# Reinforcement Learning

PhD course @ "Dottorato di ingegneria industriale e dell'informazione"
Trieste, 2024

Simone Silvetti
Research & Development

# Who am I?

**Simone Silvetti**  (silvetti@esteco.com)

➔   Studied mathematics in Rome

➔   Phd in Computer Science @ Udine

➔   Currently working in ESTECO

application of quantitative formal methods and machine learning techniques to Verification and Model-based Testing of Complex Systems

# Who am I?

**Simone Silvetti**  (silvetti@esteco.com)

➔ Studied mathematics in Rome

➔ Phd in Computer Science @ Udine

application of quantitative formal methods and machine learning techniques to Verification and Model-based Testing of Complex Systems

➔ Currently working in ESTECO

Numerical Methods Group

multi-objective optimization algorithms, machine learning, object-oriented programming

# Who am I?

**Simone Silvetti**  (silvetti@esteco.com)

➔   Studied mathematics in Rome

➔   Phd in Computer Science @ Udine

application of quantitative formal methods and machine learning techniques to Verification and Model-based Testing of Complex Systems

➔   Currently working in ESTECO

Numerical Methods Group

multi-objective optimization algorithms, machine learning, object-oriented programming

Research and Development

process mining, research projects related to technology and domains useful for ESTECO products

# Who am I?

**Simone Silvetti**  (silvetti@esteco.com)

➔    Studied mathematics in Rome

➔    Phd in Computer Science @ Udine

application of quantitative formal methods and machine learning techniques to Verification and Model-based Testing of Complex Systems

➔    Currently working in ESTECO

Numerical Methods Group

multi-objective optimization algorithms, machine learning, object-oriented programming

I worked on "Inverse Reinforcement Learning" applied to autonomous driving

Research and Development

process mining, research projects related to technology and domains useful for ESTECO products

# Who are you?



Dottorato in INGEGNERIA INDUSTRIALE E DELL'INFORMAZIONE

## Do you know Python? Numpy, Scipy?
13 responses



| | |
|---|---|
| I know Python and Numpy… | 1 (7.7%) |
| I know Python. I have a go… | 1 (7.7%) |
| I know only Python | 1 (7.7%) |
| No | 2 (15.4%) |
| Yes | 5 (38.5%) |
| Yes, I used some time ago… | 1 (7.7%) |
| Yes. (working/research ex… | 1 (7.7%) |
| pyhton and n… | 1 (7.7%) |

## During your studies have you participated in courses of Reinforcement Learning? If yes, which topics have you covered?
13 responses

No

Only partially

I did not partecipate to any course.

I have never participated in a course about Reinforcement Learning.

I have never participated at any course of bayesian optimization

no

Foundations

## Will you follow the "Learning-based Controllers and the Reality Gap" course?
5 responses



- ● Yes
- ● No

60%
40%

© ESTECO SpA

# Reference

A book from Sutton et al.



Reinforcement
Learning

An Introduction
second edition

Richard S. Sutton and Andrew G. Barto

Free available [here](http://incompleteideas.net/book/the-book-2nd.html)!

http://incompleteideas.net/book/the-book-2nd.html

# Reference

A book from Sutton et al.

Reinforcement
Learning

An Introduction
second edition

Richard S. Sutton and Andrew G. Barto

Free available [here](http://incompleteideas.net/book/the-book-2nd.html)!

http://incompleteideas.net/book/the-book-2nd.html



YouTube

**Google DeepMind**

@Google_DeepMind · 482K subscribers · 186 videos

Artificial intelligence could be one of humanity's most useful inventions. Google DeepMind …  ＞

🔔 Subscribed ⌄

Home   Videos   Shorts   Live   Podcasts   **Playlists**   Community   🔍

Created playlists                                                    ＝ Sort by

9 videos — Inside Google DeepMind — View full playlist
1 video — Visualising AI — View full playlist
6 videos — Scholarships | AI by you — View full playlist
4 videos — Unfolded: Meet the scientists using AlphaFold — View full playlist
5 videos — Life at DeepMind — View full playlist
8 videos — The story of AlphaFold — View full playlist

43 videos — Learning resources — View full playlist
8 videos — Talks | AI for science — View full playlist
10 episodes — DeepMind: The Podcast - Season 2 — View full podcast
9 episodes — DeepMind: The Podcast - Season 1 — View full podcast
13 videos — DeepMind x UCL | Deep Learning Lecture Series 2021 — View full playlist
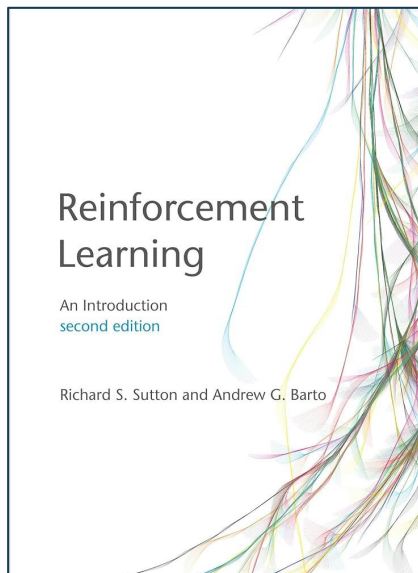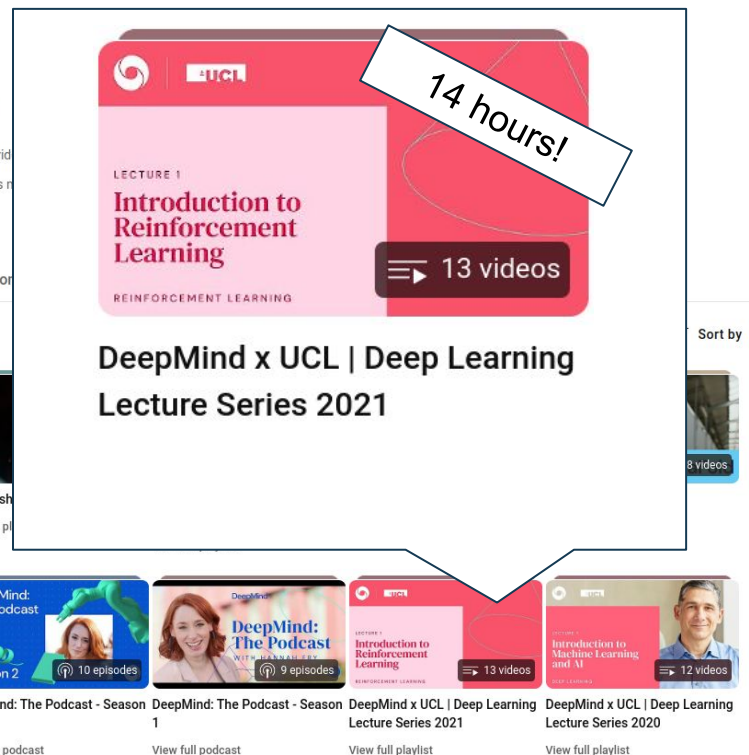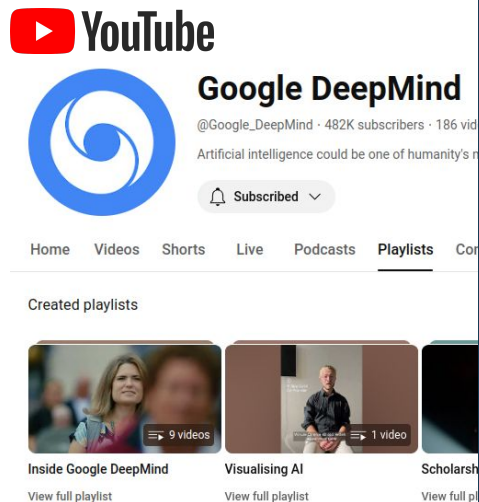12 videos — DeepMind x UCL | Deep Learning Lecture Series 2020 — View full playlist
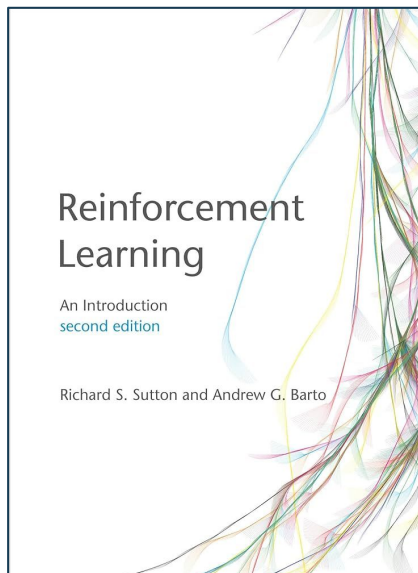
© ESTECO SpA

# Reference

A book from Sutton et al.

Reinforcement Learning

An Introduction
second edition

Richard S. Sutton and Andrew G. Barto

Free available [here](http://incompleteideas.net/book/the-book-2nd.html)!

http://incompleteideas.net/book/the-book-2nd.html

▶ YouTube

**Google DeepMind**

@Google_DeepMind · 482K subscribers · 186 vid...

