

#### ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA Interdepartmental research centre Alma Mater Research Institute for Human-centered Artificial Intelligence

# A Brief Introduction to Reinforcement Learning

Michele Lombardi

DISI/ALMA-AI – University of Bologna

### **Three Words for Three Types of Learning**

#### **Supervised Learning**

... is about learning from someone





ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA Interdepartmental Research centre Alma Mater Research Institute for Human-centered Artificial Intelligence

### **Three Words for Three Types of Learning**

#### **Supervised Learning**

... is about learning from someone

#### **Unsupervised Learning**

... is about looking at patterns





ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA Interdepartmental Research centre Alma Mater Research Institute for Human-centered Artificial Intelligence

### **Three Words for Three Types of Learning**

#### **Supervised Learning**

... is about learning from someone

#### **Unsupervised Learning**

... is about looking at patterns

#### **Reinforcement Learning**

... is about learning by doing





ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA Interdepartmental Research centre Alma Mater Research Institute for Uman-centered Artificial Intelligenci

#### Learning how to play Atari games



THE STORE

ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA Interdepartmental research centre Alma Mater Research institute for Human-centered Artificial Intelligence

https://www.youtube.com/watch?v=eG1Ed8PTJ18

#### Beating world masters at Go



https://www.youtube.com/watch?v=\_OV0Hlj8Fb8

ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA Interdepartmental research centre Alma Mater Research Institute for Human-centered Artificial Intelligence

#### **Beating professionals at team-based esports**

Human View	AI View					
	3.006	-1.386	-0.4695	0.883	1	0.84
	-0.3154	-0.5425	-0.5	0.866	0	0.82
	3.11	-1.36	-0.9336	0.3584	1	0.78
	-2.324	2.863	0.9746	0.225	0	0.86
	3.037	-1.361	-0.7773	0.6294	1	0.82
	-1.387	2.951	0.988	0.1565	0	0.74
	3.023	-0.9395	0.05234	-0.9985	0	0.66
	2.951	-0.5747	0.01746	1	0	0.72
	2.963	-1.303	0.3906	0.9204	0	0.68
	2.834	-3.164	0.01746	-1	0	0.68
	3.127	-1.368	0.6562	0.755	1	0.55
	3.088	-1.366	0.4695	0.883	0	0.55
A CONTRACT MARKED	2.984	-1.398	-0.225	0.9746	1	0.55
	3.037	-1.391	0.788	0.6157	0	0.55
	3.076	-1.438	0.883	0.4695	0	0.55
	-2.412	2.846	0.996	0.08716	1	0.3



https://www.twitch.tv/videos/410533063?t=01h32m02s

ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA Interdepartmental Research centre Alma mater Research Institute for Human-centered Artificial Intelligence

#### **Robotics**



https://www.youtube.com/watch?v=Q4bMcUk6pcw&feature=emb\_logo



ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA Interdepartmental research centre Alma mater research institute for Human-centered artificial intelligence

## **Ingredients of Reinforcement Learning**

### In RL there's no dataset, but an environment to play with

At each step t:

- We observe the environment state
- We perform an action

#### Then:

- We receive a reward
- The environment moves to the next state



• I.e. what action to choose, based on the observed state





ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA Interdepartmental research centre Alma Mater Research Institute for JMAN-Centered Artificial Intelligenc

### A Simple Example

### A food inventory control problem

- Producing yields one 🍑 in the next step
- Producing costs one

- Orders require one or more
- Orders always deplete the inventory
- Un-met orders cost three 💰 💰 💰
- Orders are not a priori known



 $\rightarrow$ 



ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA Interdepartmental research centre Alma Mater Research Institute for JMAN-Centered Artificial Intelligens

### **A Simple Example**

### A possible representation

- Steps = opportunities to produce
- Actions = whether to produce or not
- State = inventory situation + current step



How do we determine an optimal policy?



-CENTERED ARTIFICIAL INTELLIGENC

#### Let's consider all possible choices





ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA Interdepartmental research centre Alma Mater Research institute for Human-centered Artificial intelligence

#### ...And then go backwards



• For the last step we know the quality of each choice...



ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA Interdepartmental research centre Alma mater research institute for Human-centered artificial intelligence

#### ...And then go backwards



- For the last step we know the quality of each choice...
- ...And therefore the value best possible outcome



ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA Interdepartmental Research centre Alma Mater Research Institute for Iuman-centered Artificial Intelligence

#### ...And then go backwards



- Then we can do the same for the qualities in the previous step
- ...and the same for the values



ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA Interdepartmental Research centre Alma Mater Research Institute for Iuman-centered Artificial Intelligence

Now we just need to pick always pick the best Q



- The Q values define something called a Q-function
- ...And this approach is known as dynamic programming



ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA Nterdepartmental research centre Alma mater research institute for Man-centered artificial institute for

#### In practice, enumeration is not viable

But we can circumvent it:

• Every optimal Q function satisfies a certain mathematical relation

state & action  

$$Q(s,a) = r(s,a) + \max_{a} Q(T(s,a),a)$$
  
reward the best Q of the next state



ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA Interdepartmental Research centre Alma Mater Research Institute for UMAN-CENTERED ARTIFICIAL INTELLIGENC

#### In practice, enumeration is not viable

But we can circumvent it:

• Every optimal Q function satisfies a certain mathematical relation





ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA Interdepartmental Research centre Alma mater Research Institute for JMAN-Centered Artificial Intelligenc

#### In practice, enumeration is not viable

But we can circumvent it:

• Every optimal Q function satisfies a certain mathematical relation



- Now the trick is starting from random Qs...
- ...Sample actions and observe rewards...
- ...And adjust the Q function until it satisfies the relation



ALMA MATER STUDIORUM Università di Bologna Interdepartmental Research Centre Alma Mater Research Institute for Man-Centered Artificial Intelligenc

## Since we sample, we can now deal with uncertainty!

Think of:

- Unpredictable orders
- Machine failures
- Uncertain power supply...







ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA Interdepartmental Research centre Alma Mater Research Institute for Iuman-centered Artificial Intelligenci

### **Deep Reinforcement Learning**

#### Storing Q values for each state & action is also impractical

...But we can replace them with a Neural Network!

- Way more scalable
- Can handle sensorial information (e.g. images)

... Of course is not as trivial as this ;-)



INTERDEPARTMENTAL RESEARCH CENTRE Alma Mater Research Institute for Human-centered Artificial Intelligenc

#### We can the evaluate some of the advantages of RL

- Has no need to generate labels
- Can adapt to dynamic conditions





ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA Interdepartmental Research centre Alma Mater Research Institute for UMAN-CENTERED ARTIFICIAL INTELLIGENC

#### We can the evaluate some of the advantages of RL

- Has no need to generate labels
- Can adapt to dynamic conditions
- Can devote absurd computational resource to learning
- Can process huge amount of data
- Can live for (virutal) hundreds of years





ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA Interdepartmental research centre Alma mater research institute for uman-centered artificial intelligenc

#### We can the evaluate some of the advantages of RL

- Has no need to generate labels
- Can adapt to dynamic conditions
- Can devote absurd computational resource to learning
- Can process huge amount of data
- Can live for (virutal) hundreds of years

My favorite:

- If we can provide an environment to play with...
- ...We can think of building a RL system





ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA Interdepartmental research centre Alma Mater Research Institute for UMAN-CENTERED ARTIFICIAL INTELLIGENC

#### There are several caveats:

RL systems need to experiment

- I.e. they must be allowed to make mistakes
- Exploration vs exploitation trade-off

A good problem understanding is necessary

- Choosing the state has an impact on the policy
- Representing states may be practically impossible (hence: approximation)
- How to represent actions? Which DRL variant?



ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA Interdepartmental research centre Alma mater research institute for UMAN-Centered Artificial Intelligenc

#### ...And of course, there are plain drawbacks

- RL systems often require very large training times
- They don't require supervision... but what about survelliance?
- RL systems can be (very) opaque!
- They often don't play along well with other decision support techniques

Can some of these spur interesting research? Of course!



ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA Interdepartmental Research Centre Alma Mater Research Institute for UMAN-CENTERED ARTIFICIAL INTELLIGENC



#### ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA Interdepartmental Research centre Alma Mater Research Institute for Human-centered Artificial Intelligence

A Brief Introduction to Reinforcement Learning

# **Questions?**

michele.lombardi2@unibo.it

www.unibo.it