

# COMP3667: Reinforcement learning practical 5

## Exploration vs exploitation with Monte Carlo methods

robert.lieck@durham.ac.uk

**This is the version *with* answers!**

[Notebook with code](#)

## 1 Overview

Welcome to the fifth reinforcement learning practical. In this practical, we will be focusing on the fundamental exploration vs exploitation problem with the aim to build intuition and develop good practices in conducting experiments. In particular, we will:

- Evaluate different  $\epsilon$ -greedy strategies for Monte Carlo reinforcement learning.
- Evaluate different GLIE strategies.
- Learn about the heuristic: “optimism in the face of uncertainty”.
- A mini competition—who can design the fastest converging and most sample efficient agent, just by tuning the  $\epsilon$ -schedule?

## Setup

For this practical, you will only need the following packages:

```
1 import numpy as np
2 import rldurham as rld
3
4 import matplotlib.pyplot as plt
5 from IPython import display as disp
```

As before, we will also use the frozen lake gym environment:

```
1 # import the frozen lake gym environment
2 env = rld.make('FrozenLake-v1', is_slippery=False)
3 rld.seed_everything(42, env)
```

## Monte Carlo methods with $\epsilon$ -greedy exploration

When tuning RL agents, it’s important to evaluate convergence across multiple experiments. You can use tools like ‘weights & biases’ for this. For now, we’ll simply keep track of our past experiments in a list and plot them.

```
1 # plotting
2 ep_reward = 0
3 reward_list = [[]]
4 plot_data = [[]]
5 plot_labels = []
6 auc = [0]
7 experiment_id = 0
```

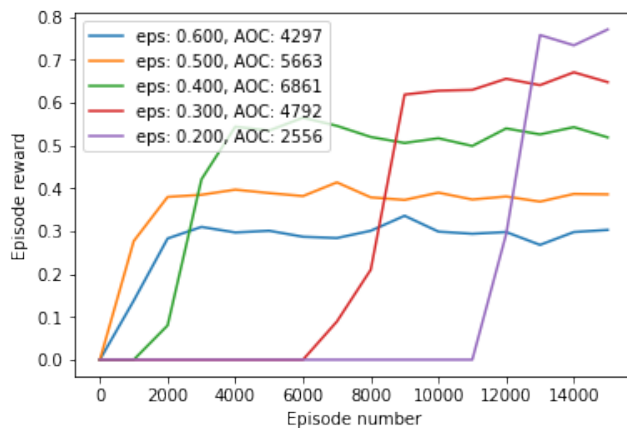
You can evaluate the following ( $\epsilon$ -greedy MC) multiple times to plot in different colours:

```
1 Q = np.zeros([env.observation_space.n, env.action_space.n])
2 n_s_a = np.zeros([env.observation_space.n, env.action_space.n])
3 num_episodes = 15000
4 epsilon = np.random.rand()*0.2+0.35 # change this each run
5 plot_labels.append("eps: {:.3f}".format(epsilon))
6
7 for episode in range(num_episodes):
8     state, info = env.reset()
9     done = False
10    results_list = []
11    result_sum = 0.0
12    while not done:
13        if np.random.rand() > epsilon:
14            action = np.argmax(Q[state, :])
15        else:
16            action = env.action_space.sample()
17        new_state, reward, term, trunc, info = env.step(action)
18        done = term or trunc
19        results_list.append((state, action))
20        result_sum += reward
21        state = new_state
22
23    for (state, action) in results_list:
24        n_s_a[state, action] += 1.0
25        alpha = 1.0 / n_s_a[state, action]
26        Q[state, action] += alpha * (result_sum - Q[state, action])
27
28    reward_list[experiment_id].append(result_sum)
29    auc[experiment_id] += int(result_sum)
30
31    if episode % 1000 == 0:
32        plot_data[experiment_id].append([episode, np.array(reward_list[experiment_id]).mean(),
33                                         np.array(reward_list[experiment_id]).std()])
34        reward_list[experiment_id] = []
35
36        for i in range(len(plot_data)):
37            plt.plot([x[0] for x in plot_data[i]], [x[1] for x in plot_data[i]],
38                    '-',
39                    # color=colors[i],
40                    label=plot_labels[i] + ', AUC: '+str(auc[i]))
41        plt.xlabel('Episode number')
42        plt.ylabel('Episode reward')
43        plt.legend()
44        disp.clear_output(wait=True)
45        plt.show()
46
47
48    env.close()
49
50    # advance to next experiment (so when we run the cell again we retain previous results)
51    experiment_id += 1
52    reward_list.append([])
53    plot_data.append([])
54    auc.append(0)
```

Take your time to study the code and relate it to what was covered in the lecture. If you run the cell multiple times, you will see it keeps track of the previous results and plots them in different colours.

## Questions

- Evaluate the above cell several times for different  $\epsilon$  values [0.6, 0.5, 0.4, 0.3, 0.2] (it should display previous experiment results on the same graph in different colours):



- What behaviour do you observe with small  $\epsilon$  vs with large  $\epsilon$  agents?

Large epsilon agents do a lot of initial exploration, so they quickly discover reward but are not good at exploiting it. In contrast, small epsilon agents exploit too early and thus get stuck in local environment minima—they take a lot longer to discover the reward. However, when they do discover the reward, they are good at exploiting it again.

- Does changing  $\epsilon$  impact the expected episode length (how long it takes to run)?

Yes in this case, when the agent learns to exploit the environment it acts greedily and reaches the terminal state (with the reward) more quickly. Low epsilon values—on average—take longer to discover the states with high reward, in this case resulting in longer episodes.

