



10-601 Introduction to Machine Learning

Machine Learning Department
School of Computer Science
Carnegie Mellon University

Q-Learning

Matt Gormley
Lecture 27
April 16, 2018

Reminders

- **Homework 7: HMMs**
 - Out: Wed, Apr 04
 - Due: Mon, Apr 16 at 11:59pm
- **Homework 8: Reinforcement Learning**
 - Out: Mon, Apr 16
 - Due: Fri, Apr 27 at 11:59pm
 - Recitation: Mon, Apr 23 (instead of lecture)

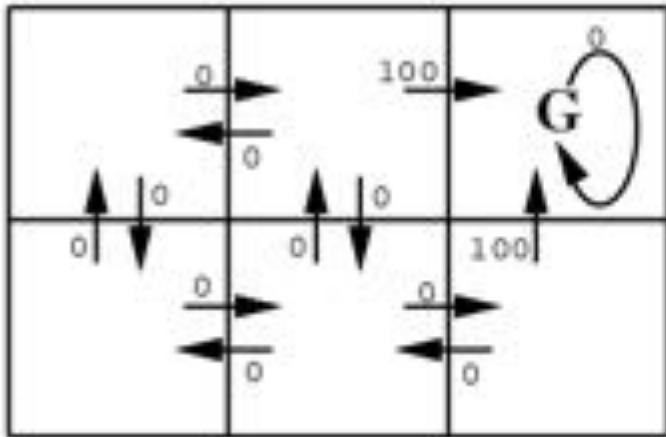
Q-LEARNING

Q-Learning

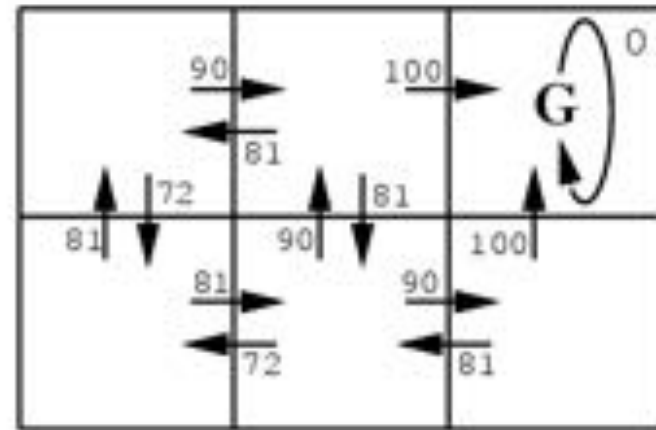
Whiteboard

- Motivation: What if we have zero knowledge of the environment?
- Q-Function: Expected Discounted Reward