Artificial intelligence could be one of humanity's ...

Subscribed ⌄

Home    Videos    Shorts    Live    Podcasts    Playlists    Co...

Created playlists                                                    Sort by

9 videos
Inside Google DeepMind
View full playlist

1 video
Visualising AI
View full playlist

Scholarsh...
View full...

14 hours!

LECTURE 1
Introduction to Reinforcement Learning
REINFORCEMENT LEARNING
13 videos

DeepMind x UCL | Deep Learning Lecture Series 2021

8 videos

43 videos
Learning resources
View full playlist

8 videos
Talks | AI for science
View full playlist

10 episodes
DeepMind: The Podcast - Season 2
View full podcast

9 episodes
DeepMind: The Podcast - Season 1
View full podcast

13 videos
DeepMind x UCL | Deep Learning Lecture Series 2021
View full playlist

12 videos
DeepMind x UCL | Deep Learning Lecture Series 2020
View full playlist

# Reference

A book from Sutton et al.

Reinforcement Learning

An Introduction
second edition

Richard S. Sutton and Andrew G. Barto

Free available [here](http://incompleteideas.net/book/the-book-2nd.html)!

▶ YouTube

ICTP
Quantitative Life Science

**ICTP Quantitative Life Sciences**

@ictpquantitativelifescienc9505 · 5.14K subscribers · 1K videos

More about this channel ›

Subscribe

Home    Videos    Live    Playlists    Community    🔍

⏸ 29 videos

2020-2021 Reinforcement Learning (QLS-RL)

40 hours!

Prof. Antonio Celani

# Introduction

What is Reinforcement Learning?

# A map

| | |
|---|---|
| **Science** | the systematic study of physical and natural world through observation, experimentation, and the testing of theories against the evidence obtained |
| **Formal Science** | uses formal systems to generate knowledge |
| **Computer Science** | is the study of computation, information and automation |
| *intelligence* · · · · · · · → | |
| **Artificial Intelligence** | enabling machines to perceive their environment and uses learning and intelligence to take actions that maximize their chances of achieving defined goals |
| *statistics & data* · · · · · · · → | |
| **Machine Learning** | development and study of **statistical algorithms** that can learn from data and **generalize** to unseen data, and thus perform tasks without explicit instructions. |
| *environment* · · · · · · · · · → | |
| **Supervised Learning**  **Unsupervised Learning**  **Reinforcement Learning** | technique that trains software to make decisions to achieve the most optimal results |

# A definition

| Reinforcement Learning | technique that trains software to make decisions to achieve the most optimal results |

# A definition

Reinforcement
Learning

technique that trains ~~software~~ to ~~make decisions~~ to ~~achieve the most optimal results~~

# A definition

Reinforcement Learning

technique that trains **agents** to ~~make decisions~~ to ~~achieve the most optimal results~~

# A definition

| Reinforcement Learning |
|---|

technique that trains **agents** to **map states into actions** to ~~achieve the most optimal results~~
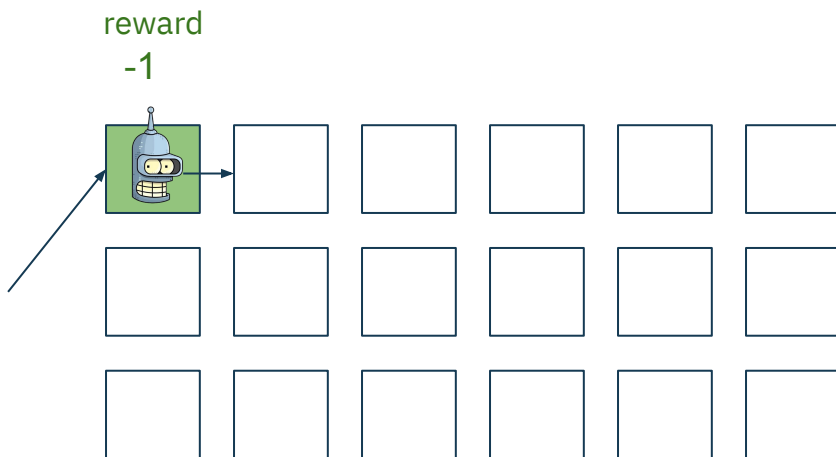
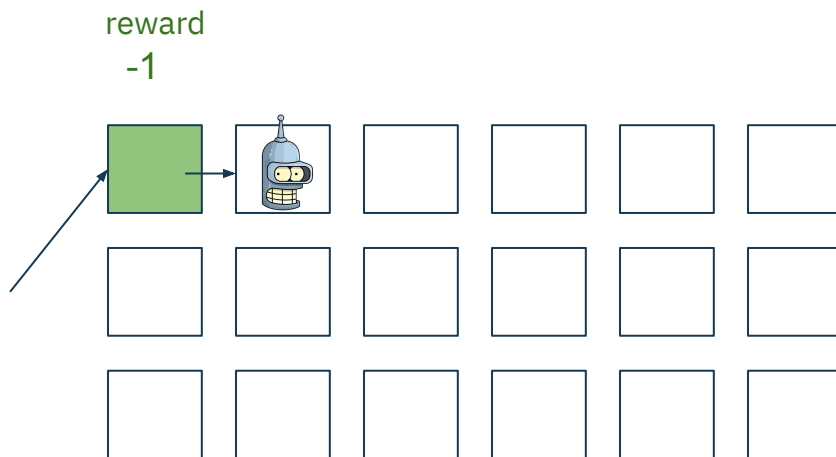# A definition

| Reinforcement Learning |
|---|

technique that trains **agents** to **map states into actions** to **maximize a cumulative reward**

# A definition

Reinforcement Learning

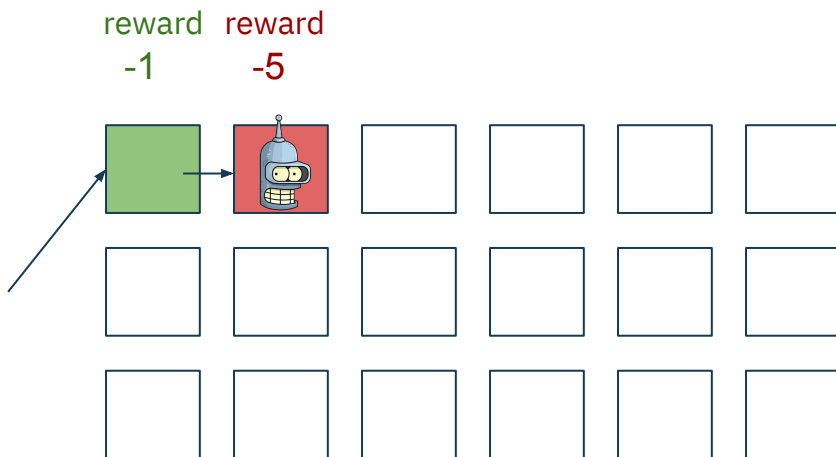technique that trains **agents** to **map states into actions** to **maximize a cumulative reward**

agent

# A definition

Reinforcement Learning

technique that trains **agents** to **map states into actions** to **maximize a** <u>**cumulative reward**</u>

Goal

fast!

agent
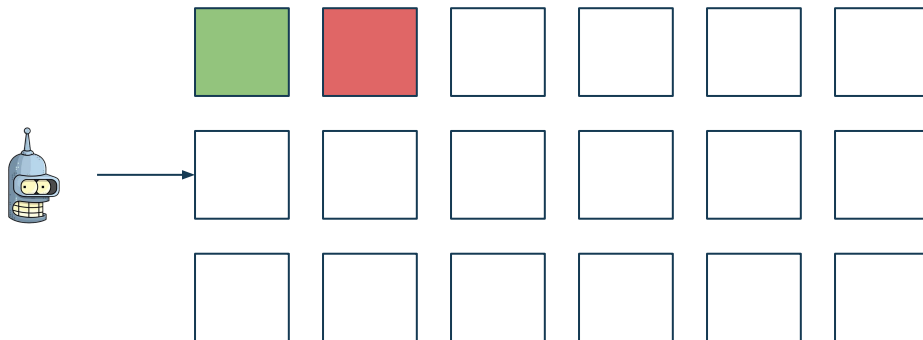
# A definition

Reinforcement
Learning

technique that trains **agents** to **map states into actions** to **maximize a cumulative reward**

actions

agent

# A definition

Reinforcement Learning — technique that trains **agents** to **map states into actions** to **maximize a cumulative reward**

