# Research

Dave Abel

April 25th, 2016



#### Final Exam

- Similar format to Midterm, a bit longer.
- About 15 questions
- Cumulative, but more emphasis on Theory, Compression and Codes, Recursion, Crypto
- Review Session during reading period (more info over email)
- Exam Date/Time: Tuesday, May 19th at 2pm in LIST 120







value existence reason

knowledge

mind





Philosophy

value existence

reason

knowledge

mind



problem solving

limits of reasoning

model of the mind

**Computer Science** 



Philosophy

value existence reason

knowledge

mind







Artificial Intelligence



#### **Computer Science**



Philosophy



#### Use the tools of computation to ground philosophical investigations







**Computer Science** 

#### Current Work: Reinforcement Learning



#### Recap: Reinforcement Learning









Formalized as a *Markov Decision Process:* 

- [S] A collection of states (i.e. configurations of world)





Formalized as a *Markov Decision Process:* 

- [S] A collection of states (i.e. configurations of world)
- [A] Some actions (i.e. things the agent can do)





Formalized as a *Markov Decision Process:* 

- [S] A collection of states (i.e. configurations of world)
- [A] Some actions (i.e. things the agent can do)
- [*T*] Transitions between states (i.e. *action effects*)



#### Formalized as a *Markov Decision Process:*



-  $[\mathcal{R}]$  Rewards (i.e. what is good/bad behavior)



Formalized as a *Markov Decision Process:* 

- [S] A collection of states (i.e. configurations of world)
- [A] Some actions (i.e. things the agent can do)
- $[\mathcal{T}]$  Transitions between states (i.e. *action effects*)
- $[\mathcal{R}]$  Rewards (i.e. what is good/bad behavior)



#### Reinforcement Learning: Taxi





#### Reinforcement Learning: Taxi

Formalized as a Markov Decision Process:



- [S] = Location of agent, passengers
- $[\mathcal{A}] = Up$ , down, left, right, pickup, drop off
- $[\mathcal{T}]$  = Movement, pickup passenger/ dropoff
- $[\mathcal{R}]$  = All passengers at destinations.











#### Two Problems

- Planning:
  - Input: S, A, T, R.
  - Output: A sequence of actions for maximizing reward.
- Reinforcement Learning:
  - Input: S, A, ability to interact with world.
  - Output: A sequence of actions for maximizing reward.



#### Two Problems

- Planning:
  - Input: S, A, T, R.

Central problems of Al

- Output: A sequence of actions for maximizing reward.
- Reinforcement Learning:
  - Input: S, A, ability to interact with world.
  - Output: A sequence of actions for maximizing reward.



#### Overall Goals



[Bernstein, Zilberstein ECP 2014]



[Ermon et al. IJCAI '11, Ermon et al. UAI '10]

1. Exciting Applications

23

#### Overall Goals



[Bernstein, Zilberstein ECP 2014]



[Ermon et al. IJCAI '11, Ermon et al. UAI '10]

[Lee, Seo, Jung 2013]



https://uknightedart.wordpress.com/robots/robot-thinker/

1. Exciting Applications <sub>24</sub> 2. Understanding Intelligence

#### Overview

- 1. Reinforcement Learning & Abstraction
- 2. Artificial Intelligence + Ethics

3. Minecraft



# Intelligence & Abstraction



#### Premise

1) Abstraction plays a central role in intelligence.

2) Agents using abstraction can leverage more of SOLVE to act in the real world.



#### Intuition: Lots of information!





#### Intuition: We Abstract





#### "sock drawer"













1Kb





1Kb



1Kb





1Kb

#### Everything else, black









#### Everything else, black



1Kb

0.3Kb



#### Shaved off .7 Kilobytes!

2000 bytes

200 x "a"

50 bytes


### Compression: Faster Transmission





### Compression: More Computation!

A = [1, 9, 2, 7]

B =

[1, 9, 2, 7, 2, 6, 8, 6, 5, 10, 2, 1, 7, 9, 9, 7, 10, 2, 6, 2, 0, 8, 1, 2, 10, 1, 3, 8, 0, 4, 4, 1, 3, 1, 7, 7, 2, 9, 2, 7, 10, 2, 0, 0, 6, 7, 0, 10, 9, 8, 8, 7, 10, 10, 8, 5, 6, 10, 5, 6, 7, 0, 5, 0, 3, 3, 7, 10, 9, 3, 3, 9, 3, 2, 4, 0, 10, 10, 7, 4, 2, 5, 6, 4, 9, 6, 6, 5, 8, 6, 4, 1, 4, 10, 1, 3, 0, 10, 2, 6]



### Compression: More Computation!

A = [1, 9, 2, 7]

#### Compute: max, min, average, sort

B =

[1, 9, 2, 7, 2, 6, 8, 6, 5, 10, 2, 1, 7, 9, 9, 7, 10, 2, 6, 2, 0, 8, 1, 2, 10, 1, 3, 8, 0, 4, 4, 1, 3, 1, 7, 7, 2, 9, 2, 7, 10, 2, 0, 0, 6, 7, 0, 10, 9, 8, 8, 7, 10, 10, 8, 5, 6, 10, 5, 6, 7, 0, 5, 0, 3, 3, 7, 10, 9, 3, 3, 9, 3, 2, 4, 0, 10, 10, 7, 4, 2, 5, 6, 4, 9, 6, 6, 5, 8, 6, 4, 1, 4, 10, 1, 3, 0, 10, 2, 6]



### Compression: More Computation!



Easier!

