/informatik & security **/fh///**
st.pölten

# Reinforcement Learning

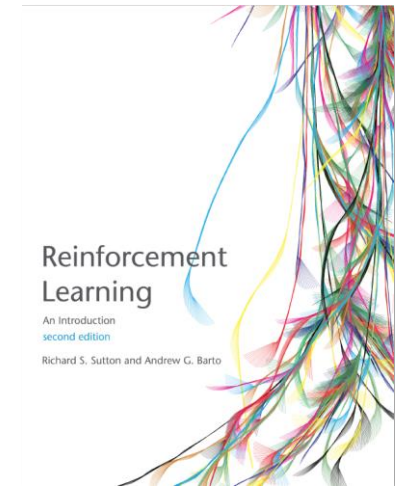Studiengang: Data Science and Business Analytics

Vortragender: Sebastian Eresheim

Ort, Datum: St. Pölten, 05.12.2021

# Organizational

Lecture is based on:

- *Reinforcement Learning: An Introduction*
  Richard Sutton, Andrew Barto
  www.incompleteideas.net/book/the-book.html

# Overview

**Machine Learning**

Supervised Learning

Unsupervised Learning

Reinforcement Learning

# Overview

Reinforcement Learning is motviated by learning:

- like a human

- via interaction (Trial and Error)
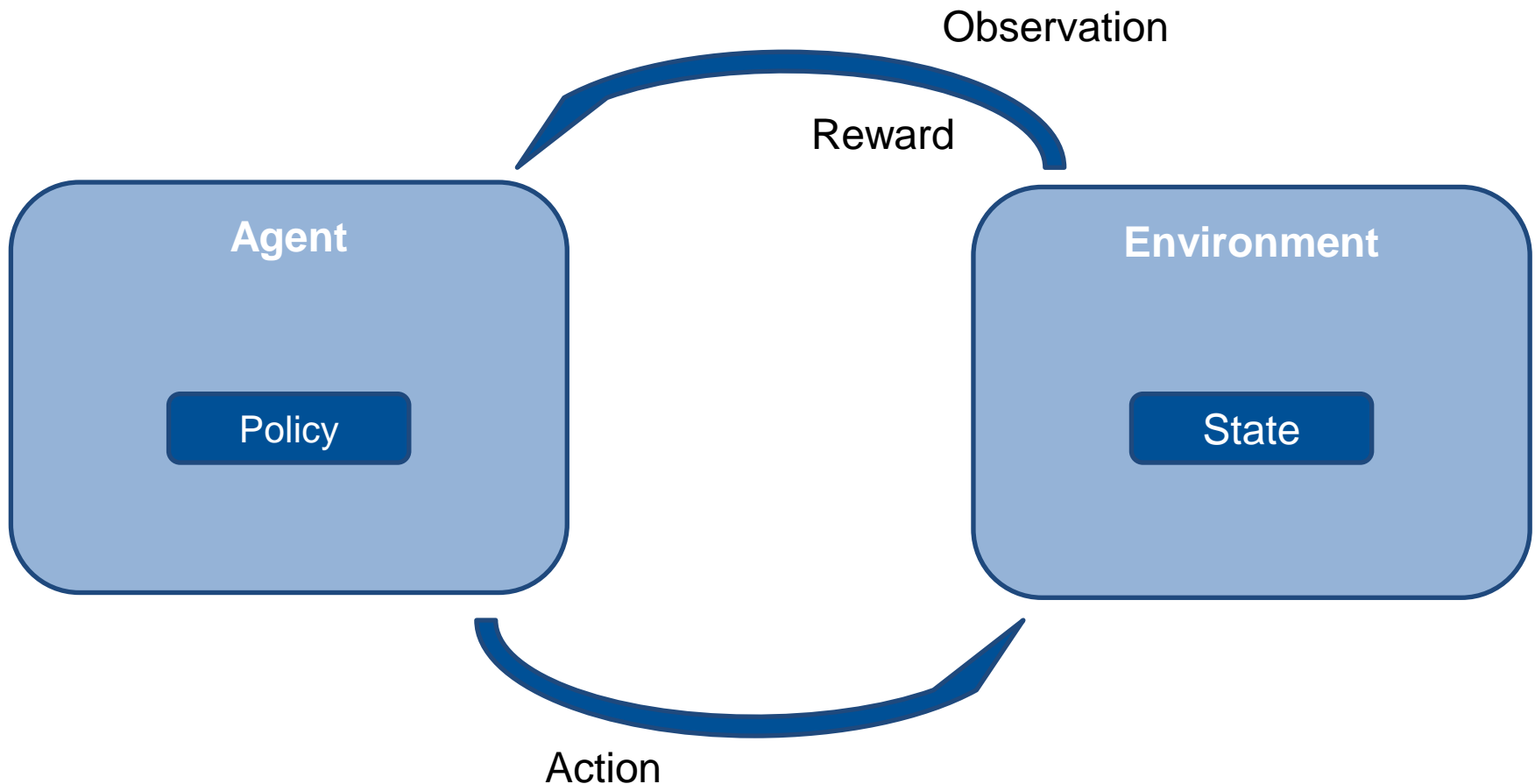
- „learning **what** to do!" not „how to do it".

