

# Artificial Intelligence

Introduction

Marc Toussaint  
University of Stuttgart  
Winter 2019/20

# What is AI?

- AI is a research field

# What is AI?

- AI is a research field
- AI research is about systems that take decisions
  - AI formalizes decision processes: interactive (decisions change world state), or passive
  - AI distinguishes between agent(s) and the external world: decision variables

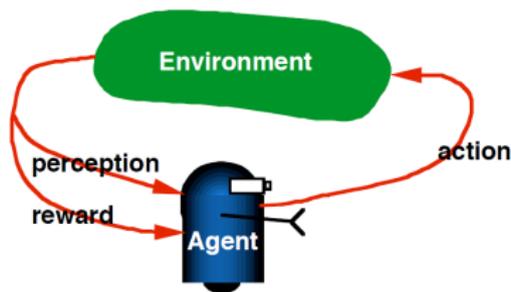
# What is AI?

- AI is a research field
- AI research is about systems that take decisions
  - AI formalizes decision processes: interactive (decisions change world state), or passive
  - AI distinguishes between agent(s) and the external world: decision variables
- ... systems that take optimal/desirable decisions
  - AI formalizes decision objectives
  - AI aims for systems to exhibit functionally *desirable behavior*

# What is AI?

- AI is a research field
- AI research is about systems that take decisions
  - AI formalizes decision processes: interactive (decisions change world state), or passive
  - AI distinguishes between agent(s) and the external world: decision variables
- ... systems that take optimal/desirable decisions
  - AI formalizes decision objectives
  - AI aims for systems to exhibit functionally *desirable behavior*
- ... systems that take optimal decisions on the basis of all available information
  - AI is about inference
  - AI is about learning

# 1) ... systems that take decisions

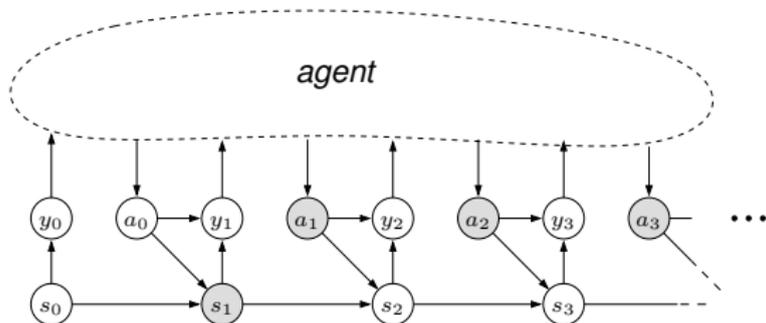


- We assume the agent is in *interaction* with a domain.
  - The world is in a state  $s_t \in \mathcal{S}$  (see below on what that means)
  - The agent senses observations  $y_t \in \mathcal{O}$
  - The agent decides on an action  $a_t \in \mathcal{A}$
  - The world transitions to a new state  $s_{t+1}$
- The *observation*  $y_t$  describes all information received by the agent (sensors, also rewards, feedback, etc) if not explicitly stated otherwise

(The technical term for this is a POMDP)

# State

- The notion of *state* is often used imprecisely
- At any time  $t$ , we assume the world is in a state  $s_t \in \mathcal{S}$
- $s_t$  is a *state description* of a domain such that future observations  $y_{t+}, t^+ > t$  are conditionally independent of all history observations  $y_{t-}, t^- < t$  given  $s_t$  and future actions  $a_{t:t+}$ :



- Notes:
  - Intuitively,  $s_t$  describes everything about the world that is “relevant”
  - Worlds do not have additional latent (hidden) variables to the state  $s_t$

## Examples

- What is a sufficient definition of *state* of a computer that you interact with?
- What is a sufficient definition of *state* for a thermostat scenario?  
(First, assume the 'room' is an isolated chamber.)
- What is a sufficient definition of *state* in an autonomous car case?

## Examples

- What is a sufficient definition of *state* of a computer that you interact with?
  - What is a sufficient definition of *state* for a thermostat scenario?  
(First, assume the 'room' is an isolated chamber.)
  - What is a sufficient definition of *state* in an autonomous car case?
- in real worlds, the exact *state* is practically not representable
- all models of domains will have to make approximating assumptions  
(e.g., about independencies)

# How can agents be formally described?

...or, what formal classes of agents do exist?

