The reinforcement learning problen

Solution techniques

Sample complexity

▲ロト ▲ □ ト ▲ □ ト ▲ □ ト ● ● の Q ()

Hands-on

Reinforcement learning tutorial

Wouter Caarls

Koroibot Summer School, September 25th, 2014

he reinforcement learning problem

Solution techniques

Sample complexity

Hands-on

What is reinforcement learning?



[R. Nagel, 2008]

The reinforcement learning problem

Solution techniques

Sample complexity

▲□▶▲□▶▲□▶▲□▶ □ のQで

Hands-on

What is reinforcement learning?

Supervised learning

Given samples of $x \rightarrow y$, find y = f(x) (regression)

Reinforcement learning

Given samples of $(x, y) \rightarrow z$, find y = f(x) that maximizes z.

Unsupervised learning

Given a dataset x, find interesting patterns (clusters, relations, etc.)

The reinforcement learning problem

Solution techniques

Sample complexity

▲□▶▲□▶▲□▶▲□▶ □ のQで

Hands-on

What is reinforcement learning?

Supervised learning

Given samples of $x \rightarrow y$, find y = f(x) (regression)

Reinforcement learning

Given samples of $(x, y) \rightarrow z$, find y = f(x) that maximizes *z*.

Unsupervised learning

Given a dataset *x*, find interesting patterns (clusters, relations, etc.)

he reinforcement learning problen

Solution techniques

Sample complexity

Hands-on

Learning by trial and error



Goal

Find actions that maximize the reward received over the lifetime of the agent.

The reinforcement learning problem

Solution techniques

Sample complexity

▲□▶▲□▶▲□▶▲□▶ □ のQで

Hands-on

Why reinforcement learning?

Specify the problem, not the solution

- Easier: desired outcome is usually known
- Still tricky: reward engineering
- Can deal with delayed rewards
 - Optimizes over all future rewards
- Model-free
 - No mapping, system identification, etc. beforehand
 - Works in any (static) environment

The reinforcement learning problem

Solution techniques

Sample complexity

▲□▶▲□▶▲□▶▲□▶ □ のQで

Hands-on

Why reinforcement learning?

- Specify the problem, not the solution
 - Easier: desired outcome is usually known
 - Still tricky: reward engineering
- Can deal with delayed rewards
 - Optimizes over all future rewards
- Model-free
 - No mapping, system identification, etc. beforehand
 - Works in any (static) environment

The reinforcement learning problem

Solution techniques

Sample complexity

▲□▶▲□▶▲□▶▲□▶ □ のQで

Hands-on

Why reinforcement learning?

- Specify the problem, not the solution
 - Easier: desired outcome is usually known
 - Still tricky: reward engineering
- Can deal with delayed rewards
 - Optimizes over all future rewards
- Model-free
 - No mapping, system identification, etc. beforehand
 - Works in any (static) environment

he reinforcement learning probler

Solution techniques

Sample complexity

Hands-on

Context of learning control



▲□▶▲圖▶▲≧▶▲≧▶ 差 のへ⊙

The reinforcement learning problem

Solution techniques

Sample complexity

▲□▶▲□▶▲□▶▲□▶ □ のQで

Hands-on

Interactive (no model)



All actions are sequentially executed in the real world.

- Sequentially
 - Observations only occur in trajectories
 - Cannot evaluate different actions in the same state
- Real world
 - Actions take real time
 - Actions must be selected as fast as possible
 - Dangerous actions can damage the system

Introduction
00000000000

he reinforcement learning problem

Solution techniques

Sample complexity

▲ロト ▲ □ ト ▲ □ ト ▲ □ ト ● ● の Q ()

Hands-on

Cost per step



Reward is available at intermediate steps and only based on the current transition.

- Intermediate steps
 - Allows us to assign credit to particular actions
 - May be delayed
- Based on current transition
 - No integrated costs at the end of a trajectory
 - No "target arrival time" unless time can be observed

ne reinforcement learning problem

Solution techniques

Sample complexity

Hands-on

Example: LEO

Main characteristics:

- '2D' (boom keeps hip axis horizontal)
- \circ \sim 50 cm
- ~1.7 kg
- 7 servo motors (Dynamixel)
- On-board computing
- Autonomous (except power)

Task:

Learn to walk



[Schuitema et al., 2010]

▲ロト ▲ □ ト ▲ □ ト ▲ □ ト ● ● の Q ()

he reinforcement learning problen

Solution techniques

Sample complexity

◆□▶ ◆□▶ ◆ □▶ ◆ □▶ ─ □ ─ のへぐ

Hands-on

Example: LEO



[Schuitema, 2010]

he reinforcement learning problem

Solution techniques

Sample complexity

Hands-on

Example: Pancake flipping



[Kormuschev et al., 2010]

Sample complexity

Hands-on

Outline

Introduction

- 2 The reinforcement learning problem
 - Markov decision processes
 - Value functions
- 3 Solution techniques
- 4 Sample complexity

5 Hands-on

▲□▶▲圖▶▲圖▶▲圖▶ ▲国 ● ●

The reinforcement learning problem

Solution techniques

Sample complexity

▲ロト ▲ □ ト ▲ □ ト ▲ □ ト ● ● の Q ()

Hands-on

Markov chain



- Next state depends only on current state
- Transitions may be stochastic, but distribution must be static given the current state

The reinforcement learning problem

Solution techniques

Sample complexity

▲□▶ ▲圖▶ ▲臣▶ ★臣▶ = 臣 = のへで

Hands-on

Markov process



Action also influences next state

The reinforcement learning problem

Solution techniques

Sample complexity

< □ > < 同 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Hands-on

Markov decision process



- Introduces optimization criterion: find actions that maximize reward
- Reward (cost) depends only on the state transition

States and actions

- Given current state and action, next state distribution is static $p(s'|s, a) = \mathcal{T}(s, a, s')$
- History-independent $p(s_n|s_1, a_1, s_2, a_2, \dots, s_{n-1}, a_{n-1}) = p(s_n|s_{n-1}, a_{n-1})$



States and actions

- Given current state and action, next state distribution is static $p(s'|s, a) = \mathscr{T}(s, a, s')$
- History-independent $p(s_n|s_1, a_1, s_2, a_2, \dots, s_{n-1}, a_{n-1}) = p(s_n|s_{n-1}, a_{n-1})$



