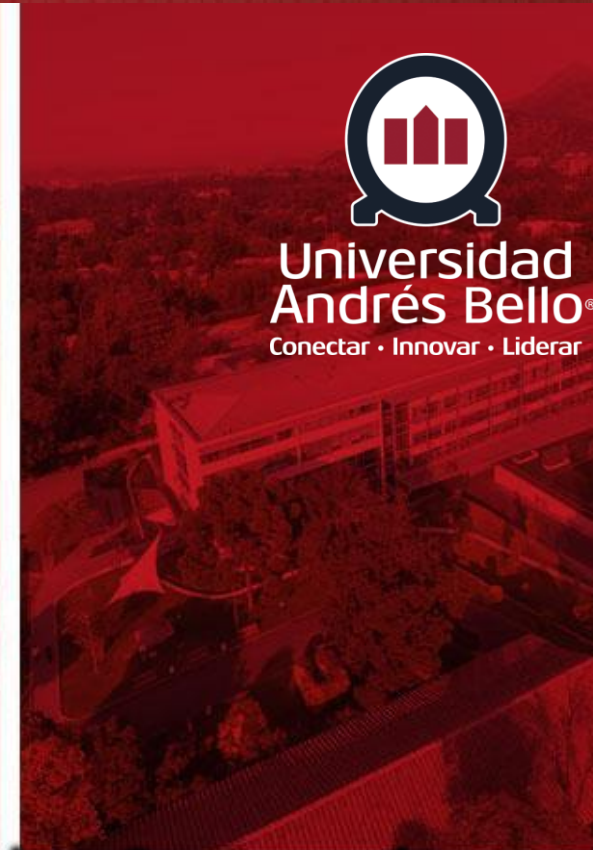
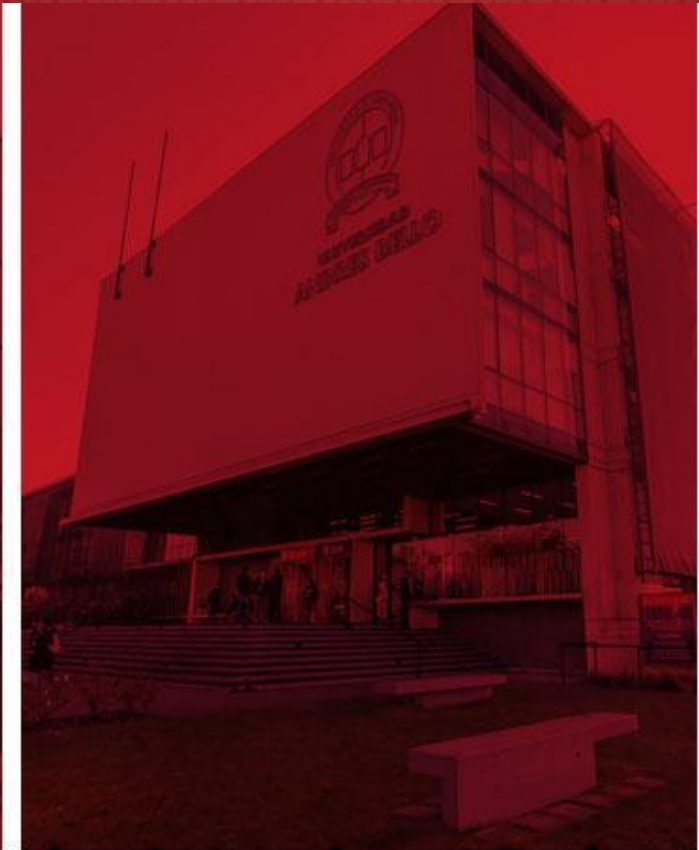


Practical Introduction to Reinforcement Learning with Gym in Python



Universidad
Andrés Bello®
Conectar · Innovar · Liderar

21st IEEE Latin American Robotics Symposium

Material available at
www.miguelsolis.info

Dr. Miguel A. Solís
Universidad Andrés Bello
November 12, 2024

LARS
2024



- Mobile Robotics
- Bio-inspired Robotics
- Intelligent Robotics
- Cognitive Robotics
- Evolutive Robotics
- Human-Robot Interaction
- Micro-Robotics
- Nano-Robotics
- Tele-Operated Robotics
- Swarm Robotics

