrace{a \sim \text{Uniform}(\{a_1 \dots a_k\})}$ is a random action (choose an action uniformly randomly). \hookrightarrow explore **with probability** ϵ

Optimistic Initial Values

if initial value
is larger than q_{π} , the agent will systematically exploit the actions
Optimistic 初期価値をもつ行動 (P1) の行動は optimistic initial value が
action agent が systematic exploitation (P2) する。agent が 2 つある理由。

$$Q_1(\text{P}) = 2.0$$

$$Q_1(\text{Y}) = 2.0$$

$$Q_1(\text{B}) = 2.0$$



u

Incremental update rule

$$\text{NewEstimate} \leftarrow \text{OldEstimate} + \text{StepSize} [\text{Target} - \text{OldEstimate}]$$

$$3.8 = 4 + \frac{1}{n} (2 - 4)$$

$$Q_{n+1} = Q_n + \frac{1}{n} \left(\underbrace{R_n - Q_n}_{\text{the new reward is our target}} \right)$$

$\frac{1}{n}$ (StepSize) produces a number between 0 to 1, (0 to \$1).

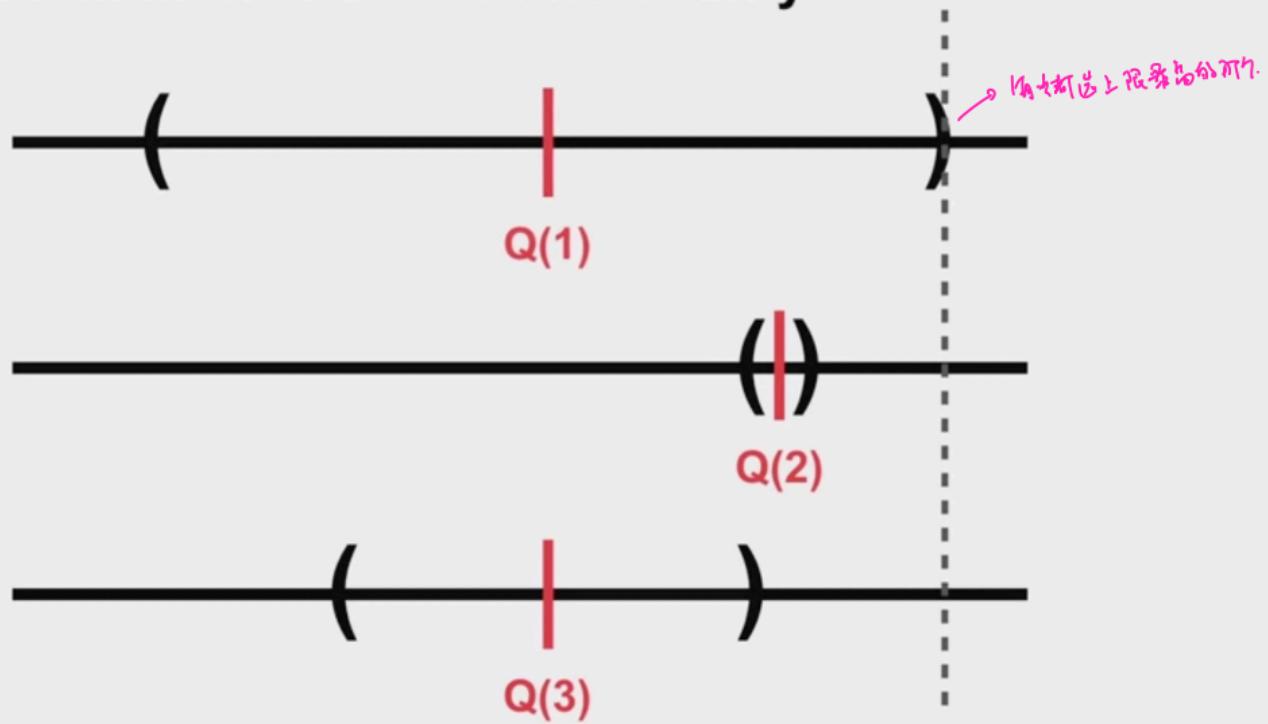
Upper-Confidence Bound (UCB) Action Selection

$$A_t \doteq \operatorname{argmax} \left[Q_t(a) + c \sqrt{\frac{\ln t}{N_t(a)}} \right]$$

Exploit Explore

use the strategy of:

Optimism in the Face of Uncertainty



Sequential Decision-Making

Latest Submission Grade 100%

1. What is the incremental rule (sample average) for action values?

1 / 1 point

$Q_{n+1} = Q_n + \frac{1}{n}[R_n - Q_n]$

$Q_{n+1} = Q_n - \frac{1}{n}[R_n - Q_n]$

$Q_{n+1} = Q_n + \frac{1}{n}[Q_n]$

$Q_{n+1} = Q_n + \frac{1}{n}[R_n + Q_n]$

✓ **Correct**

Correct! At each time step the agent moves its prediction in the direction of the error by the step size (here $1/n$).