Sample complexity

Hands-on

States and actions

- Given current state and action, next state distribution is static $p(s'|s, a) = \mathcal{T}(s, a, s')$
- History-independent $p(s_n|s_1, a_1, s_2, a_2, \dots, s_{n-1}, a_{n-1}) = p(s_n|s_{n-1}, a_{n-1})$



Sample complexity

Hands-on

States and actions

- Given current state and action, next state distribution is static $p(s'|s, a) = \mathcal{T}(s, a, s')$
- History-independent $p(s_n|s_1, a_1, s_2, a_2, \dots, s_{n-1}, a_{n-1}) = p(s_n|s_{n-1}, a_{n-1})$



States and actions

- Given current state and action, next state distribution is static $p(s'|s, a) = \mathscr{T}(s, a, s')$
- History-independent $p(s_n|s_1, a_1, s_2, a_2, \dots, s_{n-1}, a_{n-1}) = p(s_n|s_{n-1}, a_{n-1})$



▲□▶▲□▶▲□▶▲□▶ □ のQで

Rewards

Determine the task

- Depend only on the current state transition p(r|s, a, s') = p(s, a, s')
- Learning goal is to maximize the expected return $R_t = r_{t+1} + r_{t+2} + \dots + r_T$
- Should reward desirable behaviors and punish undesirable ones
- In practice, often given by designer instead of environment

Define the problem, not the solution (reward engineering).

The reinforcement learning problem

Solution techniques

Sample complexity

▲□▶▲□▶▲□▶▲□▶ □ のQで

Hands-on

Episodic and continuing tasks

- Some tasks are naturally episodic: Reaching, grasping, standing up, kicking, etc.
- Others are continuing: Holding, balancing, walking, etc.
- Continuing tasks (T = ∞) may have infinite returns
- Introduce discounted return:

$$R_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \cdots$$
$$= \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$

• Discount factor $\gamma \in [0,1]$ determines task horizon

Control policy

Specifies which action to take in each state

- Due to Markov property, only has to depend on current state $\pi(a|s) = f(s)$
- Optimal policy \u03c6^{*} maximizes expected return

$$\pi^* = \operatorname*{arg\,max}_{\pi} E_{\pi} \{ R_t \}$$

Optimal subproblems

In an acyclic MDP, the optimal policy for a state *s* is independent of the optimal policy for the preceeding states, leading to recursive solutions.

▲□▶▲□▶▲□▶▲□▶ □ のQで

Control policy

Specifies which action to take in each state

- Due to Markov property, only has to depend on current state $\pi(a|s) = f(s)$
- Optimal policy \u03c0^{*} maximizes expected return

$$\pi^* = \operatorname*{arg\,max}_{\pi} E_{\pi} \{ R_t \}$$

Optimal subproblems

In an acyclic MDP, the optimal policy for a state *s* is independent of the optimal policy for the preceeding states, leading to recursive solutions.

The reinforcement learning problem

Solution techniques

Sample complexity

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

Hands-on

Control scheme



The reinforcement learning problem

Solution techniques

Sample complexity

▲□▶ ▲圖▶ ▲臣▶ ★臣▶ = 臣 = のへで

Hands-on

Control scheme



The reinforcement learning problem

Solution techniques

Sample complexity

Hands-on

Robot task modeled as an MDP

- States: $s = [\theta, \dot{\theta}]$
- Actions:

 $a = \tau$

- Transition function: $s' = s + \int_0^t \operatorname{eom}(s, a) dt + \mathcal{N}(\mu, \sigma)$
- Rewards:
 - $r = -sQs^T Pa^2$
- Policy:

 $a = \pi(s)$



< □ > < 同 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Introduction
0000000000

The reinforcement learning problem

Solution techniques

Sample complexity

Hands-on

Value functions

Learning goal: find optimal policy π^* that maximizes the expected return R_t

 $\pi^* = \operatorname*{arg\,max}_{\pi} E_{\pi} \{ R_t \}$

Define a value function V^{π} that stores the expected return for each state *s* under a certain policy π :

$$egin{aligned} & \chi^{\pi}(s) = E_{\pi} \left\{ R_t | s_t = s
ight\} \ &= E_{\pi} \left\{ \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | s_t = s
ight\} \end{aligned}$$

then

$$\pi^* = \operatorname*{arg\,max}_{\pi} V^{\pi}$$

 $V^*(s) = V^{\pi^*}(s) = \operatorname*{max}_{\pi} V^{\pi}(s)$

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ ▲ □ ● の Q ()

Introduction
0000000000

The reinforcement learning problem

Solution techniques

Sample complexity

▲□▶▲□▶▲□▶▲□▶ □ のQで

Hands-on

Value functions

Learning goal: find optimal policy π^* that maximizes the expected return R_t

$$\pi^* = \operatorname*{arg\,max}_{\pi} E_{\pi} \{ R_t \}$$

Define a value function V^{π} that stores the expected return for each state *s* under a certain policy π :

$$egin{aligned} &\mathcal{V}^{\pi}(s) = \mathcal{E}_{\pi}\left\{\mathcal{R}_{t}|s_{t}=s
ight\} \ &= \mathcal{E}_{\pi}\left\{\sum_{k=0}^{\infty}\gamma^{k}r_{t+k+1}|s_{t}=s
ight\} \end{aligned}$$

then

$$\pi^* = rg\max_{\pi} V^{\pi}$$

 $V^*(s) = V^{\pi^*}(s) = \max_{\pi} V^{\pi}(s)$

Introduction
0000000000

The reinforcement learning problem

Solution techniques

Sample complexity

▲□▶▲□▶▲□▶▲□▶ □ のQで

Hands-on

Value functions

Learning goal: find optimal policy π^* that maximizes the expected return R_t

$$\pi^* = \operatorname*{arg\,max}_{\pi} E_{\pi} \{ R_t \}$$

Define a value function V^{π} that stores the expected return for each state *s* under a certain policy π :