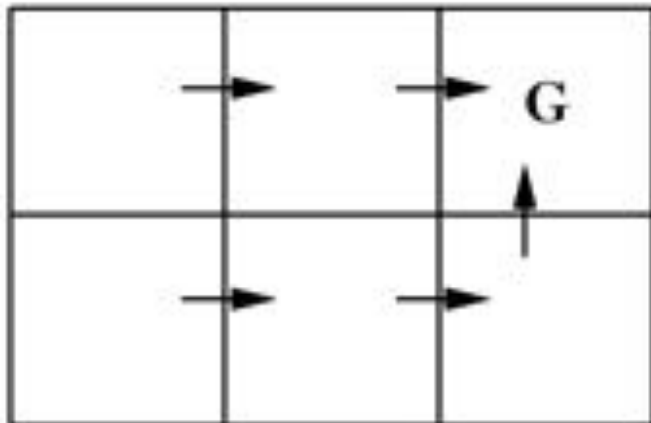
Example: Robot Localization



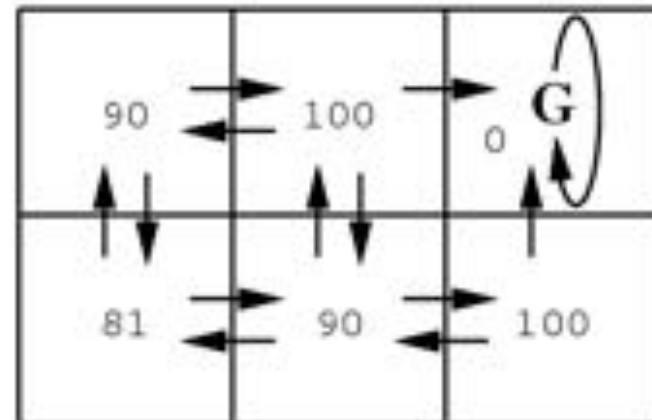
$r(s, a)$ (immediate reward) values



$Q(s, a)$ values



One optimal policy



$V^*(s)$ values

Q-Learning

Whiteboard

- Q-Learning Algorithm
 - Case 1: Deterministic Environment
 - Case 2: Nondeterministic Environment
- Convergence Properties
- Exploration Insensitivity
- Ex: Re-ordering Experiences
- ϵ -greedy Strategy

DEEP RL EXAMPLES

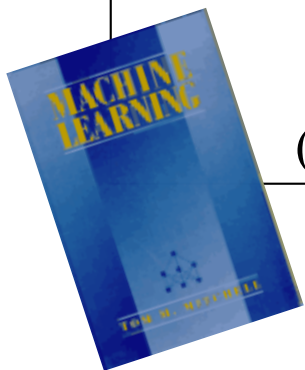
TD Gammon → Alpha Go

Learning to beat the masters at board games

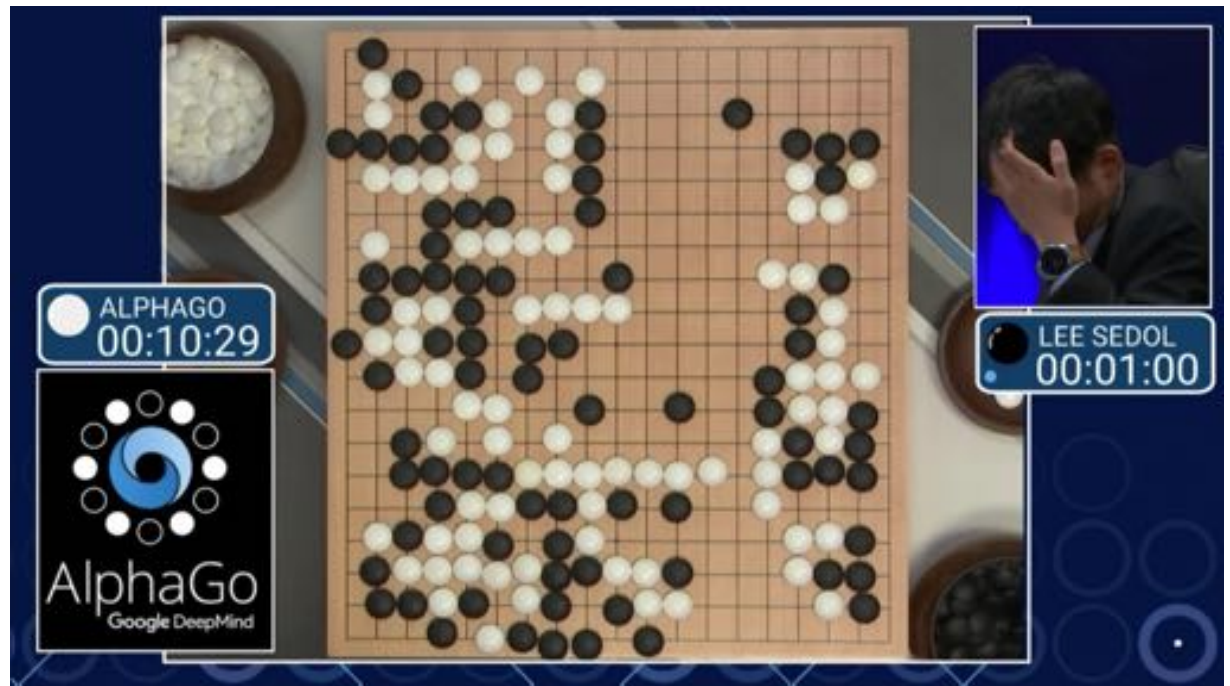
THEN

“...the world’s top computer program for backgammon, TD-GAMMON (Tesauro, 1992, 1995), learned its strategy by playing over one million practice games against itself...”

