

# Week 2

Wednesday, February 23, 2022

4:22 PM

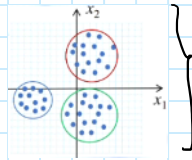
## Reinforcement Learning (Poi Refuereb)

- ↳ learning by prizes/punishment
- ↳ It is an Artificial Learning type.  
also known as  
Machine Learning

### 3 types:

#### 1. Non-supervised Learning

- The problems of this kind must be defined in a space of States



This is a simple state space,  
a space with 2 variables ( $x_1, x_2$ )

- ↳ Regularly, the problem is a series of points in this state space.  
i.e. one point is an instance of the problem.

given a new point,  
the machine is expected to classify it

- The objective of NS Learning is to group these points in the state space, forming CLUSTERS (Groups).

- ↳ This groups come from the assumption that there exists a distance that separates the instances of the problem

The machine finds out on its own the classes.

- The data is not previously classified or labeled, thus the machine gives the data a group by their similarities

#### 2. Supervised Learning

- Once again the problem is inside a state space, and we got the same space.
- The difference with these problems is that we provide the input to the machine for it to learn, and this input is already classified inside groups  
↳ instances that belong to a known class

=

The machine learns this by building a curve inside

the state space  
↓  
that curve separates the data

The machine thus learns what does an entity of class  $c_i$  look like.

- We expect that, if we got an unknown point (test), the machine is able to classify it correctly.

- Thus, by the data being classified, we say it is Labeled data.

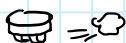
#### 3. Reinforcement Learning:

- In this case, there is not a series of instances given to the machine,  
↳ unlike the NS and S learning

instead, what we want is the machine to learn to perform a series of actions in the world

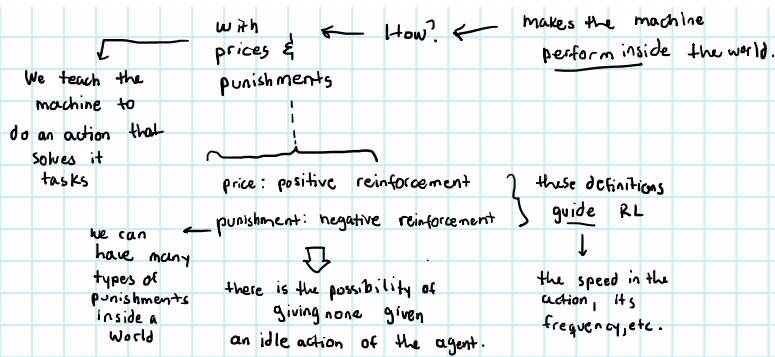
- ↳ in order to solve a PROBLEM.

i.e. a sweeper robot:




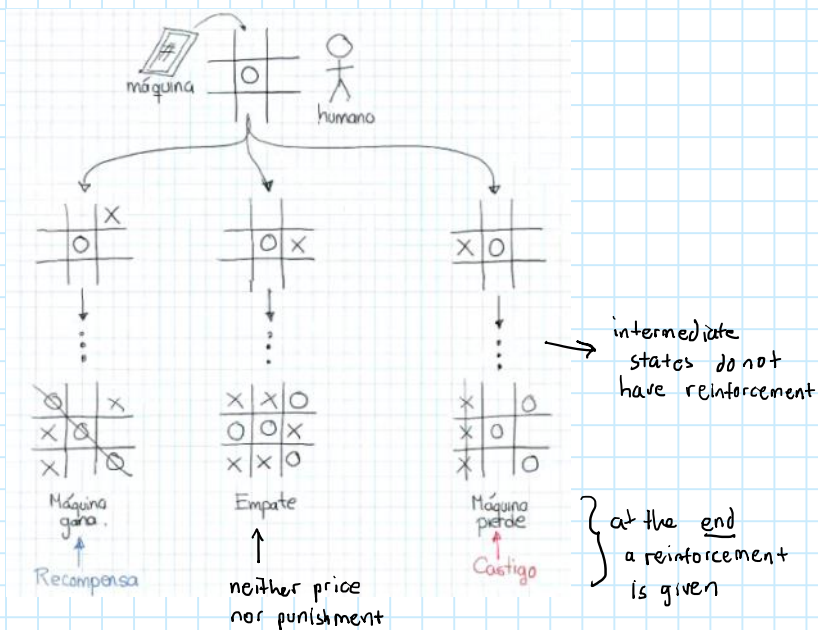
} the NS and S Learning recognise things in the world

↓  
Reinforcement Learning



- o In order to do RL, we need to define:
  - a World
  - a series of actions, available to the agent
  - the Reinforcement (price/punishment)

- o Another example is: the game 



- no feedback at any time.
- NS: not classified
- a) data: NON-LABELED, also known as NON STRUCTURED data.
  - b) Function: Grouping data (clustering) as defined entities  
Reduce dimensionality, since high dimensionalities slow down ML algorithms
  - c) Objective: find a structure in the data, give it a label

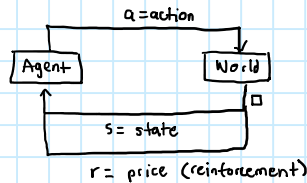
- feedback: we tell the class of data
- S:
- a) data: LABELED, also known as STRUCTURED data
  - b) Function: classification or regressions
  - c) Objective: prediction (for new data), since the class already exists.

