

# CS 378: Autonomous Intelligent Robotics FRI-II

#### Instructor: Jivko Sinapov

http://www.cs.utexas.edu/~jsinapov/teaching/cs378\_fall2016/

# Reinforcement Learning (Part 2)



observation

#### Announcements

# Volunteers needed for robot study

Sign up sheet here:

https://docs.google.com/spreadsheets/d/1Gr 2GqlPt8kdTJwlZ3FerxU0J8oIoGt0pEbeR37iCHqY/ edit#gid=0

Further details will be made available on Canvas via an announcement

# FAI Talk this Friday

#### "Turning Assistive Machines into Assistive Robots" Brenna Argall Northwestern University

#### Friday, Sept. 9<sup>th</sup>, 11 am @ GDC 6.302

[ https://www.cs.utexas.edu/~ai-lab/fai/ ] or google "fai ut cs"

# **Robotics Seminar Series Talk**

"Learning from and about humans using an autonomous multi-robot mobile platform" Jivko Sinapov

#### **UT** Austin

Wed., Sept. 7<sup>th</sup>, 3 pm @ GDC 5.302

# **Robot Training**

- Sign up for a robot training session next week at: https://docs.google.com/spreadsheets/d/1kz6QMPa-xkdFQNyV 0Biif913JKf73GIUGx9bVSLo\_R8/edit?usp=sharing
- Link will be posted as Announcement on Canvas

# Preliminary Project "Presentations"

- Date: September 13th
- Form groups of 2-3 prior to the date
- Be prepared to talk about 2-3 project ideas for 5-10 minutes
- Email me your group info, i.e., who is in it

# **Project Ideas**

# Project Idea: Improve the robot's grasping ability

- Currently, the robot does not "remember" which grasps succeeded and which failed
- If the robot were to log the context of the grasp (e.g., the position of the gripper relative to the object's point cloud) and the outcome, it could incrementally learn a model to predict the outcome given the context
- The robot's current grasping software is described in http://wiki.ros.org/agile\_grasp

# Project Idea: Object Handover

- Currently, the arm can let go of an object or close its fingers upon sufficient contact using haptic feedback
- Can you make it so that it can move towards an object held by a human and grasp it based on visual and haptic feedback?

# Project Idea: Learning about objects from humans

- The robot is currently able to grasp an object from a table and navigate to an office
- Can we use the GUI to ask humans questions about objects and store this information in a database that can be used for learning recognition models?

# Project Idea: Large-Scale 3D object mapping

 Can we combine 3D Plane detection and Clustering to detect and map objects in the environment?





# Project Idea: Learning an object manipulation skill

• Example: Pressing a button

## **Project Idea: Enhance Virtour**

www.cs.utexas.edu/~larg/bwi\_virtour

# Reinforcement Learning (Part 2)



observation

# Markov Decision Process (MDP)



# Markov Decision Process (MDP)



The reward and state-transition observed at time *t* after picking action *a* in state *s* is independent of anything that happened before time *t* 



# Maze World



- Rewards: -1 per time-step
- Actions: N, E, S, W
- States: Agent's location

# Maze World



- Rewards: -1 per time-step
- Actions: N, E, S, W
- States: Agent's location

#### State Representation: Factored vs. Tabula Rasa

## Maze Example: Policy



Arrows represent policy  $\pi(s)$  for each state s

## Maze Example: Value Function



Numbers represent value  $v_{\pi}(s)$  of each state s

## Maze Example: Policy



Arrows represent policy  $\pi(s)$  for each state s

# Maze Example: Model



- Agent may have an internal model of the environment
- Dynamics: how actions change the state
- Rewards: how much reward from each state
- The model may be imperfect
- Grid layout represents transition model  $\mathcal{P}^{a}_{ss'}$
- Numbers represent immediate reward R<sup>a</sup><sub>s</sub> from each state s (same for all a)