# Overview

In reinforcement learning an agent learns how to map actions to observations in order to maximise a numerical reward, which it receives from an environment.
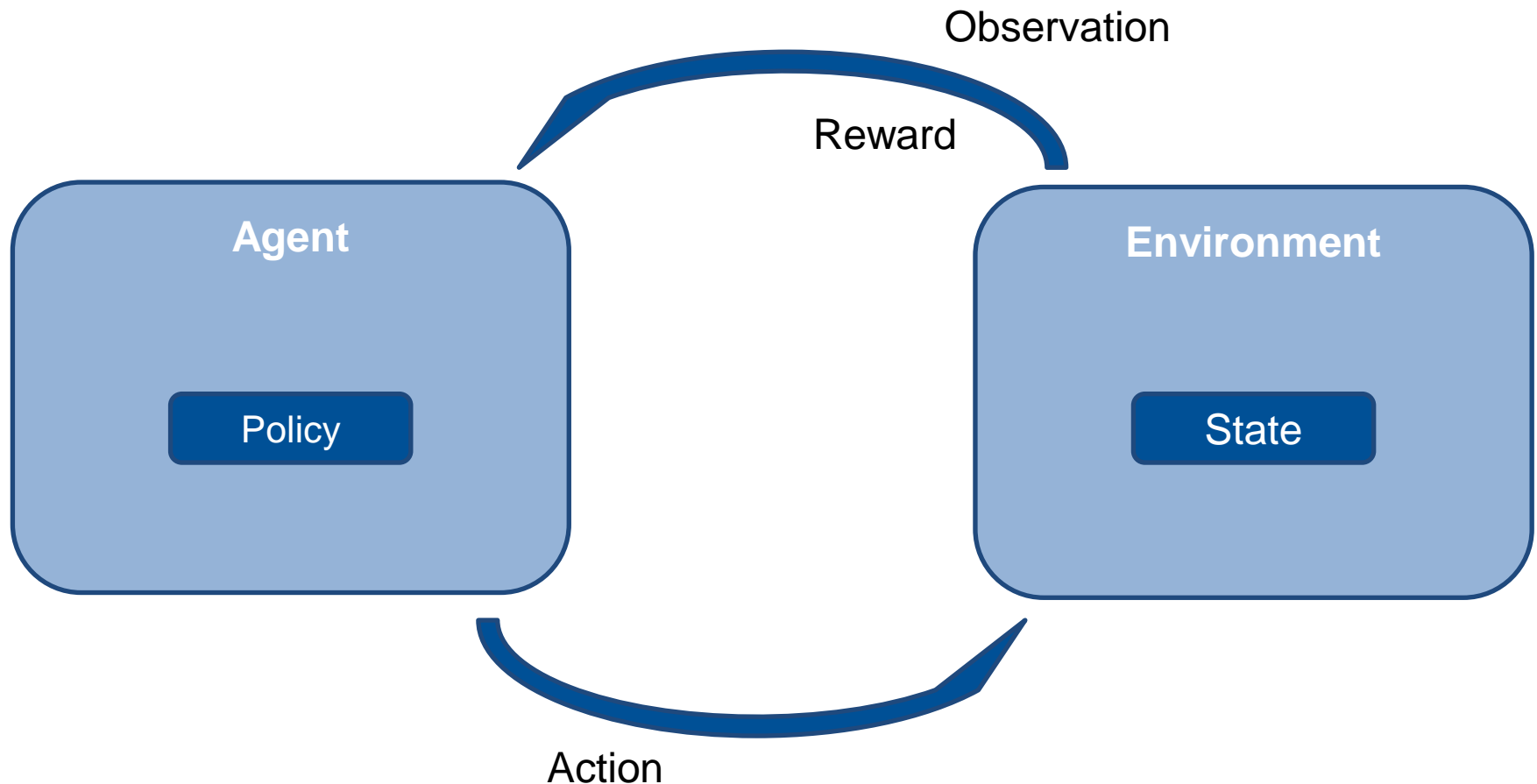
# Overview