states

actions

agent

# A definition

Reinforcement
Learning

technique that trains **agents** to **map states into actions** to **maximize a cumulative reward**
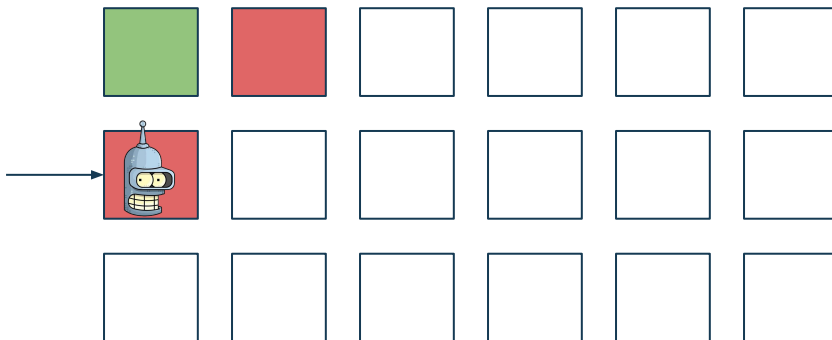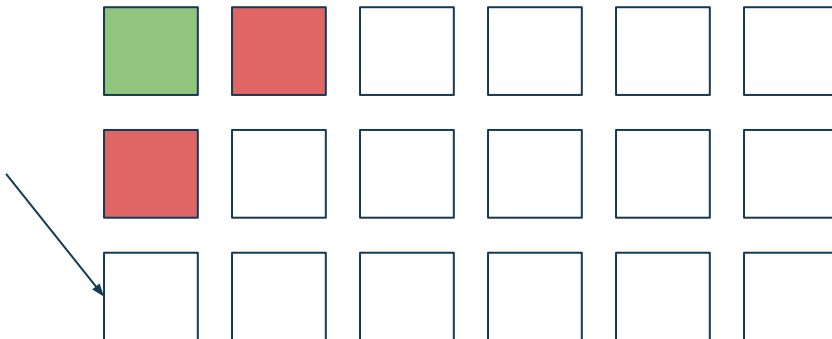
states

action

agent

# A definition

Reinforcement
Learning

technique that trains **agents** to **map states into actions** to **maximize a cumulative reward**

states

action

# A definition

Reinforcement Learning

technique that trains **agents** to **map states into actions** to **maximize a cumulative reward**

reward
-1

# A definition

Reinforcement
Learning

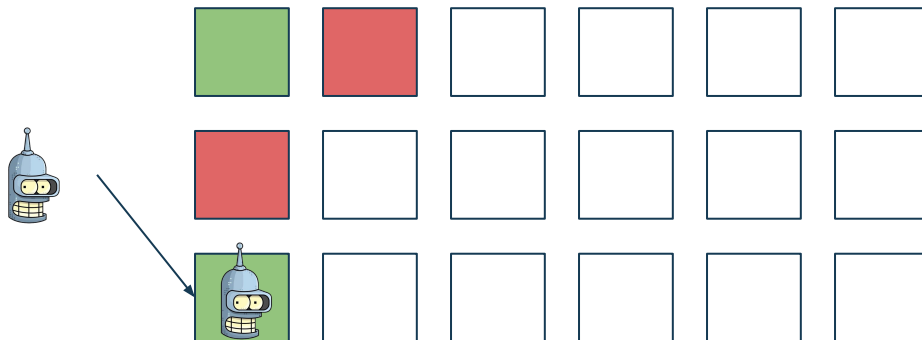technique that trains **agents** to **map states into actions** to **maximize a cumulative reward**

reward
-1

# A definition

technique that trains **agents** to **map states into actions** to **maximize a cumulative reward**
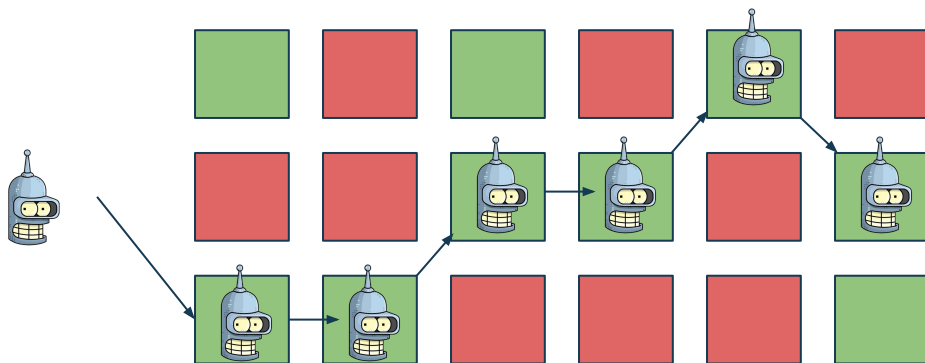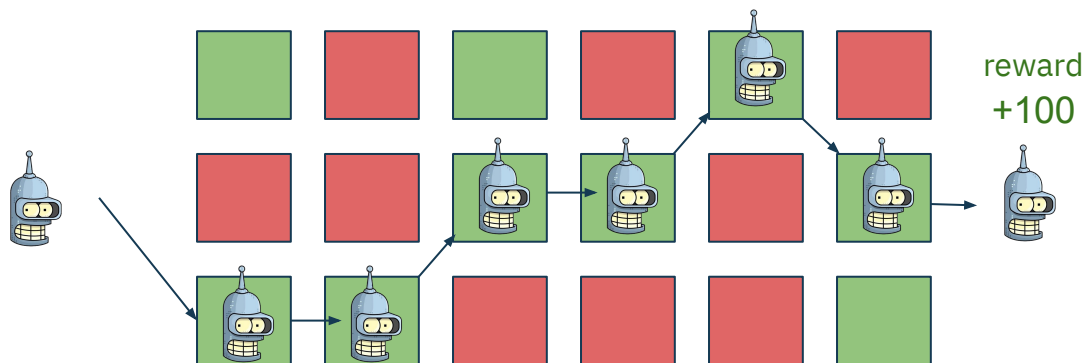


reward
-1

# A definition

Reinforcement Learning

technique that trains **agents** to **map states into actions** to **maximize a cumulative reward**
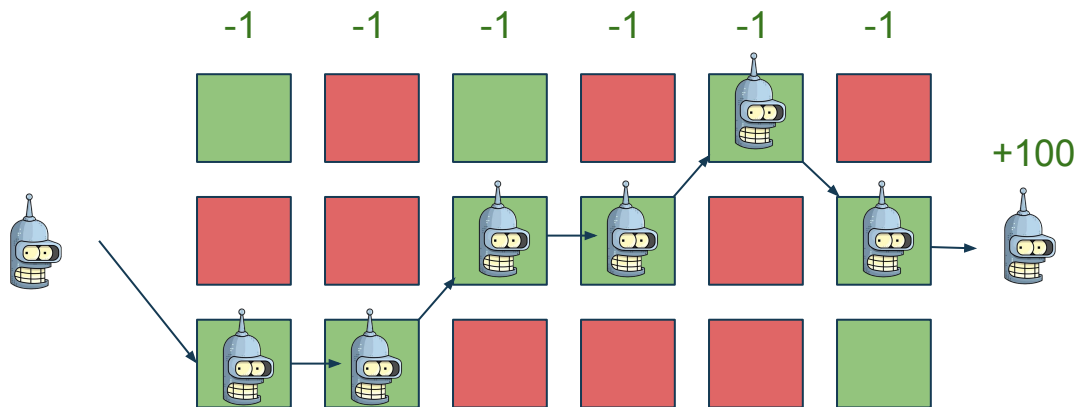
Cumulative Reward

-6

reward
-1

reward
-5

# A definition

Reinforcement Learning

technique that trains **agents** to **map states into actions** to **maximize a cumulative reward**

# A definition

| Reinforcement Learning |
|---|

technique that trains **agents** to **map states into actions** to **maximize a cumulative reward**

| Cumulative Reward |
|---|
| -5 |

reward
-5

# A definition

Reinforcement Learning — technique that trains **agents** to **map states into actions** to **maximize a cumulative reward**

# A definition

Reinforcement Learning

technique that trains **agents** to **map states into actions** to **maximize a cumulative reward**

Cumulative Reward

-1

# A definition

Reinforcement Learning

technique that trains **agents** to **map states into actions** to **maximize a cumulative reward**

Cumulative Reward

-6

# A definition

Reinforcement Learning

technique that trains **agents** to **map states into actions** to **maximize a cumulative reward**

