Reinforcement Learning

Bartlett

ntroduction

 Theory

Algorithms

Question:

Bootcamp 6: Reinforcement Learning



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

April 22, 2016

Overview



Reinforcement Learning

> Guss & Bartlett

. . .

miroducti

...---,

Algorithm

Question

- 1 Introduction
- 2 Theory
- 3 Algorithms
- 4 Questions



Reinforcement Learning

> Guss & Bartlett

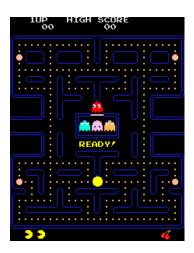
Introduction

Theory

Algorithm

Question

How would you solve pacman with machine learning?





Reinforcement Learning

> Guss & Bartlett

Introduction

Theory

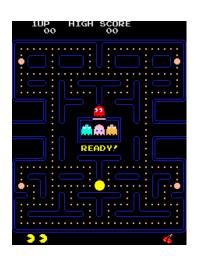
Algorithm

Question

How would you solve pacman with machine learning?

Find a model which takes screen pixels to actions:

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





Reinforcement Learning

> Guss & Bartlett

Introduction

Algorithm

Theory

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?





Reinforcement Learning

Guss & Bartlett

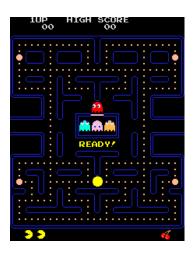
Introduction

Theory

Algorithm

Questions





Solution: Reinforcement Learning



Reinforcement Learning

Guss &

Introduction

Theory

Algorithm

Questio

 Supervised learning is not the most general formulation of learning.



Solution: Reinforcement Learning



Reinforcement Learning

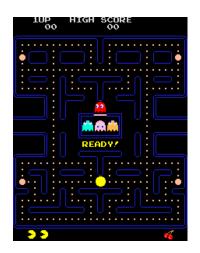
Guss &

Introduction

Theory

Algorithm

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



Solution: Reinforcement Learning



Reinforcement Learning

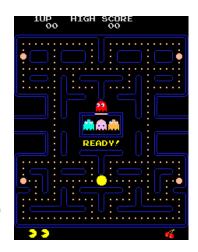
Introduction

Can we make algorithms which improve with crude reward signals?

Machine learning without explicit objective functions



Reinforcement Learning (RL)



The Core Idea



Reinforcement Learning

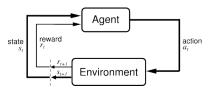
> Guss & Bartlett

troduction

Theory

Algorithm:

Questior



- 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

Dartiett

Introductio

Theory

Algorithm

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

- $lue{1}$ State space, ${\cal S}$
- 2 Action space, A
- f 3 Reward space, $\cal R$
- 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

Bartlett

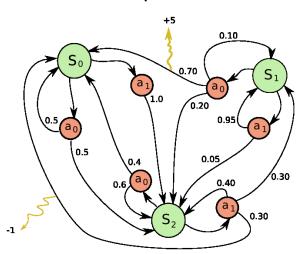
ntroductio

Theory

Algorithm

Question

Example MDP



Pacman as an MDP



Reinforcement Learning

> Guss & Bartlett

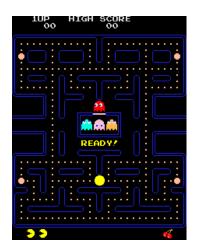
Introduction

Theory

Algorithm

Ŭ

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



Policies/Agents



Reinforcement Learning

> Guss & Bartlett

ntroduction

Theory

Algorithm

Questio

Two different types of agents

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

Policy Trajectories

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

Value under a policy



Reinforcement Learning

Bartlet

Introduction

Theory

Algorithms

Question

The **state value** is a function of a given state for an agent π defined as

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

- $oldsymbol{1}$ γ is the discount factor
- \mathbf{Z} $\pi(s_n)$ is the action the agent π makes after seeing state s_n .
- $r(s_n, \pi(s_n))$ is the reward the agent gets from taking that action.