$$egin{aligned} &\mathcal{V}^{\pi}(s) = \mathcal{E}_{\pi}\left\{\mathcal{R}_{t}|s_{t}=s
ight\} \ &= \mathcal{E}_{\pi}\left\{\sum_{k=0}^{\infty}\gamma^{k}r_{t+k+1}|s_{t}=s
ight\} \end{aligned}$$

then

$$\pi^* = rg\max_{\pi} V^{\pi}$$

 $V^*(s) = V^{\pi^*}(s) = \max_{\pi} V^{\pi}(s)$

The reinforcement learning problem

V

Solution techniques

Sample complexity

▲□▶ ▲圖▶ ▲臣▶ ★臣▶ = 臣 = のへで

Hands-on

Recursive definition

$$\pi(s) = E_{\pi} \left\{ \sum_{k=0}^{\infty} \gamma^{k} r_{t+k+1} | s_{t} = s \right\}$$
$$= E_{\pi} \left\{ r_{t+1} + \sum_{k=1}^{\infty} \gamma^{k} r_{t+k+1} | s_{t} = s \right\}$$
$$= E_{\pi} \left\{ r_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^{k} r_{t+k+2} | s_{t} = s \right\}$$
$$= E_{\pi} \left\{ r_{t+1} + \gamma V^{\pi}(s_{t+1}) | s_{t} = s \right\}$$

Allows us to write the optimal policy in terms of the optimal value function

$$\pi^*(s) = \arg\max_{a} E_{\pi} \{ r_{t+1} + \gamma V^*(s_{t+1}) | s_t = s, a_t = a \}$$

The reinforcement learning problem

Solution techniques

Sample complexity

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ ▲ □ ● ● ● ●

Hands-on

Recursive definition

$$V^{\pi}(s) = E_{\pi} \left\{ \sum_{k=0}^{\infty} \gamma^{k} r_{t+k+1} | s_{t} = s \right\}$$
$$= E_{\pi} \left\{ r_{t+1} + \sum_{k=1}^{\infty} \gamma^{k} r_{t+k+1} | s_{t} = s \right\}$$
$$= E_{\pi} \left\{ r_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^{k} r_{t+k+2} | s_{t} = s \right\}$$
$$= E_{\pi} \left\{ r_{t+1} + \gamma V^{\pi}(s_{t+1}) | s_{t} = s \right\}$$

Allows us to write the optimal policy in terms of the optimal value function

$$\pi^*(s) = \arg\max_{a} E_{\pi} \{ r_{t+1} + \gamma V^*(s_{t+1}) | s_t = s, a_t = a \}$$

The reinforcement learning problem

Solution techniques

Sample complexity

▲ロト ▲ □ ト ▲ □ ト ▲ □ ト ● ● の Q ()

Hands-on

Recursive definition

$$V^{\pi}(s) = E_{\pi} \left\{ \sum_{k=0}^{\infty} \gamma^{k} r_{t+k+1} | s_{t} = s \right\}$$
$$= E_{\pi} \left\{ r_{t+1} + \sum_{k=1}^{\infty} \gamma^{k} r_{t+k+1} | s_{t} = s \right\}$$
$$= E_{\pi} \left\{ r_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^{k} r_{t+k+2} | s_{t} = s \right\}$$
$$= E_{\pi} \left\{ r_{t+1} + \gamma V^{\pi}(s_{t+1}) | s_{t} = s \right\}$$

Allows us to write the optimal policy in terms of the optimal value function

$$\pi^*(s) = \operatorname*{arg\,max}_{a} E_{\pi} \{ r_{t+1} + \gamma V^*(s_{t+1}) | s_t = s, a_t = a \}$$
The reinforcement learning problem

Solution techniques

Sample complexity

▲ロト ▲ □ ト ▲ □ ト ▲ □ ト ● ● の Q ()

Hands-on

Recursive definition

$$V^{\pi}(s) = E_{\pi} \left\{ \sum_{k=0}^{\infty} \gamma^{k} r_{t+k+1} | s_{t} = s \right\}$$
$$= E_{\pi} \left\{ r_{t+1} + \sum_{k=1}^{\infty} \gamma^{k} r_{t+k+1} | s_{t} = s \right\}$$
$$= E_{\pi} \left\{ r_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^{k} r_{t+k+2} | s_{t} = s \right\}$$
$$= E_{\pi} \left\{ r_{t+1} + \gamma V^{\pi}(s_{t+1}) | s_{t} = s \right\}$$

Allows us to write the optimal policy in terms of the optimal value function

$$\pi^*(s) = \arg\max_{a} E_{\pi} \{ r_{t+1} + \gamma V^*(s_{t+1}) | s_t = s, a_t = a \}$$

The reinforcement learning problem

V

Solution techniques

Sample complexity

Hands-on

Recursive definition

$$egin{aligned} \pi(s) &= E_{\pi} \left\{ \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | s_t = s
ight\} \ &= E_{\pi} \left\{ r_{t+1} + \sum_{k=1}^{\infty} \gamma^k r_{t+k+1} | s_t = s
ight\} \ &= E_{\pi} \left\{ r_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} | s_t = s
ight\} \ &= E_{\pi} \left\{ r_{t+1} + \gamma V^{\pi}(s_{t+1}) | s_t = s
ight\} \end{aligned}$$

Allows us to write the optimal policy in terms of the optimal value function

$$\pi^*(s) = \operatorname*{arg\,max}_{a} E_{\pi} \{ r_{t+1} + \gamma V^*(s_{t+1}) | s_t = s, a_t = a \}$$

▲□▶▲圖▶▲≣▶▲≣▶ ≣ のQ@

The reinforcement learning problem

Solution techniques

Sample complexity

Hands-on

Recursive definition

$$egin{aligned} V^{\pi}(s) &= E_{\pi} \left\{ \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | s_t = s
ight\} \ &= E_{\pi} \left\{ r_{t+1} + \sum_{k=1}^{\infty} \gamma^k r_{t+k+1} | s_t = s
ight\} \ &= E_{\pi} \left\{ r_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} | s_t = s
ight\} \ &= E_{\pi} \left\{ r_{t+1} + \gamma V^{\pi}(s_{t+1}) | s_t = s
ight\} \end{aligned}$$

Allows us to write the optimal policy in terms of the optimal value function

$$\pi^*(s) = \arg\max_{a} E_{\pi} \{ r_{t+1} + \gamma V^*(s_{t+1}) | s_t = s, a_t = a \}$$

