

Reinforcement Learning

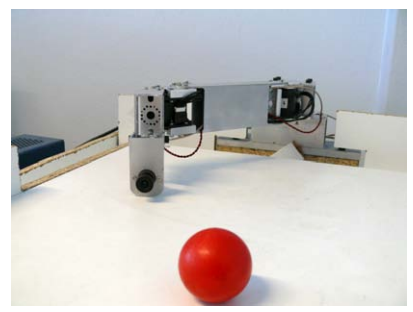
Part I: The Classical Setting

Lucian Buşoniu, Jelmer van Ast, Robert Babuška

Knowledge-Based Control Systems
2010-03-01

Demo: RL for a robot goalkeeper

Learn how to catch ball, using video camera image



Outline

- 1 Reinforcement learning basics
- 2 Algorithms
- 3 Accelerating RL

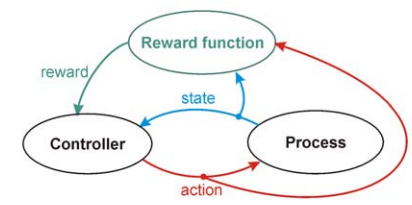
- 1 Reinforcement learning basics
 - Introduction
 - Elements of RL
 - RL solution
- 2 Algorithms
- 3 Accelerating RL

Why learning?

- Learning** can find solutions that:
- 1 cannot be found in advance
 - problem too complex (e.g., controlling highly nonlinear systems)
 - problem not fully known beforehand (e.g., robotic exploration of extraterrestrial planets)
 - 2 steadily improve
 - 3 adapt to time-varying environments

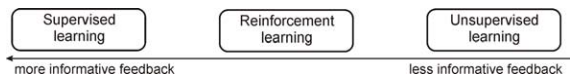
Essential for any **intelligent** system

Principle of RL



- Interact with a system through **states** and **actions**
- Receive **rewards** as performance feedback
- Inspired by human and animal learning

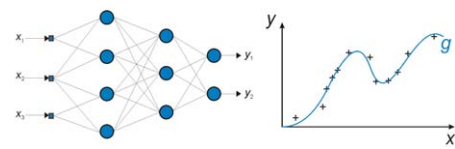
RL on the Machine Learning spectrum



Spectrum: Supervised learning



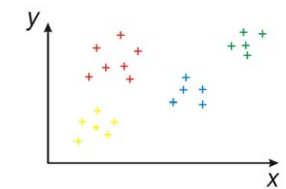
- For each input sample x , **correct output** y is known
- Infer input-output relationship $y \approx g(x)$
- Example: **neural networks**



Spectrum: Unsupervised learning



- Only input samples x available – **no outputs**
- Find patterns in the data
- Example: **clustering**



Reinforcement learning basics Algorithms Accelerating RL

Introduction

Spectrum: Reinforcement learning

- Correct outputs not available, **only rewards**
- Find optimal control behavior

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

Introduction

Reinforcement learning = Control

Reinforcement learning is about **control**:
optimal, adaptive, and model-free

This presentation: **classical RL** – discrete states and actions

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

Elements of RL

- Reinforcement learning basics
 - Introduction
 - Elements of RL
 - RL solution
- Algorithms
- Accelerating RL

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

Elements of RL

A simple cleaning robot example

- Cleaning robot in a 1-D world
- Either pick up trash (reward +5) or power pack (reward +1)
- After picking up item, episode terminates

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

Elements of RL

Cleaning robot: State & action

- Robot in given **state** x (cell)
- and takes **action** u (e.g., move right)

- State space** $X = \{0, 1, 2, 3, 4, 5\}$
- Action space** $U = \{-1, 1\} = \{\text{left, right}\}$

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

Elements of RL

Cleaning robot: Transition & reward

- Robot reaches **next state** x'
- and receives **reward** r = quality of transition (here, +5 for collecting trash)

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

Elements of RL

Cleaning robot: Transition & reward functions

- Transition function** (process behavior):

$$x' = f(x, u) = \begin{cases} x & \text{if } x \text{ is terminal (0 or 5)} \\ x + u & \text{otherwise} \end{cases}$$
- Reward function** (immediate performance):

$$r = \rho(x, u) = \begin{cases} 1 & \text{if } x = 1 \text{ and } u = -1 \text{ (powerpack)} \\ 5 & \text{if } x = 4 \text{ and } u = 1 \text{ (trash)} \\ 0 & \text{otherwise} \end{cases}$$
- Note**: Terminal states cannot be left & do not accumulate rewards!

