

# Bootcamp 6: Reinforcement Learning



William H. Guss, James Bartlett  
{wguss, james}@ml.berkeley.edu  
Machine Learning at Berkeley

April 22, 2016

## 1 Introduction

## 2 Theory

## 3 Algorithms

## 4 Questions

# Problem: ML for Pacman.



Reinforcement  
Learning

Guss &  
Bartlett

Introduction

Theory

Algorithms

Questions

*How would you solve pacman  
with machine learning?*



# Problem: ML for Pacman.



Reinforcement  
Learning

Guss &  
Bartlett

Introduction

Theory

Algorithms

Questions

*How would you solve pacman  
with machine learning?*

**Find a model which takes  
screen pixels to actions:**

$$\pi_{\theta} : s_t \mapsto a_t.$$



# Problem: ML for Pacman.



Reinforcement  
Learning

Guss &  
Bartlett

Introduction

Theory

Algorithms

Questions

*How would you solve pacman  
with machine learning?*

**Find a model which takes  
screen pixels to actions:**

$$\pi_{\theta} : s_t \mapsto a_t.$$

*What is your loss function?  
Data?*



# Problem: ML for Pacman.



Reinforcement  
Learning

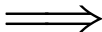
Guss &  
Bartlett

Introduction

Theory

Algorithms

Questions



# Solution: Reinforcement Learning



Reinforcement  
Learning

Guss &  
Bartlett

Introduction

Theory

Algorithms

Questions

- Supervised learning is *not* the most general formulation of learning.



# Solution: Reinforcement Learning



Reinforcement  
Learning

Guss &  
Bartlett

Introduction

Theory

Algorithms

Questions

- Supervised learning is *not* the most general formulation of learning.
- Humans learn through reward and penalty





# Solution: Reinforcement Learning



Reinforcement  
Learning

Guss &  
Bartlett

Introduction

Theory

Algorithms

Questions

- Can we make algorithms which improve with crude reward signals?

**Machine learning without  
explicit objective functions**



**Reinforcement Learning (RL)**



# The Core Idea



Reinforcement  
Learning

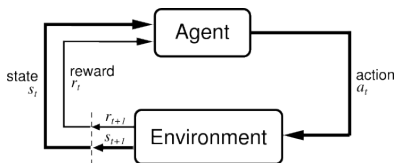
Guss &  
Bartlett

Introduction

Theory

Algorithms

Questions



- Models (agents) take action  $a_t$  in some environment.
- Environment provides state  $s_t$ , reward  $r_t$ .
- Models learn to maximize reward  $r_t$ ,  $\forall t$ .

# Markov Decision Process (MDP)



Reinforcement  
Learning

Guss &  
Bartlett

Introduction

Theory

Algorithms

Questions

Environment,  $E = (\mathcal{S}, \mathcal{A}, \mathcal{R}, \rho, r)$ .

- 1 State space,  $\mathcal{S}$
- 2 Action space,  $\mathcal{A}$
- 3 Reward space,  $\mathcal{R}$
- 4 Transition distribution,  $\rho(s' \mid s, a)$ . Given a previous state  $s$  and action  $a$ , environment gives  $s'$ .
- 5 Reward function  $r(s, a) \in \mathcal{R}$ .

**Markov Property:**  $\rho(s' \mid s, a)$  depends only on  $s, a$  not previous states!

# Markov Decision Process (MDP)



Reinforcement  
Learning

Guss &  
Bartlett

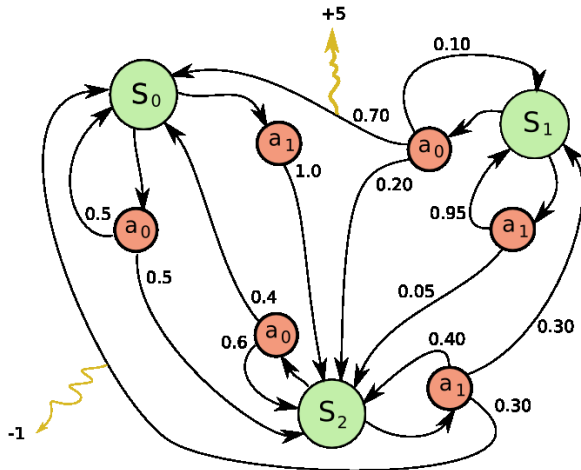
Introduction

Theory

Algorithms

Questions

## Example MDP



# Pacman as an MDP



Reinforcement  
Learning

Guss &  
Bartlett

Introduction

Theory

Algorithms

Questions

- $\mathcal{S} = \mathbb{R}^{256 \times 256}$ , images as state space.
- $\mathcal{A} = \{\uparrow, \downarrow, \rightarrow, \leftarrow\}$ , joystick as action space.
- $r(s_t, a_t)$  = change in score.
- $\rho(s_{t+1} \mid s_t, a_t)$  = next frame of game after joystick action  $a_t$ .