The reinforcement learning problem

Solution techniques

Sample complexity

Hands-on

State-action value functions

Stores expected return for each state-action combination

$$egin{aligned} \mathcal{Q}^{\pi}(s,a) &= \mathcal{E}_{\pi} \left\{ \mathcal{R}_{t} | s_{t} = s, a_{t} = a
ight\} \ &= \mathcal{E}_{\pi} \left\{ \sum_{k=0}^{\infty} \gamma^{k} r_{t+k+1} | s_{t} = s, a_{t} = a
ight\} \ &= \mathcal{E}_{\pi} \left\{ r_{t+1} + \gamma \mathcal{Q}^{\pi}(s_{t+1}, a_{t+1}) | s_{t} = s, a_{t} = a
ight\} \end{aligned}$$

Allows us to find optimal actions without calculating expectation

$$Q^*(s, a) = \max_{\pi} Q^{\pi}(s, a)$$
$$\pi^*(s) = \arg\max_{a} Q^*(s, a)$$

▲ロト▲舂▶▲目▶▲目▶ 目 のへで

The reinforcement learning problem

Solution techniques

Sample complexity

▲□▶▲□▶▲□▶▲□▶ □ のQで

Hands-on

State-action value functions

Stores expected return for each state-action combination

$$Q^{\pi}(s, a) = E_{\pi} \{ R_t | s_t = s, a_t = a \}$$

= $E_{\pi} \left\{ \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | s_t = s, a_t = a \right\}$
= $E_{\pi} \{ r_{t+1} + \gamma Q^{\pi}(s_{t+1}, a_{t+1}) | s_t = s, a_t = a \}$

Allows us to find optimal actions without calculating expectation

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

 $\pi^*(s) = rgmax_a Q^*(s,a)$

Introduction
0000000000

Solution techniques

Sample complexity

Hands-on



Introduction
0000000000

Solution techniques

Sample complexity

Hands-on



Introduction
0000000000

Solution techniques

Sample complexity

Hands-on



Introduction
0000000000

Solution techniques

Sample complexity

Hands-on



Introduction
0000000000

Solution techniques

Sample complexity

Hands-on



Introduction
0000000000

Solution techniques

Sample complexity

Hands-on



The reinforcement learning problem

Solution techniques

Sample complexity

Hands-on



The reinforcement learning problem

Solution techniques

Sample complexity

Hands-on



The reinforcement learning problem

Solution techniques

Sample complexity

Hands-on



The reinforcement learning problem

Solution techniques

Sample complexity

Hands-on



The reinforcement learning problem

Solution techniques

Sample complexity

Hands-on



The reinforcement learning problem

Solution techniques

Sample complexity

▲□▶▲□▶▲□▶▲□▶ □ のQで

Hands-on

Value function for pendulum swing-up



- System starts in down position, goal is to stabilize in the up position
- Can apply torque to the rotation axis

▲□▶▲□▶▲□▶▲□▶ □ のQで

Control

Can now evaluate value functions for arbitrary policies, but how do we find the optimal policy?

• Optimal policy is greedy with respect to optimal value function $\pi^*(s) = \arg \max_a Q^*(s, a)$

• Let
$$\pi'(s) = \operatorname{arg\,max}_a Q^{\pi}(s,a)$$
, then $Q^{\pi'} \geq Q^{\pi}$

Leads to iterative dynamic programming solutions

he reinforcement learning problen

Solution techniques

Sample complexity

▲□▶▲□▶▲□▶▲□▶ □ のQで

Hands-on

Outline

Introduction

2 The reinforcement learning problem

3 Solution techniques

- Dynamic programming
- Temporal difference learning
- Policy search
- Sample complexity

5 Hands-on

Introduction
0000000000

Solution techniques

Sample complexity

Hands-on

Overview



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ ─臣 ─のへ(?)

Solution techniques

Sample complexity

▲□▶▲□▶▲□▶▲□▶ □ のQで

Hands-on

Dynamic programming

Solves problems by breaking them down into simpler subproblems

• For MDPs, the solution to a future state is a subproblem to the solution for the current state

Follow the policy that is greedy with respect to the optimal value function, which satisfies the Bellman equation

$$egin{aligned} V^*(s) &= \max_a E_{s'} \left\{
ho(s,a,s') + \gamma V^*(s')
ight\} \ &= \max_a \sum_{s'}
ho(s'|s,a) \left(
ho(s,a,s') + \gamma V^*(s')
ight) \end{aligned}$$

Sample complexity

▲ロト ▲ □ ト ▲ □ ト ▲ □ ト ● ● の Q ()

Hands-on

Value iteration

Elementary solution method: iteratively approximate optimal value function

```
function VALUEITERATION
    t \leftarrow 0
    V_0(s) \leftarrow 0 for all s
    repeat
         for each s do
              V_{t+1}(s) \leftarrow \max_{a \sum s'} p(s'|s,a) (\rho(s,a,s') + \gamma V_t(s'))
         end for
         t \leftarrow t + 1
    until ||V_t - V_{t-1}|| < \varepsilon
    return V
end function
```

Bootstrapping method: updates estimates based on estimates

he reinforcement learning problen

Solution techniques

Sample complexity

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ ▲ □ ● ● ● ●

Hands-on

Ditching the model



- Observations only occur in trajectories
- Cannot evaluate different actions in the same state

he reinforcement learning problen

Solution techniques

Sample complexity

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ ▲ □ ● ● ● ●

Hands-on

Ditching the model



- Observations only occur in trajectories
- Cannot evaluate different actions in the same state

ne reinforcement learning problem

Solution techniques

Sample complexity

(日)

Hands-on

Ditching the model



- Observations only occur in trajectories
- Cannot evaluate different actions in the same state