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

Elements of RL

Markov decision process

- State space X
- Action space U
- Transition function $x' = f(x, u)$
- Reward function $r = \rho(x, u)$

... form a **Markov decision process**

Note: stochastic formulation possible

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

Elements of RL

Policy

- Policy** h : mapping from x to u (state feedback)
- Determines controller behavior

Example: $h(0) = *$ (terminal state, action is irrelevant), $h(1) = -1$, $h(2) = 1$, $h(3) = 1$, $h(4) = 1$, $h(5) = *$

TU Delft

- 1 Reinforcement learning basics
 - Introduction
 - Elements of RL
 - RL solution
- 2 Algorithms
- 3 Accelerating RL

Learning goal

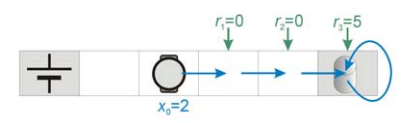
Find h that maximizes **discounted return**:

$$R^h(x_0) = \sum_{k=0}^{\infty} \gamma^k r_{k+1} = \sum_{k=0}^{\infty} \gamma^k \rho(x_k, h(x_k))$$

from any x_0

- Discount factor $\gamma \in [0, 1)$:
- induces a "pseudo-horizon" for optimization
 - bounds infinite sum
 - encodes increasing uncertainty about the future
 - helps convergence of algorithms

Cleaning robot: Return



Assume h always goes right

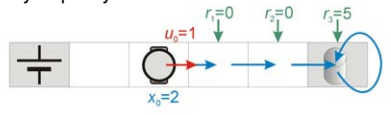
$$R^h(2) = \gamma^0 r_1 + \gamma^1 r_2 + \gamma^2 r_3 + \gamma^3 0 + \gamma^4 0 + \dots = \gamma^2 \cdot 5$$

Because x_3 is terminal, all remaining rewards are 0

Q-function

- **Q-function** of policy h :

$$Q^h(x_0, u_0) = \rho(x_0, u_0) + \gamma R^h(x_1)$$
 (return after taking u_0 in x_0 and then following h)
- Simply fix the first action in the sequence, independently of policy



- Why Q-function? Useful to choose actions (later)

Bellman equation

- Develop Q-function one step ahead:

$$Q^h(x_0, u_0) = \rho(x_0, u_0) + \gamma R^h(x_1)$$

$$= \rho(x_0, u_0) + \gamma [\rho(x_1, h(x_1)) + \gamma R^h(x_2)]$$

$$= \rho(x_0, u_0) + \gamma Q^h(x_1, h(x_1))$$
 Also, $x_1 = f(x_0, u_0)$

⇒ **Bellman equation for Q^h**

$$Q^h(x, u) = \rho(x, u) + \gamma Q^h(f(x, u), h(f(x, u)))$$

Optimal solution

- **Optimal Q-function:**

$$Q^* = \max_h Q^h$$
- ⇒ Greedy policy in Q^* :

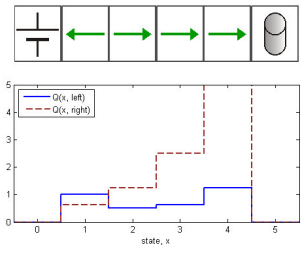
$$h^*(x) = \arg \max_u Q^*(x, u)$$
 is **optimal** (achieves maximal returns)

Bellman optimality equation (for Q^*)

$$Q^*(x, u) = \rho(x, u) + \gamma \max_{u'} Q^*(f(x, u), u')$$

Cleaning robot: Optimal solution

Discount factor $\gamma = 0.5$



- 1 Reinforcement learning basics
- 2 Algorithms
 - Taxonomy
 - Q-learning
 - SARSA
- 3 Accelerating RL

Types of algorithms

- By model knowledge
- 1 **Model-based** – dynamic programming
 f, ρ known
 - 2 **Model-free** – proper reinforcement learning
 f, ρ unknown, only transition data (x, u, x', r) available
 - 3 **Model-learning RL**
estimate f and ρ from transition data

Reinforcement learning basics Algorithms Accelerating RL

Taxonomy

Types of algorithms (cont'd)

By level of interaction

- Offline
data collected in advance
- Online
controller learns by interacting with the process

By path to optimal solution

- Off-policy
find Q^* , use it to compute h^*
- On-policy
find Q^h , improve h , repeat

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

Taxonomy

Algorithms in this lecture

Online model-free reinforcement learning:

Off-policy	On-policy
Q-learning	SARSA

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

Q-learning

- Reinforcement learning basics
- Algorithms
 - Taxonomy
 - Q-learning
 - SARSA
- Accelerating RL

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

Q-learning

Off-policy online RL: Q-learning

Recall off-policy: find Q^* , use it to compute h^*

- Take Bellman optimality equation at some (x, u) :

$$Q^*(x, u) = \rho(x, u) + \gamma \max_{u'} Q^*(f(x, u), u')$$
- Turn into **iterative update**:

$$Q(x, u) \leftarrow \rho(x, u) + \gamma \max_{u'} Q(f(x, u), u')$$
- Instead of model f, ρ , use **transition sample** $(x_k, u_k, x_{k+1}, r_{k+1})$ at each step k :

$$Q(x_k, u_k) \leftarrow r_{k+1} + \gamma \max_{u'} Q(x_{k+1}, u')$$
 Note: $x_{k+1} = f(x_k, u_k), r_{k+1} = \rho(x_k, u_k)$