In reinforcement learning an **agent** learns how to map <u>actions</u> to <u>observations</u> in order to maximise a numerical <u>reward</u>, which it receives from an **environment**.
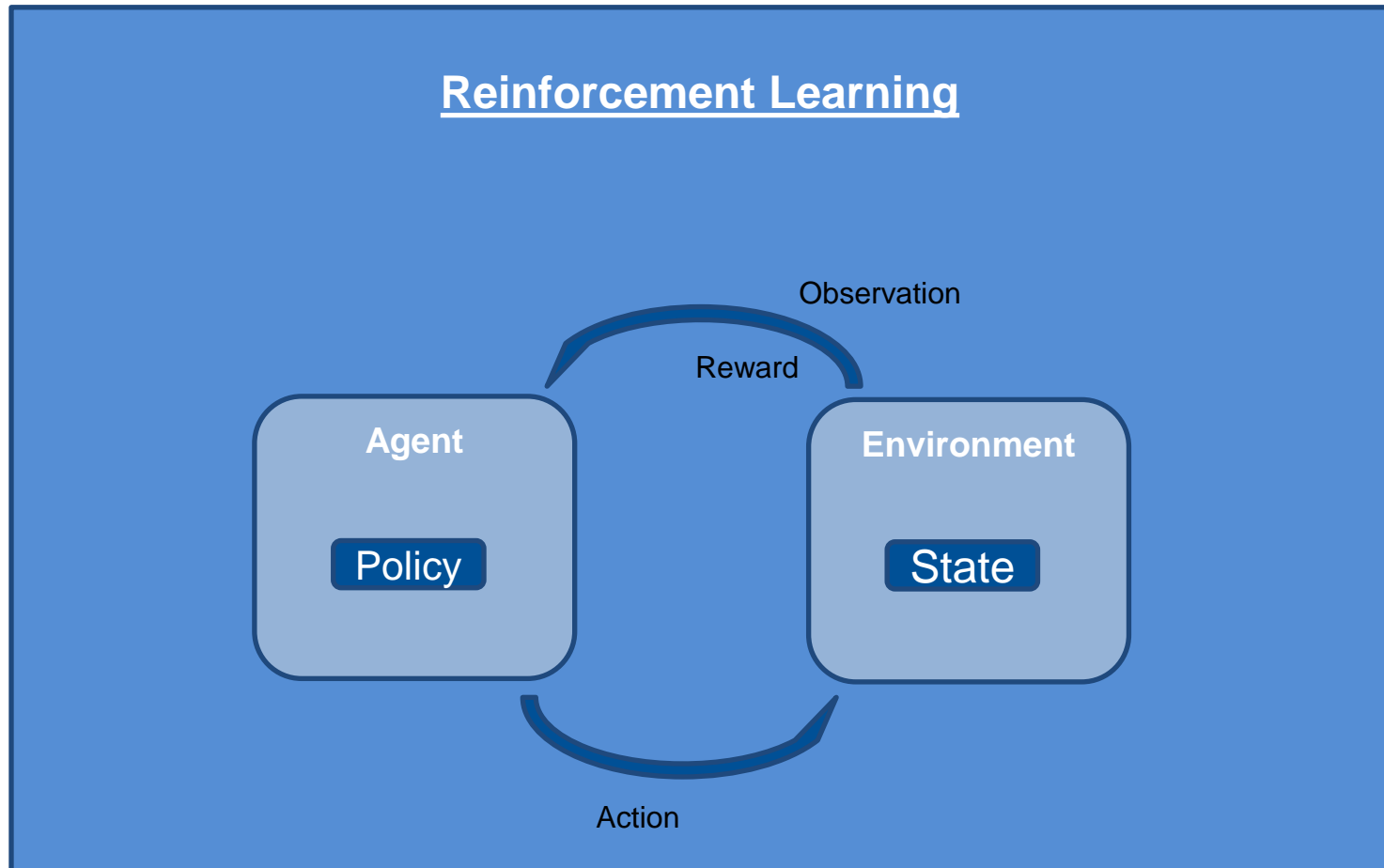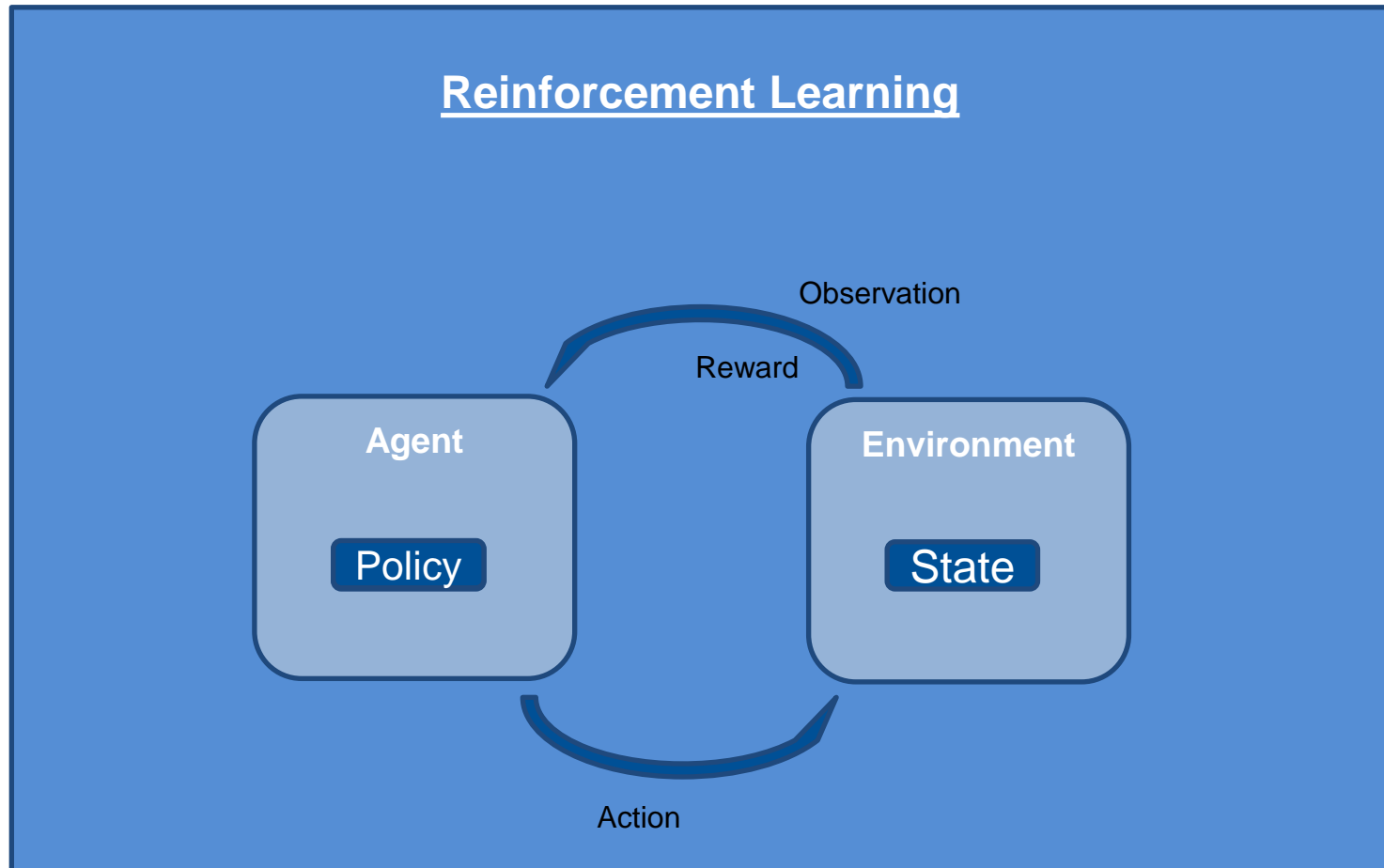
# Preliminaries

