# Data-driven and Learning-based Control

## 4th hands-on session

Erica Salvato

UNIVERSITÀ
DEGLI STUDI
DI TRIESTE

## Table of Contents

▶ A brief recap

▶ 4th hands-on session

We can classify Reinforcement Learning approaches in:

- **Value-function methods:**
  The policy is implicitly defined via $V(x^{(k)})$ or $Q(x^{(k)}, u^{(k)}) \rightarrow$ **Critic**

- **Policy optimization methods:**
  The policy is a parameterized function whose weights are learned in order to maximize the expected cumulative discounted reward$\rightarrow$ **Actor**

- **Actor-critic methods:**
  Merging the two ideas by guiding the actor's learning on the basis of the critic's estimated return

# Value-function methods taxonomy

1 A brief recap

We can classify Value-function methods:

- depending on the **state set representation**:

  — Tabular

  — Function approximation

- depending on **how the policy is used in the evaluation/improvement steps**:

  — On-policy

  — Off-policy

# Value-function methods taxonomy

We can classify Value-function methods:

- depending on the **state set representation**:

    — Tabular

    — Function approximation

- depending on **how the policy is used in the evaluation/improvement steps**:

    — On-policy

    — Off-policy

Value function methods assume to work with a discrete compact action set, namely

$$\mathcal{U} = \{u_1, u_2, \ldots, u_h\}$$

## On-policy vs. Off-policy

- On-policy means that the evaluation is based on the **current policy**:

$$Q\left(x^{(k)}, u^{(k)}\right) = Q\left(x^{(k)}, u^{(k)}\right) + \alpha \left[r^{(k+1)} + \gamma Q\left(x^{(k+1)}, u^{(k+1)}\right) - Q\left(x^{(k)}, u^{(k)}\right)\right]$$

- Off-policy means that the evaluation is based on the **greedy policy**:

$$Q\left(x^{(k)}, u^{(k)}\right) = Q\left(x^{(k)}, u^{(k)}\right) + \alpha \left[r^{(k+1)} + \gamma \max Q\left(x^{(k+1)}, u\right) - Q\left(x^{(k)}, u^{(k)}\right)\right]$$

**Initialization.** $Q \, \forall x \in \mathcal{X} \, u \in \mathcal{U}$

**Repeat** (for each episode)

— Set $x^{(0)}$
— Select an action $u^{(k)}$ with $\epsilon$-greedy policy
— Repeat for each step of the episode until the terminal condition is met
  ○ Perform the action $u^{(k)}$, observe $x^{(k+1)}$ and $r^{(k+1)}$
  ○ Select an action $u^{(k+1)}$ with $\epsilon$-greedy policy
  ○ $Q\left(x^{(k)}, u^{(k)}\right) = Q\left(x^{(k)}, u^{(k)}\right) + \alpha \left[r^{(k+1)} + \gamma Q\left(x^{(k+1)}, u^{(k+1)}\right) - Q\left(x^{(k)}, u^{(k)}\right)\right]$
  ○ Update $x^{(k)} = x^{(k+1)}$, $u^{(k)} = u^{(k+1)}$

## Q-Learning algorithm

1 A brief recap

**Initialization.** $Q \, \forall x \in \mathcal{X} \, u \in \mathcal{U}$

**Repeat** (for each episode)

— Set $x^{(0)}$

— Repeat for each step of the episode until the terminal condition is met

 ○ Select an action $u^{(k)}$ with $\epsilon$-greedy policy

 ○ Perform the action $u^{(k)}$, observe $x^{(k+1)}$ and $r^{(k+1)}$

 ○ $Q\left(x^{(k)}, u^{(k)}\right) = Q\left(x^{(k)}, u^{(k)}\right) + \alpha \left[r^{(k+1)} + \gamma \max Q\left(x^{(k+1)}, u\right) - Q\left(x^{(k)}, u^{(k)}\right)\right]$

 ○ Update $x^{(k)} = x^{(k+1)}$,

- Tabular means that the $Q\left(x^{(k)}, u^{(k)}\right)$ representation is a lookup table, i.e., one entry for every state/action pair.