Value under a policy



Reinforcement Learning

Bartlet

Introduction

Theory

A.1. 2.1

Aigorithms

Questions

The **state-action value** for an agent π 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, π^* is one that take action

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

Problems in Reinforcement Learning



Reinforcement Learning

Bartlett

Introductio

Theory

Algorithms

Policy Optimization: maximize the expected reward with respect to a policy π ;

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

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



Reinforcement Learning

Bartlett

Introduction

Theory

Algorithms

Assorted Algorithms

We'll go over:

- Behavioral Cloning
- Q-Learning
- Policy Iteration

Learn at home:

- Value iteration
- Temporal Difference Methods
- Inverse Reinforcement Learning.



Reinforcement Learning

> Guss & Bartlett

Introduction

Theory

Algorithms

Juestions

Behavioral Cloning: Supervised learning in MDPs using and expert agent expert π^* !



Reinforcement Learning

Bartlett

Introduction

Tilcory

Algorithms

Questic

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 π_{θ} find θ^* st

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

where \mathcal{L} is some loss function.



Reinforcement Learning

Bartlett

Introduction

Algorithms

Questic

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 π_{θ} find θ^* st

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

where \mathcal{L} is some loss function.

Show, don't tell!



Reinforcement Learning

Dartiett

Introduction

Tricory

Algorithms

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 π_{θ} find θ^* st

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

where \mathcal{L} is some loss function.

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



Reinforcement Learning

> Guss & Bartlett

ntroduction

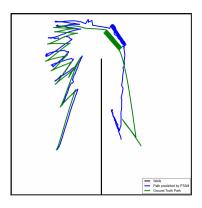
T1.

Algorithms

Ŭ

Issue: Compounding Error

Given some irreducible error $\epsilon = 0.001$





Reinforcement Learning

> Guss & Bartlet

Algorithms

Issue: Distribution Mismatch

• States expert dataset \mathcal{D} generated by π^* have different distribution than those generated by π_{θ} .

 \implies No self correction.





Reinforcement Learning

> Guss & Bartlet

ntroduction

Theory

Algorithms

Issue: Distribution Mismatch

States expert dataset \mathcal{D} generated by π^* have different distribution than those generated by π_{θ} .

⇒ No self correction.

Solution: DAgger.

- Do BC on \mathcal{D} and generate E_0 states generated by π_{θ} .
- Label E_0 with expert level actions and add to \mathcal{D} .





Reinforcement Learning

> Guss & Bartlett

ntroduction

THEOLY

Algorithms

Questio

Goals of Q-learning

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



Reinforcement Learning

> Guss & Bartlett

ntroduction

 Theory

Algorithms

Question

Goals of Q-learning

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



Reinforcement Learning

> Guss & Bartlett

miroductic

.

Algorithms

Question:

How do we define best?



Reinforcement Learning

> Guss & Bartlett

ntroductior

Tricory

Algorithms

Questior

How do we define best?

Given some state s_t , the **best** agent, π^* is one that takes action

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



Reinforcement Learning

Bartlett

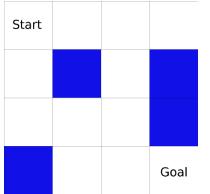
Introductio

Theory

Algorithms

Questions

An example: Frozen Lake Problem





Reinforcement Learning

> Guss & Bartlett

ntroductio

Theory

Algorithms

Ougetier

- 100 reward for reaching the goal
- $\blacksquare 0$ otherwise

How do we keep track of this long term reward?



Reinforcement Learning

> Guss & Bartlett

ntroductio

i neory

Algorithms

Questio

lacksquare 100 reward for reaching the goal

• 0 otherwise

How do we keep track of this long term reward?

Q function



Reinforcement Learning

Bartlett

ntroductio

Theory

Algorithms

Questions

How do we actually calculate the ${\it Q}$ function?



Reinforcement Learning

> Guss & Bartlett

troduction

_.