- What is the AUC (area under the curve) measuring here?

In this case, it's the total reward accumulated by the agent up to that point in training. Generally speaking, a fast-converging and sample efficient agent will have a higher AUC than others—although it depends on how many episodes you evaluate this over.

Now change the environment to the more difficult 8x8 frozen lake environment:

```
1 env = rld.make('FrozenLake8x8-v1', is_slippery=False)
2 rld.seed_everything(42, env)
```

## Question

- Experiment with  $\epsilon$  values in this new environment. What do you observe in contrast to the previous environment?

The rewards are more sparse in the new environment state space. Therefore you need to do much more exploration to discover them (typically needing epsilon around 0.9). Then, naturally, the agent is very poor at exploiting the new environment with this much exploration.

## 2 Evaluating GLIE strategies

A learning policy is considered GLIE (Greedy at the Limit with Infinite Exploration) if these two conditions are met:

- The agent continues to explore everything:

$$\lim_{k \rightarrow \infty} N_k(s, a) = \infty,$$

where  $N_k(s, a)$  is the number of times an action  $a$  was taken in state  $s$  at step  $k$ .

2. The policy converges on a greedy policy:

$$\lim_{k \rightarrow \infty} \pi_k(a|s) = \mathbf{1}(a = \arg \max_{a' \in \mathcal{A}} Q_k(s, a')).$$

In other words, after training for ever and ever, the final policy will be greedy with respect to  $Q$ .

## Exercise

- Return to the previous simple environment and design a GLIE policy that achieves an AUC > 12,000. It helps to design these in `desmos.com`—make sure to have the y-axis maximum at 1, and the x-axis maximum at 15,000 for this environment.

There are many solutions for this, one that works quite well is:

```
epsilon = min(1.0, 100.0/(episode+1))
```

## Competition

Now return to the more difficult 8x8 frozen lake environment and try to optimise the AUC.

- Compete with those around you (not allowed more than 15,000 episodes). When you are satisfied with your best schedule, evaluate it 5 times and write down the mean AUC across those 5 trials. Is it stable or does it only sometimes work?
- Now, lets try the ‘hardcore’ version of this environment. Enable `is_slippery=True` and try to design a schedule that performs well over 5 trials of this slippery 8x8 environment. Again, write down the mean AUC across those 5 trials. Is it stable or does it only sometimes work? Why?

The stochastic environment has higher variance and thus more evaluation is needed before optimising the policy.

- Who can design the best schedule for the deterministic and the slippery environments respectively? What are the challenges here?

## 3 Optimism in the face of uncertainty

A guiding principle to address the exploration-exploitation dilemma is “optimism in the face of uncertainty”. The basic idea is that if you are very certain, you can safely exploit, but if you are uncertain about the value of an action you should act “optimistically”, so basically “assume rewards everywhere” (unless you know better), which leads to exploration. This principle can be mathematically formulated in different ways and it is even possible to prove certain guarantees about convergence rates (all good policies of this kind are also GLIE, of course). One of the most simple and popular policies is called UCB1<sup>1</sup> (upper confidence bounds) and it defines an optimistic (upper confidence) bound on the estimated action value (average return) as

$$B_{(s,a)} = \underbrace{\widehat{R}_{(s,a)}}_{\text{unbiased value estimate}} + \underbrace{2 C_p \sqrt{\frac{2 \log n_s}{n_{(s,a)}}}}_{\text{bonus for little explored actions}}, \quad (1)$$

where  $\widehat{R}_{(s,a)}$  is the average return of action  $a$  in state  $s$ ;  $n_s$  is the number of visits to state  $s$ ;  $n_{(s,a)}$  is the number of times action  $a$  was taken in state  $s$ ; the returns are assumed to be in  $[0, 1]$ ; and the constant  $C_p > 0$  controls exploration. Actions are then chosen by maximising this optimistic upper bound instead of the unbiased value estimate

$$a^* = \arg \max_a B_{(s,a)}. \quad (2)$$

You can try out this policy by making some simple modifications to your code

<sup>1</sup>Auer P, Cesa-Bianchi N, Fischer P (2002) Finite-time Analysis of the Multiarmed Bandit Problem. Machine Learning 47:235–256. <https://doi.org/10.1023/A:1013689704352>; Kocsis L, Szepesvári C (2006) Bandit based monte-carlo planning. In: Machine Learning: ECML 2006. Springer, pp 282–293

```

1  ...
2  # define exploration constant and plot label
3  c = .1
4  plot_labels.append(f"UCB1: c={c}")
5  ...
6
7  for episode in range(num_episodes+1):
8      ...
9      while not done:
10         if np.any(n_s_a[state] == 0):
11             # first, try every action at least once
12             action = np.random.choice(np.argwhere(n_s_a[state] == 0).flatten())
13         else:
14             # then, use upper bounds to choose actions
15             R_s_a = Q[state, :]
16             n_s_a_ = n_s_a[state]
17             n_s = n_s_a_.sum()
18             B_s_a = R_s_a + 2 * c * np.sqrt(2 * np.log(n_s) / n_s_a_)
19             action = np.random.choice(np.argwhere(B_s_a==np.max(B_s_a)).flatten())
20         ...

```

## Exercise

Play around with the constant  $C_p$  and see how the UCB1 policy performs in comparison to the static  $\epsilon$ -greedy policies and your best scheduled  $\epsilon$ -greedy policy. What differences can you observe?

For the UCB1 policy, it may take longer before it discovers a useful policy (compared to a good scheduled  $\epsilon$ -greedy policy), but it then exploits very efficiently.