# Reinforcement Learning
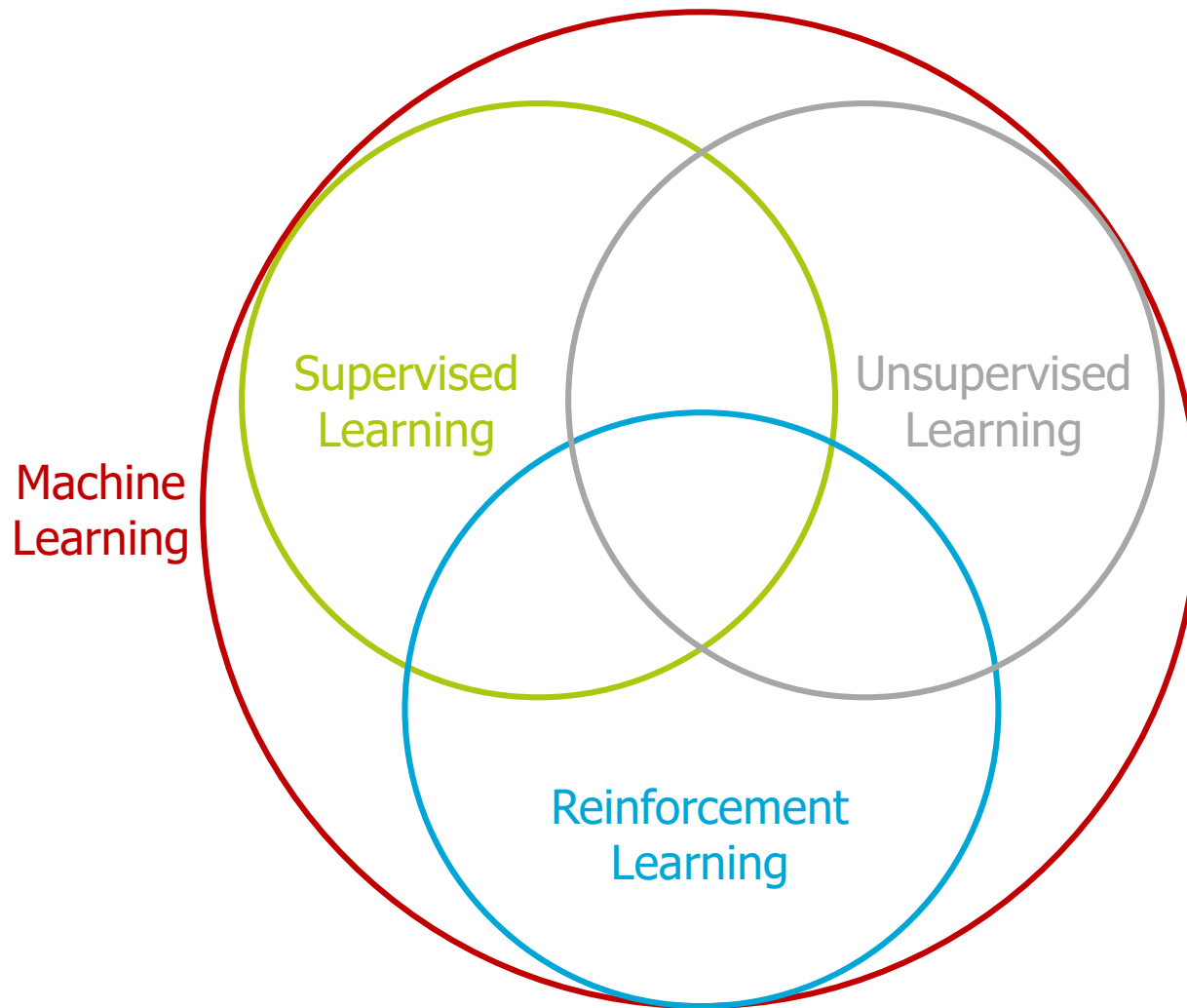
Jens Kober

j.kober@tudelft.nl

# Outline

- Introduction
- Notation & basic concepts
- Types of RL algorithms
  - Optimal control
  - Value function methods
  - Policy search
- Summary & outlook

TUDelft

# Types of Machine Learning

- Supervised learning
- Unsupervised learning
- Semi-supervised learning
- Reinforcement learning

➢Definitions?
➢Differences?
➢Examples?