Observation

Reward

**Agent**

Policy

**Environment**

State

Action

# Preliminaries

Observation

Reward

**Agent**

Policy

**Environment**

State

Action

# Preliminaries

**Reinforcement Learning**

Observation

Reward

| Agent | Environment |
|---|---|
| Policy | State |

Action

# Preliminaries

**Reinforcement Learning**

Observation

Reward

**Agent**

Policy

**Environment**

State

Action

# Preliminaries

# Preliminaries

Tabular Methods

Approximate Solution Methods
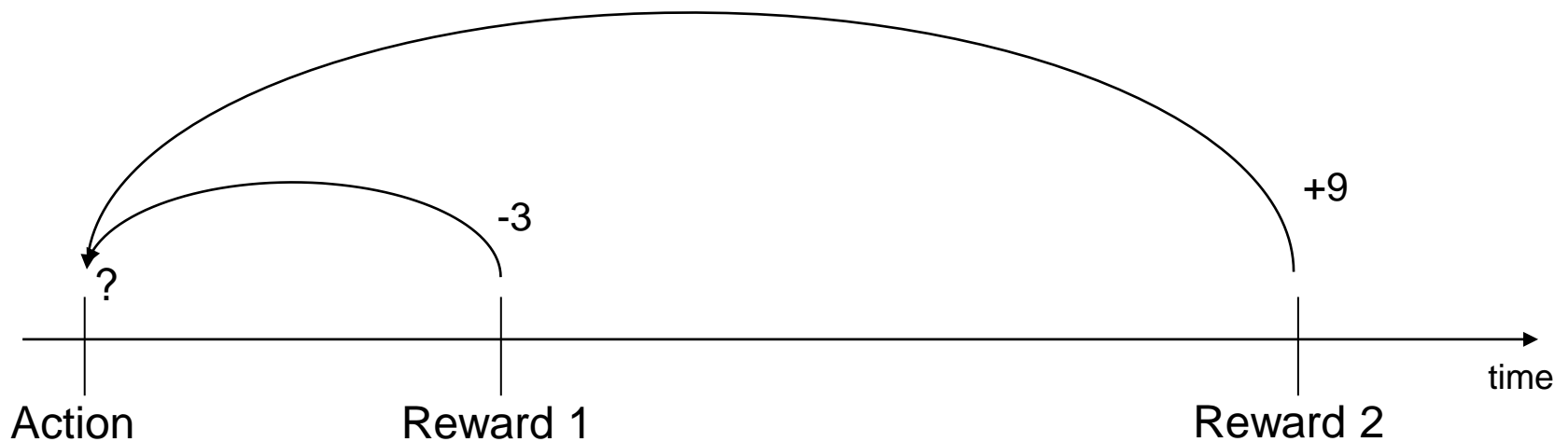
**Direct Policy Methods**
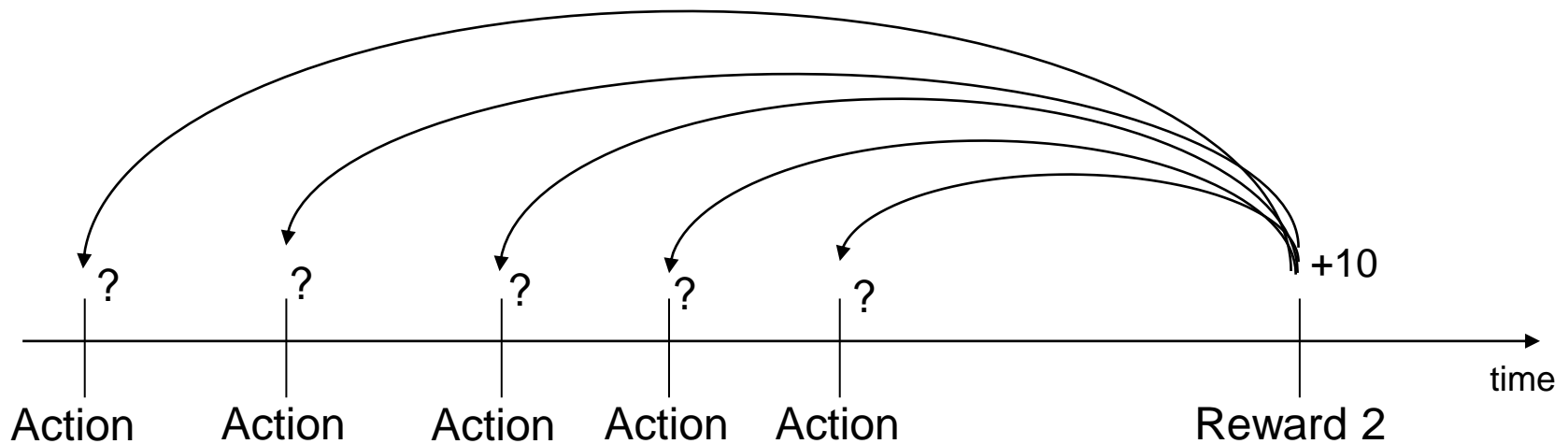
# Overview - Problems