- Basic alternative agent models:

- The agent maps  $y_t \mapsto a_t$   
(**stimulus-response** mapping.. non-optimal)
- The agent stores all previous observations and maps

$$f : y_{0:t}, a_{0:t-1} \mapsto a_t$$

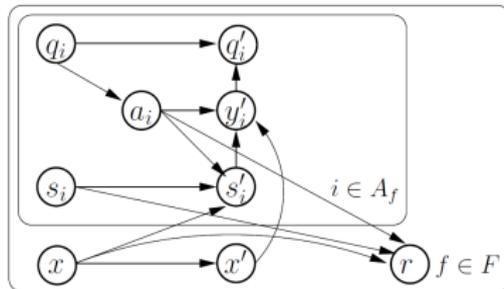
$f$  is called **agent function**. This is the most general model, including the others as special cases.

- The agent stores only the recent history and maps  
 $y_{t-k:t}, a_{t-k:t-1} \mapsto a_t$  (crude, but may be a good heuristic)
- The agent is some machine with its own **internal state**  $n_t$ , e.g., a computer, a finite state machine, a brain... The agent maps  $(n_{t-1}, y_t) \mapsto n_t$  (internal state update) and  $n_t \mapsto a_t$
- The agent maintains a full probability distribution (**belief**)  $b_t(s_t)$  over the state, maps  $(b_{t-1}, y_t) \mapsto b_t$  (Bayesian belief update), and  $b_t \mapsto a_t$



## Multi-agent domain models

(The technical term for this is a Decentralized POMDPs)



(from Kumar et al., IJCAI 2011)

- This is a special type (simplification) of a general DEC-POMDP
- Generally, this level of description is very general, but NEXP-hard  
Approximate methods can yield very good results, though

## 2) ... take optimal/desirable decisions

- A cognitive scientist or psychologist: “Why are you AI people always so obsessed with optimization? Humans are not optimal!”

## 2) ... take optimal/desirable decisions

- A cognitive scientist or psychologist: “Why are you AI people always so obsessed with optimization? Humans are not optimal!”
- That’s a total misunderstanding of what “being optimal” means.
- Optimization principles are a means to describe systems:
  - Feynman’s “unworldliness measure” objective function
  - Everything can be cast optimal – under *some* objective
  - Optimality principles are just a scientific means of formally describing systems and their behaviors (esp. in physics, economy, ... and AI)
  - Toussaint, Ritter & Brock: *The Optimization Route to Robotics – and Alternatives*. Künstliche Intelligenz, 2015

## 2) ... take optimal/desirable decisions

- A cognitive scientist or psychologist: “Why are you AI people always so obsessed with optimization? Humans are not optimal!”
- That’s a total misunderstanding of what “being optimal” means.
- Optimization principles are a means to describe systems:
  - Feynman’s “unworldliness measure” objective function
  - Everything can be cast optimal – under *some* objective
  - Optimality principles are just a scientific means of formally describing systems and their behaviors (esp. in physics, economy, ... and AI)
  - Toussaint, Ritter & Brock: *The Optimization Route to Robotics – and Alternatives*. Künstliche Intelligenz, 2015
- Vocabulary:
  - Utility, typically “expected future accumulated discounted reward”
  - A “rational” agent: taking decisions to maximize utility

# What are interesting objectives?

- Learn to control all degrees of freedom of the environment that are controllable
  - DOFs are mechanical/kinematics DOFs, objects, light/temperature, mood of humans
  - This objective is generic: no preferences, not limits
  - Implies to actively go exploring and finding controllable DOFs
  - Acting to Learning (instead of 'Learning to Act' for a fixed task)
  - Related notions in other fields: (*Bayesian*) *Experimental Design*, *Active Learning*, curiosity, intrinsic motivation
- At time  $T$ , the system will be given a random task (e.g., random goal configuration of DOFs); the objective then is to reach it as quickly as possible

## More on objectives

- The value alignment dilemma
- What are objectives that describe things like “creativity”, “empathy”, “morale”, etc?
- Coming up with objective functions that imply desired behavior is a core part of AI research

**3) ... on the basis of all available information**

# AI = ML?

- In large parts, ML is: (let's call this  $ML^0$ )

*Fitting a function  $f : x \mapsto y$  to given data  $D = \{(x_i, y_i)\}_{i=1}^n$*