TUDelft

# Types of Machine Learning

# Reinforcement Learning

- Learning by trial & error


http://hiit-blog.dailyhiit.com

- Learning by rewards & punishments


edwinandrewlove.blogspot.com


http://pixabay.com

# Outline

- Introduction
- Notation & basic concepts
- Types of RL algorithms
  - Optimal control
  - Value function methods
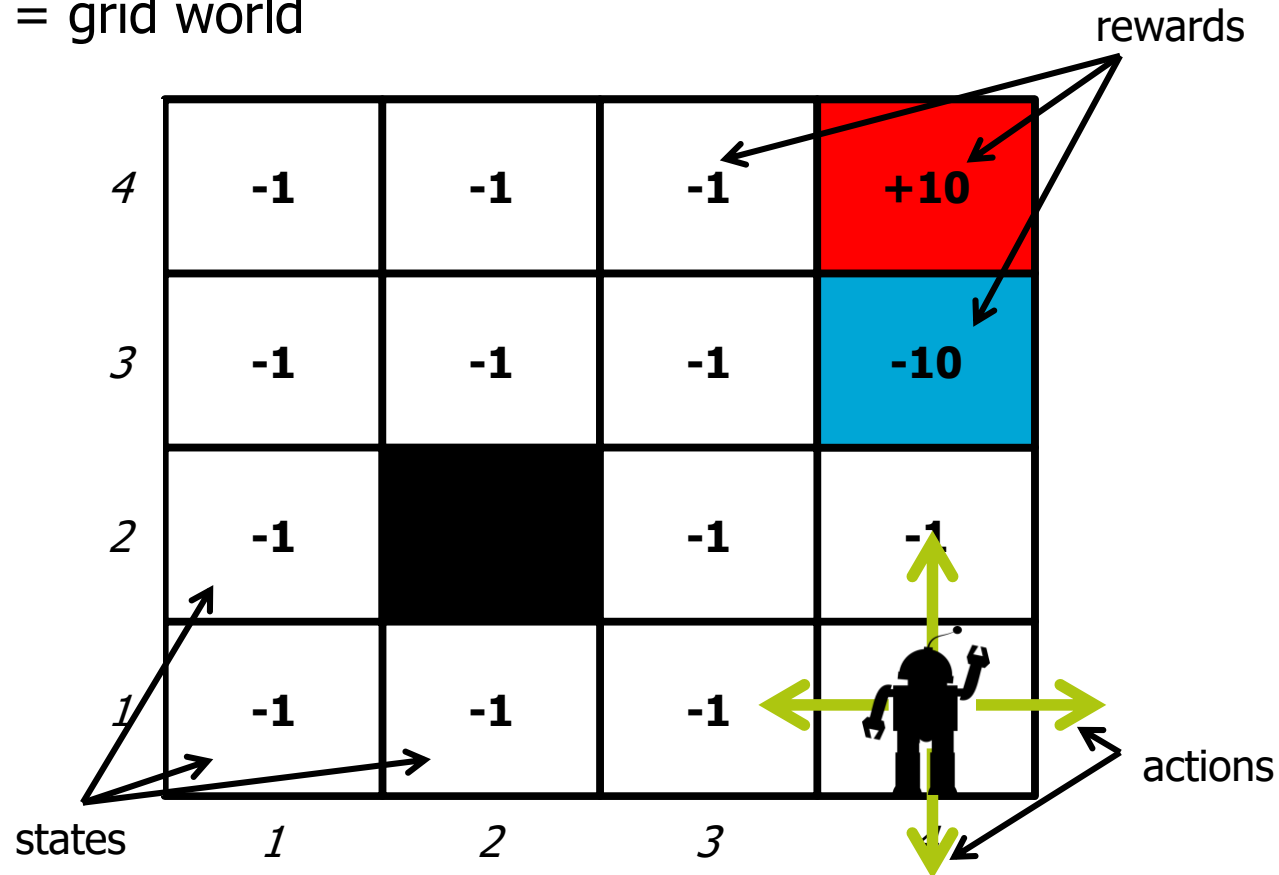  - Policy search
- Summary & outlook

