

# Reinforcement Learning Lab

## Lesson 3: Monte Carlo Reinforcement Learning Methods

Gabriele Roncolato and Alberto Castellini

University of Verona  
*email: gabriele.roncolato@univr.it*

Academic Year 2024-25



UNIVERSITÀ  
di **VERONA**

Dipartimento  
di **INFORMATICA**

# Environment Setup

The first step for the setup of the laboratory environment is to update the repository and load the **miniconda** environment.

## Safe Procedure

Always back up the previous lessons' solutions before executing the repository update.

- Update the repository of the lab:

```
cd RL-Lab  
git stash  
git pull  
git stash pop
```

- Activate the *miniconda* environment:

```
conda activate rl-lab
```

# Today Assignment

In today's lesson, we implement the **On Policy Monte Carlo Control** algorithm in Python. In particular, the file to complete is:

---

`RL-Lab/lessons/lesson_3_code.py`

---

Inside the file, two functions are partially implemented. The objective of this lesson is to complete them.

- **`def on_policy_mc_epsilon_soft()`**
- **`def on_policy_mc_exploring_starts()`**

Expected results can be found in:

---

`RL-Lab/results/lesson_3_results.txt`

---

# Pseudocode - On Policy Monte Carlo

**On-policy first-visit MC control (for  $\epsilon$ -soft policies), estimates  $\pi \approx \pi_*$**

Algorithm parameter: small  $\epsilon > 0$

Initialize:

$\pi \leftarrow$  an arbitrary  $\epsilon$ -soft policy

$Q(s, a) \in \mathbb{R}$  (arbitrarily), for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$

$Returns(s, a) \leftarrow$  empty list, for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$

Repeat forever (for each episode):

Generate an episode following  $\pi$ :  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode,  $t = T-1, T-2, \dots, 0$ :

$G \leftarrow \gamma G + R_{t+1}$

Unless the pair  $S_t, A_t$  appears in  $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$ :

Append  $G$  to  $Returns(S_t, A_t)$

$Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$

$A^* \leftarrow \arg\max_a Q(S_t, a)$  (with ties broken arbitrarily)

For all  $a \in \mathcal{A}(S_t)$ :

$$\pi(a|S_t) \leftarrow \begin{cases} 1 - \epsilon + \epsilon/|\mathcal{A}(S_t)| & \text{if } a = A^* \\ \epsilon/|\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$$

**Figure:** Pseudocode for the on-policy monte carlo control algorithm for  $\epsilon$ -soft policies, the implementation is from the Sutton and Barto book *Reinforcement Learning: An Introduction*

# Pseudocode - On Policy Monte Carlo

## Monte Carlo ES (Exploring Starts), for estimating $\pi \approx \pi_*$

Initialize:

$\pi(s) \in \mathcal{A}(s)$  (arbitrarily), for all  $s \in \mathcal{S}$

$Q(s, a) \in \mathbb{R}$  (arbitrarily), for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$

$Returns(s, a) \leftarrow$  empty list, for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$

Loop forever (for each episode):

Choose  $S_0 \in \mathcal{S}, A_0 \in \mathcal{A}(S_0)$  randomly such that all pairs have probability  $> 0$

Generate an episode from  $S_0, A_0$ , following  $\pi$ :  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode,  $t = T-1, T-2, \dots, 0$ :

$G \leftarrow \gamma G + R_{t+1}$

Unless the pair  $S_t, A_t$  appears in  $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$ :

Append  $G$  to  $Returns(S_t, A_t)$

$Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$

$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a)$

**Figure:** Pseudocode for the on-policy monte carlo control algorithm with exploring starts, the implementation is from the Sutton and Barto book *Reinforcement Learning: An Introduction*

# Assignment Notes

Today's assignment is based on the same environment as the first lesson (*DangerousGridWorld*). The suggested assignment's solution uses the `sample_episode()` function. Consult the first tutorial for more information.

## First Visit vs Every Visit

The given pseudocode is for the first visit version. However, the most straightforward every-visit approach works for the *DangerousGridWorld* environment. The suggestion is to use the every visit approach, which does not require the check: *unless the pair  $S_t, A_t$  appears in ...* (6<sup>th</sup> line of pseudocode).

## Results Disclaimer

Given the (high) stochasticity of the method, the obtained results may differ from those suggested. The crucial requirement is to obtain a policy that reaches the goal position.