The reinforcement learning problem

Solution techniques

Sample complexity

▲□▶▲□▶▲□▶▲□▶ □ のQで

Hands-on

Temporal difference learning

Value iteration requires known transition model p(s'|s, a). Temporal difference learning only samples it:

function TDLEARNING(π)

$$egin{array}{l} s \leftarrow s_0 \ V^\pi(s) \leftarrow 0 ext{ for all } s \end{array}$$

repeat

$$\begin{array}{l} \text{Sample } s' \sim p(s'|s,\pi(s)) \text{ and reward } r = \rho(s,a,s') \\ V^{\pi}(s) \leftarrow V^{\pi}(s) + \alpha \left(r + \gamma V^{\pi}(s') - V^{\pi}(s)\right) \\ s \leftarrow s' \end{array}$$

until convergence

return V

end function

lpha is a learning rate that acts as an exponential moving average filter

The reinforcement learning problem

Solution techniques

Sample complexity

▲□▶▲□▶▲□▶▲□▶ □ のQで

Hands-on

Temporal difference learning

Value iteration requires known transition model p(s'|s, a). Temporal difference learning only samples it:

function TDLEARNING(π)

$$s \leftarrow s_0
onumber V^\pi(s) \leftarrow 0$$
 for all s

repeat

Sample $s' \sim p(s'|s, \pi(s))$ and reward $r = \rho(s, a, s')$ $V^{\pi}(s) \leftarrow V^{\pi}(s) + \alpha (r + \gamma V^{\pi}(s') - V^{\pi}(s))$ $s \leftarrow s'$ Temporal difference error

until convergence

return V

end function

lpha is a learning rate that acts as an exponential moving average filter

The reinforcement learning problem

Solution techniques

Sample complexity

Hands-on

Q-learning: Temporal difference control

TD learning estimates state-value function $V^{\pi}(s)$. For control, we must estimate state-action value function Q(s, a) instead. Q-learning estimates the optimal state-action value function $Q^*(s, a)$:

function QLEARNING

```
s \leftarrow s_0
    Q(s, a) \leftarrow 0 for all s, a
    repeat
         a \leftarrow \arg \max_{a} Q(s, a)
         Sample s' \sim p(s'|s, a) and reward r = p(s, a, s')
         Q(s,a) \leftarrow Q(s,a) + \alpha (r + \gamma \max_{a'} Q(s',a') - Q(s,a))
         s \leftarrow s'
    until convergence
    return Q
end function
```

Introduction
0000000000

.

Exploration

To guarantee convergence, every state has to be sampled an infinite number of times.

- In dynamic programming, this is guaranteed by sweeping the entire state space.
- When only trajectory sampling is available, requires exploration
- Use a stochastic exploration policy derived from Q(s, a), e.g. (ε-greedy)

$$\pi(a|s) = \left\{egin{array}{cc} 1 - arepsilon + rac{arepsilon}{|A|} & ext{if} a = ext{arg} \max_{a'} Q(s,a') \ rac{arepsilon}{|A|} & ext{otherwise} \end{array}
ight.$$

Exploration-exploitation trade-off

Choose between exploiting current knowledge (refining current known best policy) or exploring new knowledge (finding better policies)

-

Introduction
0000000000

Exploration

To guarantee convergence, every state has to be sampled an infinite number of times.

- In dynamic programming, this is guaranteed by sweeping the entire state space.
- When only trajectory sampling is available, requires exploration
- Use a stochastic exploration policy derived from Q(s, a), e.g. (ε-greedy)

$$\pi(a|s) = \left\{egin{array}{cc} 1 - arepsilon + rac{arepsilon}{|A|} & ext{if} a = ext{arg} \max_{a'} Q(s,a') \ rac{arepsilon}{|A|} & ext{otherwise} \end{array}
ight.$$

Exploration-exploitation trade-off

Choose between exploiting current knowledge (refining current known best policy) or exploring new knowledge (finding better policies)

The reinforcement learning problem

Solution techniques

Sample complexity

▲□▶▲□▶▲□▶▲□▶ □ のQで

Hands-on

SARSA: on-policy TD control

Q-learning finds the value function of the optimal policy π^* while following a stochastic exploratory policy π . SARSA finds the value function of π itself:

$$Q(s,a) \leftarrow Q(s,a) + \alpha \left(r + \gamma Q(s',a') - Q(s,a)\right)$$

- Better convergence properties when using function approximation
- Safer during learning



Introduction
0000000000

Solution techniques

Sample complexity

▲ロト ▲ □ ト ▲ □ ト ▲ □ ト ● ● の Q ()

Hands-on

Batch techniques

Q-learning and SARSA treat each sample only once, but that is inefficient. We can instead use all samples $(s_i, a_i) \rightarrow (r_i, s'_i)$ simultaneously (fitted Q-iteration)



The reinforcement learning problem

Solution techniques

Sample complexity

▲□▶▲□▶▲□▶▲□▶ □ のQで

Hands-on

Fitted Q iteration

function FITTEDQITERATION $k \leftarrow 0, \hat{Q}_k \leftarrow 0$ repeat for each sample *i* do $\pi(s_i') = \operatorname{arg\,max}_a \hat{Q}_k(s_i', a)$ $\hat{Q}_{k+1}(s_i, a_i) \leftarrow r_i + \gamma \hat{Q}_k(s'_i, \pi(s'_i))$ end for $k \leftarrow k+1$ until $\|\hat{Q}_k - \hat{Q}_{k-1}\| < \varepsilon$ return \hat{Q}_k end function

 \hat{Q} can be approximated by any supervised learning algorithm.

The reinforcement learning problem

Solution techniques

Sample complexity

Hands-on

Least squares approximation

We know that for every sample $(s_i, a_i) o (r_i, s_i')$, \hat{Q}^{π} has to satisfy

$$\hat{Q}^{\pi}(s_i, a_i) = r_i + \gamma \hat{Q}^{\pi}(s'_i, \pi(s'_i))$$

Given features $\phi(s, a)$ and a linear function approximator $\hat{Q}(s, a) = \phi(s, a)^T \theta$ with parameters θ , every sample *i* induces a constraint

$$\phi(s_i, a_i)^T \theta = r_i + \gamma \phi(s'_i, \pi(s'_i))^T \theta, \text{such that}$$
$$(\phi(s_i, a_i)^T - \gamma \phi(s'_i, \pi(s'_i))^T) \theta = r_i$$

In this case, we can estimate \hat{Q}^{π} directly using least squares fitting, although we still need to iterate $\pi(s) \leftarrow \max_a \phi(s, a)^T \theta$ to get \hat{Q}^* .

◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ ◆ ○ ◆ ○ ◆ ○ ◆ ○ ◆

The reinforcement learning problem

Solution techniques

Sample complexity

◆□▶ ◆□▶ ◆□▶ ◆□▶ → □ ◆○◆

Hands-on

Least squares approximation

We know that for every sample $(s_i, a_i) \rightarrow (r_i, s'_i)$, \hat{Q}^{π} has to satisfy

$$\hat{Q}^{\pi}(s_i,a_i) = r_i + \gamma \hat{Q}^{\pi}(s'_i,\pi(s'_i))$$

Given features $\phi(s, a)$ and a linear function approximator $\hat{Q}(s, a) = \phi(s, a)^T \theta$ with parameters θ , every sample *i* induces a constraint

$$\phi(s_i, a_i)^T \theta = r_i + \gamma \phi(s'_i, \pi(s'_i))^T \theta$$
, such that
 $(\phi(s_i, a_i)^T - \gamma \phi(s'_i, \pi(s'_i))^T) \theta = r_i$

In this case, we can estimate \hat{Q}^{π} directly using least squares fitting, although we still need to iterate $\pi(s) \leftarrow \max_a \phi(s, a)^T \theta$ to get \hat{Q}^* .

The reinforcement learning problem

Solution techniques

Sample complexity

Hands-on

Least squares approximation

We know that for every sample $(s_i, a_i)
ightarrow (r_i, s_i')$, \hat{Q}^{π} has to satisfy

$$\hat{Q}^{\pi}(s_i,a_i) = r_i + \gamma \hat{Q}^{\pi}(s'_i,\pi(s'_i))$$

Given features $\phi(s, a)$ and a linear function approximator $\hat{Q}(s, a) = \phi(s, a)^T \theta$ with parameters θ , every sample *i* induces a constraint

$$\phi(s_i, a_i)^T \theta = r_i + \gamma \phi(s'_i, \pi(s'_i))^T \theta$$
, such that
 $(\phi(s_i, a_i)^T - \gamma \phi(s'_i, \pi(s'_i))^T) \theta = r_i$

In this case, we can estimate \hat{Q}^{π} directly using least squares fitting, although we still need to iterate $\pi(s) \leftarrow \max_a \phi(s, a)^T \theta$ to get \hat{Q}^* .
he reinforcement learning problen

Solution techniques

Sample complexity

Hands-on

Direct policy search

Why use a value function?



Policy search works directly on the parameters θ of a parameterized policy $\pi(s; \theta)$.

he reinforcement learning problem

Solution techniques

Sample complexity

Hands-on

Direct policy search

Why use a value function?



Policy search works directly on the parameters θ of a parameterized policy $\pi(s; \theta)$.

Solution techniques

Sample complexity

Hands-on

Policy parameterization

What form should the policy take?

- Low-level
 - PD?
 - Dynamic movement primitives?
 - Radial basis functions?
 - Neural network?
- High-level
 - Via point placement?
 - Center of mass trajectory?
 - etc.



Solution techniques

Sample complexity

Hands-on

Policy parameterization

What form should the policy take?

- Low-level
 - PD?
 - Dynamic movement primitives?
 - Radial basis functions?
 - Neural network?
- High-level
 - Via point placement?
 - Center of mass trajectory?
 - etc.



Solution techniques

Sample complexity

Hands-on

Policy parameterization

What form should the policy take?

- Low-level
 - PD?
 - Dynamic movement primitives?
 - Radial basis functions?
 - Neural network?
- High-level
 - Via point placement?
 - Center of mass trajectory?
 - etc.



Solution techniques

Sample complexity

Hands-on

Policy parameterization

What form should the policy take?

- Low-level
 - PD?
 - Dynamic movement primitives?
 - Radial basis functions?
 - Neural network?
- High-level
 - Via point placement?
 - Center of mass trajectory?
 - etc.



Solution techniques

Sample complexity

Hands-on

Policy parameterization

What form should the policy take?

- Low-level
 - PD?
 - Dynamic movement primitives?
 - Radial basis functions?
 - Neural network?
- High-level
 - Via point placement?
 - Center of mass trajectory?
 - etc.



Solution techniques

Sample complexity

Hands-on

Policy parameterization

What form should the policy take?

- Low-level
 - PD?
 - Dynamic movement primitives?
 - Radial basis functions?
 - Neural network?
- High-level
 - Via point placement?
 - Center of mass trajectory?
 - etc.



The reinforcement learning problem

Solution techniques

Sample complexity

▲□▶▲□▶▲□▶▲□▶ □ のQで

Hands-on

Black-box Pl²

Black-box optimization routine, competitive with state-of-the-art Markov-based schemes.

function BLACKBOXPI²

 $g \leftarrow 0$ Initialize $heta_0$ randomly

repeat

for k in K do $\varepsilon_k = \mathcal{N}(0, \Sigma), R_k = \text{rollout}(\theta_a + \varepsilon_k)$ end for for k in K do $\mathbf{W}_{k} = \frac{e^{\frac{1}{K}R_{k}}}{\sum_{k=1}^{K}e^{\frac{1}{K}R_{k}}}$ end for $\theta_{a+1} = \theta_a + \sum_{k=1}^{K} w_k (\theta_a + \varepsilon_k)$ $g \leftarrow g + 1$ until converged end function

The reinforcement learning problem

Solution techniques

Sample complexity

▲□▶▲□▶▲□▶▲□▶ □ のQで

Hands-on

Policy search versus TD learning

Advantages

- Smoother policies
- Parameter space could be smoother
- Parameter space is generally smaller

Disadvantages

- Parameterization requires system knowledge
- Only locally optimal (may get stuck in local minima)
- Rollouts are noisy

he reinforcement learning problem

Solution techniques

Sample complexity

▲□▶▲□▶▲□▶▲□▶ □ のQで

Hands-on

Outline



2 The reinforcement learning problem

3 Solution techniques

4

Sample complexity

- Eligibility traces
- Value function approximation
- Indirect reinforcement learning
- Local bias



Introduction	The reinforcement learning problem	Solution techniques	Sample complexity	Hands-on 00000000
Sample	complexity			

Sample complexity is the number of samples required to learn the task (interaction time).