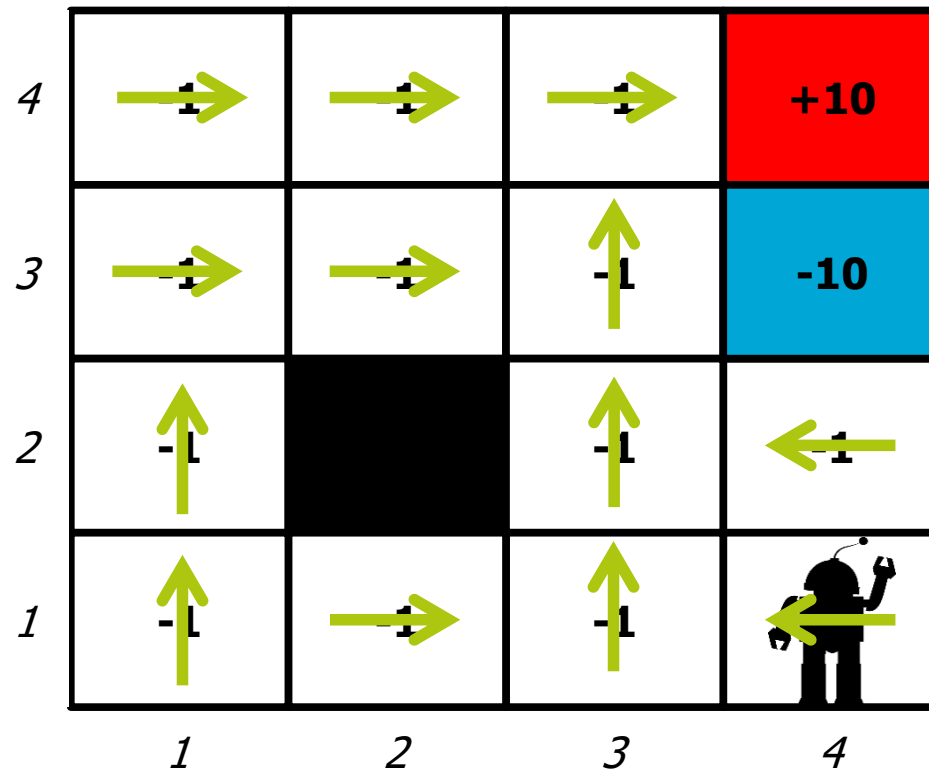TUDelft
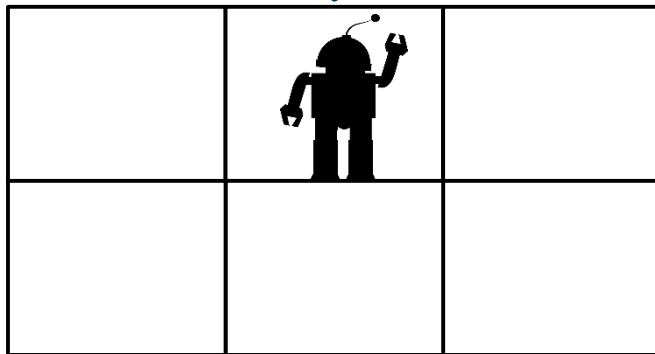
# Classical RL Example

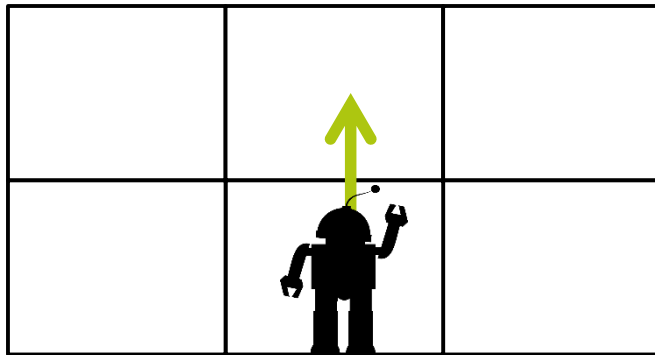- Maze

# Classical RL Example

- Maze = grid world

# Classical RL Example

- Strategy = policy

# Classical RL Example

- Deterministic grid world

- Stochastic grid world



transition probabilities

60%

20%   20%

# Formal Notation

state $s_t$

policy $\pi$

Agent

state $s_t$

reward $r_t$ $\boxed{\text{-1}}$

action $a_t$ ↑

$r_{t+1}$

Environment

$s_{t+1}$

# Formal Notation



- $s \in \mathcal{S}$ set of states (discrete or continuous)

- $a \in \mathcal{A}$ set of actions (discrete or continuous)