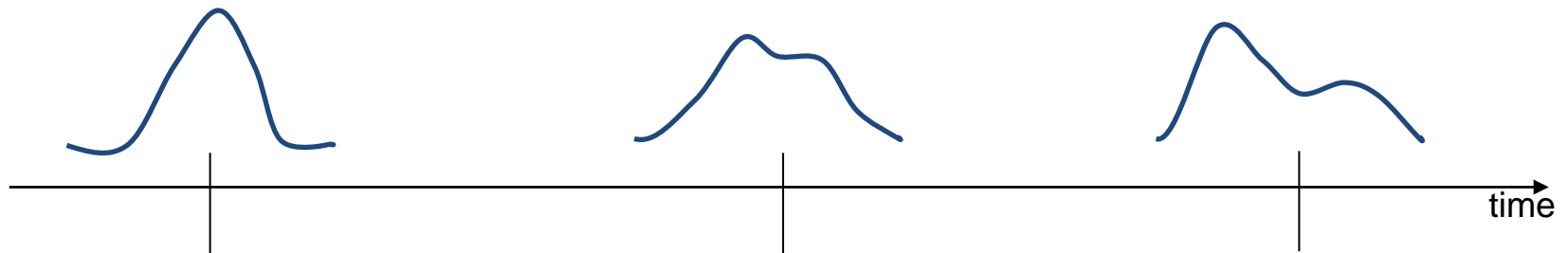Exploitation  vs  Exploration

# Overview - Problems

Delayed Reward

# Overview - Problems

Credit Assignment Problem

# Overview - Problems

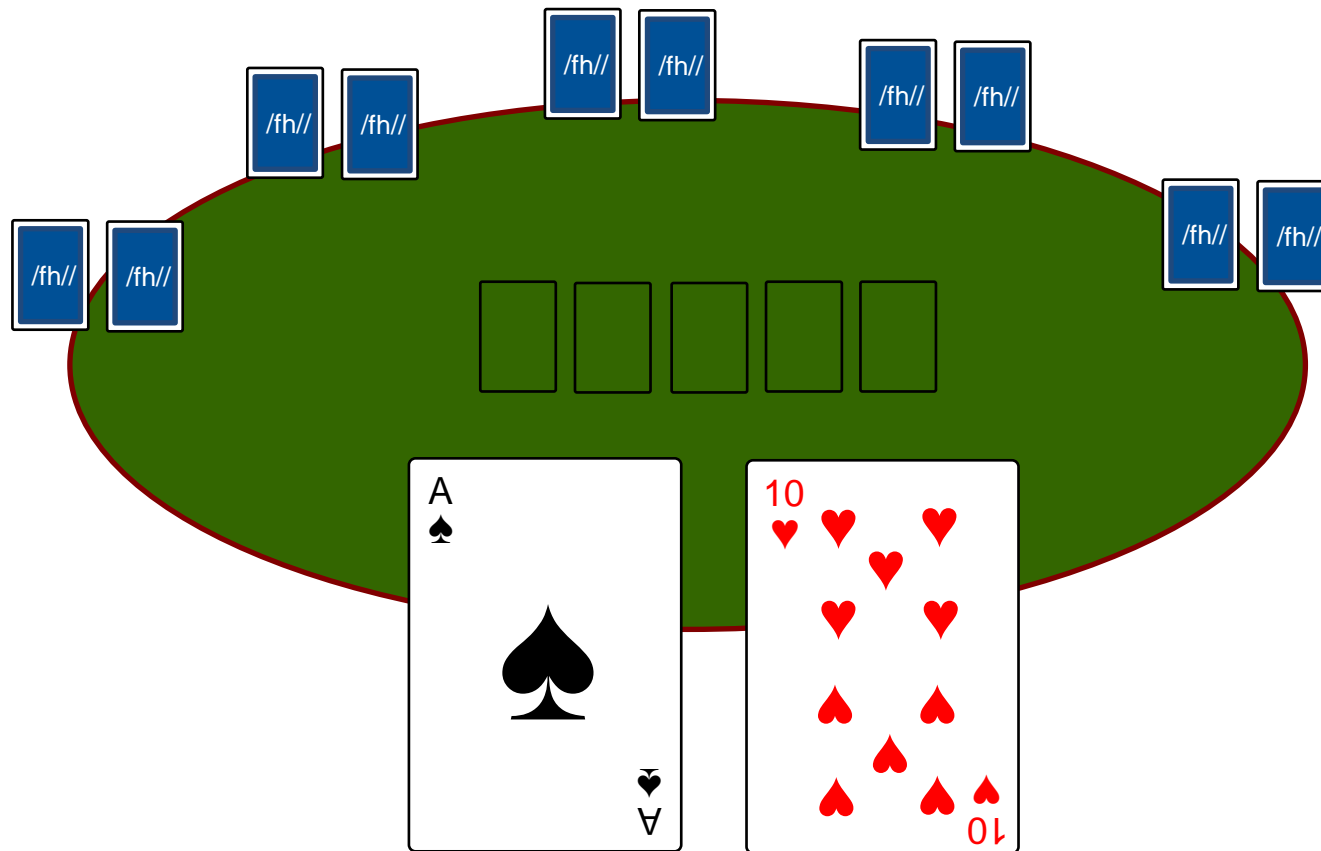Non-Stationarity

# Overview - Problems

Partial Observability

# Overview - Problems

Multi-Agent



Image: https://deepmind.com/blog/article/capture-the-flag-science

# Milestones of RL

*/ informatik & security* **/fh///**
st.pölten

Checkers Player
*Arthur Samuel 1959*

TD-Algorithm
*Richard Sutton*

Atari 2600
*DeepMind 2013*

1950    1960    1970    1980    1990    2000    2010    2020

Dynamic Programming
*Richard Bellmann 1953*

Actor-Critic Architectur
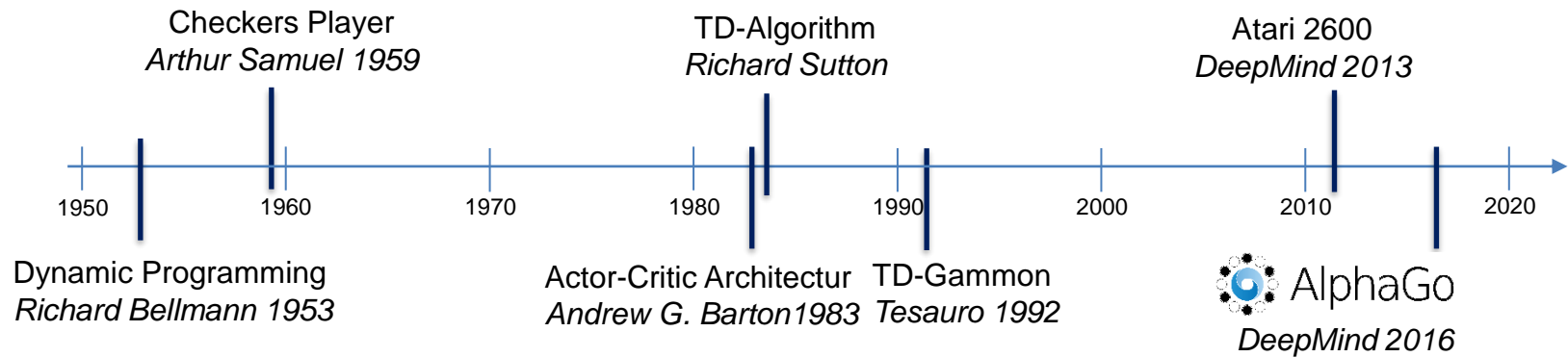*Andrew G. Barton1983*