Action and Perception

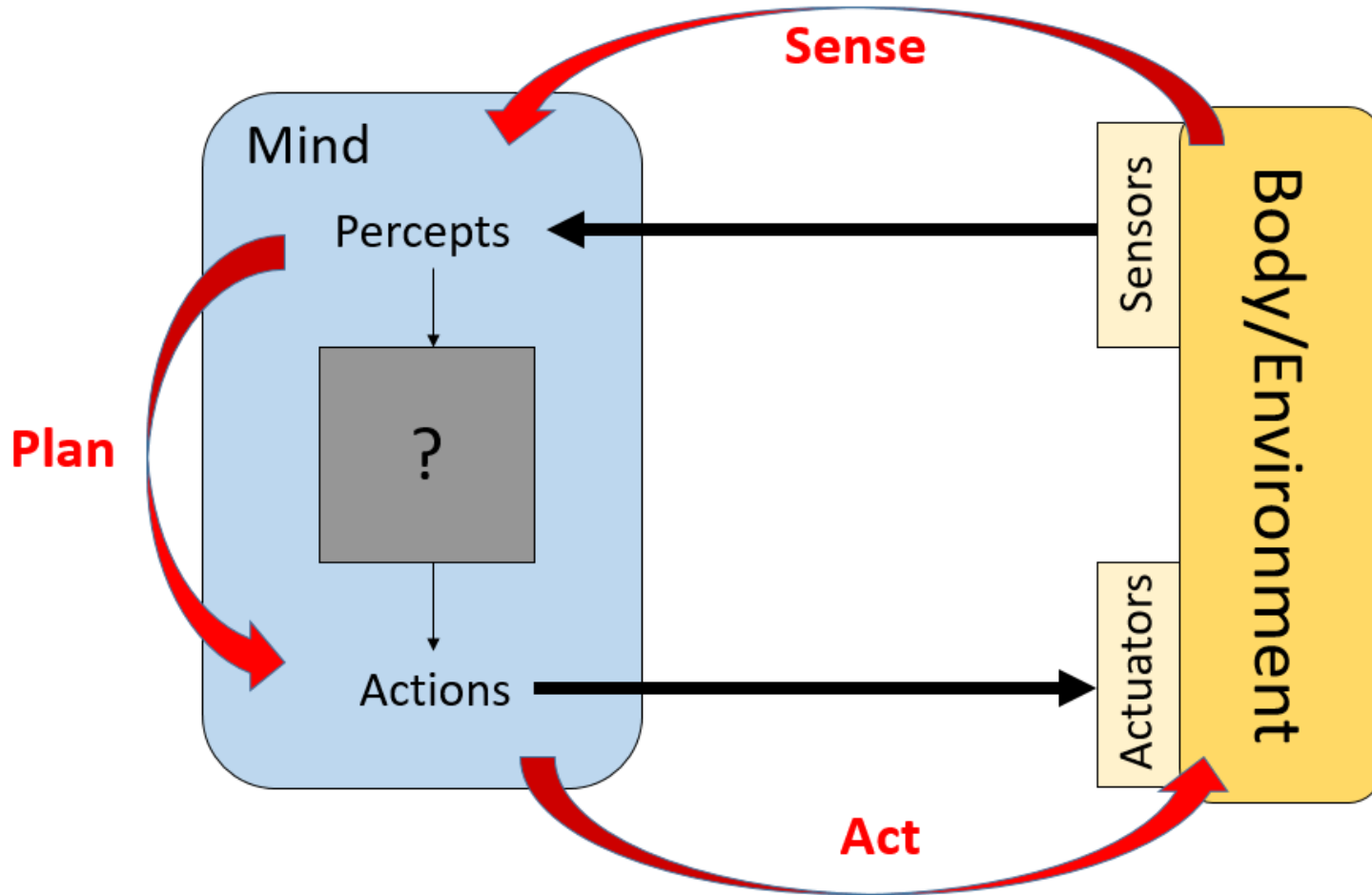


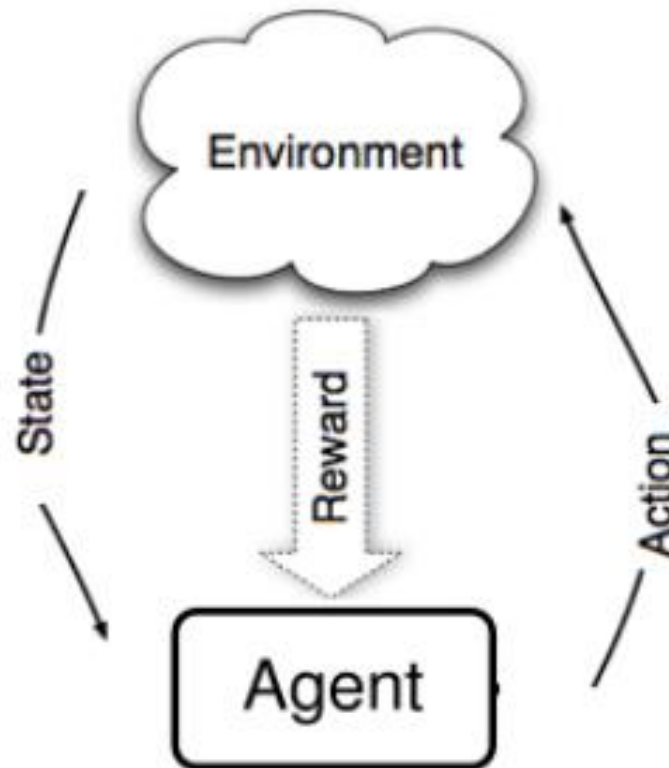
Image extracted from “Robotic Systems”, Kris Hauser, University of Illinois at Urbana-Champaign.

Ways of interacting between agent/environment:

- reactive: simplest way of decisión making through finite states machine (if A when X, then execute B leaving on Y)

Ways of interacting between agent/environment:

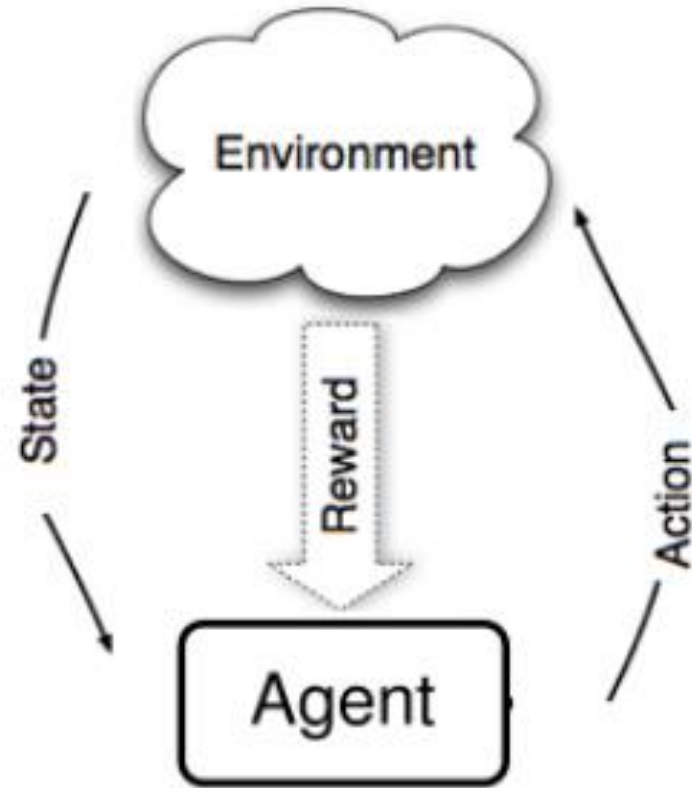
- intelligent: behavior shaping through heuristics
- incremental learning: behavior shaping through patterns recognition and prediction of outcomes from the environment



- Learning through interactions with environment.
- Applications
 - Sequential decision problems
 - Adaptive systems

Edward L. Thorndike – Animal Intelligence: An experimental study of the associate processes in animals (1898).