- $r$ reward $\mathcal{R}_{ss'}^{a} = E\left\{ r_{t+1} \mid a_t = a, s_t = s, s_{t+1} = s' \right\}$

- $\pi(s) = a$ policy

Goal: find $\pi^*$ maximizing return $R$

# Return $R$

- Finite Horizon $H$

$$R_t = r_{t+1} + r_{t+2} + \ldots + r_{t+H} = \sum_{h=1}^{H} r_{t+h}$$

- Discounted $\quad 0 \le \gamma < 1$

$$R_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+2} + \ldots = \sum_{h=0}^{\infty} \gamma^h r_{t+h+1}$$

- Average

$$R_t = \lim_{H \to \infty} \left( \frac{1}{H} \sum_{h=1}^{H} r_{t+h} \right)$$

# Markov Property

- Transition probability

$$P\left(s_{t+1} \mid a_t, s_t, a_{t-1}, s_{t-1}, \ldots, a_1, s_1\right)$$

➢Markov property

$$\mathcal{P}_{ss'}^a = P\left(s_{t+1} = s' \mid s_t = s, a_t = a\right)$$

- Markov decision process (MDP)

$$\mathcal{S}, \mathcal{A}, \mathcal{P}_{ss'}^a, \mathcal{R}_{ss'}^a, \gamma$$

# Bellman Principle of Optimality

"An optimal policy has the property that whatever the initial state and initial decision are, the remaining decisions must constitute an optimal policy with regard to the state resulting from the first decision."

Bellman, 1957

- Dynamic Programming
- Optimal Control

# Value Functions

- State value function

$$V^{\pi}(s) = E^{\pi}\left\{\sum_{h=0}^{\infty} \gamma^h r_{t+h+1} \mid s_t = s\right\} = \sum_{s' \in S} \mathcal{P}_{ss'}^{\pi}\left[\mathcal{R}_{ss'}^{\pi} + \gamma V^{\pi}(s')\right]$$

$$V^*(s) = \max_{a \in \mathcal{A}}\left(\sum_{s' \in S} \mathcal{P}_{ss'}^{a}\left[\mathcal{R}_{ss'}^{a} + \gamma V^*(s')\right]\right)$$

- State-action value function

$$Q^{\pi}(s,a) = \sum_{s' \in S} \mathcal{P}_{ss'}^{a}\left[\mathcal{R}_{ss'}^{a} + \gamma V^{\pi}(s')\right]$$

$$Q^*(s,a) = \sum_{s' \in S} \mathcal{P}_{ss'}^{a}\left[\mathcal{R}_{ss'}^{a} + \gamma V^*(s')\right]$$

$$= \sum_{s' \in S} \mathcal{P}_{ss'}^{a}\left[\mathcal{R}_{ss'}^{a} + \gamma\left(\max_{a' \in \mathcal{A}} Q^*(s',a')\right)\right]$$

**T̃U**Delft

# Optimal Policy

- State-action value function

$$Q^*(s,a) = \sum_{s' \in S} \mathcal{P}_{ss'}^a \left[ \mathcal{R}_{ss'}^a + \gamma \left( \max_{a' \in \mathcal{A}} Q^*(s',a') \right) \right]$$

➢ $\pi^*(s) = \arg\max_{a \in \mathcal{A}} \left( Q^*(s,a) \right)$

- State value function

$$V^*(s) = \max_{a \in \mathcal{A}} \left( \sum_{s' \in S} \mathcal{P}_{ss'}^a \left[ \mathcal{R}_{ss'}^a + \gamma V^*(s') \right] \right)$$

➢ $\pi^*(s) = \arg\max_{a \in \mathcal{A}} \left( \sum_{s' \in S} \mathcal{P}_{ss'}^a \left[ \mathcal{R}_{ss'}^a + \gamma V^*(s') \right] \right)$

$\tilde{T}UDelft$

# Optimal State-Action Value Function

- max. 10 steps, stop at fields $\pm$ 10, stay at walls, deterministic, $\gamma = 1$ , reward: -1 per step, $\pm 10$ for reaching states

# Outline

- Introduction
- Notation & basic concepts
- Types of RL algorithms
  - Optimal control
  - Value function methods
  - Policy search
- Summary & outlook