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

Q-learning

Q-learning (cont'd)

- Instead of model, use **transition sample** $(x_k, u_k, x_{k+1}, r_{k+1})$ at each step k :

$$Q(x_k, u_k) \leftarrow r_{k+1} + \gamma \max_{u'} Q(x_{k+1}, u')$$
 Note: $x_{k+1} = f(x_k, u_k), r_{k+1} = \rho(x_k, u_k)$
- Finally, make update **incremental**:

$$Q(x_k, u_k) \leftarrow Q(x_k, u_k) + \alpha_k [r_{k+1} + \gamma \max_{u'} Q(x_{k+1}, u') - Q(x_k, u_k)]$$

$$\alpha_k \in (0, 1] \text{ learning rate}$$

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

Q-learning

Complete Q-learning algorithm

Q-learning

```

for every trial do
  initialize  $x_0$ 
  repeat for each step  $k$ 
    take action  $u_k$ 
    measure  $x_{k+1}$ , receive  $r_{k+1}$ 
     $Q(x_k, u_k) \leftarrow Q(x_k, u_k) + \alpha_k [r_{k+1} + \gamma \max_{u'} Q(x_{k+1}, u') - Q(x_k, u_k)]$ 
  until terminal state
end for
    
```

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

Q-learning

Exploration-exploitation tradeoff

- Essential condition for convergence to Q^* : all (x, u) pairs must be visited infinitely often
- ⇒ **Exploration** necessary: sometimes, choose actions randomly
- Exploitation** of current knowledge is also necessary: sometimes, choose actions greedily:

$$u_k = \arg \max_{\bar{u}} Q(x_k, \bar{u})$$

Exploration-exploitation tradeoff crucial for performance of online RL

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

Q-learning

Exploration-exploitation: ϵ -greedy strategy

- Simple solution: **ϵ -greedy**

$$u_k = \begin{cases} \arg \max_{\bar{u}} Q(x_k, \bar{u}) & \text{with probability } (1 - \epsilon_k) \\ \text{a random action} & \text{with probability } \epsilon_k \end{cases}$$
- Exploration probability $\epsilon_k \in (0, 1)$ is usually decreased over time

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

Q-learning

Cleaning robot: Q-learning demo

Parameters: $\alpha = 0.2, \epsilon = 0.3$ (constant)
 $x_0 = 2$ or 3 (randomly)

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

SARSA

- 1 Reinforcement learning basics
- 2 Algorithms
 - Taxonomy
 - Q-learning
 - SARSA
- 3 Accelerating RL

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

SARSA

On-policy online RL: SARSA

Recall on-policy: **find** Q^h , improve h , repeat

Similar to Q-learning:

- 1 Take Bellman equation for Q^h , at some (x, u) :

$$Q^h(x, u) = \rho(x, u) + \gamma Q^h(f(x, u), h(f(x, u)))$$
- 2 Turn into iterative update:

$$Q(x, u) \leftarrow \rho(x, u) + \gamma Q(f(x, u), h(f(x, u)))$$
- 3 Use sample $(x_k, u_k, r_{k+1}, x_{k+1}, u_{k+1})$ at each step k :

$$Q(x_k, u_k) \leftarrow r_{k+1} + \gamma Q(x_{k+1}, u_{k+1})$$

Note: $u_{k+1} = h(f(x_k, u_k))$, $h =$ policy being followed

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

SARSA

SARSA (cont'd)

- 4 Use sample $(x_k, u_k, r_{k+1}, x_{k+1}, u_{k+1})$ at each step k :

$$Q(x_k, u_k) \leftarrow r_{k+1} + \gamma Q(x_{k+1}, u_{k+1})$$

Note: $u_{k+1} = h(f(x_k, u_k))$, $h =$ policy being followed
- 4 Make update incremental:

$$Q(x_k, u_k) \leftarrow Q(x_k, u_k) + \alpha_k \cdot [r_{k+1} + \gamma Q(x_{k+1}, u_{k+1}) - Q(x_k, u_k)]$$

$(x_k, u_k, r_{k+1}, x_{k+1}, u_{k+1}) =$
(State, Action, Reward, State, Action) = SARSA