Animal modifies its behavior according to trial-and-error interactions with the environment.



- Agent perceives current state, s .
- Agent executes action, a , as output.
- Agent receives a reward, r , as reinforcement signal.

- A reinforcement learning problema is formally formulated through an MDP.
- MDP: Markov Decision Process.

A given state s_k comes from a (first-order) Markov Process if and only if:

$$Pr\{s_{k+1}|s_k\} = Pr\{s_{k+1}|s_1, \dots, s_k\}.$$

A reinforcement learning problema, formulated as an MDP is given by the tuple (S, A, T, R) where:

- S : set of states
- A : set of actions
- $T: S \times A \times S \rightarrow [0, 1]$ (transition function, unknown)
- $R: S \times A \times S \rightarrow \mathbb{R}$ (rewards function)
- $\pi: S \rightarrow A$ (policy)

State on time step k:

$$s_k \in \mathcal{S}$$

Action executed on time step k:

$$a_k \in \mathcal{A}$$

State-transition function:

$$s_{k+1} = T(s_k, a_k)$$

$$r_k = R(s_k, a_k, s_{k+1})$$

- $r_k > 0$
- $r_k = 0$
- $r_k < 0$

Policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$

A policy is said to be optimal, if it maximizes the long-term reward

Value function:

$$V^\pi(s_k) = r_k + \gamma r_{k+1} + \gamma^2 r_{k+2} + \dots$$
$$= \sum_{i=0}^{\infty} \gamma^i r_{k+i}.$$

Constraints:

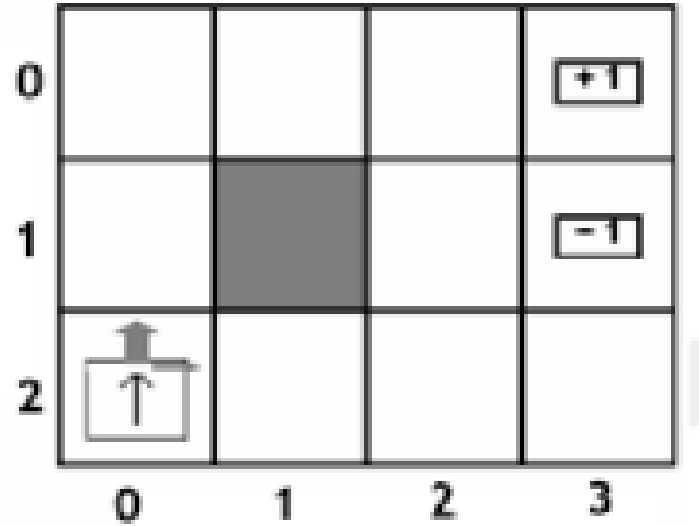
$$0 \leq \gamma < 1.$$

r_k bounded

A policy π^* is optimal if value function for that policy is optimal:

$$V^*(s) = V^{\pi^*}(s) \geq V^\pi(s) \quad \forall s, \pi$$

Toy example: Grid World



- States: location within the grid.
- Actions: up, left, right, down.
- Rewards: $+1$, -1 , -0.1

- Value Iteration
- Policy Iteration
- Temporal Difference
 - Q-Learning
 - SARSA

It consists on k iterations, until $V_k(s) - V_{k-1}(s)$ is small enough (according to a given tolerance), with updates given by:

$$V_k(s) = \max_a \sum_{s'} Pr\{s', s, a\} \cdot (R(s', s, a) + \gamma V_{k-1}(s'))$$

Toy example

0			$+1$	
1			-1	
2	\uparrow			
	0	1	2	3

- After 1 iteration:

0	-0.1	-0.1	1	$+1$
1	-0.19		-0.1	-1
2	-0.1	-0.1	-0.19	-0.1
	0	1	2	3

- After 2 iterations:

0	0.62	0.8	1	$+1$
1	0.46		0.8	-1
2	-0.27	-0.19	0.62	-0.27
	0	1	2	3

Toy example

- After 2 iterations:

0	0.62	0.8	1	<input type="checkbox"/>
1	0.46		0.8	<input type="checkbox"/>
2	0.31	0.46	0.62	0.46
	0	1	2	3

- On convergence:

0	0.62	0.8	1	<input type="checkbox"/>
1	0.46		0.8	<input type="checkbox"/>
2	0.31	0.46	0.62	0.46
	0	1	2	3

Just as $V(s)$ gives some value related to state, $Q(s,a)$ gives some value to taking a certain action on such state.