Cumulative Reward

94

reward
+100

# On reward

Goal and reward coherence

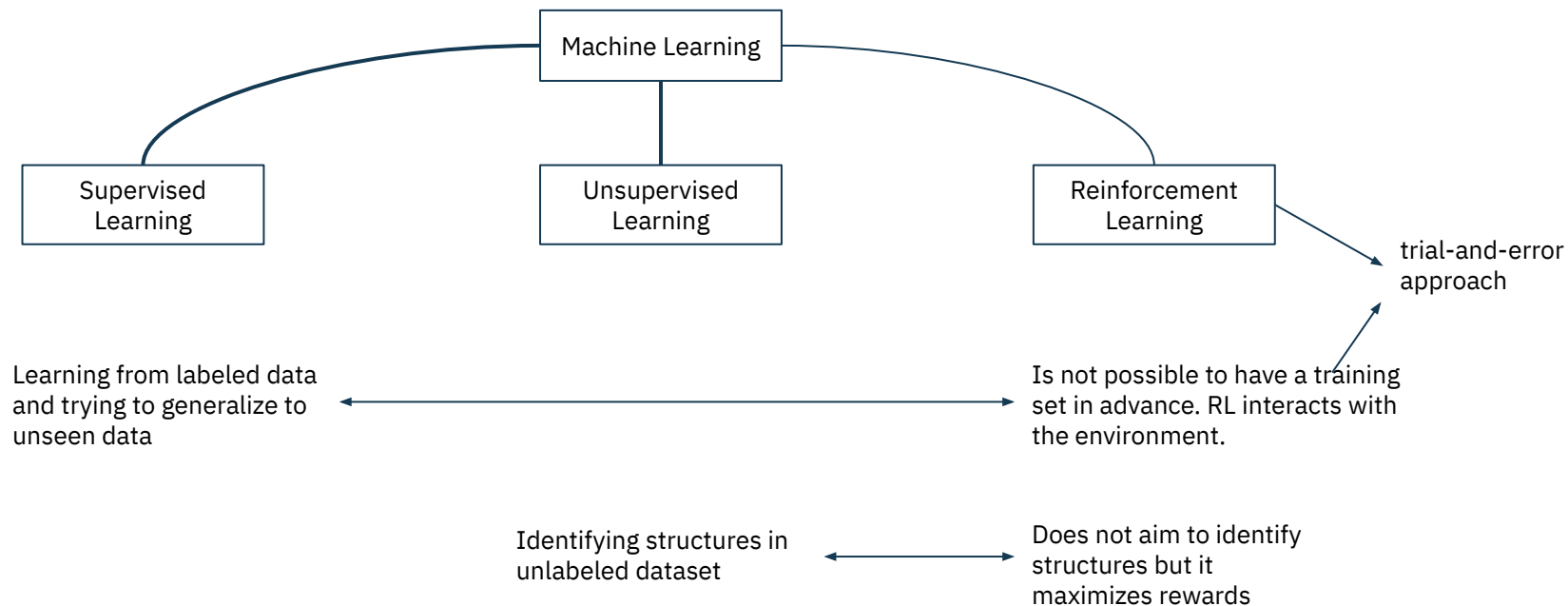we want the agent goes as fast as possible from A to B. We need to choose an appropriate reward signal!
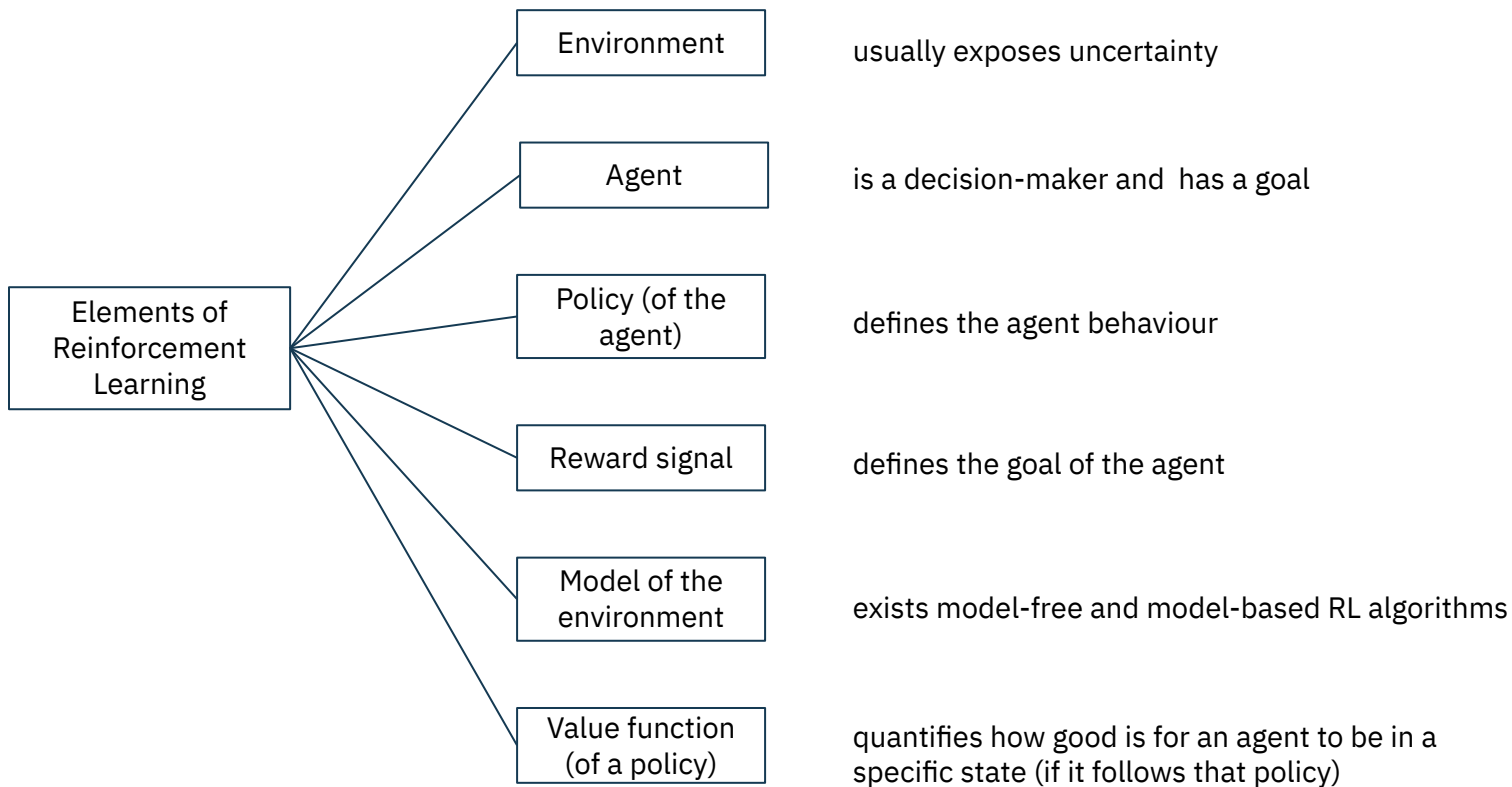
Cumulative Reward

94

-1   -1   -1   -1   -1   -1

+100

# On reward

we want the agent goes as fast as possible from A to B. We need to choose an appropriate reward signal!
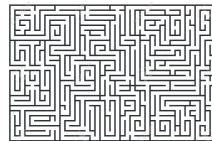
Cumulative Reward

94

Why negative values?

| -1 | -1 | -1 | -1 | -1 | -1 |

+100

# A definition



Machine Learning

Supervised Learning

Unsupervised Learning

Reinforcement Learning

trial-and-error approach

Learning from labeled data and trying to generalize to unseen data

Is not possible to have a training set in advance. RL interacts with the environment.

Identifying structures in unlabeled dataset

Does not aim to identify structures but it maximizes rewards

# Elements of RL

Mathematical definition

# List of the ingredients

```
                        ┌─────────────────┐
                        │   Environment   │        usually exposes uncertainty
                        └─────────────────┘

                        ┌─────────────────┐
                        │      Agent       │        is a decision-maker and  has a goal
                        └─────────────────┘

┌─────────────────┐     ┌─────────────────┐
│  Elements of    │     │  Policy (of the │        defines the agent behaviour
│ Reinforcement   │─────│     agent)      │
│    Learning     │     └─────────────────┘
└─────────────────┘
                        ┌─────────────────┐
                        │  Reward signal  │        defines the goal of the agent
                        └─────────────────┘

                        ┌─────────────────┐
                        │  Model of the   │        exists model-free and model-based RL algorithms
                        │  environment    │
                        └─────────────────┘

                        ┌─────────────────┐
                        │ Value function  │        quantifies how good is for an agent to be in a
                        │  (of a policy)  │        specific state (if it follows that policy)
                        └─────────────────┘
```

# Observability


environment


agent

# Observability



environment

action

agent

# Observability



environment

action

interpreter

reward

observation of the state

agent

© ESTECO SpA

# Observability



environment

interpreter

action

reward

observation of the state

agent

model of the
environment

# Observability



environment

action

interpreter

reward

observation of the state

agent

knowledge

model of the
environment

# Observability

RL = learning + prediction + controlling

Building a model of
the environment

Knowing the
cumulative reward
I'll get following a
policy

Discovering the
best action

# Knowledge of the environment

The two axes of knowledge

# Knowledge of the environment

The two axes of knowledge

Pure planning
problem

Markov
Decision
Process
(MDP)

observability

knowledge of the model

Empirical
knowledge ⟷ Epistemic
knowledge

© ESTECO SpA

# Knowledge of the environment

The two axes of knowledge

Pure planning problem

Markov
Decision
Process
(MDP)

inference

Partially
Observable
MDP

Planning with
uncertainty

We have model but
some params are
unknown

observability

knowledge of the model

Empirical
knowledge

Epistemic
knowledge

© ESTECO SpA

# Knowledge of the environment

## The two axes of knowledge

Pure planning
problem

Model free

Markov
Decision
Process
(MDP)