TUDelft

# Types of RL Algorithms



Prior Knowledge

Data
$\{(s, a, s', r)\}$

Transition & Reward
$\mathcal{P}^a_{ss'}, \mathcal{R}^a_{ss'}$

Value Function
$V^*$

Value Function
$V^*$

Policy
$\pi^*$

Policy
$\pi^*$

Policy
$\pi^*$

Optimal Control

Value Function Methods

Policy Search Methods

TUDelft

# Optimal Control

- **Model-based:** requires models of reward and transition functions
- Fast (no real world interactions required)
- Models can be also be learned
- RL often takes advantage of model errors
- Examples:
  - Policy Iteration
  - Value Iteration

TUDelft

# Value Iteration

- Bellman optimality equation

$$Q^*(s,a) = \sum_{s' \in S} \mathcal{P}^a_{ss'} \left[ \mathcal{R}^a_{ss'} + \gamma \left( \max_{a' \in \mathcal{A}} Q^*(s',a') \right) \right]$$

➤ Turn into an iterative update

**Q-Iteration**

**repeat** at each iteration $i$
    **for all** $s, a$ **do**

$$Q_{i+1}(s,a) \leftarrow \sum_{s' \in S} \mathcal{P}^a_{ss'} \left[ \mathcal{R}^a_{ss'} + \gamma \left( \max_{a' \in \mathcal{A}} Q_i(s',a') \right) \right]$$

    **end for**
**until** convergence to $Q^*$

➤ Once $Q^*$ available $\pi^*(s) = \arg \max_{a \in \mathcal{A}} \left( Q^*(s,a) \right)$

# 1ˢᵗ Iteration

$Q_0$

| | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| **4** | 0 <br> 0 **-1** 0 <br> 0 | 0 <br> 0 **-1** 0 <br> 0 | 0 <br> 0 **-1** 0 <br> 0 | **+10** |
| **3** | 0 <br> 0 **-1** 0 <br> 0 | 0 <br> 0 **-1** 0 <br> 0 | 0 <br> 0 **-1** 0 <br> 0 | **-10** |
| **2** | 0 <br> 0 **-1** 0 <br> 0 | ■ | 0 <br> 0 **-1** 0 <br> 0 | 0 <br> 0 **-1** 0 <br> 0 |
| **1** | 0 <br> 0 **-1** 0 <br> 0 | 0 <br> 0 **-1** 0 <br> 0 | 0 <br> 0 **-1** 0 <br> 0 | 0 <br> 0 **-1** 0 <br> 0 |

$Q_1$

| | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| **4** | -1 <br> -1 **-1** -1 <br> -1 | -1 <br> -1 **-1** -1 <br> -1 | -1 <br> -1 **-1** 9 <br> -1 | **+10** |
| **3** | -1 <br> -1 **-1** -1 <br> -1 | -1 <br> -1 **-1** -1 <br> -1 | -1 <br> -1 **-1** -11 <br> -1 | **-10** |
| **2** | -1 <br> -1 **-1** -1 <br> -1 | ■ | -1 <br> -1 **-1** -1 <br> -1 | -11 <br> -1 **-1** -1 <br> -1 |
| **1** | -1 <br> -1 **-1** -1 <br> -1 | -1 <br> -1 **-1** -1 <br> -1 | -1 <br> -1 **-1** -1 <br> -1 | -1 <br> -1 **-1** -1 <br> -1 |

# 2ⁿᵈ Iteration



$Q_1$

| | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 4 | -1 **-1** -1 / -1 | -1 **-1** -1 / -1 | -1 **-1** 9 / -1 | **+10** |
| 3 | -1 **-1** -1 / -1 | -1 **-1** -1 / -1 | -1 **-1** -11 / -1 | **-10** |
| 2 | -1 **-1** -1 / -1 | | -1 **-1** -1 / -1 | -11 / -1 **-1** -1 / -1 |
| 1 | -1 **-1** -1 / -1 | -1 **-1** -1 / -1 | -1 **-1** -1 / -1 | -1 **-1** -1 / -1 |

$Q_2$

| | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 4 | -2 **-1** -2 / -2 | -2 **-1** 8 / -2 | -1 **-1** 9 / -1 | **+10** |
| 3 | -2 **-1** -2 / -2 | -2 **-1** -2 / -2 | -2 **-1** -11 / -2 | **-10** |
| 2 | -2 **-1** -2 / -2 | | -2 **-1** -2 / -2 | -11 / -2 **-1** -2 / -2 |
| 1 | -2 **-1** -2 / -2 | -2 **-1** -2 / -2 | -2 **-1** -2 / -2 | -2 **-1** -2 / -2 |

