

Master's degree in Control System Engineering






Reinforcement Learning LAB 1

k-armed Bandits



Github repo

 ShangtongZhang / reinforcement-learning-an-introduction Public

 **Code**  Issues **11**  Pull requests  Actions  Projects  Security  Insights









 master ▾

 **1** branch  **0** tags



ShangtongZhang Update README.md

c7cc538 on May 11  **314** commits

	chapter01	pythonic for chap01	3 years ago
	chapter02	Improved the display of the Greek alphabets.	2 years ago
	chapter03	add state labels	2 years ago
	chapter04	Fix a plotting issue, thanks @QuangTran4810	3 years ago
	chapter05	Update the axis limit	2 years ago
	chapter06	Update random_walk.py	16 months ago
	chapter07	Pythonic edits	3 years ago
	chapter08	Merge pull request #138 from vinnik-dmitry07/patch-1	2 years ago

<https://github.com/ShangtongZhang/reinforcement-learning-an-introduction>



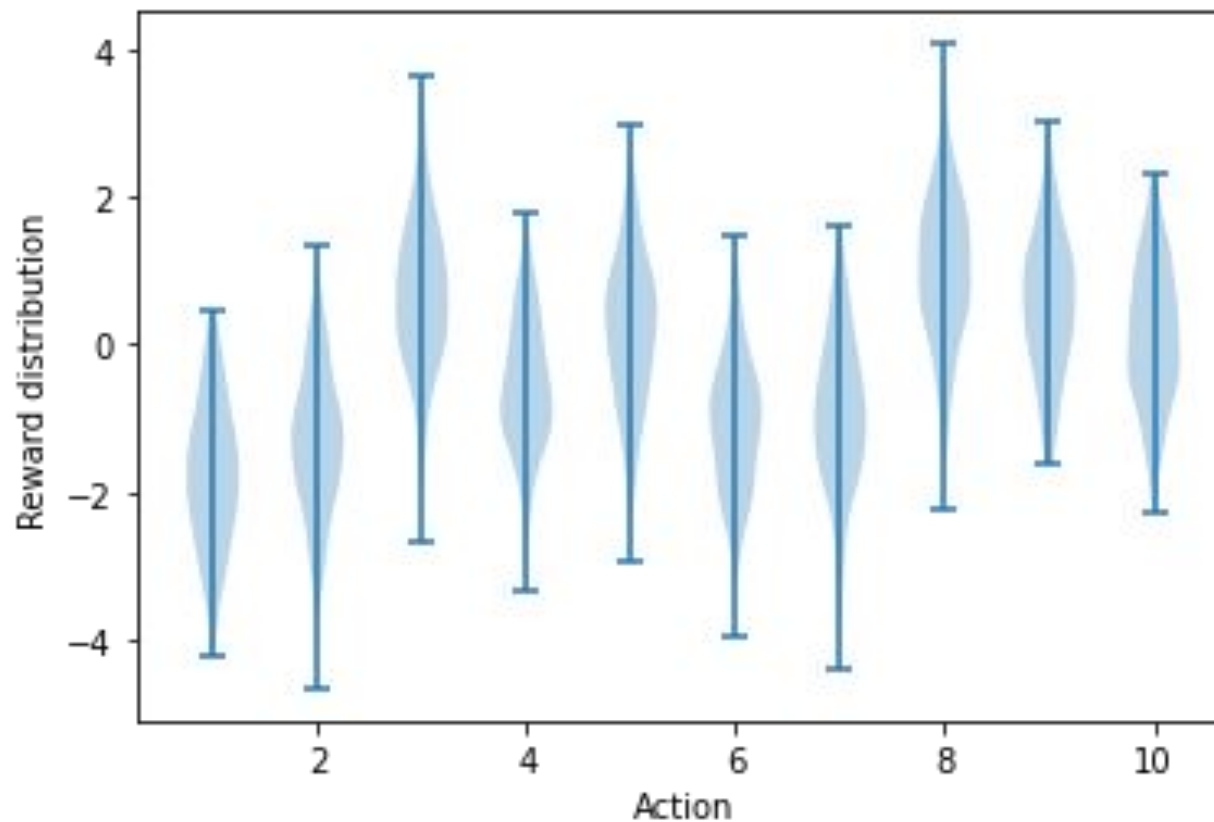
Bandit Class



```
class Bandit:
```

```
# @k_arm: # of arms
```

10 arms in our simulations!