(Mitchell, 1997)

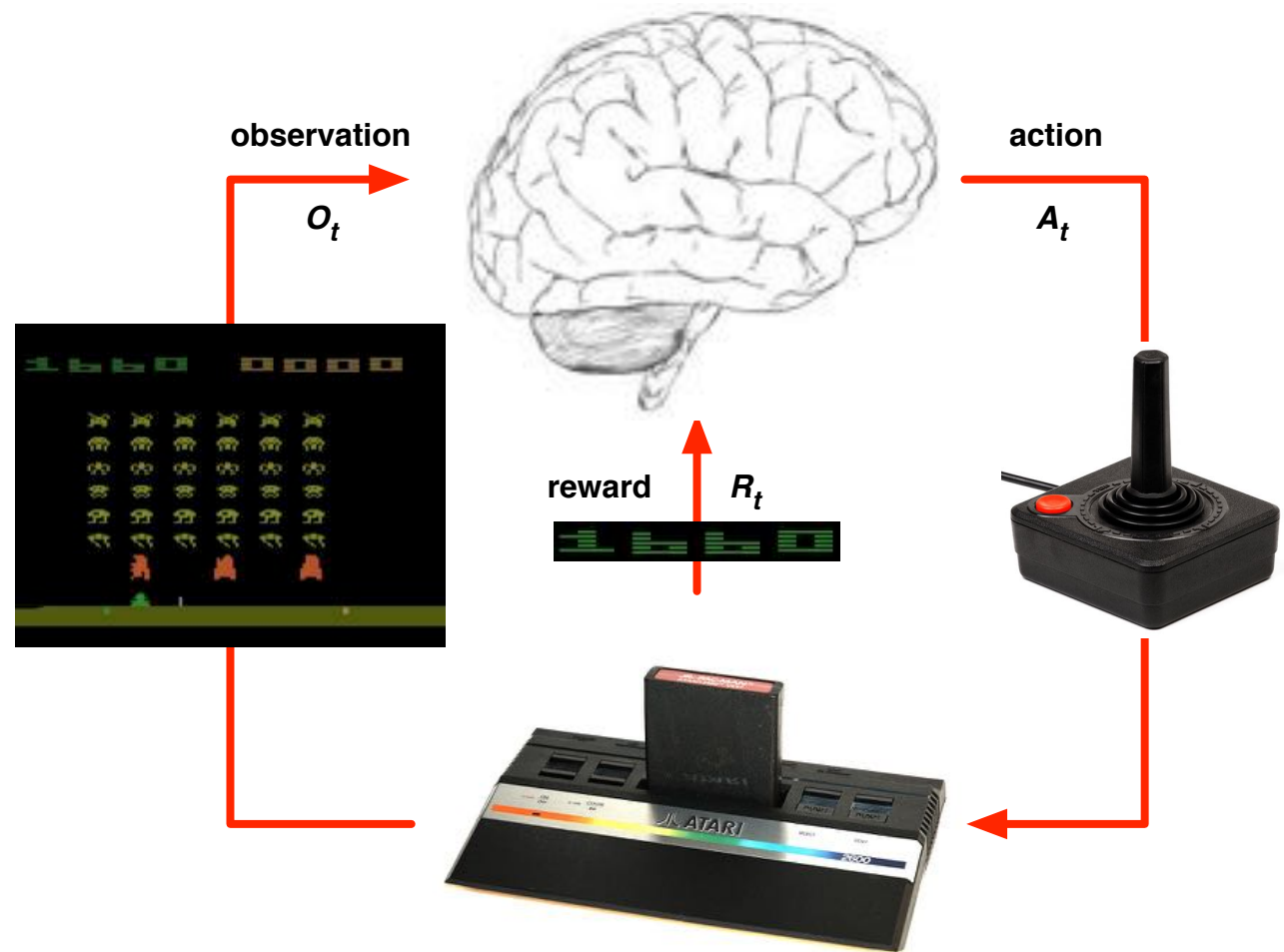


NOW



Playing Atari with Deep RL

- Setup: RL system observes the pixels on the screen
- It receives rewards as the game score
- Actions decide how to move the joystick / buttons