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

SARSA

Complete SARSA algorithm

SARSA

for every trial do

initialize x_0 , choose initial action u_0

repeat for each step k

apply u_k , measure x_{k+1} , receive r_{k+1}

choose **next** action u_{k+1}

$Q(x_k, u_k) \leftarrow Q(x_k, u_k) + \alpha_k \cdot [r_{k+1} + \gamma Q(x_{k+1}, u_{k+1}) - Q(x_k, u_k)]$

until terminal state

end for

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

SARSA

Exploration-exploitation in SARSA

- For convergence—besides infinite exploration—SARSA requires **policy to eventually become greedy**
- E.g., ϵ -greedy

$$u_k = \begin{cases} \arg \max_{\bar{u}} Q(x_k, \bar{u}) & \text{with probability } (1 - \epsilon_k) \\ \text{a random action} & \text{with probability } \epsilon_k \end{cases}$$

with $\lim_{k \rightarrow \infty} \epsilon_k = 0$
- Greedy actions \Rightarrow policy implicitly improved!
 (Recall on-policy: find Q^h , improve h , repeat)

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

SARSA

Cleaning robot: SARSA demo

Parameters like Q-learning: $\alpha = 0.2$, $\epsilon = 0.3$ (constant)
 $x_0 = 2$ or 3 (randomly)

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

- 1 Reinforcement learning basics
- 2 Algorithms
- 3 Accelerating RL
 - Eligibility traces
 - Experience replay

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

Accelerating RL

In practice, transition data costs:

- time
- profits (suboptimal performance due to exploration)
- process wear & tear

Fast RL = **use data efficiently**

(computational demands are secondary)

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

Eligibility traces

- Leave decaying **trace** along state-action trajectory:

- $\lambda \in [0, 1]$ decay rate, γ discount factor
- Implementation:

$e(x, u) \leftarrow 0$ for all x, u

for each step k **do**

$e(x, u) \leftarrow \lambda\gamma e(x, u)$ for all x, u

$e(x_k, u_k) \leftarrow 1$

end for

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

Eligibility traces

Q(λ)-learning

- Recall basic Q-learning only updates $Q(x_k, u_k)$:

$$Q(x_k, u_k) \leftarrow Q(x_k, u_k) + \alpha_k \cdot [r_{k+1} + \gamma \max_{u'} Q(x_{k+1}, u') - Q(x_k, u_k)]$$
- Q(λ)-learning updates **all eligible pairs**:

$$Q(x, u) \leftarrow Q(x, u) + \alpha_k \cdot e(x, u) \cdot [r_{k+1} + \gamma \max_{u'} Q(x_{k+1}, u') - Q(x_k, u_k)] \text{ for all } x, u$$
- Note: exploratory actions break causality
 \Rightarrow reset eligibility trace to 0

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

Eligibility traces

Complete Q(λ)-learning algorithm

Q(λ)-learning

```

for every trial do
  e(x, u) ← 0 for all x, u
  initialize x0
  repeat for each step k
    take action uk
    measure xk+1, receive rk+1
    if uk exploratory then e(x, u) ← 0 for all x, u
    else e(x, u) ← λγe(x, u) for all x, u
    end if
    e(xk, uk) ← 1
    Q(x, u) ← Q(x, u) + αk · e(x, u) · [rk+1 + γ maxu' Q(xk+1, u') - Q(xk, uk)] for all x, u
  until terminal state
end for
  
```

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

Eligibility traces

Cleaning robot: Q(λ)-learning demo

Parameters: $\alpha = 0.2$, $\varepsilon = 0.3$ (like basic Q-learning), $\lambda = 0.5$
 $x_0 = 2$ or 3 (randomly)

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

Eligibility traces

SARSA(λ)

Similar to Q-learning:

- Basic SARSA:

$$Q(x_k, u_k) \leftarrow Q(x_k, u_k) + \alpha_k \cdot [r_{k+1} + \gamma Q(x_{k+1}, u_{k+1}) - Q(x_k, u_k)]$$
- SARSA(λ)-learning:

$$Q(x, u) \leftarrow Q(x, u) + \alpha_k \cdot e(x, u) \cdot [r_{k+1} + \gamma Q(x_{k+1}, u_{k+1}) - Q(x_k, u_k)] \text{ for all } x, u$$
- SARSA on-policy, including exploration
 \Rightarrow exploratory actions not a problem

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

Eligibility traces

Complete SARSA(λ) algorithm

SARSA(λ)

```

for every trial do
  e(x, u) ← 0 for all x, u
  initialize x0, choose initial action u0
  repeat for each step k
    apply uk, measure xk+1, receive rk+1
    choose next action uk+1
    e(x, u) ← λγe(x, u) for all x, u
    e(xk, uk) ← 1
    Q(x, u) ← Q(x, u) + αk · e(x, u) · [rk+1 + γ Q(xk+1, uk+1) - Q(xk, uk)] for all x, u
  until terminal state
end for
  
```

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

Eligibility traces

Cleaning robot: SARSA(λ) demo

Parameters: $\alpha = 0.2$, $\varepsilon = 0.3$ (like basic SARSA), $\lambda = 0.5$
 $x_0 = 2$ or 3 (randomly)

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

Eligibility traces

Effects of eligibility trace

- Accelerates learning: fewer trials to convergence
- However: too large λ can make algorithm settle on suboptimal solution!