Bandit Class



```
class Bandit:
```

```
# @k_arm: # of arms
```

- Allows to simulate interaction with the environment

- Implements:

- ϵ -greedy policy

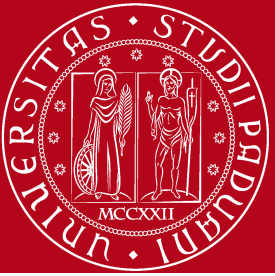
$$A_t = \begin{cases} \text{best action with } P = 1 - \epsilon \\ \text{uniform random } P = \epsilon \end{cases}$$

- UCB

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

- Gradient bandit

$$\Pr\{A_t = a\} = \frac{e^{H_t(a)}}{\sum_{b=1}^k e^{H_t(b)}}$$



Bandit Class



```
class Bandit:
```

```
# @k_arm: # of arms
```

```
# @epsilon: probability for exploration in epsilon-greedy algorithm
```

$$A_t = \begin{cases} \text{best action with } P = 1 - \epsilon \\ \text{uniform random } P = \epsilon \end{cases}$$



Bandit Class

```
class Bandit:  
  
    # @k_arm: # of arms  
  
    # @epsilon: probability for exploration in epsilon-greedy algorithm  
  
    # @initial: initial estimation for each action
```

$$Q_{t_0}(a) \quad \forall a \in A$$



Bandit Class



```
class Bandit:
```

```
# @k_arm: # of arms
```

```
# @epsilon: probability for exploration in epsilon-greedy algorithm
```

```
# @initial: initial estimation for each action
```

```
# @step_size: constant step size for updating estimations
```

$$\begin{cases} H_{t+1}(A_t) \leftarrow H_t(A_t) + \alpha (R_t - \bar{R}_t) (1 - \pi_t(A_t)) \\ H_{t+1}(a) \leftarrow H_t(a) - \alpha (R_t - \bar{R}_t) \pi_t(a) \quad \text{for all } a \neq A_t \end{cases}$$

$$Q_{n+1} \leftarrow Q_n + \alpha (R_n - Q_n)$$



Bandit Class



```
class Bandit:  
  
    # @k_arm: # of arms  
  
    # @epsilon: probability for exploration in epsilon-greedy algorithm  
  
    # @initial: initial estimation for each action  
  
    # @step_size: constant step size for updating estimations  
  
    # @sample_averages: if True, use sample averages to update estimations
```

$$Q_{n+1} \leftarrow Q_n + \alpha (R_n - Q_n)$$



Bandit Class



```
class Bandit:
```

```
# @k_arm: # of arms
```

```
# @epsilon: probability for exploration in epsilon-greedy algorithm
```

```
# @initial: initial estimation for each action
```

```
# @step_size: constant step size for updating estimations
```

```
# @sample_averages: if True, use sample averages to update estimations
```

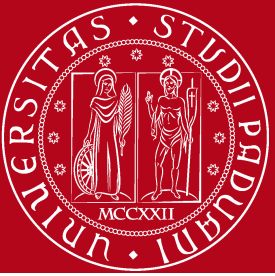
```
# @UCB_param: if not None, use UCB algorithm to select action
```

```
# @gradient: if True, use gradient based bandit algorithm
```

**Algorithm
Selection**

$$\operatorname{argmax}_a [Q_t(a) + c \sqrt{\frac{\ln t}{N_t(a)}}]$$

$$Pr\{A_t = a\} = \frac{e^{H_t(a)}}{\sum_{b=1}^k e^{H_t(b)}}$$



Bandit Class



class Bandit:

@k_arm: # of arms

@epsilon: probability for exploration in epsilon-greedy algorithm

@initial: initial estimation for each action

@step_size: constant step size for updating estimations

@sample_averages: if True, use sample averages to update estimations

@UCB_param: if not None, use UCB algorithm to select action

@gradient: if True, use gradient based bandit algorithm

@gradient_baseline: if True, use average reward baseline for gradient bandit

$$\begin{cases} H_{t+1}(A_t) \leftarrow H_t(A_t) + \alpha (R_t - \bar{R}_t) (1 - \pi_t(A_t)) \\ H_{t+1}(a) \leftarrow H_t(a) - \alpha (R_t - \bar{R}_t) \pi_t(a) \quad \text{for all } a \neq A_t \end{cases}$$



Simulation of Bandits



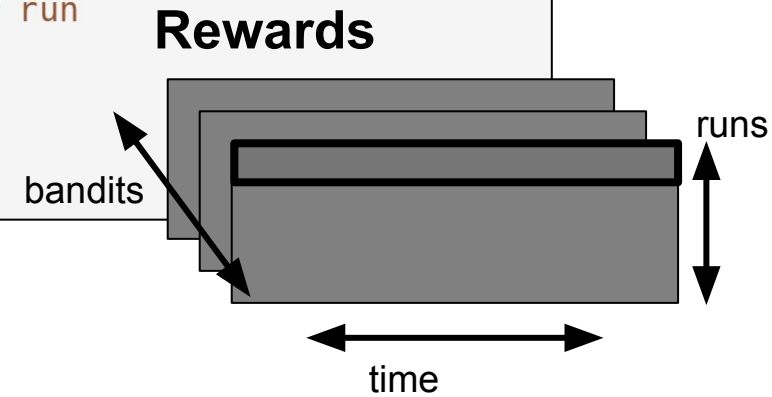
```
def simulate(runs, time, bandits):  
    # @runs: number of times each badit is simulated  
    # @time: interaction time (number of steps) for each run  
    # @bandits: list of Bandit objects
```



Simulation of Bandits

```
def simulate(runs, time, bandits):  
    # @runs: number of times each bandit is simulated  
    # @time: interaction time (number of steps) for each run  
    # @bandits: list of Bandit objects
```

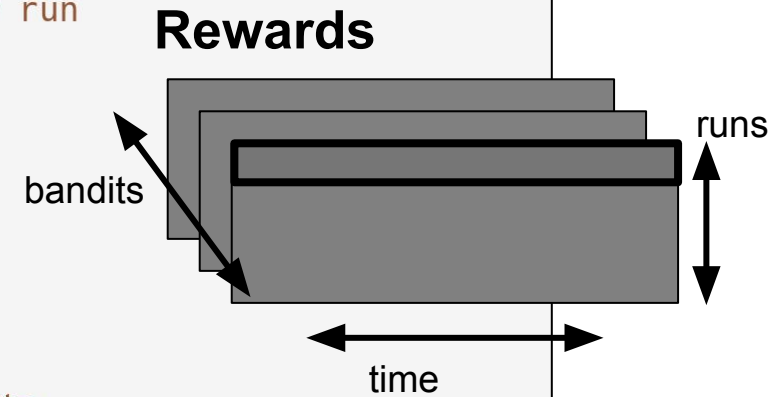
```
    rewards = np.zeros((len(bandits), runs, time))  
    best_action_counts = np.zeros(rewards.shape)
```