Trial-and-error
approaches

observability

inference

Full RL

Partially
Observable
MDP

Planning with
uncertainty

We have model but
some params are
unknown

knowledge of the model

Empirical
knowledge

⟷

Epistemic
knowledge

© ESTECO SpA

# Knowledge of the environment

The two axes of knowledge

Markovian process
only matter knowledge
of the actual state

Pure planning
problem

Model free

Trial-and-error
approaches

Markov
Decision
Process
(MDP)

observability

inference

Full RL

Partially
Observable
MDP

Planning with
uncertainty

We have model but
some params are
unknown

knowledge of the model

Empirical
knowledge

Epistemic
knowledge

# (finite) Markov Decision Process



| trajectory | $S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, R_3, \ldots$ |

| dynamics | $p(s', r \mid s, a) \;\doteq\; \Pr\{S_t = s', R_t = r \mid S_{t-1} = s, A_{t-1} = a\}$ |

$$\sum_{s' \in \mathcal{S}} \sum_{r \in \mathcal{R}} p(s', r \mid s, a) = 1, \;\; \text{for all } s \in \mathcal{S}, a \in \mathcal{A}(s)$$

Perfect knowledge
of the model

# (finite) Markov Decision Process

| | |
|---|---|
| trajectory | $S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, R_3, \ldots$ |
| dynamics | $p(s', r \,|\, s, a) \;\doteq\; \Pr\{S_t = s', R_t = r \mid S_{t-1} = s, A_{t-1} = a\}$ |

| | |
|---|---|
| state-transition probability | $p(s' \,|\, s, a) \;\doteq\; \Pr\{S_t = s' \mid S_{t-1} = s, A_{t-1} = a\} \;=\; \sum_{r \in \mathcal{R}} p(s', r \,|\, s, a)$ |
| expected reward (I) | $r(s, a) \;\doteq\; \mathbb{E}[R_t \mid S_{t-1} = s, A_{t-1} = a] \;=\; \sum_{r \in \mathcal{R}} r \sum_{s' \in \mathcal{S}} p(s', r \,|\, s, a)$ |
| expected reward (II) | $r(s, a, s') \;\doteq\; \mathbb{E}[R_t \mid S_{t-1} = s, A_{t-1} = a, S_t = s'] \;=\; \sum_{r \in \mathcal{R}} r \, \frac{p(s', r \,|\, s, a)}{p(s' \,|\, s, a)}$ |

# Reward signal

I have my goal

**Reward hypothesis:** that all of what we mean by <u>goals and purposes</u> can be well thought of as the <u>maximization of the expected value of the cumulative sum of a received scalar signal (called reward).</u>

| Reward | $R_{t+1}$ |

$$G_t \doteq R_{t+1} + R_{t+2} + R_{t+3} + \cdots + R_T$$

Return

$$G_t \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

Discounted Return

# Reward signal

I have my goal

**Reward hypothesis:** that all of what we mean by <u>goals and purposes</u> can be well thought of as the <u>maximization of the expected value of the cumulative sum of a received scalar signal (called reward).</u>

| Reward | $R_{t+1}$ |
|---|---|

Short-term view

| Return | $G_t \doteq R_{t+1} + R_{t+2} + R_{t+3} + \cdots + R_T$ |
|---|---|

| Discounted Return | $G_t \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$ |
|---|---|

Long-term view

# Reward signal

I have my goal

**Reward hypothesis:** that all of what we mean by <u>goals and purposes</u> can be well thought of as the <u>maximization of the expected value of the cumulative sum of a received scalar signal (called reward).</u>

Reward $R_{t+1}$

Short-term view

Return $G_t \doteq R_{t+1} + R_{t+2} + R_{t+3} + \cdots + R_T$

Discounted Return $G_t \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$

Long-term view

$G_t = R_{t+1} + \gamma G_{t+1}$  **RECURSIVE DEFINITION**

# Policy

| Policy |
|---|

is a mapping from states to probabilities of selecting each possible action

$$\pi : \mathcal{S} \times \mathcal{A} \to [0, 1]$$

If we are at time t, $\pi(a|s)$ is the probability of having $A_t = a \wedge S_t = s$

can be *deterministic*

# Value function

| Value Function | is a function that quantify how good is to be on a state and follows a specific policy |

$$v_\pi : \mathcal{S} \to \mathbb{R}$$

| state-value function | $v_\pi(s) \;\dot{=}\; \mathbb{E}_\pi[G_t \mid S_t = s] \;=\; \mathbb{E}_\pi\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \;\middle|\; S_t = s\right], \text{ for all } s \in \mathcal{S}$ |

| action-value function | $q_\pi(s, a) \;\dot{=}\; \mathbb{E}_\pi[G_t \mid S_t = s, A_t = a] \;=\; \mathbb{E}_\pi\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \;\middle|\; S_t = s, A_t = a\right]$ |

# Solving a RL problem

find a policy that achieves the maximum reward over the long run

optimal policy $\qquad \pi_* \succeq \pi \quad \forall \pi \in$ policies

# Solving a RL problem

find a policy that achieves the maximum reward over the long run

optimal policy

$$\pi_* \succeq \pi \quad \forall \pi \in \text{policies}$$

$$\pi' \succeq \pi \iff \forall s \in \mathcal{S}, \ v_{\pi'}(s) \geq v_{\pi}(s)$$

# Solving a RL problem

find a policy that achieves the maximum reward over the long run

| optimal policy |

$$\pi_* \succeq \pi \quad \forall \pi \in \text{policies}$$

$$\pi' \succeq \pi \iff \forall s \in \mathcal{S}, \ v_{\pi'}(s) \geq v_\pi(s)$$

| optimal state-value function |

$$v_*(s) \doteq \max_\pi v_\pi(s)$$

| optimal action-value function |

$$q_*(s, a) \doteq \max_\pi q_\pi(s, a)$$

# Solving a RL problem

find a policy that achieves the maximum reward over the long run

optimal policy

$$\pi_* \succeq \pi \quad \forall \pi \in \text{policies}$$

$$\pi' \succeq \pi \iff \forall s \in \mathcal{S}, \; v_{\pi'}(s) \geq v_\pi(s)$$

optimal state-value function

$$v_*(s) \doteq \max_\pi v_\pi(s)$$

optimal action-value function

$$q_*(s, a) \doteq \max_\pi q_\pi(s, a)$$

$$q_*(s, a) = \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a]$$

# Dynamic Programming

How to solve MDP problems

# Dynamic Programming
## Mr. Richard Ernest Bellman



Bellman, 1950s

Algorithm paradigm useful to solve a specific class of <u>problems</u> that can be decomposed in <u>sub-problems</u> in <u>recursive</u> way

# Dynamic Programming

## In the RL context

Collection of algorithms that can be used to compute optimal policies given a perfect model of the environment as a MDP.

Bellman, 1950s

**Key idea:** use value function to organize and structure the search of optimal policies

Consistency relation of state-value function

$$
\begin{aligned}
v_\pi(s) &\doteq \mathbb{E}_\pi[G_t \mid S_t = s] \\
&= \mathbb{E}_\pi[R_{t+1} + \gamma G_{t+1} \mid S_t = s] \\
&= \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r | s, a)\Big[r + \gamma \mathbb{E}_\pi[G_{t+1} | S_{t+1} = s']\Big] \\
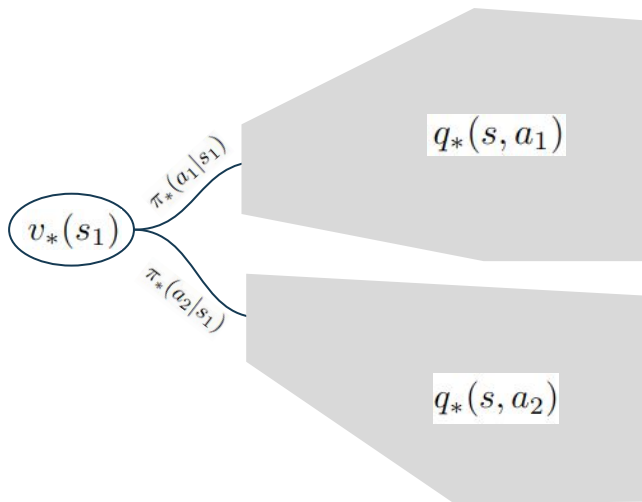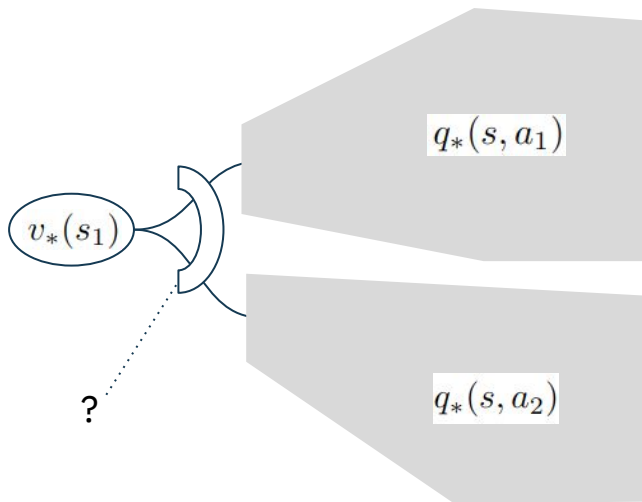&= \sum_a \pi(a|s) \sum_{s',r} p(s', r | s, a)\Big[r + \gamma v_\pi(s')\Big], \quad \text{for all } s \in \mathcal{S}.
\end{aligned}
$$

