Reinforcement Learning

Bartlett

ntroduction

Theory

Algorithm

Practice

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

Dartiett

Introductio

Theory

Algorit

Practic

Question

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



Reinforcement Learning

> Guss & Bartlett

Introduction

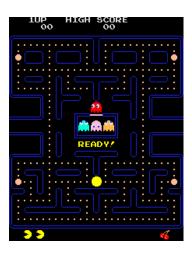
Theory

Algorith

D.

Question

How would you solve pacman with machine learning?





Reinforcement Learning

> Guss & Bartlett

Introduction

Theory

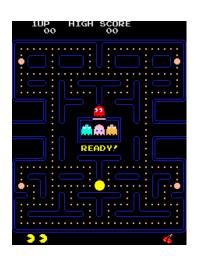
Algorithr

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

Theony

Algorith

Б. ..

Ouestie

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 &

Introduction

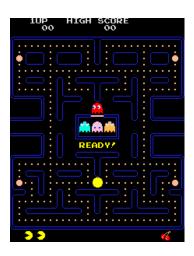
Theory

. ...

Dractice

Question





Solution: Reinforcement Learning



Reinforcement Learning

Guss &

Introduction

Theory

0

Dractica

Question

 Supervised learning is not the most general formulation of learning.



Solution: Reinforcement Learning



Reinforcement Learning

Guss &

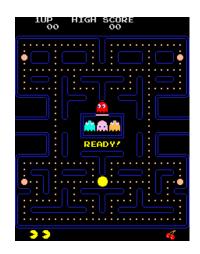
Introduction

Theory

Algorith

Question

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



Solution: Reinforcement Learning



Reinforcement Learning

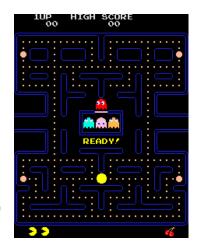
Introduction

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

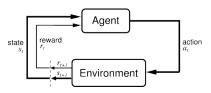
troduction

Theory

. ...

Dractice

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

Dartiett

Introductio

Theory

Algorithr

Practic

Question

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

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

> Guss & Bartlett

ntroductio

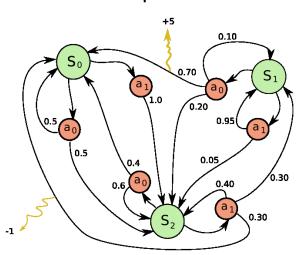
Theory

Algorith

Practice

Question

Example MDP



Pacman as an MDP



Reinforcement Learning

> Guss & Bartlett

ntroduction

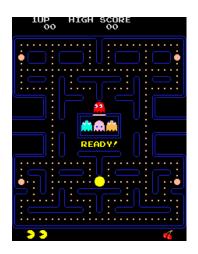
Theory

Algorith

Б. ...

Question

- $S = \mathbb{R}^{256 \times 256}$, 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

ntroductio

Theory

Algorithm

Question

Two different types of agents

- Deterministic policy $a = \pi(s)$ acts in E.
- 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

Theory

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

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

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

Value under a policy



Reinforcement Learning

Bartlett

Introductio

Theory

Algorithm

Question

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

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

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

$$s_t \xrightarrow{\pi} a_0 \xrightarrow{\rho,r} s_7, r_7 \xrightarrow{\pi} a_3 \xrightarrow{\rho,r} \cdots$$

Value under a policy



Reinforcement Learning

Guss & Bartlett

Introductio

Theory

Practice

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

> Guss & Bartlett

Introductio

Theory

, tigoritiii

Practice

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

I heory

Algorithms

Practice

Questio

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

ntroduction

Theory

Algorithms

Б. ..

Question

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



Reinforcement Learning

> Guss & Bartlett

Introductior

Theory

Algorithms

Question

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

> Guss & Bartlett

Introduction

THEOLY

Algorithms

Practice

Questions

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

.

Algorithms

Dractica

Practice

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.

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

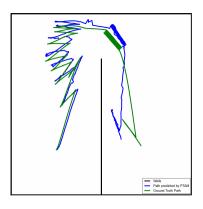


Reinforcement Learning

Algorithms

Issue: Compounding Error

Given some irreducible error $\epsilon = 0.001$





Reinforcement Learning

> Guss & Bartlet

ntroduction

Theory

Algorithms

Ouestion

Issue: Distribution Mismatch

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

⇒ No self correction.





Reinforcement Learning

> Guss & Bartlet

Introduction

Algorithms

Donath

. ..

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

I heory

Algorithms

Practio

Question

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

- 1 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

Bartlet

Introductio

Theory

Algorithms

Practio

Question

How do we define best?



Reinforcement Learning

> Guss & Bartlett

ntroduction

 Theory

Algorithms

200

Question

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

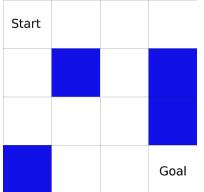
ntroductio

 Theory

Algorithms

Tactice

An example: Frozen Lake Problem





Reinforcement Learning

> Guss & Bartlett

ntroduction

 Theory

Algorithms

A

■ 100 reward for reaching the goal

lacksquare 0 otherwise

How do we keep track of this long term reward?