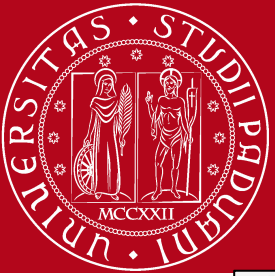




Simulation of Bandits

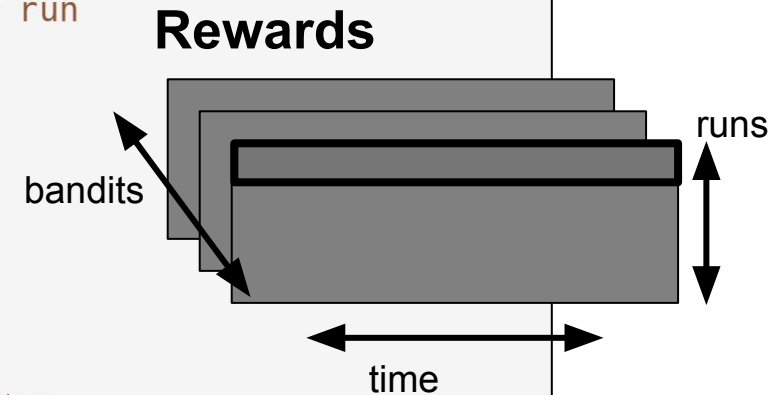
```
def simulate(runs, time, bandits):  
    # @runs: number of times each badit is simulated  
    # @time: interaction time (number of steps) for each run  
    # @bandits: list of Bandit objects  
  
    rewards = np.zeros((len(bandits), runs, time))  
    best_action_counts = np.zeros(rewards.shape)  
  
    for i, bandit in enumerate(bandits):  
        for r in range(runs):  
  
            bandit.reset()    # start from initial state
```





Simulation of Bandits

```
def simulate(runs, time, bandits):  
    # @runs: number of times each badit is simulated  
    # @time: interaction time (number of steps) for each run  
    # @bandits: list of Bandit objects  
  
    rewards = np.zeros((len(bandits), runs, time))  
    best_action_counts = np.zeros(rewards.shape)  
  
    for i, bandit in enumerate(bandits):  
        for r in range(runs):  
  
            bandit.reset()    # start from initial state  
  
            for t in range(time):  
  
                action = bandit.act()          # decide the action  
                reward = bandit.step(action)    # get reward  
                rewards[i, r, t] = reward       # save reward
```





Simulation of Bandits

```
def simulate(runs, time, bandits):
    # @runs: number of times each badit is simulated
    # @time: interaction time (number of steps) for each run
    # @bandits: list of Bandit objects

    rewards = np.zeros((len(bandits), runs, time))
    best_action_counts = np.zeros(rewards.shape)

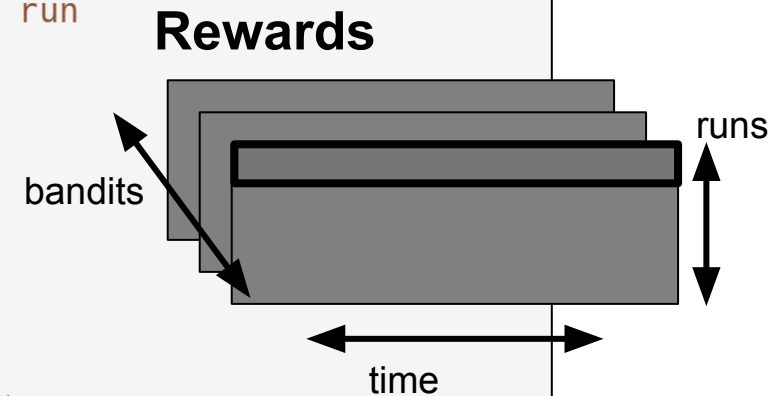
    for i, bandit in enumerate(bandits):
        for r in range(runs):

            bandit.reset()    # start from initial state

            for t in range(time):

                action = bandit.act()          # decide the action
                reward = bandit.step(action)    # get reward
                rewards[i, r, t] = reward       # save reward

                if action == bandit.best_action:  # Count best action
                    best_action_counts[i, r, t] = 1  # for stats
```





Simulation of Bandits

```
def simulate(runs, time, bandits):
    # @runs: number of times each badit is simulated
    # @time: interaction time (number of steps) for each run
    # @bandits: list of Bandit objects

    rewards = np.zeros((len(bandits), runs, time))
    best_action_counts = np.zeros(rewards.shape)

    for i, bandit in enumerate(bandits):
        for r in range(runs):

            bandit.reset()    # start from initial state

            for t in range(time):

                action = bandit.act()          # decide the action
                reward = bandit.step(action)    # get reward
                rewards[i, r, t] = reward       # save reward

                if action == bandit.best_action:  # Count best action
                    best_action_counts[i, r, t] = 1  # for stats

    # Averages over runs
    mean_best_action_counts = best_action_counts.mean(axis=1)
    mean_rewards = rewards.mean(axis=1)
    return mean_best_action_counts, mean_rewards
```