TD-Gammon
*Tesauro 1992*

# Atari Games

**Starting out - 10 minutes of training**

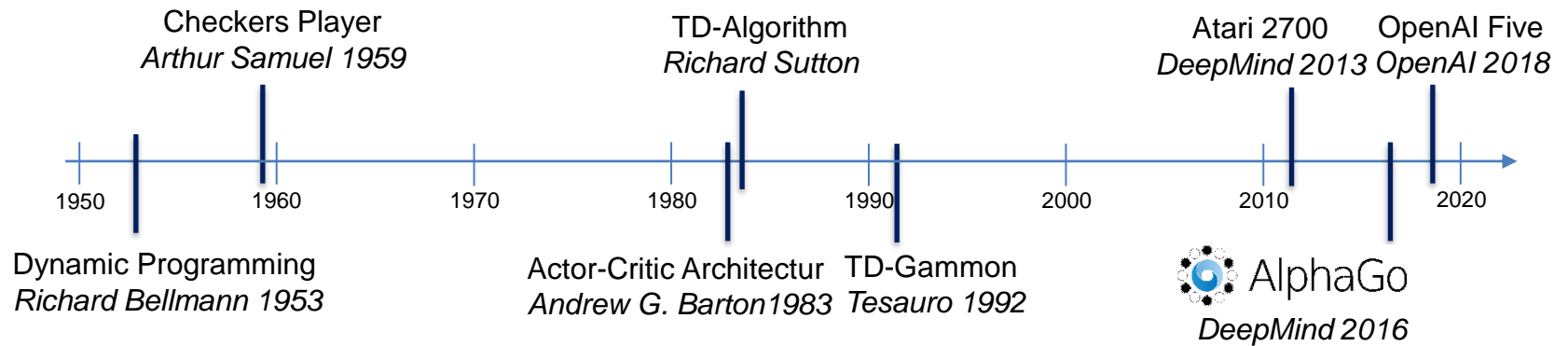**The algorithm tries to hit the ball back, but it is yet too clumsy to manage.**

Source: Two Minute Papers https://www.youtube.com/watch?v=V1eYniJ0Rnk

# Milestones of RL

/informatik & security **/fh///**
st. pölten

Checkers Player
*Arthur Samuel 1959*

TD-Algorithm
*Richard Sutton*

Atari 2600
*DeepMind 2013*

1950    1960    1970    1980    1990    2000    2010    2020

Dynamic Programming
*Richard Bellmann 1953*

Actor-Critic Architectur
*Andrew G. Barton1983*

TD-Gammon
*Tesauro 1992*

AlphaGo

*DeepMind 2016*

# AlphaGo – AlphaZero

Image: https://cdn-images-1.medium.com/max/1600/0*J2HH4TDqGCuKbufS.png

# Milestones of RL

Checkers Player
*Arthur Samuel 1959*

TD-Algorithm
*Richard Sutton*

Atari 2700
*DeepMind 2013*

OpenAI Five
*OpenAI 2018*

1950    1960    1970    1980    1990    2000    2010    2020

Dynamic Programming
*Richard Bellmann 1953*

Actor-Critic Architectur
*Andrew G. Barton1983*

TD-Gammon
*Tesauro 1992*

AlphaGo
*DeepMind 2016*

# Open AI 5



Source: Youtube https://www.youtube.com/watch?v=yBEidvm_tZQ

# Milestones of RL

Checkers Player
*Arthur Samuel 1959*

TD-Algorithm
*Richard Sutton*

Atari 2700
*DeepMind 2013*

OpenAI Five
*OpenAI 2018*

1950    1960    1970    1980    1990    2000    2010    2020

Dynamic Programming
*Richard Bellmann 1953*

Actor-Critic Architectur
*Andrew G. Barton1983*

TD-Gammon
*Tesauro 1992*

AlphaGo
*DeepMind 2016*

AlphaStar
*DeepMind 2019*

# AlphaStar



Source: Youtube https://www.youtube.com/watch?v=cUTMhmVh1qs

# Milestones of RL

Checkers Player
*Arthur Samuel 1959*

TD-Algorithm
*Richard Sutton*

Atari 2700
*DeepMind 2013*

OpenAI Five
*OpenAI 2018*

1950    1960    1970    1980    1990    2000    2010    2020

Dynamic Programming
*Richard Bellmann 1953*

Actor-Critic Architectur
*Andrew G. Barton1983*

TD-Gammon
*Tesauro 1992*

AlphaGo
*DeepMind 2016*

AlphaStar
*DeepMind 2019*

# Multi-armed Bandits

Exploitation vs Exploration

# INTERACTIVE PART I

# Multi-armed Bandits

**n** choices of **k** different actions

After each choice: numerical reward from <u>stationary distribution</u>

Objective: maximize total reward over n choices

# Multi-armed Bandits

$A_t \ldots$ Action at timestep t

$R_t \ldots$ Reward at timestep t

optimal action-value function:

(for the multi-armed bandit problem)

$$q_*(a) := \mathbb{E}[R_t | A_t = a]$$

We usually don't know that value!

$q_t(a) \ldots$ Estimate of $q_*(a)$ at timestep t

# INTERACTIVE PART II

# Multi-armed Bandits

Estimator of action-value function:

$$q_t(a) := \frac{\text{sum of rewards when } a \text{ taken before } t}{\text{number of times } a \text{ taken before } t} = \frac{\sum_{i=1}^{t-1} R_i \cdot \mathbb{1}_{[A_i=a]}}{\sum_{i=1}^{t-1} \mathbb{1}_{[A_i=a]}}$$

Indicator function:

$$\mathbb{1}_{[A_i=a]}(a) := \begin{cases} 1 & \text{if } A_i = a \\ 0 & \text{if } A_i \neq a \end{cases}$$

complicated way of saying:
average of rewards

# Action Selection Methods

_Random_ action selection method:


Randomly select one of the possible actions.
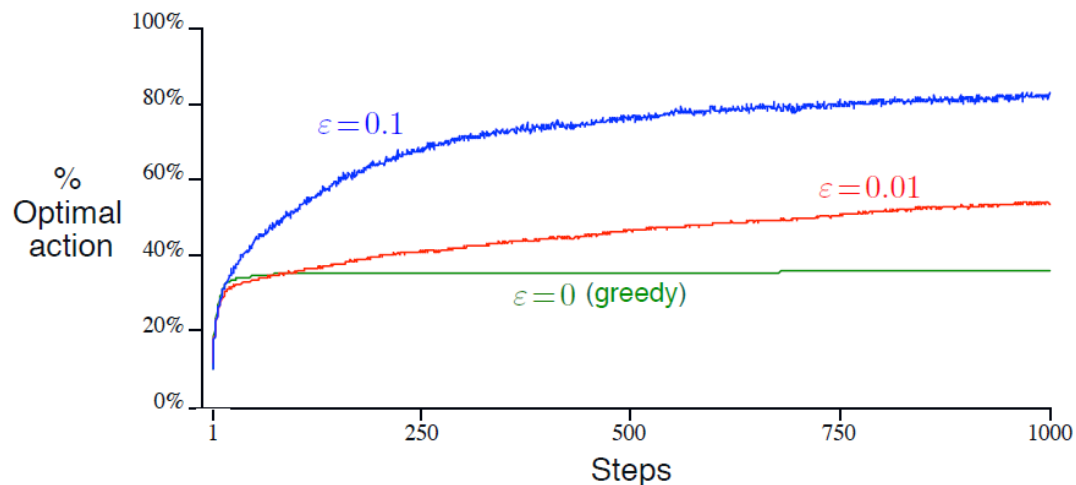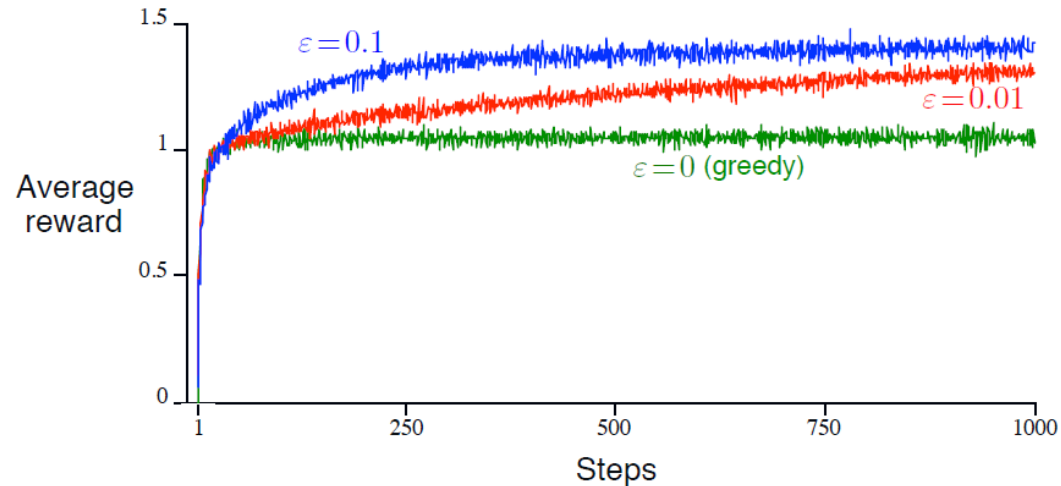
# Action Selection Methods

Graph of a probability mass function (pmf) of the distribution:

A                                       20 %

B                                         20 %

C                                         20 %

D                                         20 %

E                                         20 %

## No exploitation - Only exploration

# Action Selection Methods

> *Greedy* action selection method:
>
> Always select the action with the highest action-value according to the action-value function.
>
> $$A_t = \operatorname*{argmax}_a q_t(a)$$

# Action Selection Methods

Graph of a probability mass function(pmf) of the distribution:



A     0 %

B     0 %

C     0 %

D     100 %

E     0 %

Only exploitation - no exploration

# Action Selection Methods

$\varepsilon$-*Greedy* action selection method:

In most $(1-\varepsilon)$ cases select the best action an in a small amount of cases select a random action.

$$A_t = \begin{cases} \underset{a}{\mathrm{argmax}}\, q_t(a) & \text{in } (100-\varepsilon) \text{ \% cases} \\ \text{random action} & \text{in } \varepsilon \text{ \% cases} \end{cases}$$

Mostly exploitation - small exploration

# Action Selection Methods

Graph of a probability mass function (pmf) of the distribution:

| | |
|---|---|
| A | 1 % |
| B | 96 % |
| C | 1 % |
| D | 1 % |
| E | 1 % |

Mostly exploitation - small exploration

# Multi-armed Bandits

# Incremental Implementation

Simplified: only one action

$$q_n = \frac{R_1 + R_2 + \cdots + R_n}{n}$$

$$
\begin{aligned}
q_n \quad &= \frac{1}{n} \sum_{i=1}^{n} R_i \\
&= \frac{1}{n} \left( R_n + \sum_{i=1}^{n-1} R_i \right) \\
&= \frac{1}{n} \left( R_n + (n-1) \frac{1}{n-1} \sum_{i=1}^{n-1} R_i \right) \\
&= \frac{1}{n} \left( R_n + (n-1) Q_{n-1} \right) \\
&= \frac{1}{n} \left( R_n + n Q_{n-1} - Q_{n-1} \right) \\
&= Q_{n-1} + \frac{1}{n} \left( R_n - Q_{n-1} \right)
\end{aligned}
$$

# Incremental Implementation

$$q_{n+1} = q_n + \frac{1}{n+1}(R_n - q_n)$$

$$NewEstimate \leftarrow OldEstimate + StepSize[\underbrace{Target - OldEstimate}_{Error}]$$