Does not include computation time (computational complexity)

イロト 不得 とうほう イヨン

3

ne reinforcement learning problem

Solution techniques

Sample complexity

▲ロト ▲ □ ト ▲ □ ト ▲ □ ト ● ● の Q ()

Hands-on

Eligibility traces

Standard TD methods update only the current state.



ne reinforcement learning problem

Solution techniques

Sample complexity

▲ロト ▲ □ ト ▲ □ ト ▲ □ ト ● ● の Q ()

Hands-on

Eligibility traces

Standard TD methods update only the current state.



ne reinforcement learning problem

Solution techniques

Sample complexity

▲ロト ▲ □ ト ▲ □ ト ▲ □ ト ● ● の Q ()

Hands-on

Eligibility traces

Standard TD methods update only the current state.



he reinforcement learning problen

Solution techniques

Sample complexity

▲ロト ▲ □ ト ▲ □ ト ▲ □ ト ● ● の Q ()

Hands-on

Eligibility traces

Standard TD methods update only the current state.



ne reinforcement learning problem

Solution techniques

Sample complexity

▲ロト ▲ □ ト ▲ □ ト ▲ □ ト ● ● の Q ()

Hands-on

Eligibility traces

Standard TD methods update only the current state.



ne reinforcement learning problem

Solution techniques

Sample complexity

▲ロト ▲ □ ト ▲ □ ト ▲ □ ト ● ● の Q ()

Hands-on

Eligibility traces

Standard TD methods update only the current state.



he reinforcement learning problen

Solution techniques

Sample complexity

▲ロト ▲ □ ト ▲ □ ト ▲ □ ト ● ● の Q ()

Hands-on

Eligibility traces

Standard TD methods update only the current state.



Introduction
0000000000

Solution techniques

Sample complexity

▲□▶▲□▶▲□▶▲□▶ □ のQで

Hands-on

Eligibility traces

Standard TD methods update only the current state.

Introduction
0000000000

Solution techniques

Sample complexity

▲□▶▲□▶▲□▶▲□▶ □ のQで

Hands-on

Eligibility traces

Standard TD methods update only the current state.

Introduction
0000000000

Solution techniques

Sample complexity

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ ▲ □ ● ● ● ●

Hands-on

Eligibility traces

Standard TD methods update only the current state.



Introduction
0000000000

Solution techniques

Sample complexity

▲ロト ▲ □ ト ▲ □ ト ▲ □ ト ● ● の Q ()

Hands-on

Eligibility traces

Standard TD methods update only the current state.





Introduction
0000000000

Solution techniques

Sample complexity

▲ロト ▲ □ ト ▲ □ ト ▲ □ ト ● ● の Q ()

Hands-on

Eligibility traces

Standard TD methods update only the current state.





Introduction
0000000000

Solution techniques

Sample complexity

▲ロト ▲ □ ト ▲ □ ト ▲ □ ト ● ● の Q ()

Hands-on

Eligibility traces

Standard TD methods update only the current state.





Solution techniques

Sample complexity

▲□▶▲□▶▲□▶▲□▶ □ のQで

Hands-on

Value function approximation

Robots work in continuous state spaces. How to discretize the value function?

- LEO: 12-dimensional state-action space. Assume each discretized in N steps → N¹² values!
- Say N = 10, float values \rightarrow 4 TB memory to be stored and filled.

Generalization

To be sample-efficient, experience must be generalized. Generalization always has a bias-variance trade-off.

Solution techniques

Sample complexity

▲□▶▲□▶▲□▶▲□▶ □ のQで

Hands-on

Value function approximation

Robots work in continuous state spaces. How to discretize the value function?

- LEO: 12-dimensional state-action space. Assume each discretized in N steps → N¹² values!
- Say N = 10, float values \rightarrow 4 TB memory to be stored and filled.

Generalization

To be sample-efficient, experience must be generalized. Generalization always has a bias-variance trade-off.

Introduction
0000000000

Solution techniques

Sample complexity

Hands-on

Tile coding

Linear function approximator using overlapping grids.



Convergence guarantees due to linearity $\hat{Q}(s) = \phi(s, a)^T \theta$.

・ロト・日本・日本・日本・日本・日本

ne reinforcement learning problem

Solution techniques

Sample complexity

▲ロト ▲ □ ト ▲ □ ト ▲ □ ト ● ● の Q ()

Hands-on

Neural networks



- Less parameters (faster generalization)
- Updates have global, non-linear effects (no convergence guarantees, *forgetting*)

Introduction
0000000000

Solution techniques

Sample complexity

イロト 不得 とうほう イヨン

Hands-on





- Local changes
- Forests provide finer granularity and variance estimation.

he reinforcement learning problem

Solution techniques

Sample complexity

▲□▶▲□▶▲□▶▲□▶ □ のQで

Hands-on

Indirect reinforcement learning

Direct reinforcement learning solves MDPs by interacting directly with the environment



Indirect reinforcement learning learns a model as an intermediary step.

he reinforcement learning problem

Solution techniques

Sample complexity

▲□▶▲□▶▲□▶▲□▶ □ のQで

Hands-on

Indirect reinforcement learning

Direct reinforcement learning solves MDPs by interacting directly with the environment



Indirect reinforcement learning learns a model as an intermediary step.

Solution techniques

Sample complexity

Hands-on

Model approximation

Generalizes in a different space than value function approximation. It approximates the transition function $p(s'|s, a) = \mathcal{T}(s, a, s')$. Since in a Markov process this is static, it is a supervised learning problem.

- Allows simulating additional transitions (mental rehearsal)
- Provides an estimate of $\frac{\partial s}{\partial a}$ for better policy updates
- Allows derivation of π from state-value function V instead of state-action value function Q.

Model learning versus system identification

Model learning is performed on-line and therefore has an exploration-exploitation trade-off. Model fidelity should be high along the trajectory of the optimal policy.

Solution techniques

Sample complexity

Hands-on

Model approximation

Generalizes in a different space than value function approximation. It approximates the transition function $p(s'|s, a) = \mathcal{T}(s, a, s')$. Since in a Markov process this is static, it is a supervised learning problem.

- Allows simulating additional transitions (mental rehearsal)
- Provides an estimate of $\frac{\partial s}{\partial a}$ for better policy updates
- Allows derivation of π from state-value function V instead of state-action value function Q.