  Consequently, it assumes to work with a discrete compact state set.

|       | $u_1$         | $u_2$         | $\ldots$ | $u_h$         |
|-------|---------------|---------------|----------|---------------|
| $x_1$ | $Q\left(x_1, u_1\right)$ | $Q\left(x_1, u_2\right)$ | $\ldots$ | $Q\left(x_1, u_h\right)$ |
| $x_2$ | $Q\left(x_2, u_1\right)$ | $Q\left(x_2, u_2\right)$ | $\ldots$ | $Q\left(x_2, u_h\right)$ |
| $\ldots$ | $\ldots$   | $\ldots$      | $\ldots$ | $\ldots$      |
| $x_l$ | $Q\left(x_l, u_1\right)$ | $Q\left(x_l, u_2\right)$ | $\ldots$ | $Q\left(x_l, u_h\right)$ |

- Function approximation means to define an approximate representation of the value function $\hat{V}(x)$ or of the action-value function $\hat{Q}(x, u)$.

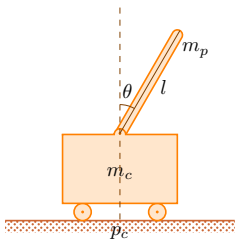$$Q(x, u) \approx \hat{Q}_\theta(x, u)$$

**Table of Contents**

► A brief recap

► 4th hands-on session

$$\ddot{\theta} = \frac{g\sin(\theta) \ + \ \cos(\theta)\left[\dfrac{-F - m_p\,l\,\dot{\theta}^2\sin(\theta)}{m_c + m_p}\right] - \dfrac{\mu_p\dot{\theta}}{m_p\,l}}{l\left[\dfrac{4}{3} - \dfrac{m_p\,cos^2(\theta)}{m_c + m_p}\right]}$$

$$\ddot{p}_c = \frac{F + m_p\,l\left[\dot{\theta}^2\,\sin(\theta) \ - \ \ddot{\theta}\,\cos(\theta)\right]}{m_c \ + \ m_p}$$

Given the cart and pole system previously defined, we would like to apply tabular Q-Learning with the goal of balancing the cart and pole.

1. Given the $\mathcal{X}, \mathcal{U}$ of the first hands-on define
   — the set of possible initial conditions of an episode
   — the terminal conditions of the episode

**Expectation:** Argue the answer

Given the cart and pole system previously defined, we would like to apply tabular Q-Learning with the goal of balancing the cart and pole.

1. Given the $\mathcal{X}, \mathcal{U}$ of the first hands-on define
   — the set of possible initial conditions of an episode
   — the terminal conditions of the episode

**Expectation:** Argue the answer

2. Provide a proper discretization of the state space and define the size of the $Q$

**Expectation:** Argue the answer

3. Create a code able to select the discretized version of the state given a measured state of the system

**Expectation:** Create a code able to convert measured states of the cart and pole system into the discretized version (Notice that we need it only to access the correct position in the table).

3. Create a code able to select the discretized version of the state given a measured state of the system

**Expectation:** Create a code able to convert measured states of the cart and pole system into the discretized version (Notice that we need it only to access the correct position in the table).

4. Given the terminal conditions of the episode create a code able to perform a fixed number of episodes of a fixed number of steps starting from random initial conditions and applying random inputs

**Expectation:** Create a code able to simulate the cart and pole for different episodes of a fixed number of steps, where each episode starts with a different initial condition randomly chosen. (Notice that terminal condition means that the episode terminates even if the number of time-steps have not been performed)

5. Create a function able to perform $\epsilon$-greedy policy for a chosen $\epsilon$

**Expectation:** Given the $\epsilon$ and the $Q$ the code should return the action to be performed following $\epsilon$-greedy policy

5. Create a function able to perform $\epsilon$-greedy policy for a chosen $\epsilon$

**Expectation:** Given the $\epsilon$ and the $Q$ the code should return the action to be performed following $\epsilon$-greedy policy

6. Create a code that applies Tabular Q-Learning algorithm

**Expectation:** Properly select the parameters $\alpha$, $\gamma$, $\epsilon$ and motivate the choice. The code should use all the results of the previous steps in order to construct the Tabular Q-Learning algorithm for balancing the cart and pole system. Plot the cumulative discounted reward

7. Copy and past your Q-learning algorithm and apply the changes needed to convert it into a SARSA algorithm

**Expectation:** The code should use all the results of the previous steps in order to construct the Tabular SARSA algorithm for balancing the cart and pole system. Plot the cumulative discounted reward

Questions' time!

UNIVERSITÀ
DEGLI STUDI
DI TRIESTE