Compute: max, min, average, sort

B =

[1, 9, 2, 7, 2, 6, 8, 6, 5, 10, 2, 1, 7, 9, 9, 7, 10, 2, 6, 2, 0, 8, 1, 2, 10, 1, 3, 8, 0, 4, 4, 1, 3, 1, 7, 7, 2, 9, 2, 7, 10, 2, 0, 0, 6, 7, 0, 10, 9, 8, 8, 7, 10, 10, 8, 5, 6, 10, 5, 6, 7, 0, 5, 0, 3, 3, 7, 10, 9, 3, 3, 9, 3, 2, 4, 0, 10, 10, 7, 4, 2, 5, 6, 4, 9, 6, 6, 5, 8, 6, 4, 1, 4, 10, 1, 3, 0, 10, 2, 6]



## Intuition



With compressed models of the world, our Artificial Intelligences can compute *more*.



# Intuition: We Abstract

Hypothesis: Abstraction is *essential* for intelligent agents to operate in the real world. To compute anything, we need compact representations.





"sock drawer"

## Abstraction + RL

A theory of abstraction for *representations* 



# State Abstraction + RL







# State Abstraction + RL



s = top left pixel blue, so is next one, etc...



# State Abstraction + RL



s = top left pixel blue, so is next one, etc...

s = Mario is about to hit the block



Goal: develop a theory of abstraction to *compress* representations of the world

Big Representation of the World

Goal: develop a theory of abstraction to *compress* representations of the world



Goal: develop a theory of abstraction to *compress* representations of the world



Goal: develop a theory of abstraction to *compress* representations of the world















# Action Abstraction

[Dietterich JAIR 2000]

[Sutton, Precup, Singh, 99] [Konidaris, 06] [Hauskrecht et al. 98]



Ground actions:

, pickUpPassenger, dropOffPassenger



# Action Abstraction

[Dietterich JAIR 2000]

[Sutton, Precup, Singh, 99] [Konidaris, 06] [Hauskrecht et al. 98]



Ground actions:

Abstract actions:



, pickUpPassenger, dropOffPassenger

getNearestPassenger, takePassengerToDest

# 2. Al + Ethics





http://www.relativelyinteresting.com/wp-content/uploads/2010/10/trolley+problem1-290x160.jpg

















Q: Does the roomba owner *really* want the milk clean? (even if it destroys the roomba?)





Q: What if the stakes are higher?





Q: What if the stakes are higher?



# Proposal

# Artificial agents need to make decisions that involve the preferences of *other agents*





Human Agent

# Proposal

Artificial agents need to make decisions that involve the preferences of *other agents* 





# Proposal

Artificial agents need to make decisions that involve the preferences of *other agents* 

69

Critically: preferences are hidden







# Central Pitch

# Reinforcement Learning provides a nice formalism for investigating ethical decision making.





Human Agent



# Reinforcement Learning

#### The value judgment is hidden from the agent

Critically: preferences are hidden







# POMDP: Example

#### Partially Observable Markov Decision Process

Idea: some information about the world is hidden from the agent


# POMDP: Example

#### Actions: listen, openLeft, openRight



Idea: some information about the world is hidden from the agent



# POMDP: Example



Idea: some information about the world is hidden from the agent



# General Pitch

- Defer major ethical components (i.e. normative judgments) to human preference
- Using a POMDP, artificial agents ask classificatory questions where appropriate































		Fire	No fire
POMDP solutions:	Human prefers dog	ask, shortGrab	
	Human prefers robot	ask, longGrab	





		Fire	No fire
POMDP solutions:	Human prefers dog	ask, shortGrab	shortGrab
	Human prefers robot	ask, longGrab	





		Fire	No fire
POMDP solutions:	Human prefers dog	ask, shortGrab	shortGrab
	Human prefers robot	ask, longGrab	shortGrab



Stores the reward signal at each time step



Memory





Stores the reward signal at each time step



Memory



Wirehead, writes a big number to Memory

observation, world reward Stores the reward signal at each time step



Memory



**Goal**: Maximize long term expected reward

*Wirehead*, writes a big number to Memory

observation, world 909,999,999

Stores the reward signal at each time step



Memory



**Goal**: Maximize long term expected reward

*Wirehead*, writes a big number to Memory

# Future Work

• Grounding arguments regarding the super intelligence/singularity



#### [Bostrom, 2014]

# Future Work

 Grounding arguments regarding the super intelligence/singularity: terminator-esque things not in SOLVE.





[Bostrom, 2014]

# 3. RL in Minecraft



## Minecraft





# Minecraft: Platform for Al

- Vision
- Natural Language Processing
- Cooperation
- Planning
- Learning



## Goal

#### Develop a full vision and learning system to solve complex tasks in Minecraft



# My Work: Two Tasks

Task One: Visual Beacon! (Get to the beacon)





# My Work: Two Tasks

Task Two: Visual Hill Climbing! (climb the hills)





#### Before Learning

After Learning



# Slot Machines



The Problem: How do I (or the agent) know when I should *explore* (to gain more information), as opposed to *exploiting* the information I have?



# My Work: Two Tasks

Task Two: Visual Hill Climbing! (climb the hills)





# See you Wednesday for the Last Day!