Playing Atari with Deep RL



Figure 1: Screen shots from five Atari 2600 Games: (*Left-to-right*) Pong, Breakout, Space Invaders, Seaquest, Beam Rider

Videos:

- Atari Breakout:

<https://www.youtube.com/watch?v=V1eYniJoRnk>

- Space Invaders:

<https://www.youtube.com/watch?v=ePvoFs9cGgU>

Playing Atari with Deep RL



Figure 1: Screen shots from five Atari 2600 Games: (Left-to-right) Pong, Breakout, Space Invaders, Seaquest, Beam Rider

	B. Rider	Breakout	Enduro	Pong	Q*bert	Seaquest	S. Invaders
Random	354	1.2	0	−20.4	157	110	179
Sarsa [3]	996	5.2	129	−19	614	665	271
Contingency [4]	1743	6	159	−17	960	723	268
DQN	4092	168	470	20	1952	1705	581
Human	7456	31	368	−3	18900	28010	3690
HNeat Best [8]	3616	52	106	19	1800	920	1720
HNeat Pixel [8]	1332	4	91	−16	1325	800	1145
DQN Best	5184	225	661	21	4500	1740	1075

Table 1: The upper table compares average total reward for various learning methods by running an ϵ -greedy policy with $\epsilon = 0.05$ for a fixed number of steps. The lower table reports results of the single best performing episode for HNeat and DQN. HNeat produces deterministic policies that always get the same score while DQN used an ϵ -greedy policy with $\epsilon = 0.05$.

Q-Learning

Whiteboard

- Approximating the Q function with a neural network
- Deep Q-Learning
- Experience Replay

Alpha Go

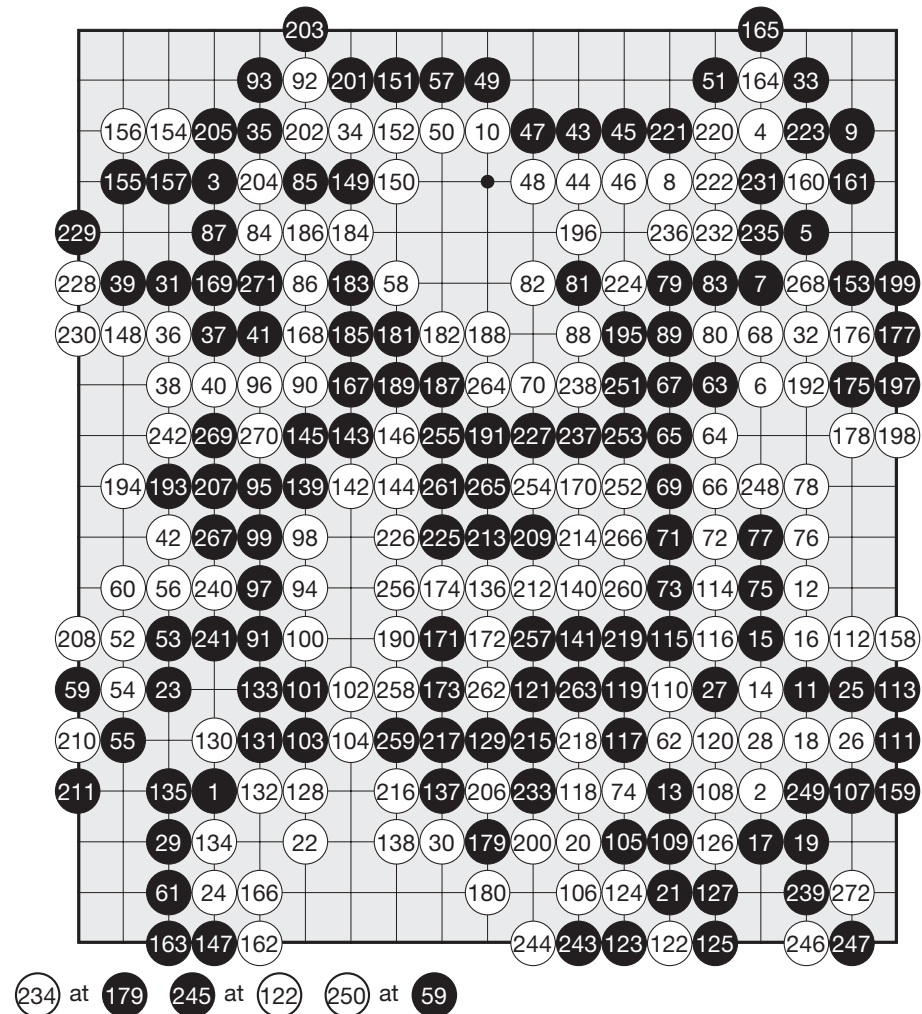
Game of Go (圍棋)