Optimal policy π^* satisfies:

$$Q^{\pi^*}(s, a) \geq Q^{\pi}(s, a) \quad \forall s, a, \pi$$

Consists on iterating over every (state,action) pair, for given hyperparameters (alpha and gamma).

Update expression:

$$Q(s_k, a_k) \leftarrow (1 - \alpha)Q(s_k, a_k) + \alpha \left(r_{k+1} + \gamma \max_a Q(s_{k+1}, a) \right)$$

Assuming $\hat{Q} = Q^*$, then optimal action for every state could be obtained by means of maximizing:

$$\pi^*(s_k) = \underset{a_k}{\operatorname{arg\,max}} Q^*(s_k, a_k)$$



Image extracted from "Introduction to Q-learning with OpenAI Gym", Gelana Tostaeva, Medium.

Toy example: Grid World

Algorithm does not iterate over the whole states space, just those visited states.

This (toy) example is episodic (when agent gets to final state, episode finish and a new episode starts with the agent starting again).

Must pay attention to local optimum.

0	der.	izq.	.	meta
1	arr.			prohib
2	arr.	izq.		
	0	1	2	3

Toy example: Grid World

Local optimum stuck:

0	der.	izq.		meta
1	arr.			prohib
2	arr.	izq.		
	0	1	2	3

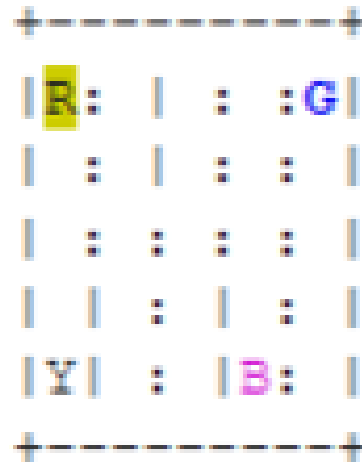
Including exploration

0	der.	der.	der.	meta
1	arr.			prohib
2	arr.	izq.		
	0	1	2	3

Taxi problem (Gym)



Actions



- 0: south
- 1: north
- 2: east (right)
- 3: west (left)
- 4: pick up passenger
- 5: leave passenger

Actions: 6

States: 500 (passenger_location, taxi_location, destination)

4 possible destinations

passenger_location: 4 possible locations for origin, or same location as the taxi

taxi_location: 25 possible locations according to the map

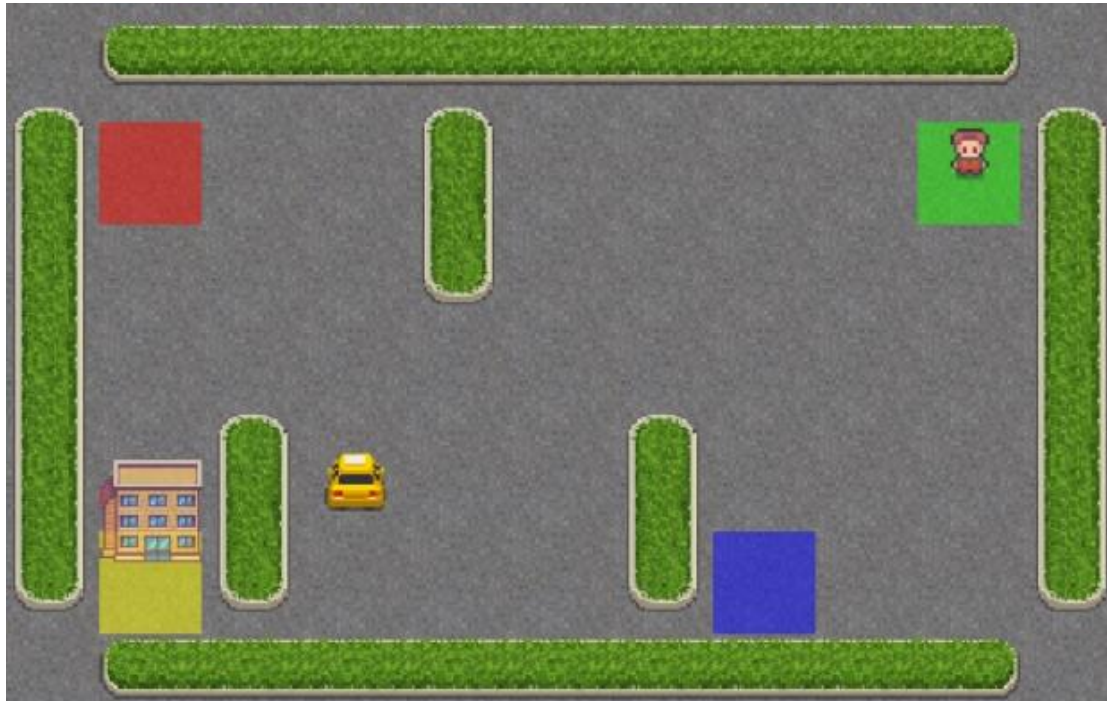
Gym verification:

env.action_space

env.observation_space

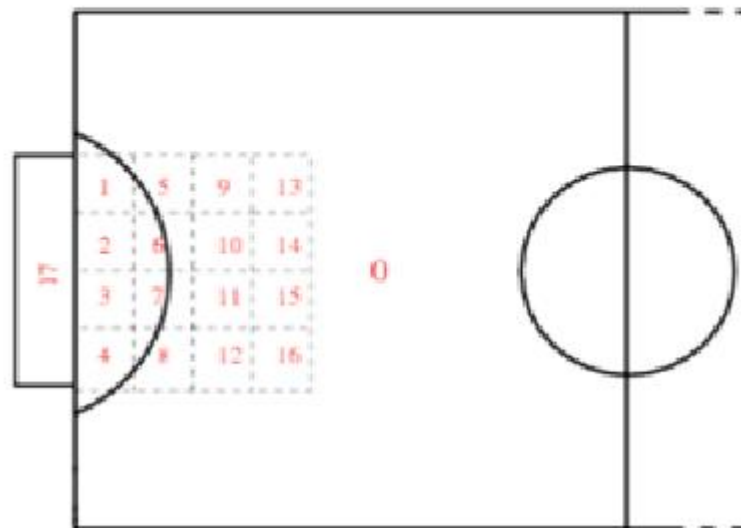
- Passenger left on correct location: 20 points
- Discount of 1 point anytime taxi moves with passenger without arriving to destination
- Discount of 10 points for leaving passenger on illegal location

12112024_Tutorial_MSolis.ipynb



Open with any Python Notebook Interpreter, or
colab.research.google.com

- ▶ G.A. Ahumada, C.J. Nettle and M.A. Solis, 'Accelerating Q-learning through Kalman Filter Estimations applied in a RoboCup SSL Simulation', **Proceedings** of the 10th IEEE Latin American Robotics Symposium, 2013.



Defensive strategy generation



State is composed by

- $\text{dist}(\text{Keeper}, \text{ball})$, $\text{dist}(\text{Taker}, \text{ball})$
- $\text{Dist}(\text{Keeper_Team_A}, \text{Keeper_Team_B})$
- $\text{Dist}(\text{Keeper_Team_A}, \text{Taker_Team_B})$
- $\text{Angle}(\text{Keeper_Team_A}, \text{Taker_Team_B})$

Ollino, F., Solis, M. A., & Allende, H. (2018). Batch reinforcement learning on a RoboCup Small Size League keepaway strategy learning problem. In 4th Congress on Robotics and Neuroscience, CRoNe 2018. CEUR-WS.

Rewards design

Delay on actions execution

Tabular representation

- affordances

(Object, Action, Effect)

Given an object and a certain action, what effect does it have?

Given an object and a desired effect, what is the required action?



- continual reinforcement learning (open-ended)

Final questions?

(material available at www.miguelsolis.info)