# AI = ML?

- In large parts, ML is: (let's call this  $ML^0$ )  
*Fitting a function  $f : x \mapsto y$  to given data  $D = \{(x_i, y_i)\}_{i=1}^n$*
- And what does that have to do with AI?
- Literally, not much:
  - The decision made by a  $ML^0$  method is only a single decision: Decide on the function  $f \in \mathcal{H}$  in the hypothesis space  $\mathcal{H}$
  - This “single-decision process” is not interactive;  $ML^0$  does not formalize/model/consider how the choice of  $f$  changes the world
  - The objective  $\mathcal{L}(f, D)$  in  $ML^0$  depends only on the given static data  $D$  and the decision  $f$ , not how  $f$  might change the world
  - “Learning” in  $ML^0$  is not an interactive process, but some method to pick  $f$  on the basis of the static  $D$ . The typically iterative optimization process is not the decision process that  $ML^0$  focusses on.

**But**, function approximation can be used to help solving AI (interactive decision process) problems:

- Building a trained/fixed  $f$  into an interacting system based on human expert knowledge

E.g., an engineer trains a NN  $f$  to recognize street signs based on a large fixed data set  $D$ . S/he builds  $f$  into a car to drive autonomously. That car certainly solves an AI problem;  $f$  continuously makes decisions that change the state of the world. But  $f$  was never optimized literally for the task of interacting with the world; it was only optimized to minimize a loss on the static data  $D$ . It is not a priori clear that minimizing the loss of  $f$  on  $D$  is related to maximizing rewards in car driving. That fact is the expertise of the engineer; it is not some AI algorithm that discovered that fact. At no place there was an AI method that “learns to drive a car”. There was only an optimization method to minimize the loss of  $f$  on  $D$ .

- Using ML in interactive decision processes

E.g., to approximate a  $Q$ -function, state evaluation function, reward function, the system dynamics, etc. Then, other AI methods can use these approximate models to actually take decisions in the interactive context.

## Summary – AI is about:

- Systems that interact with the environment
  - We distinguish between 'system' and 'environment' (cf. embodiment)
  - We just introduced basic models of interaction
  - A core part of AI research is to develop formal models for interaction
- Systems that aim to manipulate their environment towards 'desired' states (optimality)
  - Optimality principles are a standard way to describe desired behaviors
  - We sketched some interesting objectives
  - Coming up with objective functions that imply desired behavior is a core part of AI research
- Systems that exploit all available data, i.e., learn and infer
  - Machine Learning plays a major role, but it is usually only a component in an embedding AI framework

# Organisation

# Vorraussetzungen für die KI Vorlesung

- Mathematik für Informatiker und Softwaretechniker
- außerdem hilfreich:
  - Algorithmen und Datenstrukturen
  - Theoretische Informatik

# Vorlesungsmaterial

- Webseite zur Vorlesung:  
<https://ipvs.informatik.uni-stuttgart.de/mlr/marc/teaching/>  
die Folien und Übungsaufgaben werden dort online gestellt
- Alle Materialien des letzten Jahres sind online
- Hauptliteratur:  
*Stuart Russell & Peter Norvig: Artificial Intelligence A Modern Approach*
  - Some slides are adopted from Stuart

# Prüfung

- Schriftliche Prüfung, 90 Minuten
- Termin zentral organisiert
- keine Hilfsmittel erlaubt
- Anmeldung: Im LSF / beim Prüfungsamt
- Prüfungszulassung:
  - 50% der Punkte der Programmieraufgaben
  - UND 50% der Votieraufgaben

# Übungen

- 8 Übungsgruppen (4 Tutoren)
- 2 Arten von Aufgaben: Coding- und Votier-Übungen
- Coding-Aufgaben: Teams von bis zu 3 Studenten geben die Coding-Aufgaben zusammen ab
- Votier-Aufgaben:
  - Zu Beginn der Übung eintragen, welche Aufgaben bearbeitet wurden/präsentiert werden können
  - Zufällige Auswahl
- Schein-Kriterium:
  - 50% der Punkte der Programmieraufgaben
  - UND 50% der Votieraufgaben

- **Registrierung**

`https://ipvs.informatik.uni-stuttgart.de/mlr/teaching/course-registration/`

Passwort: ki2019