- 19x19 **board**
- Players alternately play black/white **stones**
- **Goal** is to fully encircle the largest region on the board
- **Simple** rules, but **extremely complex** game play

Game 1

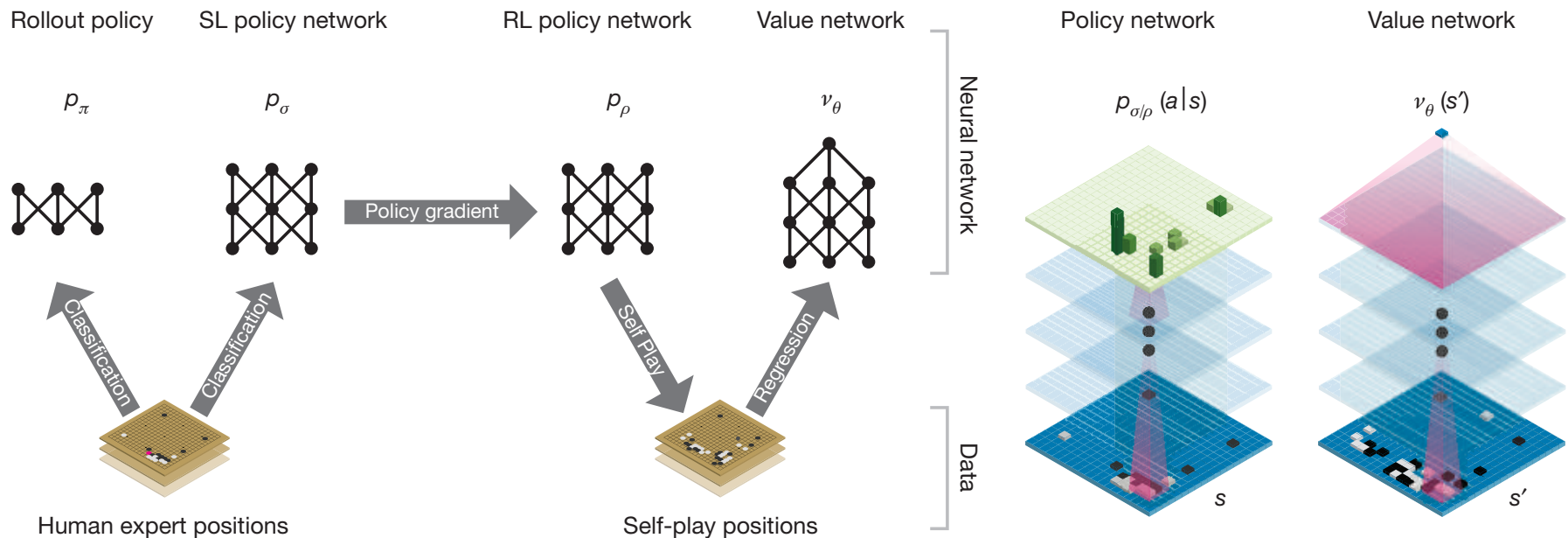
Fan Hui (Black), AlphaGo (White)