# 3rd Iteration

## $Q_2$

| | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| **4** | -2 **-1** -2 <br> -2 <br> -2 | -2 **-1** 8 <br> -2 <br> -2 | -1 **-1** 9 <br> 8 <br> -1 | **+10** |
| **3** | -2 **-1** -2 <br> -2 <br> -2 | -2 **-1** -2 <br> -2 <br> -2 | -2 **-1** -11 <br> 8 <br> -2 | **-10** |
| **2** | -2 **-1** -2 <br> -2 <br> -2 | ■ | -2 **-1** -2 <br> -2 <br> -2 | -2 **-1** -2 <br> -11 <br> -2 |
| **1** | -2 **-1** -2 <br> -2 <br> -2 | -2 **-1** -2 <br> -2 <br> -2 | -2 **-1** -2 <br> -2 <br> -2 | -2 **-1** -2 <br> -2 <br> -2 |

## $Q_3$

| | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| **4** | -3 **-1** 7 <br> -3 <br> -3 | -3 **-1** 8 <br> 7 <br> -3 | 7 **-1** 9 <br> 8 <br> 7 | **+10** |
| **3** | -3 **-1** -3 <br> -3 <br> -3 | -3 **-1** 7 <br> 7 <br> -3 | -2 **-1** -11 <br> 8 <br> -2 | **-10** |
| **2** | -3 **-1** -3 <br> -3 <br> -3 | ■ | -3 **-1** -3 <br> 7 <br> -3 | -3 **-1** -3 <br> -11 <br> -3 |
| **1** | -3 **-1** -3 <br> -3 <br> -3 | -3 **-1** -3 <br> -3 <br> -3 | -3 **-1** -3 <br> -3 <br> -3 | -3 **-1** -3 <br> -3 <br> -3 |

https://youtu.be/VCdxqn0fcnE

# Types of RL Algorithms



Prior Knowledge

Data
$\{(s, a, s', r)\}$

Transition & Reward
$\mathcal{P}^a_{ss'}, \mathcal{R}^a_{ss'}$

Value Function
$V^*$

Value Function
$V^*$

Policy
$\pi^*$

Policy
$\pi^*$

Policy
$\pi^*$

Optimal Control

Value Function Methods

Policy Search Methods

TUDelft

# Value Function Methods

- Model-free: (implicitly) learns models of reward and transition functions
- Global optimum (if everything is explored)
- Policy has global dependence
- Examples:
  - Monte Carlo
  - TD(λ)
  - Q-learning
  - SARSA

# Temporal Difference Learning

- Use sample based estimate to update the value function

$$Q(s,a) \leftarrow (1-\alpha)Q(s,a) + \alpha \left[ r + \gamma \max_{a' \in \mathcal{A}} Q(s',a') \right]$$

learning rate

- Equivalently

$$Q(s,a) \leftarrow Q(s,a) + \alpha \left[ r + \gamma \max_{a' \in \mathcal{A}} Q(s',a') - Q(s,a) \right]$$

temporal difference error

# Exploration

- Which action $a$ to pick?
- Assume $Q$ is optimal (greedy action)

$$a \leftarrow \arg\max_{a \in \mathcal{A}} \left( Q^*(s,a) \right)$$

- But $Q$ is not optimal (yet)!
- ➤ We need to explore

$\varepsilon$ -greedy policy:
- With probability $\varepsilon$ pick a random action (exploration)

$$a \leftarrow \mathrm{rand}(\mathcal{A})$$

- With probability $(1-\varepsilon)$ pick the greedy action (exploitation)

$$a \leftarrow \arg\max_{a \in \mathcal{A}} \left( Q^*(s,a) \right)$$

**TU**Delft

# Q-Learning

**Q-Learning**

**loop**

    observe state $s$

$$a \leftarrow \arg \max_{a \in \mathcal{A}} \left( Q(s,a) \right)$$

    with probability $\varepsilon, a \leftarrow \text{rand}(\mathcal{A})$

    apply $a$, observe $r$ and $s'$

$$Q(s,a) \leftarrow Q(s,a) + \alpha \left[ r + \gamma \max_{a' \in \mathcal{A}} Q(s',a') - Q(s,a) \right]$$
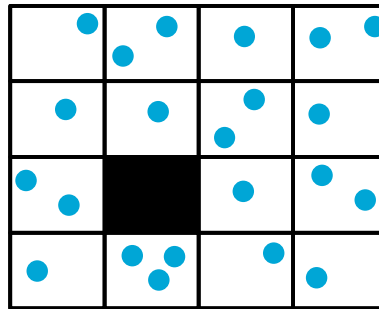
**until** convergence