Algorithms

```

How do we actually calculate the Q function? The Bellman Equation.



Reinforcement Learning

> Guss & Bartlett

ntroductio

Theory

Algorithms

Questions

How do we actually calculate the  ${\cal Q}$  function?

The Bellman Equation.

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



Reinforcement Learning

Bartlett

ntroduction

i neory

Algorithms

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.



Reinforcement Learning

> Guss & Bartlett

Introduction

Theory

Algorithms

Questions

#### An example update for Frozen Lake.

Suppose our stored  ${\cal Q}$  table looks like so:

| Up | Down | Left | Right |
|----|------|------|-------|
| 0  | 65   | 0    | 40    |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 50 | 75   | 30   | 20    |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |



Reinforcement Learning

Bartiett

ntroduction

Algorithms

Questions

An example update for Frozen Lake.

Then suppose our agent moves down from the starting square



Reinforcement Learning

> Guss & Bartlett

ntroduction

\_.

Algorithms

Questior

#### An example update for Frozen Lake.

Then we update using the Bellman equation.

$$Q(s_{t+1}, a_{t+1}) = Q(s_t, a_t) + \alpha(r_t + \gamma(\max_a Q(s_t, a) - Q(s_t, a_t)))$$

| Up | Down | Left | Right |
|----|------|------|-------|
| 0  | 65   | 0    | 40    |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 50 | 75   | 30   | 20    |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |



Reinforcement Learning

Bartlett

Introduction

Theory

Algorithms

Questions

#### An example update for Frozen Lake.

The table now looks like so:

| Up | Down | Left | Right |
|----|------|------|-------|
| 0  | 70   | 0    | 40    |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 50 | 75   | 30   | 20    |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |
| 0  | 0    | 0    | 0     |



Reinforcement Learning

Bartlet

ntroduction

Algorithms

**Policy Iteration:** Given access to the MDP, use policy evaluation to iteratively serach for better policies!

■ Choose a policy at random,  $\pi$ .



Reinforcement Learning

Algorithms

**Policy Iteration:** Given access to the MDP, use policy evaluation to iteratively serach for better policies!

- Choose a policy at random,  $\pi$ .
- Alternate between
  - Evaluate policy  $\pi \to V^{\pi}$ .



Reinforcement Learning

Bartlet

ntroduction

Theory

Algorithms

**Policy Iteration:** Given access to the MDP, use policy evaluation to iteratively serach for better policies!

- Choose a policy at random,  $\pi$ .
- Alternate between
  - Evaluate policy  $\pi \to V^{\pi}$ .
  - Set new policy to be greedy policy for  $V^{\pi}$

$$\pi(s) := \operatorname*{argmax}_{a} \mathbb{E} \left[ R(s, a) + \gamma V^{\pi}(s') \right]$$
$$:= \operatorname*{argmax}_{a} Q^{\pi}(s, a)$$



Reinforcement Learning

Algorithms

**Policy Iteration:** Given access to the MDP, use policy evaluation to iteratively serach for better policies!

- Choose a policy at random,  $\pi$ .
- Alternate between
  - Evaluate policy  $\pi \to V^{\pi}$ .
  - Set new policy to be greedy policy for  $V^{\pi}$

$$\pi(s) := \operatorname*{argmax}_{a} \mathbb{E} \left[ R(s, a) + \gamma V^{\pi}(s') \right]$$
$$:= \operatorname*{argmax}_{a} Q^{\pi}(s, a)$$

• Learn  $Q^{\pi}$  using Q-learning without  $\operatorname{argmax}$ .



Reinforcement Learning

> Guss & Bartlett

atroduction

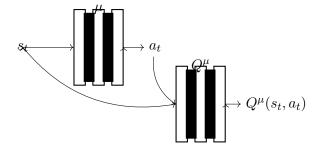
Theony

Algorithms

Algorithm

#### **Deep Determisitic Policy Gradient**

- **1** Actor neural network  $\mu: \mathcal{S} \to \mathcal{A}$
- 2 Critic network  $Q^{\mu}: \mathcal{S} \times \mathcal{A} \to \mathbb{R}$
- 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

ntroduction

Theory

Algorithms

Questions

# Questions?