 Back
Sequential Decision-Making

Graded Quiz • 45 min

Due Jan 11, 10:59 PM PST

2. Equation 2.5 (from the SB textbook, 2nd edition) is a key update rule we will use throughout the Specialization. We discussed this equation extensively in [video](#). This exercise will give you a better hands-on feel for how it works. The blue line is the target that we might estimate with equation 2.5. The red line is our estimate plotted over time.

1 / 1 point

$$q_{n+1} = q_n + \alpha_n [R_n - q_n]$$

Given the estimate update in red, what do you think was the value of the step size parameter we used to update the estimate on each time step?



$$\therefore Q_{n+1} = Q_n + \frac{1}{n} [R_n - Q_n]$$

比如假设 $Q_{n+1} = 3$, $Q_n = 2$, 故 R_n (now required) 为 estimate = 3 时的 target 值, 即 4.

- 1.0 $\therefore 3 = 2 + \frac{1}{n} [4 - 2]$
 $\therefore 1 = \frac{1}{n} [2]$
 $\therefore n = 0.5$
- 1/2 $\therefore \frac{1}{2} = 0.5$
 $\therefore StepSize = 0.5$
- 1/8 $\therefore 1/8 = 0.5$
- $1 / (t - 1)$

 Correct

Correct! We can see that the estimate is updated by about half of what the prediction error is.

 Back Sequential Decision-Making

Graded Quiz • 45 min

Due Jan 11, 10:59 PM PST

1 / 1 point

3. Equation 2.5 (from the SB textbook, 2nd edition) is a key update rule we will use throughout the Specialization. We discussed this equation extensively in [video](#). This exercise will give you a better hands-on feel for how it works. The blue line is the target that we might estimate with equation 2.5. The red line is our estimate plotted over time.

$$q_{n+1} = q_n + \alpha_n [R_n - q_n]$$

Given the estimate update in red, what do you think was the value of the step size parameter we used to update the estimate on each time step?



$$\therefore Q_{n+1} = Q_n + \alpha_n [R_n - Q_n]$$

∴ 从图设 $Q_{n+1} = 0.9, Q_n = 0.5$, 故 R_n (new reward) 的 estimate = 0.92726767 Target 值, 即 4.

若有 $0.9 = 0.5 + \frac{1}{n} [4 - 0.5]$

$$0.4 = \frac{1}{n} [3.5]$$

$$\frac{0.4}{3.5} = \frac{1}{n}$$

$$\therefore \frac{1}{n} = \frac{1}{8.75}$$

$$\therefore \alpha_n = \frac{1}{8.75}$$

$$\therefore \text{stepsize} = \frac{1}{8}$$

1 / (t - 1)

1.0

1/2

1/8

 **Correct**

Correct! We can see that the estimate is updated by $\frac{1}{8}$ of the prediction error at each time step.

← Back Sequential Decision-Making

Graded Quiz • 45 min

Due Jan 11, 10:59 PM PST

4. Equation 2.5 (from the SB textbook, 2nd edition) is a key update rule we will use throughout the Specialization. We discussed this equation extensively in [video](#). This exercise will give you a better hands-on feel for how it works. The blue line is the target that we might estimate with equation 2.5. The red line is our estimate plotted over time.

1 / 1 point

$$q_{n+1} = q_n + \alpha_n [R_n - q_n]$$

Given the estimate update in red, what do you think was the value of the step size parameter we used to update the estimate on each time step?



$$\therefore Q_{n+1} = Q_n + \alpha_n [R_n - Q_n]$$

∴ 既然 $Q_{n+1} = 4$, $Q_n = 4$, 故 R_n (now required) 为 estimate = 4 时的 target 值, 即 4.

1.0

$$\text{即 } 4 = 4 + \frac{1}{n} [4 - 4]$$

$$0 = 4 + \frac{1}{n} [0]$$

1/8

$$\therefore \frac{1}{n} = 1$$

1/2

$$\therefore \alpha_n = 1$$

1 / (t - 1)

 **Correct**

Correct! The estimate is updated to what the previous target was.

[← Back](#) Sequential Decision-Making

Graded Quiz • 45 min

Due Jan 11, 10:59 PM PST

5. Equation 2.5 (from the SB textbook, 2nd edition) is a key update rule we will use throughout the Specialization. We discussed this equation extensively in [video](#). This exercise will give you a better hands-on feel for how it works. The blue line is the target that we might estimate with equation 2.5. The red line is our estimate plotted over time.

1 / 1 point

$$q_{n+1} = q_n + \alpha_n [R_n - q_n]$$

Given the estimate update in red, what do you think was the value of the step size parameter we used to update the estimate on each time step?

