# AI Lab - Session 3

## Markov Decision Process

Riccardo Sartea

University of Verona
Department of Computer Science

May $24^{th}$ 2019

UNIVERSITÀ
di **VERONA**

Dipartimento
di **INFORMATICA**

# Start Your Working Environment

Start the previously installed (Session 1) conda environment *ai-lab*

Listing 1: Upgrade and spin up

```
cd ai-lab
git pull
conda remove --name myenv --all
conda env create -f ai-lab-environment.yml
conda activate ai-lab
jupyter notebook
```

# NumPy

## What is it

*NumPy is the fundamental package for scientific computing with Python. It contains among other things:*

- *a powerful N-dimensional array object*
- *sophisticated (broadcasting) functions*
- *tools for integrating C/C++ and Fortran code*
- *useful linear algebra, Fourier transform, and random number capabilities*

## What is it for

Fast array manipulation and mathematical operations. Think of it as a MATLAB like environment for Python: try to speed up the computations writing code in a vectorial fashion.

## Where to find it

http://www.numpy.org

## Tutorial

To open the tutorial navigate with your browser to:
*session3/session3_tutorial.ipynb*

## Assignments

- Your assignments for this session are at: *session3/session3_mdp.ipynb*. You will be required to implement value iteration and policy iteration algorithms
- In the following you can find pseudocodes for such algorithms

## Value Iteration

**Input:** $environment\ [T, R, A, S]$, $\gamma, \delta, maxiters,$
**Output:** $policy$ - state/action mapping

1: $V \leftarrow [0, ..., 0]$             ▷ Null vector of length $|S|$
2: $iter \leftarrow 0$
3: **repeat**             ▷ Compute Bellman Equation
4:      $V' \leftarrow V$
5:      $iter \leftarrow iter + 1$
6:      **for each** $s$ **in** $S$ **do**
7:          $V_s \leftarrow \max_{a \in A_s} \sum_{s' \in S} T(s, a, s')(R(s, a, s') + \gamma V_{s'})$

8: **until** $max(|V - V'|) < \delta$ **or** $iter = maxiters$
9: $\pi \leftarrow [0, ..., 0]$             ▷ Null vector of length $|S|$
10: **for each** $s$ **in** $S$ **do**             ▷ Extract policy
11:      $\pi_s \leftarrow \operatorname*{argmax}_{a \in A_s} \sum_{s' \in S} T(s, a, s')(R(s, a, s') + \gamma V_{s'})$

12: **return** $\pi$

## Policy Iteration

**Input:** $environment\ [T, R, A, S],\ \gamma, \delta, vmaxiters, pmaxiters,$
**Output:** $policy$ - state/action mapping

1: $V \leftarrow [0, ..., 0]$           ▷ Null vector of length $|S|$
2: $\pi \leftarrow [0, ..., 0]$           ▷ Initial policy length $|S|$
3: $piter \leftarrow 0$
4: **repeat**
5:     $\pi' \leftarrow \pi$
6:     $piter \leftarrow piter + 1$
7:     $viter \leftarrow 0$
8:     **repeat**           ▷ Evaluate Policy
9:        $V' \leftarrow V$
10:       $viter \leftarrow viter + 1$
11:       **for each** $s$ **in** $S$ **do**
12:          $V_s \leftarrow \sum_{s' \in S} T(s, \pi_s, s')(R(s, \pi_s, s') + \gamma V_{s'})$
13:     **until** $max(|V - V'|) < \delta$ **or** $viter = vmaxiters$
14:     **for each** $s$ **in** $S$ **do**           ▷ Improve policy
15:        $\pi_s \leftarrow \underset{a \in A_s}{\operatorname{argmax}} \sum_{s' \in S} T(s, a, s')(R(s, a, s') + \gamma V_{s'})$
16: **until** $\pi = \pi'$ **or** $piter = pmaxiters$
17: **return** $\pi$