# Dynamic Programming

Towards the Bellman Equation

Consistency relation
of state-value
function

$$v_\pi(s) = \sum_a \pi(a|s) \sum_{s',r} p(s', r|s, a) \left[ r + \gamma v_\pi(s') \right]$$

# Dynamic Programming

Towards the Bellman Equation

Consistency relation
of state-value
function

$$v_\pi(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) \left[ r + \gamma v_\pi(s') \right]$$

# Dynamic Programming

Towards the Bellman Equation

Consistency relation
of state-value
function

$$v_\pi(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) \left[ r + \gamma v_\pi(s') \right]$$

# Dynamic Programming

Towards the Bellman Equation

Consistency relation of state-value function

$$v_\pi(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) \left[ r + \gamma v_\pi(s') \right]$$



$q_\pi(s, a_1)$

$v_\pi(s_1)$

$\pi(a_1|s_1)$

$\pi(a_2|s_1)$

$q_\pi(s, a_2)$

# Dynamic Programming

Towards the Bellman Equation

Consistency relation
of state-value
function

$$v_\pi(s) = \sum_a \pi(a|s) q_\pi(s, a)$$

$q_\pi(s, a_1)$

$v_\pi(s_1)$

$\pi(a_1|s_1)$

$\pi(a_2|s_1)$

$q_\pi(s, a_2)$

# Dynamic Programming

Towards the Bellman Equation

Consistency relation of state-value function

$$v_\pi(s) = \sum_a \pi(a|s) q_\pi(s, a)$$



$q_\pi(s, a_1)$

$v_\pi(s_1)$

$\pi(a_1|s_1)$

$\pi(a_2|s_1)$

$q_\pi(s, a_2)$

© ESTECO SpA

# Dynamic Programming

Towards the Bellman Equation

What about the optimal policy and the optimal state-value function?

Consistency relation of state-value function

$$v_\pi(s) = \sum_a \pi(a|s) q_\pi(s, a)$$

It's an average

$v_\pi(s_1)$

$\pi(a_1|s_1)$

$\pi(a_2|s_1)$

$q_\pi(s, a_1)$

$q_\pi(s, a_2)$

# Dynamic Programming
## Towards the Bellman Equation

What about the optimal policy and the optimal state-value function?

Consistency relation of state-value function

$$v_\pi(s) = \sum_a \pi(a|s) q_\pi(s, a)$$

It's an average

The optimal policy is a policy so it should satisfy the consistency relation

$v_\pi(s_1)$

$\pi(a_1|s_1)$

$\pi(a_2|s_1)$

$q_\pi(s, a_1)$

$q_\pi(s, a_2)$

# Dynamic Programming
Towards the Bellman Equation

What about the optimal policy and the optimal state-value function?

Consistency relation of state-value function

$$v_*(s) = \sum_a \pi_*(a|s) q_*(s, a)$$

It's an average

The optimal policy is a policy so it should satisfy the consistency relation

$v_*(s_1)$

$\pi_*(a_1|s_1)$

$\pi_*(a_2|s_1)$

$q_*(s, a_1)$

$q_*(s, a_2)$

# Dynamic Programming

Towards the Bellman Equation

What about the optimal policy and the optimal state-value function?

Consistency relation of state-value function

$$v_*(s) = \sum_a \pi_*(a|s) q_*(s, a)$$

It's an average

The optimal policy is a policy so it should satisfy the consistency relation

The optimal policy is *optimal*

$q_*(s, a_1)$

$v_*(s_1)$

$\pi_*(a_1|s_1)$

$\pi_*(a_2|s_1)$

$q_*(s, a_2)$

# Dynamic Programming
Towards the Bellman Equation

Consistency relation of state-value function

$$v_*(s) = \sum_a \pi_*(a|s) q_*(s, a)$$

What about the optimal policy and the optimal state-value function?

It's an average

The optimal policy is a policy so it should satisfy the consistency relation

The optimal policy is *optimal*



$q_*(s, a_1)$

$v_*(s_1)$

$q_*(s, a_2)$

?

# Dynamic Programming

Bellman Equation

Bellman equation

$$v_*(s) = \max_a q_*(s, a)$$



$q_*(s, a_1)$

$v_*(s_1)$

max

$q_*(s, a_2)$

# Dynamic Programming

Bellman Equation

Bellman equation

$$v_*(s) = \max_a q_*(s, a)$$

$$v_*(s) = \max_a \sum_{s', r} p(s', r | s, a)[r + \gamma v_*(s')]$$

# Dynamic Programming

Bellman Equation

$$v_*(s) = \max_a \sum_{s',r} p(s',r|s,a)[r + \gamma v_*(s')]$$

Bellman equation

$$q_*(s,a) = \sum_{s',r} p(s',r|s,a)[r + \gamma \max_a q_*(s',a')]$$

# Dynamic Programming

How to find the optimal policy?

$$\pi \xrightarrow{\text{Iterative procedure}} \pi^*$$

# Dynamic Programming

How to find the optimal policy?

Consistency relation of the
state-value function

$\pi$ ———— Policy evaluation ————→ $v_\pi$ ————————→ $\pi'$

# Dynamic Programming

How to find the optimal policy?

Consistency relation of the
state-value function

Bellman intuition

$$\pi \xrightarrow{\text{Policy evaluation}} v_\pi \xrightarrow{\text{Policy improvement}} \pi'$$

# Dynamic Programming

How to find the optimal policy?

Bellman intuition

Consistency relation of the
state-value function

$$\pi \xrightarrow{\text{Policy evaluation}} v_\pi \xrightarrow{\text{Policy improvement}} \pi'$$

Policy Iteration

# Dynamic Programming

How to find the optimal policy?

Bellman intuition

Consistency relation of the
state-value function

$$\pi \xrightarrow{\quad \text{Policy evaluation} \quad} v_\pi \xrightarrow{\quad \text{Policy improvement} \quad} \pi'$$

Policy Iteration

$$\pi_0 \xrightarrow{\ \text{E}\ } v_{\pi_0} \xrightarrow{\ \text{I}\ } \pi_1 \xrightarrow{\ \text{E}\ } v_{\pi_1} \xrightarrow{\ \text{I}\ } \pi_2 \xrightarrow{\ \text{E}\ } \cdots \xrightarrow{\ \text{I}\ } \pi_* \xrightarrow{\ \text{E}\ } v_*$$

Does it converge? Yes

# Dynamic Programming

## Policy evaluation

| Consistency relation of state-value function |
|:---:|

$$v_\pi(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)\left[r + \gamma v_\pi(s')\right]$$

# Dynamic Programming

## Policy evaluation

| Consistency relation of state-value function | $$v_\pi(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) \left[ r + \gamma v_\pi(s') \right]$$ |

| Iterative policy evaluation | $$v_{k+1}(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)[r + \gamma v_k(s')]$$ |

# Dynamic Programming

Policy evaluation

| Consistency relation of state-value function |
|---|

$$v_\pi(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)\left[r + \gamma v_\pi(s')\right]$$

| Iterative policy evaluation |
|---|

$$v_{k+1}(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)[r + \gamma v_k(s')]$$

2 ways of updating: *in-place* vs *two arrays version*

Faster, depends on ordering of update



propagation

# Policy Evaluation

Algorithm

**Iterative Policy Evaluation, for estimating $V \approx v_\pi$**

Input $\pi$, the policy to be evaluated
Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation
Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(terminal) = 0$

Loop:
$\quad \Delta \leftarrow 0$
$\quad$ Loop for each $s \in \mathcal{S}$:
$\qquad v \leftarrow V(s)$
$\qquad V(s) \leftarrow \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) \big[ r + \gamma V(s') \big]$
$\qquad \Delta \leftarrow \max(\Delta, |v - V(s)|)$
until $\Delta < \theta$

# Policy Evaluation

Algorithm

**Iterative Policy Evaluation, for estimating $V \approx v_\pi$**

Input $\pi$, the policy to be evaluated
Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation
Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(terminal) = 0$

Loop:
$\quad \Delta \leftarrow 0$
$\quad$ Loop for each $s \in \mathcal{S}$:
$\quad\quad v \leftarrow V(s)$
$\quad\quad V(s) \leftarrow \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) \big[ r + \gamma V(s') \big]$ consistency relation
$\quad\quad \Delta \leftarrow \max(\Delta, |v - V(s)|)$
until $\Delta < \theta$

# Policy Evaluation

Algorithm

**Iterative Policy Evaluation, for estimating $V \approx v_\pi$**

Input $\pi$, the policy to be evaluated
Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation
Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(terminal) = 0$