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

Experience replay

- Reinforcement learning basics
- Algorithms
- Accelerating RL
 - Eligibility traces
 - Experience replay

TU Delft

Reinforcement learning basics Algorithms Accelerating RL

Experience replay

Experience replay (ER)

- Store each transition sample $(x_k, u_k, x_{k+1}, r_{k+1})$ into a database
- At every step, **replay** N transitions from the database
- Improvement: replay most informative samples first: **prioritized sweeping**

TU Delft

Experience replay
ER Q-learning

Q-learning with experience replay

```

for every trial do
  initialize  $x_0$ 
  repeat for each step  $k$ 
    take action  $u_k$ 
    measure  $x_{k+1}$ , receive  $r_{k+1}$ 
     $Q(x_k, u_k) \leftarrow Q(x_k, u_k) + \alpha_k \cdot [r_{k+1} + \gamma \max_{u'} Q(x_{k+1}, u') - Q(x_k, u_k)]$ 
    add  $(x_k, u_k, x_{k+1}, r_{k+1})$  to database
    ReplayExperience
  until terminal state
end for
    
```

Experience replay
ER Q-learning (cont'd)

ReplayExperience

```

loop  $N$  times
  retrieve a sample  $(x, u, x', r)$  from database
   $Q(x, u) \leftarrow Q(x, u) + \alpha \cdot [r + \gamma \max_{u'} Q(x', u') - Q(x, u)]$ 
end loop
    
```

Summary and outlook
Summary

- **Reinforcement learning** = optimal, adaptive, model-free control
- Principle: reward signal as performance feedback
- Inspired from human and animal learning, but solid mathematical foundation
- Classical RL: small, discrete X and U (this presentation)

Summary and outlook
A final look at the algorithms

	Off-policy	On-policy
Basic RL	Q-learning Param: $\gamma, \alpha_k, \epsilon_k$	SARSA Param: $\gamma, \alpha_k, \epsilon_k$
RL with eligibility traces	$Q(\lambda)$ -learning Param: $\gamma, \alpha_k, \epsilon_k, \lambda$	SARSA(λ) Param: $\gamma, \alpha_k, \epsilon_k, \lambda$

Typical parameter values:

- γ 0.9 or larger
- α_k under 0.5 or diminishing schedule
- ϵ_k around 0.1 or diminishing schedule
- λ between 0.5 and 0.9

Summary and outlook
Outlook

- Other algorithms: actor-critic, model-learning, policy search, etc.
- Continuous X, U :
Part II – RL using function approximation
- State not fully measurable: "partially observable Markov decision process"
- RL for distributed (multi-agent) control

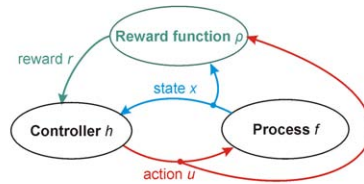
Reinforcement Learning

Part II: Approximate RL for Continuous-Space Control

Lucian Buşoniu, Jelmer van Ast, Robert Babuška

Knowledge-Based Control Systems
2010-03-03

Principle of RL



- Interact with a system through **states** and **actions**
- Receive **rewards** as performance feedback

This presentation: **approximate RL**
– continuous states & actions

Outline

- 1 Introduction
- 2 Classical offline algorithms
- 3 Approximate algorithms for continuous spaces

Recall: Solution of the RL problem

- **Q-function** Q^h of policy h
- **Optimal Q-function** $Q^* = \max_h Q^h$
Satisfies Bellman optimality equation:

$$Q^*(x, u) = \rho(x, u) + \gamma \max_{u'} Q^*(f(x, u), u')$$
- **Optimal policy** h^* – greedy in Q^* :

$$h^*(x) = \arg \max_u Q^*(x, u)$$

Why approximation?

- Classical RL – tabular representation of Q-functions:
separate Q-value for each x and u
 - But in real-life control, X, U **continuous!**
Tabular representation impossible
- ⇒ need to **approximate the Q-function**

Recall: Types of algorithms

By model knowledge

- 1 Model-based – f and ρ known (dynamic programming)
- 2 Model-free – no f and ρ , only transition data (RL)
- 3 Model-learning – estimate f and ρ from transition data

By level of interaction

- 1 Offline – data collected in advance
- 2 Online – learn by interacting with the process

By path to optimal solution

- 1 Off-policy – find Q^* , use it to compute h^*
- 2 On-policy – find Q^h , improve h , repeat

Algorithms considered

	Off-policy	On-policy
(previous lecture) Classical online RL	Q-learning	SARSA
(this lecture) Classical offline DP	Q-iteration	Policy iteration
Approximate offline	fuzzy Q-iteration	approximate policy iteration
Approximate online	approximate Q-learning	online approximate policy iteration