Reinforcement Learning

> Guss & Bartlett

ntroduction

Theory

Algorithms

Б. ..

Ouestion

■ 100 reward for reaching the goal

■ 0 otherwise

How do we keep track of this long term reward?

Q function



Reinforcement Learning

Dartiett

ntroduction

Theory

Algorithms

Practic

Question

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



Reinforcement Learning

Bartlet

troduction

Algorithms

Practice

Question

How do we actually calculate the ${\cal Q}$ function? The Bellman Equation.



Reinforcement Learning

> Guss & Bartlett

ntroduction

Theory

Algorithms

Question

How do we actually calculate the Q function?

The Bellman Equation.

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



Reinforcement Learning

> Guss & Bartlett

ntroduction

Theory

Algorithms

Dractica

Question

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

Bartlett

Introduction

Theory

Algorithms

Practice

Question

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

Bartlett

Introductior

Theory

Algorithms

Practice

Question

An example update for Frozen Lake.

Then suppose our agent moves down from the starting square



Reinforcement Learning

Bartlett

Introduction

Theory

Algorithms

D.......

Practice

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

D.

Question

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

Bartlett

ntroduction

Theory

Algorithms

Practice

• Choose a policy at random, π .

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



Reinforcement Learning

Bartlett

Introduction

THEOLY

Algorithms

Practice

Question

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

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



Reinforcement Learning

Bartlett

Introduction

Algorithms

Practice

Question

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

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

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



Reinforcement Learning

Bartlett

Introduction

Algorithms

.

Practice

Question

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

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

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

Note: Learn Q^{π} using Q-learning without argmax .

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



Reinforcement Learning

Bartlett

ntroduction

Algorithms

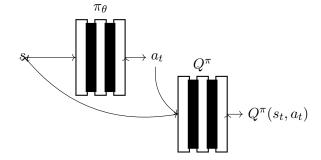
Б. ...

Practice

Question

Example: Deep Determisitic Policy Gradient

- **1** Actor neural network $\pi_{\theta}: \mathcal{S} \to \mathcal{A}$
- 2 Critic network $Q^{\pi}: \mathcal{S} \times \mathcal{A} \to \mathbb{R}$
- 3 Performance of π is $Q^{\pi}(s_t, \pi(s_t))$. Maximize performance! $\nabla_{\theta}Q^{\pi}(s_t, a_t) = \nabla_aQ^{\pi}(s_t, a) \cdot \nabla_{\theta}\pi_{\theta}(s_t)$





Reinforcement Learning

> Guss & Bartlett

itroduction

Theory

Algorithms

Practic

Question

and many more...

Papers and links to descriptions of other algorithms will be posted on github.com/mlberkeley/bootcamp.

Reinforcement Learning

> Guss & Bartlett

ntroduction

 Theory

Algorithms

Practice

Question

Practical Reinforcement Learning

Process



Reinforcement Learning

Dartiett

Introductio

Theory

Algorith

Practice

Question

- 1 Clearly define the environment.
 - What is state space, S? Action space A? Low dimensions is better!
 - Use tools like OpenAl Gym to make a simulator.
- 2 Choose a good reward function.
 - Don't make sparse reward functions:
 - +1 for winning vs. guiding agent to win.
 - Penalize for time constraints!
 - Break the problem in to smaller goals with rewards r_1, r_2, \ldots Then $r(s, a) = \sum_i w_i r_i(s, a)$.

Process



Reinforcement Learning

Dartiett

Introduction

Theory

Practice

Question

3 Choose a model.

- Deep reinforcement learning is king!
 - Discrete $A \rightarrow DQN, DDQN$
 - Continuous $A \rightarrow DDPG, NAF, A3C, TRPO$.
- Less hyperparameters better.
- Don't rule models out too early.
- 4 Train the model against the environment.
 - Choose a good exploration strategy.
 - RL can take a long time (days months), let models train to completion.
 - Do a sweep on hyperparameters.
 - Overfitting (to optimal policy) is good.
 - (You will fail here many times, go to step 3.)
- 5 Profit

$\mathsf{ProTips}^{tm}$



Reinforcement Learning

Bartlett

Thereadiction

Algorith

Practice

Question

- Fastest way to learn RL: Implement landmark papers in tensorflow, Atari DQN, DPG, ...
- Write unit tests before you train!
- Use an autodiff framework: Tensorflow, PyTorch, ...
- Avoid pixel data at all costs: CNNs/CV algorithms make training very slow!
 - Test RL models on featurized data in simualtion, then move to pixel based models for real data.

Good luck!

Reinforcement Learning

> Guss & Bartlett

ntroduction

Theory

Algorithm

Practice

Questions

Questions?