Loop:
    $\Delta \leftarrow 0$
    Loop for each $s \in \mathcal{S}$:
        $v \leftarrow V(s)$
        $V(s) \leftarrow \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)\big[r + \gamma V(s')\big]$    consistency relation
        $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
until $\Delta < \theta$    Stability of state-value function

# Policy Evaluation

Example



$$R_t = -1$$
on all transitions

Uniform Policy

actions

non-terminal state

terminal state

# Policy Evaluation

Example - 1st iteration



$$v_\pi(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)[r + \gamma v_\pi(s')]$$

# Policy Evaluation

Example - 1st iteration



$$v_\pi(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)[r + \gamma v_\pi(s')]$$

-1          0

# Policy Evaluation

Example - 1st iteration

| 0.0 | 0.0 | 0.0 | 0.0 |
|-----|-----|-----|-----|
| 0.0 | 0.0 | 0.0 | 0.0 |
| 0.0 | 0.0 | 0.0 | 0.0 |
| 0.0 | 0.0 | 0.0 | 0.0 |

$\Longrightarrow$

| 0.0 | -1.0 | -1.0 | -1.0 |
|-----|------|------|------|
| -1.0 | -1.0 | -1.0 | -1.0 |
| -1.0 | -1.0 | -1.0 | -1.0 |
| -1.0 | -1.0 | -1.0 | 0.0 |

$$v_\pi(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)[r + \gamma v_\pi(s')]$$

# Policy Evaluation

Example - 1st iteration



$$v_\pi(s) = \sum_a \pi(a|s) \sum_{s',r} p(s', r|s, a)[r + \gamma v_\pi(s')]$$

# Policy Evaluation

Example - 2nd iteration

← -1/3 ⊕

| 0.0 | -1.0 | -1.0 | -1.0 |
|------|------|------|------|
| -1.0 | -1.0 | -1.0 | -1.0 |
| -1.0 | -1.0 | -1.0 | -1.0 |
| -1.0 | -1.0 | -1.0 | 0.0 |

⟹

| 0.0 | ? |   |   |
|-----|---|---|---|
| ?   |   |   |   |
|     |   |   | ? |
|     |   | ? | 0.0 |

$$v_\pi(s) = \sum_a \underline{\pi(a|s)} \sum_{s',r} \underline{p(s',r|s,a)}[\underline{r} + \gamma \underline{v_\pi(s')}]$$

1/3      1    -1    0

# Policy Evaluation

Example - 2nd iteration



$$v_\pi(s) = \sum_a \underline{\pi(a|s)} \sum_{s',r} \underline{p(s',r|s,a)}[\underline{r} + \underline{\gamma v_\pi(s')}]$$

1/3       1    -1    -1

# Policy Evaluation

Example - 2nd iteration

←    -1/3   ✛

↓    -2/3   ✛

→    -2/3   ✛

-1.7

| 0.0 | -1.0 | -1.0 | -1.0 |
|------|------|------|------|
| -1.0 | -1.0 | -1.0 | -1.0 |
| -1.0 | -1.0 | -1.0 | -1.0 |
| -1.0 | -1.0 | -1.0 | 0.0 |

⟹

| 0.0 | ? | | |
|------|---|---|---|
| ? | | | |
| | | | ? |
| | | ? | 0.0 |

$$v_\pi(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)[r + \gamma v_\pi(s')]$$

© ESTECO SpA

# Policy Evaluation

Example - 2nd iteration

-1/3 ✛

-2/3 ✛

-2/3 ✛

-1.7

| 0.0 | -1.0 | -1.0 | -1.0 |
|------|------|------|------|
| -1.0 | -1.0 | -1.0 | -1.0 |
| -1.0 | -1.0 | -1.0 | -1.0 |
| -1.0 | -1.0 | -1.0 | 0.0 |

⟹

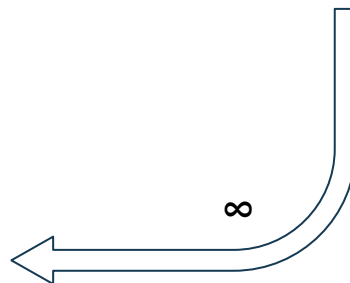| 0.0 | -1.7 | | |
|------|------|------|------|
| -1.7 | | | |
| | | | -1.7 |
| | | -1.7 | 0.0 |

$$v_\pi(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)[r + \gamma v_\pi(s')]$$

# Policy Evaluation

Example - 2nd iteration

| 0.0 | -1.0 | -1.0 | -1.0 |
|------|------|------|------|
| -1.0 | -1.0 | -1.0 | -1.0 |
| -1.0 | -1.0 | -1.0 | -1.0 |
| -1.0 | -1.0 | -1.0 | 0.0 |

$\Longrightarrow$

| 0.0 | -1.7 | -2.0 | -2.0 |
|------|------|------|------|
| -1.7 | -2.0 | -2.0 | -2.0 |
| -2.0 | -2.0 | -2.0 | -1.7 |
| -2.0 | -2.0 | -1.7 | 0.0 |

$$v_\pi(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)[r + \gamma v_\pi(s')]$$

# Policy Evaluation

Example - until the end

| 0.0 | 0.0 | 0.0 | 0.0 |
|---|---|---|---|
| 0.0 | 0.0 | 0.0 | 0.0 |
| 0.0 | 0.0 | 0.0 | 0.0 |
| 0.0 | 0.0 | 0.0 | 0.0 |

1 →

| 0.0 | -1.0 | -1.0 | -1.0 |
|---|---|---|---|
| -1.0 | -1.0 | -1.0 | -1.0 |
| -1.0 | -1.0 | -1.0 | -1.0 |
| -1.0 | -1.0 | -1.0 | 0.0 |

2 →

| 0.0 | -1.7 | -2.0 | -2.0 |
|---|---|---|---|
| -1.7 | -2.0 | -2.0 | -2.0 |
| -2.0 | -2.0 | -2.0 | -1.7 |
| -2.0 | -2.0 | -1.7 | 0.0 |

3 →

| 0.0 | -2.4 | -2.9 | -3.0 |
|---|---|---|---|
| -2.4 | -2.9 | -3.0 | -2.9 |
| -2.9 | -3.0 | -2.9 | -2.4 |
| -3.0 | -2.9 | -2.4 | 0.0 |

∞

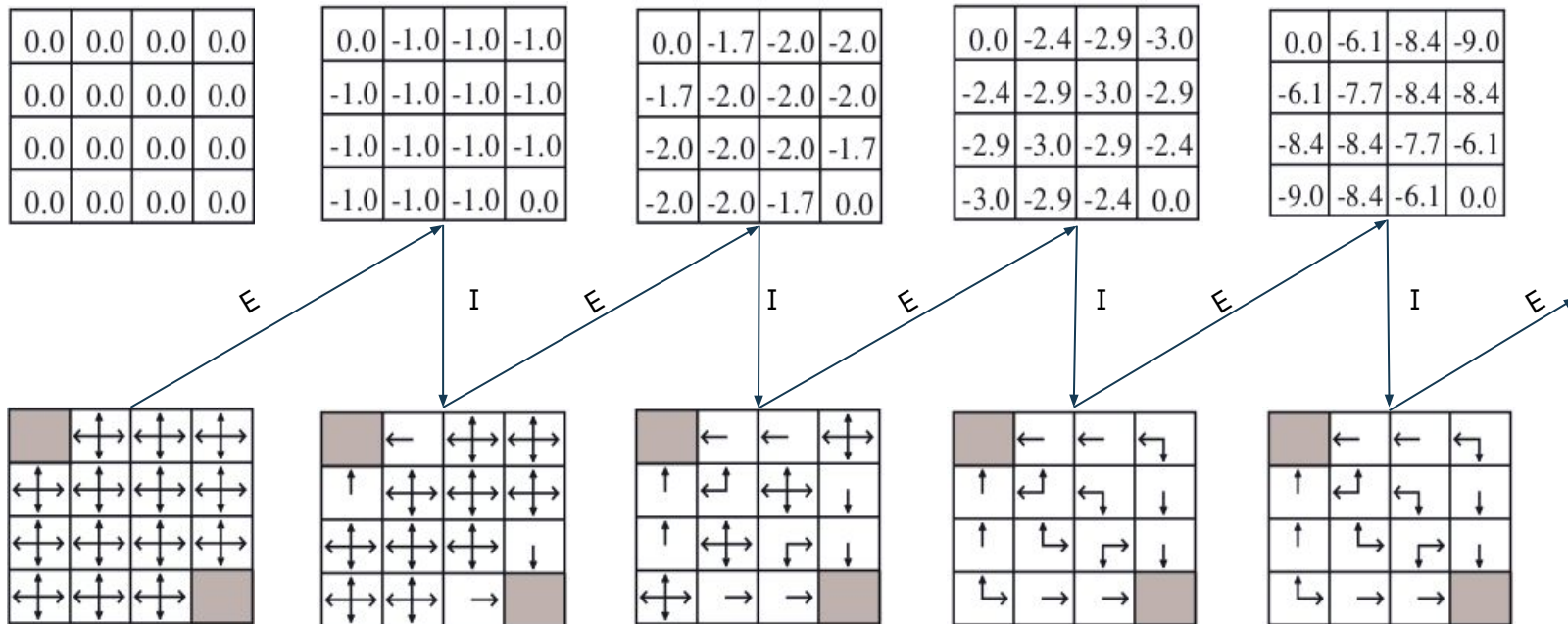| 0.0 | -14. | -20. | -22. |
|---|---|---|---|
| -14. | -18. | -20. | -20. |
| -20. | -20. | -18. | -14. |
| -22. | -20. | -14. | 0.0 |