Technical focus: Q-iteration & fuzzy Q-iteration

- 1 Introduction
- 2 Classical offline algorithms
- 3 Approximate algorithms for continuous spaces

Classical offline algorithms

	Off-policy	On-policy
Classical offline DP	Q-iteration	Policy iteration
Approximate offline	fuzzy Q-iteration	approximate policy iteration
Approximate online	approximate Q-learning	online approximate policy iteration

Classical offline algorithms
Offline, off-policy: Q-iteration

- Turn Bellman optimality equation:

$$Q^*(x, u) = \rho(x, u) + \gamma \max_{u'} Q^*(f(x, u), u')$$
 into an **iterative update**:

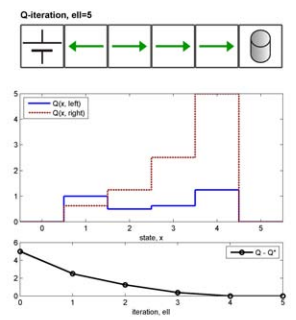
Q-iteration

repeat at each iteration ℓ
 for all x, u do
 $Q_{\ell+1}(x, u) \leftarrow \rho(x, u) + \gamma \max_{u'} Q_{\ell}(f(x, u), u')$
 end for
 until convergence to Q^*

- Once Q^* available: $h^*(x) = \arg \max_u Q^*(x, u)$

Classical offline algorithms
Cleaning robot: Q-iteration demo

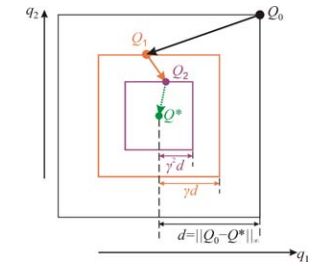
Discount factor: $\gamma = 0.5$



Classical offline algorithms
Q-iteration convergence

- Each update is a contraction with factor γ :

$$\|Q_{\ell+1} - Q^*\|_{\infty} \leq \gamma \|Q_{\ell} - Q^*\|_{\infty}$$
 ⇒ Q-iteration **monotonically converges** to Q^*



Classical offline algorithms
Offline, on-policy: Policy iteration

- Recall **on-policy**: find Q^h , improve h , repeat

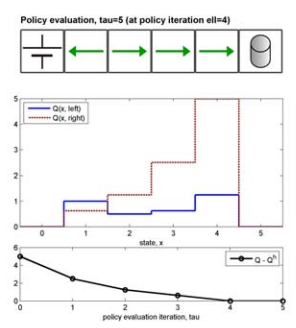
Policy iteration

starting from an initial policy
 repeat at each iteration ℓ
 policy evaluation: find $Q^{h_{\ell}}$
 policy improvement: $h_{\ell+1}(x) \leftarrow \arg \max_u Q^{h_{\ell}}(x, u)$
 until convergence to h^*

- Policy evaluation: iterative, from Bellman equation for Q^h (like Q-iteration)

Classical offline algorithms
Cleaning robot: Policy iteration demo

Initial policy: always go left



Classical offline algorithms

- 1 Introduction
- 2 Classical offline algorithms
- 3 **Approximate algorithms for continuous spaces**

Approximate algorithms for continuous spaces

	Off-policy	On-policy
Classical offline DP	Q-iteration	Policy iteration
Approximate offline	fuzzy Q-iteration	approximate policy iteration
Approximate online	approximate Q-learning	online approximate policy iteration

Approximating the Q-function

- 1 Introduction
- 2 Classical offline algorithms
- 3 **Approximate algorithms for continuous spaces**
 - Approximating the Q-function
 - Fuzzy Q-iteration
 - Approximate policy iteration
 - Online approximate policy iteration
 - Approximate Q-learning

Approximating the Q-function
Q-function approximation

- In real-life control, X, U continuous
 ⇒ **approximate Q-function** \hat{Q} must be used
- Usually, policy not approximated
 Greedy in \hat{Q} , computed on demand for given x :

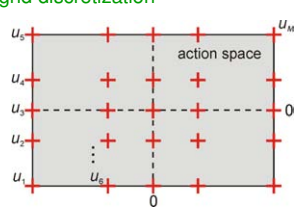
$$h(x) = \arg \max_u \hat{Q}(x, u)$$
- Approximator must ensure **efficient arg max solution**

Introduction Classical offline algorithms Approximate algorithms for continuous spaces

Approximating the Q-function

Approximating over the action space

- Approximator must ensure efficient “arg max” solution
- ⇒ Typically: **action discretization**
- Choose M discrete actions $u_1, \dots, u_M \in U$
Solve “arg max” by explicit enumeration
- Example: **grid discretization**



TU Delft

Introduction Classical offline algorithms Approximate algorithms for continuous spaces