### Sparse vs. Dense Reward

# Notation and Problem Formulation

• Overview of notation in TEXPLORE paper

# Notation

Set of States: SSet of Actions: ATransition Function:

 $\mathcal{P}: \mathcal{S} \times \mathcal{A} \mapsto \Pi(\mathcal{S})$ Reward Function:

 $\mathcal{R}:\mathcal{S} imes\mathcal{A}\mapsto\mathbb{R}$ 

### **Action-Value Function**

$$Q^*(s,a) = \mathcal{R}(s,a) + \gamma \sum_{s'} \mathcal{P}(s'|s,a) \max_{a'} Q^*(s',a')$$

### **Action-Value Function**



The reward received after taking action *a* in state *s* 

#### **Action-Value Function**

$$Q^*(s,a) = \mathcal{R}(s,a) + \gamma \sum_{s'} \mathcal{P}(s'|s,a) \max_{a'} Q^*(s',a')$$

Common algorithms to learn the action-value function include Q-Learning and SARSA

The policy consists of always taking the action that maximize the action-value function

# Q-Learning Grid World Example

https://www-s.acm.illinois.edu/sigart/docs/QLearning.pdf

# RL in a nutshell

Algorithm 1 Sequential Model-Based Architecture

1: Input: S, A	$\triangleright$ S: state space, A: action space
2: Initialize $M$ to empty model	
3: Initialize policy $\pi$ randomly	
4: Initialize s to a starting state in the MDP	
5: loop	
6: Choose $a \leftarrow \pi(s)$	
7: Take action $a$ , observe $r$ , $s'$	
8: $M \Rightarrow \text{UPDATE-MODEL}(\langle s, a, s', r \rangle)$	▷ Update model M with new experience
9: $\pi \leftarrow \text{PLAN-POLICY}(M)$	Exact planning on updated model
10: $s \leftarrow s'$	
11: end loop	

# RL in a nutshell



# Q-Learning

Guest Slides

### Pac Man Example

# Linear Function Approximator of Q\*

$$Q^{*}(s, a) = \mathcal{R}(s, a) + \gamma \sum_{s'} \mathcal{P}(s'|s, a) \max_{a'} Q^{*}(s', a')$$

Φ(s,a) = **x** where **x** is an *n*-dimensional feature vector

$$Q^{*}(\Phi(s,a)) = W_{1}^{*} X_{1} + W_{2}^{*} X_{2} + ... + W_{n}^{*} X_{n}$$

#### How does Pac-Man "see" the world?



#### How does Pac-Man "see" the world?



#### How does Pac-Man "see" the world?



# Linear Function Approximator of Q\*

$$Q^{*}(s, a) = \mathcal{R}(s, a) + \gamma \sum_{s'} \mathcal{P}(s'|s, a) \max_{a'} Q^{*}(s', a')$$

Φ(s,a) = **x** where **x** is an *n*-dimensional feature vector

$$Q^{*}(\Phi(s,a)) =$$

$$W_1 * X_1 + W_2 * X_2 + ... + W_n * X_n$$

The task now is to find the optimal weight vector w

# Can RL learn directly from images?

- Yes it can:
- http://karpathy.github.io/2016/05/31/rl/



# Video on Updating a NN's Weights

Neural Networks Demystified [Part 3: Gradient Descent]

https://www.youtube.com/watch?v=5u0jaA3qAGk

# Video of TAMER

http://labcast.media.mit.edu/?p=300

# Using RL: Essential Steps

1) Specify the state space or the state-action space

- Are the states and/or actions discrete or continous?
- 2) Specify the reward function
- If you have control over this, dense reward is better than sparse reward

3) Specify the environment (e.g., a simulator or perhaps the real world)

4) Pick your favorite RL algorithm that can handle the state and action representation

#### Resources

- BURLAP: Java RL Library: http://burlap.cs.brown.edu/
- Reinforcement Learning: An Introduction http://people.inf.elte.hu/lorincz/Files/RL\_2 006/SuttonBook.pdf