$\tilde{T}U$Delft

# Properties of Discussed Methods

- Advantages
  - Generic framework
  - Very few assumptions
  - Guaranteed to converge to optimum

- Disadvantages
  - Can take lot of iterations
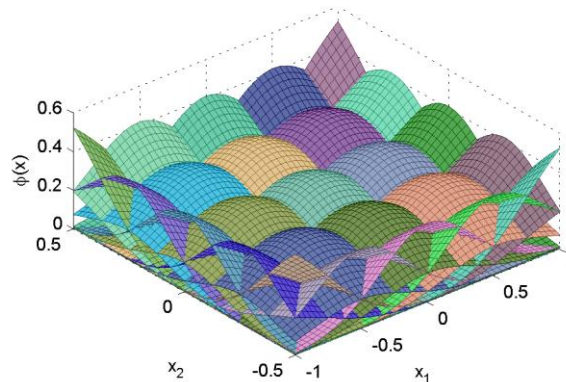  - Infeasible for continuous states/actions

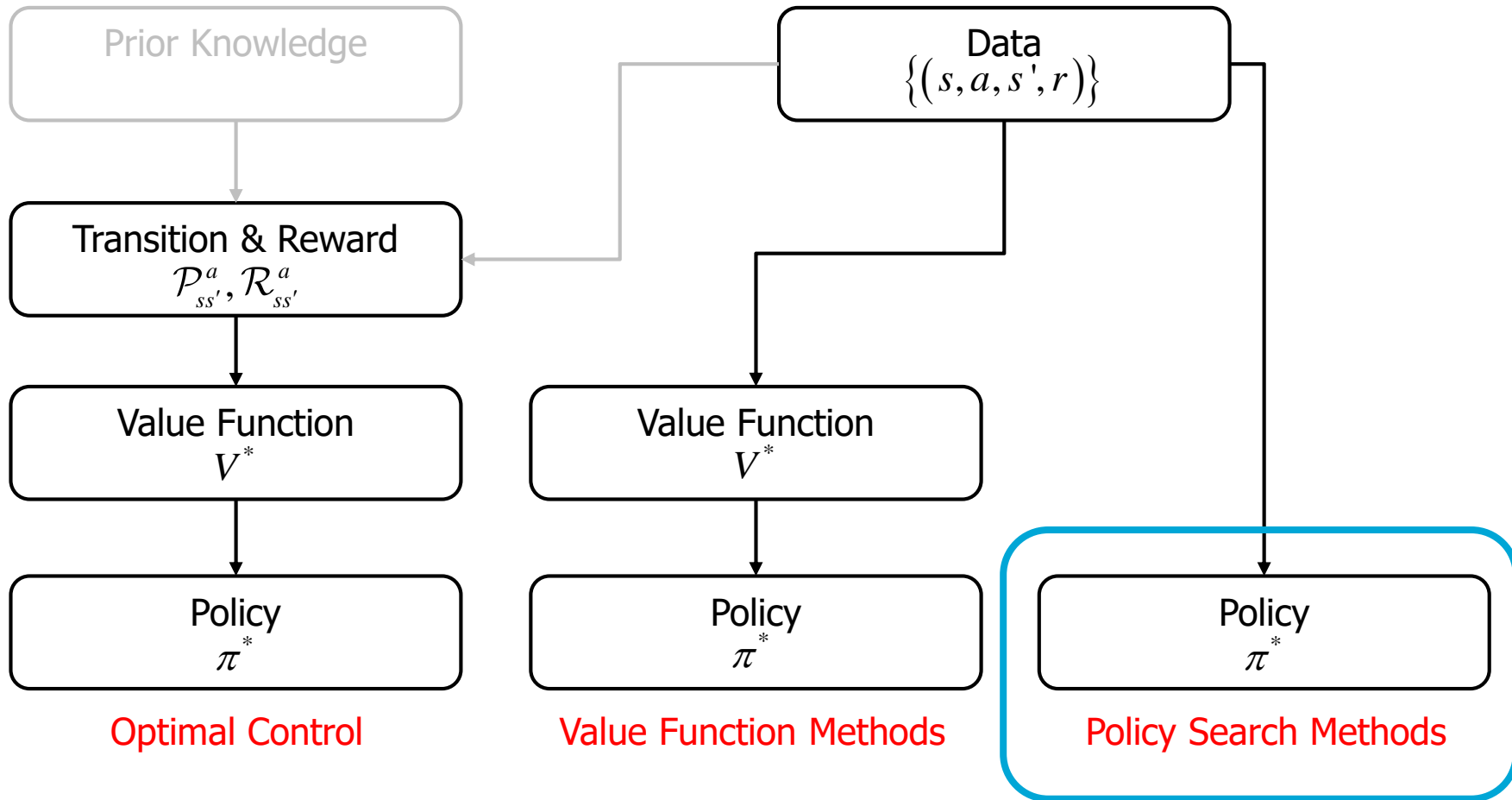TUDelft

# Dealing with Continuous States/Actions

- Discretization



- Function approximation
  - Value function
  - Policy

https://youtu.be/nM1HTp_P3lY

# Types of RL Algorithms



Prior Knowledge

Data
$\{(s, a, s', r)\}$

Transition & Reward
$\mathcal{P}_{ss'}^a, \mathcal{R}_{ss'}^a$

Value Function
$V^*$

Value Function
$V^*$

Policy
$\pi^*$

Policy
$\pi^*$

Policy
$\pi^*$

Optimal Control

Value Function Methods

Policy Search Methods

TUDelft

# Policy Search Methods

- <span style="color:red">Model-free</span>
- Local optimum
- Usually parametrized policies

- Examples:
  - Gradient-based
  - Expectation-maximization inspired