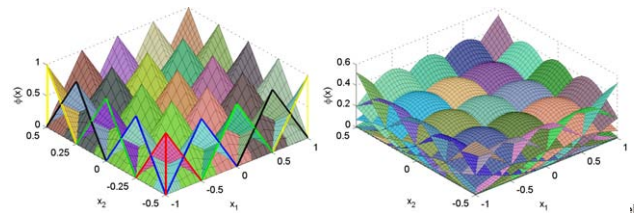
Approximating the Q-function

Approximating over the state space

- Typically: **basis functions**

$$\phi_1, \dots, \phi_N : X \rightarrow [0, \infty)$$

- Usually normalized: $\sum_i \phi_i(x) = 1$
- E.g., **fuzzy approximation, RBF network approximation**



TU Delft

Introduction Classical offline algorithms Approximate algorithms for continuous spaces

Approximating the Q-function

Linear Q-function approximation

Given:

- N basis functions ϕ_1, \dots, ϕ_N
- M discrete actions u_1, \dots, u_M

Store:

- $N \times M$ matrix of **parameters** θ
(one for each pair basis function–discrete action)

Approximate Q-function

$$\hat{Q}^\theta(x, u_j) = \sum_{i=1}^N \phi_i(x) \theta_{i,j} = [\phi_1(x) \dots \phi_N(x)] \begin{bmatrix} \theta_{1,j} \\ \vdots \\ \theta_{N,j} \end{bmatrix}$$

TU Delft

Introduction Classical offline algorithms Approximate algorithms for continuous spaces

Fuzzy Q-iteration

Offline, off-policy: Fuzzy Q-iteration

	Off-policy	On-policy
Classical offline DP	Q-iteration	Policy iteration
Approximate offline	fuzzy Q-iteration	approximate policy iteration
Approximate online	approximate Q-learning	online approximate policy iteration

(DP)

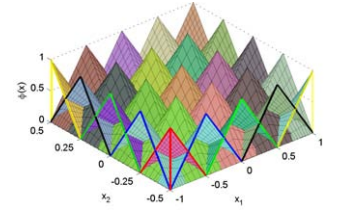
TU Delft

Introduction Classical offline algorithms Approximate algorithms for continuous spaces

Fuzzy Q-iteration

Fuzzy approximator

- Basis functions: **pyramidal membership functions (MFs)**
= cross-product of triangular MFs



- Each MF i has core (center) x_i
- $\theta_{i,j}$ can be seen as $\hat{Q}(x_i, u_j)$

TU Delft

Introduction Classical offline algorithms Approximate algorithms for continuous spaces

Fuzzy Q-iteration

Fuzzy Q-iteration

Recall classical Q-iteration:

```
repeat at each iteration ℓ
  for all x, u do
    Q_{ℓ+1}(x, u) = ρ(x, u) + γ max_{u'} Q_ℓ(f(x, u), u')
```

end for

until convergence

Fuzzy Q-iteration

```
repeat at each iteration ℓ
  for all cores x_i, discrete actions u_j do
    θ_{ℓ+1,i,j} = ρ(x_i, u_j) + γ max_{u'} Q_ℓ(f(x_i, u_j), u_j)
```

end for

until convergence

TU Delft

Introduction Classical offline algorithms Approximate algorithms for continuous spaces

Fuzzy Q-iteration

Fuzzy Q-iteration policy

- Recall optimal policy:

$$h^*(x) = \arg \max_u Q^*(x, u)$$

- Fuzzy Q-iteration policy:

$$\hat{h}^*(x) = \arg \max_{u_j, j=1, \dots, M} \hat{Q}^{\theta^*}(x, u_j)$$


(θ^* = converged parameter matrix)

TU Delft

Introduction Classical offline algorithms Approximate algorithms for continuous spaces

Fuzzy Q-iteration

Example: Inverted pendulum swing-up



- $x = [\text{angle } \alpha, \text{velocity } \dot{\alpha}]^T$
- $u = \text{voltage}$
- $\rho(x, u) = -x^T \begin{bmatrix} 5 & 0 \\ 0 & 0.1 \end{bmatrix} x - u^T 1 u$
- Discount factor $\gamma = 0.98$

- **Goal:** stabilize pointing up
- Insufficient actuation ⇒ need to swing back & forth

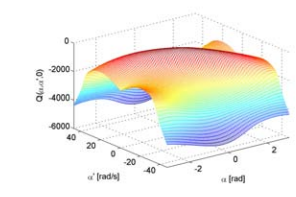
TU Delft

Introduction Classical offline algorithms Approximate algorithms for continuous spaces

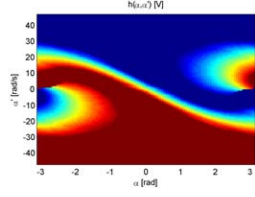
Fuzzy Q-iteration

Inverted pendulum: Near-optimal solution

Left: Q-function for $u = 0$



Right: policy



Replay

TU Delft

Introduction Classical offline algorithms Approximate algorithms for continuous spaces