# Policy Evaluation

Example - until the end

# Policy Evaluation

Example - until the end

# Policy Evaluation

Example - until the end

# Policy Improvement

How to find better policies

$$v_\pi \xrightarrow{\text{Policy improvement}} \pi'$$

| 0.0 | -14. | -20. | -22. |
|-----|------|------|------|
| -14. | -18. | -20. | -20. |
| -20. | -20. | -18. | -14. |
| -22. | -20. | -14. | 0.0 |

Policy improvement theorem

$$q_\pi(s, \pi'(s)) \geq v_\pi(s), \forall s \in \mathcal{S} \Rightarrow v_{\pi'}(s) \geq v_\pi(s), \forall s \in \mathcal{S}$$

Greedy policy approach

$$
\begin{aligned}
\pi'(s) &\doteq \arg\max_a q_\pi(s, a) \\
&= \arg\max_a \mathbb{E}[R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s, A_t = a] \\
&= \arg\max_a \sum_{s', r} p(s', r \mid s, a)\left[r + \gamma v_\pi(s')\right],
\end{aligned}
$$

# Policy Iteration

Example

# Policy Iteration

Example

propagation effect



policy convergence
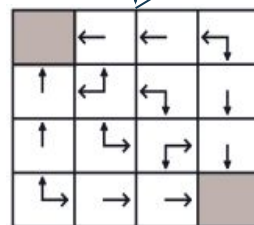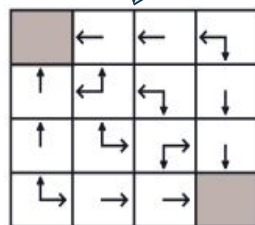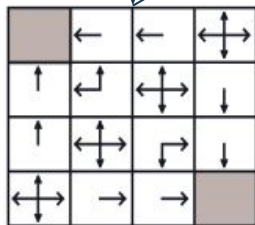
# Policy Iteration
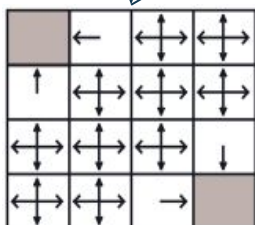
Example



propagation effect

policy convergence

# Policy Iteration

**Policy Iteration (using iterative policy evaluation) for estimating $\pi \approx \pi_*$**

1. Initialization
   $V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in \mathcal{S}$

2. Policy Evaluation
   Loop:
   $\quad \Delta \leftarrow 0$
   $\quad$ Loop for each $s \in \mathcal{S}$:
   $\quad\quad v \leftarrow V(s)$
   $\quad\quad V(s) \leftarrow \sum_{s',r} p(s',r \,|\, s, \pi(s)) \big[ r + \gamma V(s') \big]$
   $\quad\quad \Delta \leftarrow \max(\Delta, |v - V(s)|)$
   until $\Delta < \theta$ (a small positive number determining the accuracy of estimation)

3. Policy Improvement
   $policy\text{-}stable \leftarrow true$
   For each $s \in \mathcal{S}$:
   $\quad old\text{-}action \leftarrow \pi(s)$
   $\quad \pi(s) \leftarrow \arg\max_a \sum_{s',r} p(s',r \,|\, s, a) \big[ r + \gamma V(s') \big]$
   $\quad$ If $old\text{-}action \neq \pi(s)$, then $policy\text{-}stable \leftarrow false$
   If $policy\text{-}stable$, then stop and return $V \approx v_*$ and $\pi \approx \pi_*$; else go to 2

# Policy Iteration

**Policy Iteration (using iterative policy evaluation) for estimating $\pi \approx \pi_*$**

1. Initialization
   $V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in \mathcal{S}$

2. Policy Evaluation
   Loop:
       $\Delta \leftarrow 0$
       Loop for each $s \in \mathcal{S}$:
           $v \leftarrow V(s)$
           $V(s) \leftarrow \sum_{s',r} p(s',r\,|\,s,\pi(s))\big[r + \gamma V(s')\big]$
           $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
   until $\Delta < \theta$ (a small positive number determining the accuracy of estimation)

3. Policy Improvement
   $policy\text{-}stable \leftarrow true$
   For each $s \in \mathcal{S}$:
       $old\text{-}action \leftarrow \pi(s)$
       $\pi(s) \leftarrow \arg\max_a \sum_{s',r} p(s',r\,|\,s,a)\big[r + \gamma V(s')\big]$
       If $old\text{-}action \neq \pi(s)$, then $policy\text{-}stable \leftarrow false$
   If $policy\text{-}stable$, then stop and return $V \approx v_*$ and $\pi \approx \pi_*$; else go to 2

# Value Iteration

Solving efficiently the Policy Iteration

**Value Iteration, for estimating $\pi \approx \pi_*$**

Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation
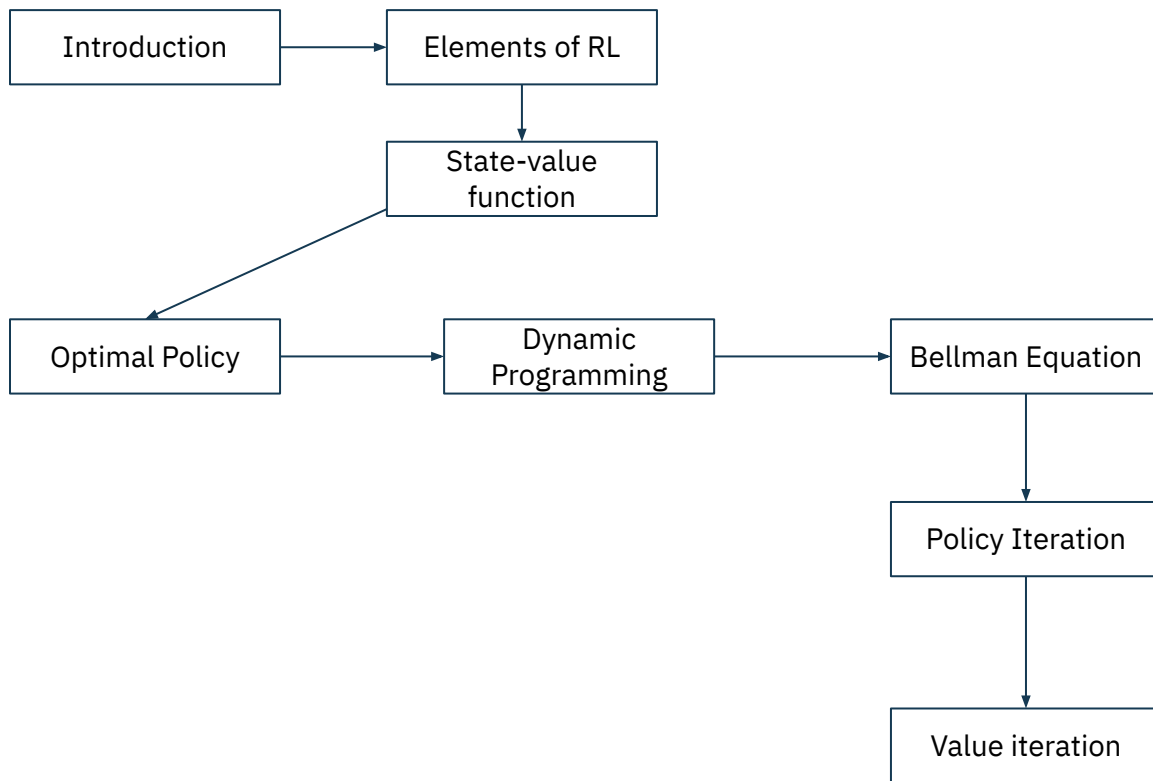Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(terminal) = 0$

Loop:
| $\quad \Delta \leftarrow 0$
| $\quad$ Loop for each $s \in \mathcal{S}$:
| $\quad\quad v \leftarrow V(s)$
| $\quad\quad V(s) \leftarrow \max_a \sum_{s',r} p(s',r|s,a)\big[r + \gamma V(s')\big]$
| $\quad\quad \Delta \leftarrow \max(\Delta, |v - V(s)|)$
until $\Delta < \theta$

Output a deterministic policy, $\pi \approx \pi_*$, such that
$\quad \pi(s) = \arg\max_a \sum_{s',r} p(s',r|s,a)\big[r + \gamma V(s')\big]$

# Recap

# Thank you

esteco.com