## Two different types of agents

- Deterministic policy  $a = \pi(s)$  acts in  $E$ .
- Stochastic policy  $a \sim \pi(a|s)$  gives a probability distribution over actions.

## Policy Trajectories

$$s_1 \xrightarrow{\pi} a_1 \xrightarrow{\rho, r} s_2, r_2 \xrightarrow{\pi} a_2 \xrightarrow{\rho, r} \dots$$

The **state value** is a function of a given state for an agent  $\pi$  defined as

$$V^\pi(s_t) = \mathbb{E} \left[ \sum_{n=t+1}^{\infty} \gamma^n r(s_n, \pi(s_n)) \right]$$

- 1  $\gamma$  is the discount factor
- 2  $\pi(s_n)$  is the action the agent  $\pi$  makes after seeing state  $s_n$ .
- 3  $r(s_n, \pi(s_n))$  is the reward the agent gets from taking that action.

The **state-action value** for an agent  $\pi$  is defined such that

$$Q^\pi(s_t, a_t) = \mathbb{E} \left[ \underbrace{r(s_t, a_t)}_{\text{reward for } a_t} + V^\pi(s_t) \right]$$

- Given some state  $s_t$ , the *best* agent,  $\pi^*$  is one that take action

$$a_t = \operatorname{argmax}_a Q(s_t, a).$$



- **Policy Optimization:** maximize the expected reward with respect to a policy  $\pi$ ;

$$\pi^* = \operatorname{argmax}_{\pi} \mathbb{E} \left[ \sum_{t=0}^{\infty} r_t \right]$$

- **Policy Evaluation:** Given some fixed policy  $\pi$  compute expected return.
  - Computing  $Q^{\pi}$ ,  $V^{\pi}$ , and other expectations on policy rollout.
  - Lets us perform policy optimization!

Reinforcement  
Learning

Guss &  
Bartlett

Introduction

Theory

**Algorithms**

Questions

# Algorithms

**Behavioral Cloning:** Supervised learning in MDPs using and expert agent expert  $\pi^*$ !

**Behavioral Cloning:** Supervised learning in MDPs using and expert agent expert  $\pi^*$ !

Given expert examples  $\mathcal{D} = (s_t, a_t = \pi^*(s_t))$  and a model  $\pi_\theta$  find  $\theta^*$  st

$$\theta^* = \underset{\theta}{\operatorname{argmin}} \mathcal{L}(a_t, \pi_\theta(s_t)).$$

where  $\mathcal{L}$  is some loss function.

**Behavioral Cloning:** Supervised learning in MDPs using and expert agent expert  $\pi^*$ !

Given expert examples  $\mathcal{D} = (s_t, a_t = \pi^*(s_t))$  and a model  $\pi_\theta$  find  $\theta^*$  st

$$\theta^* = \underset{\theta}{\operatorname{argmin}} \mathcal{L}(a_t, \pi_\theta(s_t)).$$

where  $\mathcal{L}$  is some loss function.

- Show, don't tell!

**Behavioral Cloning:** Supervised learning in MDPs using and expert agent expert  $\pi^*$ !

Given expert examples  $\mathcal{D} = (s_t, a_t = \pi^*(s_t))$  and a model  $\pi_\theta$  find  $\theta^*$  st

$$\theta^* = \underset{\theta}{\operatorname{argmin}} \mathcal{L}(a_t, \pi_\theta(s_t)).$$

where  $\mathcal{L}$  is some loss function.

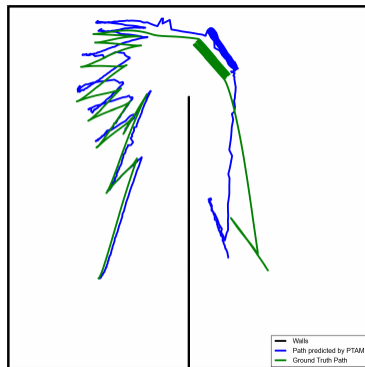
- Show, don't tell!
- No complicated machinery, just standard ML.

## Issue: Compounding Error

Given some irreducible error

$$\epsilon = 0.001$$

- $\mathcal{L}(a_0, \pi_\theta(s_0)) = \epsilon$
- $\mathcal{L}(a_1, \pi_\theta(s_1)) = 2\epsilon$
- $\mathcal{L}(a_2, \pi_\theta(s_2)) = 3\epsilon$
- $\mathcal{L}(a_3, \pi_\theta(s_3)) = 4\epsilon$
- $\mathcal{L}(a_4, \pi_\theta(s_4)) = 5\epsilon$



## Issue: Distribution Mismatch

- States expert dataset  $\mathcal{D}$  generated by  $\pi^*$  have different distribution than those generated by  $\pi_\theta$ .  
⇒ No self correction.