Fuzzy Q-iteration

Inverted pendulum: Fuzzy Q-iteration demo

MFs: 41×21 equidistant grid
 Discretization: 5 actions, logarithmically spaced around 0

TU Delft

Introduction Classical offline algorithms Approximate algorithms for continuous spaces

Fuzzy Q-iteration

Inverted pendulum: Fuzzy Q-iteration demo

Demo

TU Delft

Introduction Classical offline algorithms Approximate algorithms for continuous spaces

Fuzzy Q-iteration

Fuzzy Q-iteration convergence

Like classical Q-iteration:

- Each update is a contraction with factor γ :

$$\|\theta_{\ell+1} - \theta^*\|_{\infty} \leq \gamma \|\theta_{\ell} - \theta^*\|_{\infty}$$

\Rightarrow **Monotonic convergence to θ^***

- θ^* leads to **near-optimal** \hat{Q}^*, \hat{h}^*

TU Delft

Introduction Classical offline algorithms Approximate algorithms for continuous spaces

Approximate policy iteration

Offline, on-policy: Approximate policy iteration

	Off-policy	On-policy
Classical offline DP	Q-iteration	Policy iteration
Approximate offline	fuzzy Q-iteration	approximate policy iteration (RL)
Approximate online	approximate Q-learning	online approximate policy iteration

TU Delft

Introduction Classical offline algorithms Approximate algorithms for continuous spaces

Approximate policy iteration

Approximate policy iteration

Recall classical policy iteration:

starting from an initial policy
 repeat at each iteration ℓ
 policy evaluation: find $Q^{h_{\ell}}$
 policy improvement: $h_{\ell+1}(x) \leftarrow \arg \max_u Q^{h_{\ell}}(x, u)$
 until convergence

Approximate policy iteration (API)

starting from an initial policy
 repeat at each iteration ℓ
 approximate policy evaluation: find θ so that $\hat{Q}^{\theta} \approx Q^{h_{\ell}}$
 policy improvement: $h_{\ell+1}(x) \leftarrow \arg \max_u \hat{Q}^{\theta}(x, u)$
 until convergence

TU Delft

Introduction Classical offline algorithms Approximate algorithms for continuous spaces

Approximate policy iteration

Inverted pendulum: API demo

Basis functions: 15×9 grid of RBFs
 Discretization: 3 equidistant actions

TU Delft

Introduction Classical offline algorithms Approximate algorithms for continuous spaces

Approximate policy iteration

Inverted pendulum: API demo

Demo

Least-squares policy iteration, $\text{ell}=9$

TU Delft

Introduction Classical offline algorithms Approximate algorithms for continuous spaces

Online approximate policy iteration

Online, on-policy: Online API

	Off-policy	On-policy
Classical offline DP	Q-iteration	Policy iteration
Approximate offline	fuzzy Q-iteration	approximate policy iteration
Approximate online	approximate Q-learning	online approximate policy iteration (RL)

TU Delft

Introduction Classical offline algorithms Approximate algorithms for continuous spaces

Online approximate policy iteration

Inverted pendulum: Online API demo

Real-time learning control
 Approximator: similar to offline API (except 11×11 RBFs)

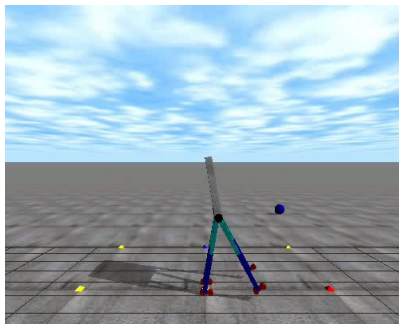
TU Delft

Approximate Q-learning
Online, off-policy: Approximate Q-learning

	Off-policy	On-policy
Classical offline DP	Q-iteration	Policy iteration
Approximate offline	fuzzy Q-iteration	approximate policy iteration
Approximate online	approximate Q-learning	online approximate policy iteration

(RL)

Approximate Q-learning
Demo: Q-learning for walking robot (Erik Schuitema)

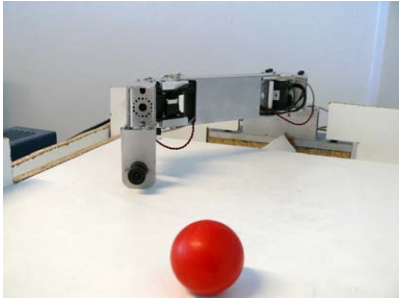


Approximate Q-learning
Recall: Experience replay

- Store each transition sample $(x_k, u_k, x_{k+1}, r_{k+1})$ into a database
- At every step, **replay** several transitions from the database

Approximate Q-learning
Demo: Q-learning for goalkeeper robot (Sander Adam)

Real-time learning control
 Employs **experience replay**



Take-home message

Approximate reinforcement learning =
Learn how to **optimally** control complex systems from scratch