AlphaGo wins by 2.5 points



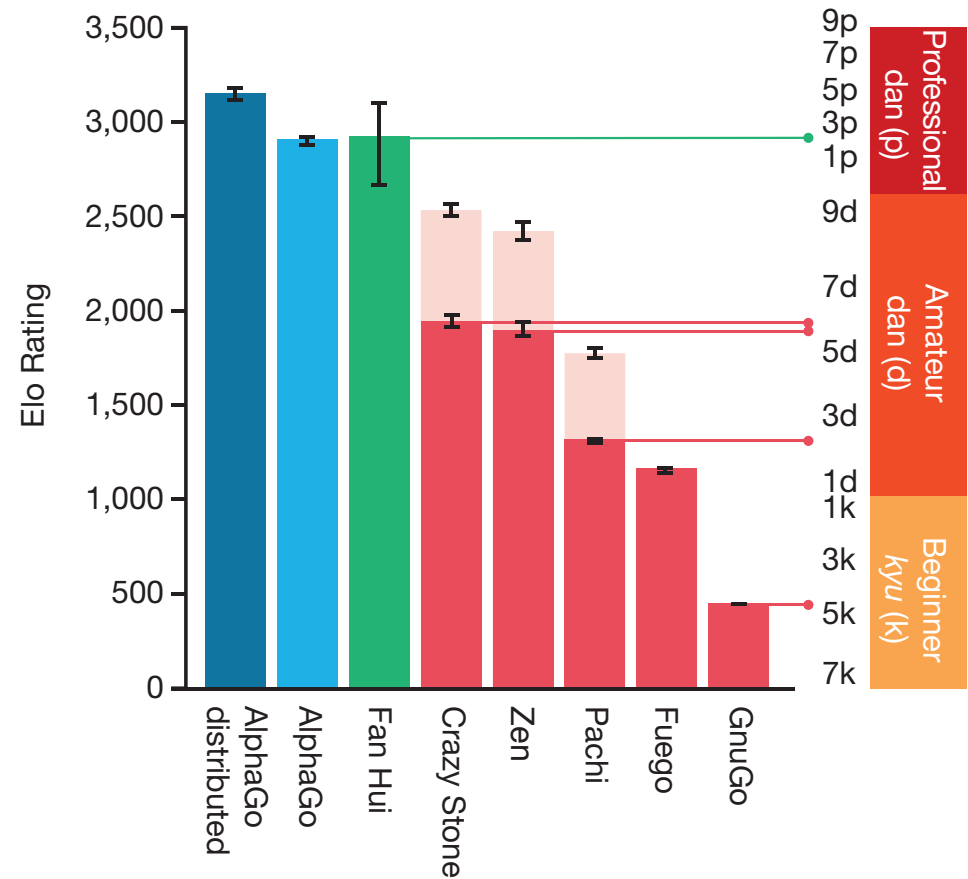
Alpha Go

- State space is too large to represent explicitly since # of sequences of moves is $O(b^d)$
 - Go: $b=250$ and $d=150$
 - Chess: $b=35$ and $d=80$
- Key idea:
 - Define a neural network to approximate the value function
 - Train by policy gradient



Alpha Go

- Results of a tournament
- From Silver et al. (2016): “a 230 point gap corresponds to a 79% probability of winning”



Learning Objectives

Reinforcement Learning: Q-Learning

You should be able to...

1. Apply Q-Learning to a real-world environment
2. Implement Q-learning
3. Identify the conditions under which the Q-learning algorithm will converge to the true value function
4. Adapt Q-learning to Deep Q-learning by employing a neural network approximation to the Q function
5. Describe the connection between Deep Q-Learning and regression