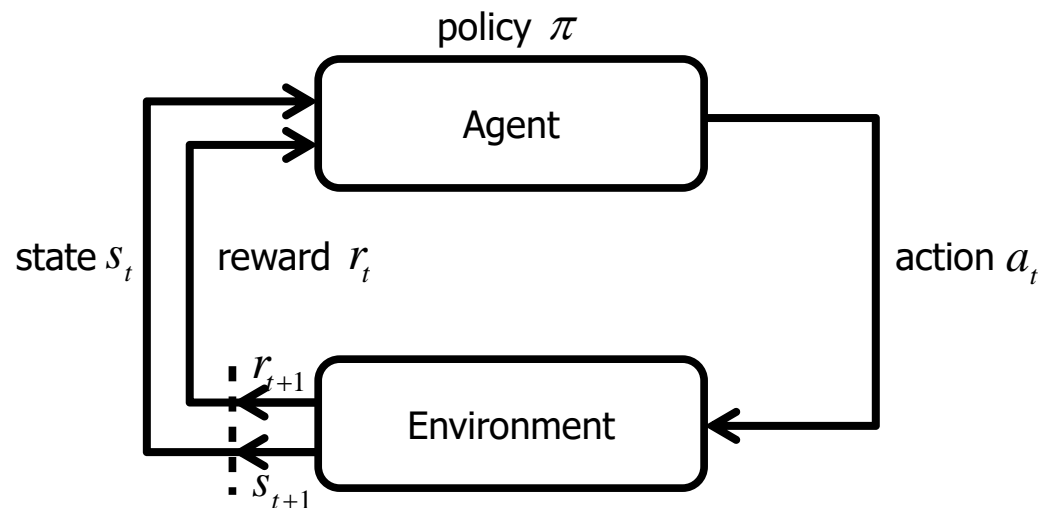  - General optimization

# Outline

- Introduction
- Notation & basic concepts
- Types of RL algorithms
  - Optimal control
  - Value function methods
  - Policy search
- Summary & outlook

TUDelft

# Summary

Reinforcement learning
- Inspired by human & animal learning
- Reward as feedback signal
- Model-free, adaptive optimal control



policy $\pi$

Agent

Environment

state $s_t$    reward $r_t$    action $a_t$

$r_{t+1}$

$s_{t+1}$

Goal: find $\pi^*$ maximizing return $R$

$\tilde{T}U$Delft

# Challenges of (Robot) RL

- Curse of dimensionality
- Exploration-exploitation trade-off
- Real-world samples
- Model uncertainty
- Goal specification

TUDelft

# Tractability Through:

- Representations
  - State and/or action discretization
  - Value function approximation
  - Pre-structured policies
- Prior knowledge
  - Demonstrations
  - Task Structure
- Models
  - Mental rehearsal