Time = 24.8 (sec)  
 Autonomy = 100.0%

Position Precision = 0.4 m  
 Speed Precision = 100.0%

Comfort = 69.2%

Drive mode = Auto  
 Distance = 452.5 (m)

## Goals of Q-learning

- 1 Approximate  $Q^{\pi^*}$ , the  $Q$  function of the optimal agent, as  $Q(s_t, a_t)$ .

## Goals of Q-learning

- 1 Approximate  $Q^{\pi^*}$ , the  $Q$  function of the optimal agent, as  $Q(s_t, a_t)$ .
- 2 Using  $Q$ , find the agent,  $\pi$ , that best approximates the optimal agent,  $\pi^*$ .

# Q-Learning (State-action Value Iteration)



Reinforcement  
Learning

Guss &  
Bartlett

Introduction

Theory

Algorithms

Questions

**How do we define best?**

## How do we define best?

Given some state  $s_t$ , the **best** agent,  $\pi^*$  is one that takes action

$$a_t = \arg \max_a Q(s_t, a).$$

# Q-Learning (State-action Value Iteration)



Reinforcement  
Learning

Guss &  
Bartlett

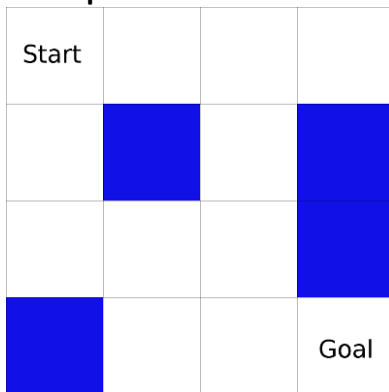
Introduction

Theory

Algorithms

Questions

## An example: Frozen Lake Problem



# Q-Learning (State-action Value Iteration)



Reinforcement  
Learning

Guss &  
Bartlett

Introduction

Theory

Algorithms

Questions

- 100 reward for reaching the goal
- 0 otherwise

**How do we keep track of this long term reward?**

# Q-Learning (State-action Value Iteration)



Reinforcement  
Learning

Guss &  
Bartlett

Introduction

Theory

Algorithms

Questions

- 100 reward for reaching the goal
- 0 otherwise

**How do we keep track of this long term reward?**

$Q$  function



# Q-Learning (State-action Value Iteration)



Reinforcement  
Learning

Guss &  
Bartlett

Introduction

Theory

Algorithms

Questions

**How do we actually calculate the  $Q$  function?**

# Q-Learning (State-action Value Iteration)



Reinforcement  
Learning

Guss &  
Bartlett

Introduction

Theory

Algorithms

Questions

**How do we actually calculate the  $Q$  function?**

**The Bellman Equation.**

# Q-Learning (State-action Value Iteration)

Reinforcement  
Learning

Guss &  
Bartlett

Introduction

Theory

Algorithms

Questions

**How do we actually calculate the  $Q$  function?**

**The Bellman Equation.**

$$Q^{\pi}(s_t, a_t) = r_t + \gamma Q^{\pi}(s_{t+1}, \pi(s_{t+1}))$$

## One $Q$ -Learning Algorithm: Tabular $Q$ -Learning

- Explore the environment
- On the way, use the Bellman equation to store a table of expected future reward ( $Q$ ) for each state-action pair.
- Use this table to pick the best possible action for any given state.

## Q-Learning (State-action Value Iteration)

## An example update for Frozen Lake.

Suppose our stored  $Q$  table looks like so:

[illegible]

# Q-Learning (State-action Value Iteration)

Reinforcement  
Learning

Guss &  
Bartlett

Introduction

Theory

Algorithms

Questions

**An example update for Frozen Lake.**

Then suppose our agent moves **down** from the starting square

## Questions

[illegible]

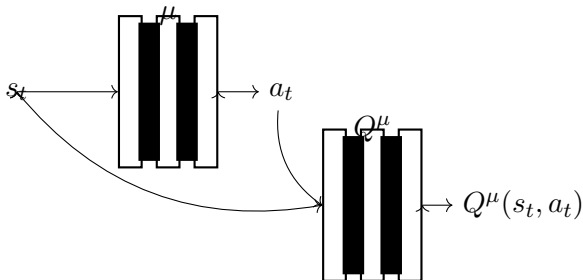




# Policy Iteration

## Deep Deterministic Policy Gradient

- 1 Actor neural network  $\mu : \mathcal{S} \rightarrow \mathcal{A}$
- 2 Critic network  $Q^\mu : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$
- 3 Performance of  $\mu$  is  $Q^\mu(s_t, \mu(s_t))$ . **Maximize performance!**  $\nabla_W Q^\mu(s_t, a_t) = \nabla_a Q^\mu(s_t, a) \cdot \nabla_W \mu(s_t)$



Reinforcement  
Learning

Guss &  
Bartlett

Introduction

Theory

Algorithms

Questions

# Questions?