Model learning versus system identification

Model learning is performed on-line and therefore has an exploration-exploitation trade-off. Model fidelity should be high along the trajectory of the optimal policy.

Solution techniques

Sample complexity

Hands-on

Model approximation

Generalizes in a different space than value function approximation. It approximates the transition function $p(s'|s, a) = \mathcal{T}(s, a, s')$. Since in a Markov process this is static, it is a supervised learning problem.

- Allows simulating additional transitions (mental rehearsal)
- Provides an estimate of $\frac{\partial s}{\partial a}$ for better policy updates
- Allows derivation of π from state-value function V instead of state-action value function Q.

Model learning versus system identification

Model learning is performed on-line and therefore has an exploration-exploitation trade-off. Model fidelity should be high along the trajectory of the optimal policy.

he reinforcement learning problem

Solution techniques

Sample complexity

Hands-on 00000000

Locally weighted regression

Very successful model approximator for robotics

- Find k nearest neighbors
 - Approximate nearest neighbor search
- Weigh according to distance

$$w(p) = e^{-\left(\frac{|p-q|_2}{h}\right)^2}$$

where *h* is the distance to the *k*th nearest neighbor

• Fit linear model using least-squares regression


Introduction

he reinforcement learning problen

Solution techniques

Sample complexity

▲□▶▲□▶▲□▶▲□▶ □ のQで

Hands-on







- Perform K simulated updates per control step
- Mix of direct and indirect reinforcement learning
- Simulated updates can be done anywhere (prioritized sweeping)



Model-learning policy search. By approximating the system dynamics with a gaussian process model, can analytically calculate $\frac{dR}{d\theta}$.



▲□▶▲□▶▲□▶▲□▶ □ のQで

[Deisenroth & Rasmussen, 2011]

Introduction	The reinforcement learning problem	Solution techniques	Sample complexity	Hands-on
PII CO				

Model-learning policy search. By approximating the system dynamics with a gaussian process model, can analytically calculate $\frac{dR}{d\theta}$.



[Deisenroth & Rasmussen, 2011]

<ロト < 同ト < 回ト < 回ト = 三日 = 三日

Introduction

he reinforcement learning problem

Solution techniques

Sample complexity

▲ロト ▲ □ ト ▲ □ ト ▲ □ ト ● ● の Q ()

Hands-on



When controlling a system, we always want to find the best action for the current state.



Therefore it makes sense to concentrate planning around the current state

Introduction

he reinforcement learning problem

Solution techniques

Sample complexity

▲ロト ▲ □ ト ▲ □ ト ▲ □ ト ● ● の Q ()

Hands-on



When controlling a system, we always want to find the best action for the current state.



Therefore it makes sense to concentrate planning around the current state

The reinforcement learning problem

Solution techniques

Sample complexity

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ ▲ □ ● ● ● ●

Hands-on

Some algorithms that use local bias

	RTDP	TEXPLORE	DYNA-2	
Model	Given	Learned	Learned	
Planning	Best-first search	Monte-carlo	TD learning	
		tree search		
		(UCT)		
Value	Updated from	Updated from	Separate for	
function	experience +	model	experience +	
	model		model	
Control	From value	From UCT	From weighted	
	function		sum of value	
			functions	

▲□▶▲□▶▲□▶▲□▶ □ のQで

Recap

- Reinforcement learning is model-free optimal control
- Value-based RL is sampled dynamic programming
 - First estimates a value function, after which the optimal policy is one-step greedy
- Policy search directly searches in the space of parameterized policies, can solve larger problems, but is only locally optimal
 - Requires suitable initialization
- Many approaches exist to make RL practical on real systems

Sample complexity

▲□▶▲□▶▲□▶▲□▶ □ のQで

Hands-on

Hands-on

Experiment with parameter settings for reinforcement learning for the pendulum swing-up problem

Manual http://wouter.caarls.org/files/kss2014_tutorial.pdf

Matlab toolbox http://wouter.caarls.org/files/kss2014_tutorial.tgz

Run pendgui from the toolbox directory

he reinforcement learning problen

Solution techniques

Sample complexity

Hands-on

Q1: Learning rate α





 $\alpha = 0.2$

 $\alpha = 0.7$

ヘロト 人間 とく ヨン く ヨン

э.

he reinforcement learning problem

Solution techniques

Sample complexity

Hands-on o●oooooc

Q2: On-policy vs off-policy learning





On-policy

Off-policy

イロト 不得 とうほう イヨン

he reinforcement learning problem

Solution techniques

Sample complexity

Hands-on ○●○○○○○○

Q2: On-policy vs off-policy learning

 $\epsilon = 0.01$



On-policy

Off-policy

イロト 不得 とうほう イヨン

he reinforcement learning problen

Solution techniques

Sample complexity

Hands-on

Q3: State space resolution





On-policy

Off-policy

イロト 不得 とうほう イヨン

he reinforcement learning problen

Solution techniques

Sample complexity

Hands-on

Q3: State space resolution





On-policy

Off-policy

・ロト ・ 何 ト ・ ヨ ト ・ ヨ ト

he reinforcement learning problen

Solution techniques

Sample complexity

Hands-on

Q4: Discount rate γ



$$\gamma = 0.97$$

 $\gamma = 0.87$

イロト イ理ト イヨト イヨト

he reinforcement learning problem

Solution techniques

Sample complexity

Hands-on

Q5: Discount rate γ (2)





$$\gamma = 0.87$$

 $\gamma = 0.57$

・ロト ・聞ト ・ヨト ・ヨト

he reinforcement learning problen

Solution techniques

Sample complexity

Hands-on

Q6: Reward function





Path rewards

Goal rewards

▲□▶▲□▶▲□▶▲□▶ □ のQで

The reinforcement learning problem

Solution techniques

Sample complexity

Hands-on

Q7: On-policy vs off-policy learning (2)

γ = 0.57, goal rewards, ε = 0.45



On-policy

Off-policy

イロト 不得 とうほう 不良 とう

3

he reinforcement learning problen

Solution techniques

Sample complexity

Hands-on

Q8: Initial value

Episodes = 2000



Initial value 0

Initial value -500

イロト 不得 とうほう イヨン

э.