- feedback: the reinforcement is the feedback
- may be only at the end, or at certain times
- RL:
- a) data: experiences from the world, that is, in which STATE it is and what REINFORCEMENT it gets.
    - What stimulates the program? → agent
      1. State: what the program can perceive
      2. Reinforcement: price/punishment, + / -
  - b) Function: associate actions to states
    - What can the agent do?  $a_0, a_1, a_2, \dots, a_n$   
actions that tell the machine what to do at every moment (politic)  
→ or state
    - \* State: the configuration of the world at an instant of time
  - c) Objective: Find the best actions: find the set of action where the

obtain the maximum positive reinforcement it can.  $\leftarrow$  maximize the price of benefit and minimize the punishment  $\leftarrow$  machine performs the best it can perform.  
 $\rightarrow$  programs must look for the most price at all times so that  
 this is the "instinct" inside the program

## Reinforcement Learning:

- This type of learning has a structure, which is the abstract structure of the problem itself. This is the structure:

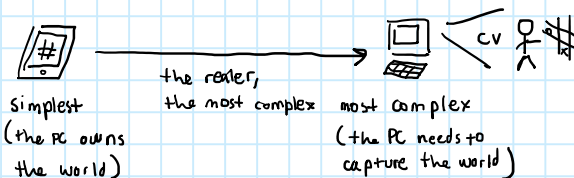


Components:

1. Agent: the program that will solve the task.
2. World/Environment

- the agent perceives both at all times
- There is an interaction between the Agent and World:
    - the agent can perceive the configuration of the world, the agent perceives the state of the world i.e. any possible board config is a world state, called  $S$ .  
 $\rightarrow$  the state ( $S$ ) is what the agent perceives of the World.
    - The other thing the agent perceives is the reinforcement ( $r$ ), which can be positive/negative. This also comes from the World.
    - Once these two are perceived, the agent performs an action ( $a$ ) to the world, depending on the state it is in.
  - = the agent and the world are in a cycle: the agent perceives and then acts
- $\rightarrow$  The purpose of RL is teaching the best action to perform at any state  $\nearrow$  maximum benefit

- The complexity of these problems can be very different:



$\rightarrow$  How can we represent the Reward (reinforcement), the State and the action?

$\rightarrow$  The representation must be something that is in the machine: numbers or strings?

- The State of the machine: it is the config of the world.

$\rightarrow$  If we use the example of # game, and we have the following state:

PC = "O"  
 Me = "X"  
 Empty = " "

we can represent the states with a mathematical object MATRIX

the board is the World

Every problem's states must be represented by any data structure we can think of

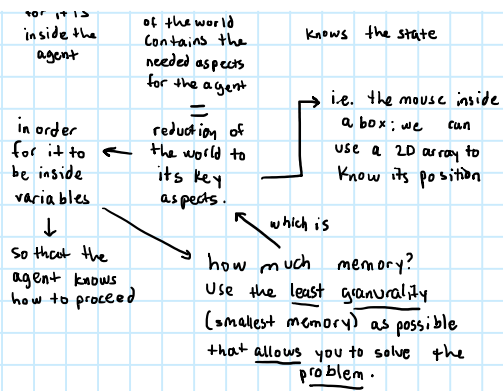
All representations of states will be variables  
 $\rightarrow$  inside RAM memory (variables)  
 $\rightarrow$  addresses of memory

$\rightarrow$  The variables can be:

- Strings "abc"
- numbers  $\rightarrow$  integers, floats
- booleanos V, F
- arrays: 2D, 3D, ...

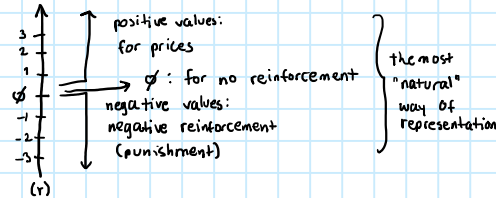
use this for representations

$\rightarrow$  the arrays are the most used for representing the World State  
 $\rightarrow$  the set of variables is an internal representation of the world (state)  
 "internal"  $\leftarrow$  the internal representation  $\leftarrow$  it is necessary, since that's how the agent



→ How to represent the reinforcement? (reward)

↳ the easiest is to use a numeric scale:



i.e. +2 if the PC wins  
-1 if no one wins  
-3 if the PC loses

the agent moves so that it gets the max benefit (the max reward)

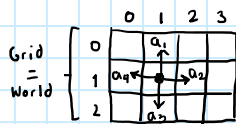
↳ avoid losing states  
↳ pursuing the winning states

"the agent's direction"

→ How do we represent the actions?

↳ a set of actions the agent has and that it can perform over the world, changing the world.

In an example where an agent moves inside a grid:



the • (agent) has  $a_1, a_2, a_3$  and  $a_4$  which involve a one-cell move, for example.

↳ the 'available' actions change every time the agent changes the state.