 1.0 1/2 1/8 1 / (t - 1) **Correct**

Correct! We can see that the estimate is updated fully to the target initially, and then over time the amount that the estimate updates is reduced. This indicates that our step size is **reducing over time**.

6. What is the exploration/exploitation tradeoff?

1 / 1 point

- The agent wants to explore to get more accurate estimates of its values. The agent also wants to exploit to get more reward. The agent cannot, however, choose to do both simultaneously.
- The agent wants to explore the environment to learn as much about it as possible about the various actions. That way once it knows every arm's true value it can choose the best one for the rest of the time.
- The agent wants to maximize the amount of reward it receives over its lifetime. To do so it needs to avoid the action it believes is worst to exploit what it knows about the environment. However to discover which arm is truly worst it needs to explore different actions which potentially will lead it to take the worst action at times.

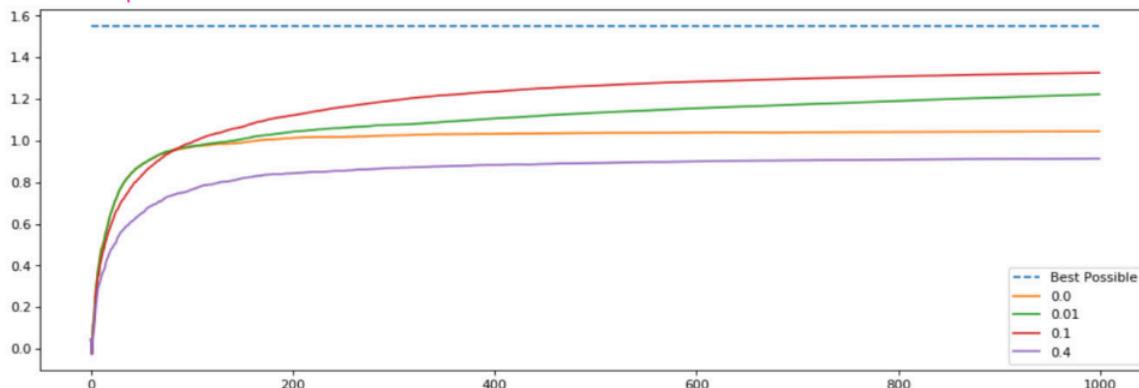
 **Correct**

Correct! The agent wants to maximize the amount of reward it receives over time, but needs to explore to find the right action.

7. Why did epsilon of 0.1 perform better over 1000 steps than epsilon of 0.01?

1 / 1 point

Explore by 0.225%



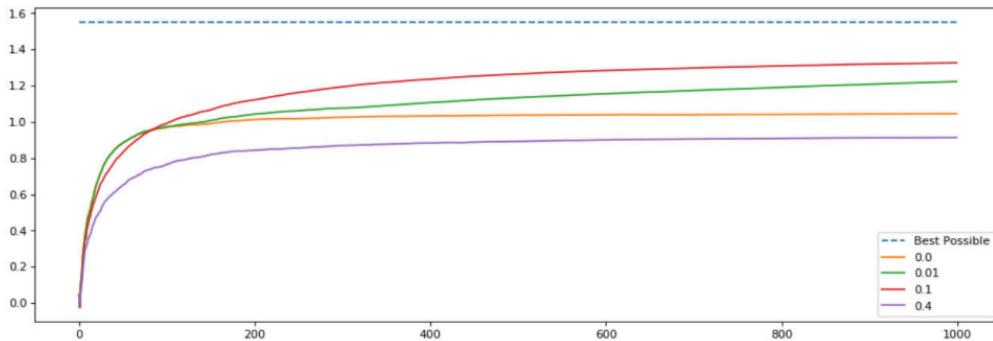
- The 0.01 agent did not explore enough. Thus it ended up selecting a suboptimal arm for longer.
- The 0.01 agent explored too much causing the arm to choose a bad action too often.
- Epsilon of 0.1 is the optimal value for epsilon in general.

Correct

Correct! The agent needs to be able to explore enough to be able to find the best arm to pull over time. Here epsilon of 0.01 does not allow for enough exploration in the time allotted.

8. If exploration is so great why did epsilon of 0.0 (a greedy agent) perform better than epsilon of 0.4?

1 / 1 point



- Epsilon of 0.4 doesn't explore often enough to find the optimal action.
- Epsilon of 0.0 is greedy, thus it will always choose the optimal arm.
- Epsilon of 0.4 explores too often that it takes many sub-optimal actions causing it to do worse over the long term.

Correct

Correct! While we want to explore to find the best arm, if we explore too much we can spend too much time choosing bad actions even when we know the correct one. In this case the action-value estimates are likely correct, however the policy does not always choose the action with the highest value.