Error

# Incremental Implementation

$$NewEstimate \leftarrow OldEstimate + StepSize[Target - OldEstimate]$$

Error

$R_1 = 1{,}93 \quad Q_1 = 0 + \dfrac{1}{1}[1{,}93 - 0] = 1{,}93$

$R_2 = 2{,}35 \quad Q_2 = 1{,}93 + \dfrac{1}{2}[2{,}35 - 1{,}93] = 2{,}14$

$R_3 = 3{,}32 \quad Q_3 = 2{,}14 + \dfrac{1}{3}[3{,}32 - 2{,}14] = 2{,}53$

# Incremental Implementation

$$NewEstimate \leftarrow OldEstimate + StepSize[Target - OldEstimate]$$

Error

$R_4 = 2{,}43 \quad Q_4 = 2{,}53 + \dfrac{1}{4}[2{,}43 - 2{,}53] = 2{,}505$

$R_5 = 3{,}28 \quad Q_5 = 2{,}505 + \dfrac{1}{5}[3{,}28 - 2{,}505] = 2{,}66$

$R_6 = 3{,}91 \quad Q_6 = 2{,}66 + \dfrac{1}{6}[3{,}91 - 2{,}66] = 2{,}868$

# Incremental Implementation

$$NewEstimate \leftarrow OldEstimate + StepSize[Target - OldEstimate]$$

Error

$R_7 = 2,34$ $\quad Q_7 = 2,868 + \dfrac{1}{7}[2,34 - 2,868] = 2,793$

$R_8 = 3,42$ $\quad Q_8 = 2,793 + \dfrac{1}{8}[3,42 - 2,793] = 2,871$

$R_9 = 2,64$ $\quad Q_9 = 2,871 + \dfrac{1}{9}[2,64 - 2,871] = 2,845$

# Incremental Implementation

$$NewEstimate \leftarrow OldEstimate + StepSize[Target - OldEstimate]$$

Error

# Optimistic initial values

# Nonstationary Problem

Constant step-size Parameter $\alpha\epsilon(0,1]$ $\qquad Q_{n+1} := Q_n + \alpha(R_n - Q_n)$

$$Q_{n+1} = Q_n + \alpha(R_n - Q_n)$$
$$= \alpha R_n + (1 - \alpha)Q_n$$
$$= \alpha R_n + (1 - \alpha)[\alpha R_{n-1} + (1 - \alpha)Q_{n-1}]$$
$$= \alpha R_n + (1 - \alpha)\,\alpha R_{n-1} + (1 - \alpha)^2 Q_{n-1}$$
$$= \alpha R_n + (1 - \alpha)\,\alpha R_{n-1} + (1 - \alpha)^2 \alpha R_{n-2} + \cdots + (1 - \alpha)^{n-1}\alpha R_1 + (1 - \alpha)^n Q_1$$

$$= (1 - \alpha)^n Q_1 + \sum_{i=1}^{n} \alpha(1 - \alpha)^{n-i} R_i$$

# Contextual Bandits

- Synonym for the full reinforcement learning problem

- Multiple situations in contrast to only one

- Action also affects the next situation not only the reward (Return)

# Multi-armed Bandits

- In which subfields can RL be devided into? What are the major differences?

- What is the mulit-armed bandit problem? What is the difference to the full RL problem?

- What is an acition-value function? How is it's estimator defined?

- What is the incremental notation of an action-value function estimator?

- Which action selecton methods exist?

# MDP

Observation

Reward

**Agent**

Policy

**Environment**

State

Action

# MDP

**Environment**

State

# MDP – Example

Tic-Tac-Toe

# MDP – Example

State Space:
$$\mathcal{S} = \{\ s_1\ ,\ s_2\ ,\ s_3\ ,\ s_4\ ,\ s_5\ , \dots\}$$

Action Space:
$$\mathcal{A} = \{\ a_1,\ a_2,\ a_3,\ a_4,\ a_5,$$
$$a_6,\ a_7,\ a_8,\ a_9\ \}$$

Reward Function:
$$\begin{cases} +1 & \text{if} \quad , \quad , \quad , \text{or} \\ -1 & \text{if} \quad , \quad , \quad , \text{or} \\ 0 & \text{else} \end{cases}$$

Time Step t

# MDP – Example



Time Step t

# MDP – Example

Time Step t

3

$s_{93}$  $s_{107}$  $s_{111}$  $s_{124}$

$a_3$  $a_9$  $a_8$

4

$s_{89}$  $s_{115}$  $s_{103}$  $s_{120}$  $s_{109}$  $s_{131}$

$a_9$  $a_3$

5

$s_{135}$  $s_{142}$

0      1      2      3      4

# MDP – Example

|   0   |     1     |   2    |     3     |   4    |
|-------|-----------|--------|-----------|--------|



$$S_0 \qquad A_0 \quad R_1 \quad S_1 \qquad A_1 \quad R_2 \quad S_2 \qquad A_2 \quad R_3 \quad S_3 \qquad A_3 \quad R_4 \quad S_4$$

# MDP – Example

$S_t, A_t, R_t$ are random variables

Like dice rolls or coin flips, they can have different results for separate executions.
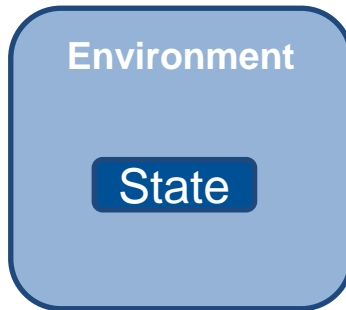
$$S_3 = s_{124} \qquad\qquad A_2 = a_2$$

$$S_3 = \ \ \qquad\qquad A_2 = $$

$$\mathbb{p}(S_{t+1} = s_{124} | S_t = s_{75}, A_t = a_4) =?$$

# MDP

# MDP – Definition

**Environment**

**State**

$\approx$

Markov Decision Process

State Space: $\mathcal{S} = \{s_1, s_2, \dots, s_N\}$

Action Space: $\mathcal{A} = \{a_1, a_2, \dots, a_M\}$

Set of Rewards: $\mathcal{R} \subset \mathbb{R}$

State-Transition
Probability Function: $p: \mathcal{S} \times \mathcal{S} \times \mathcal{A} \rightarrow [0,1]$

Reward Function: $R: \mathcal{R} \times \mathcal{S} \times \mathcal{A} \rightarrow [0,1]$

# MDP – Transition Function

State-Transition Probability Function: $p: \mathcal{S} \times \mathcal{S} \times \mathcal{A} \rightarrow [0,1]$

$$p(s', s, a) \coloneqq \mathbb{p}(S_{t+1} = s' | S_t = s, A_t = a)$$

- Defines how the <u>environment</u> behaves

- Function is usually not known

- Tic-tac-toe: contains rules of the games, as well as the opponents behavior

# MDP – Reward Function
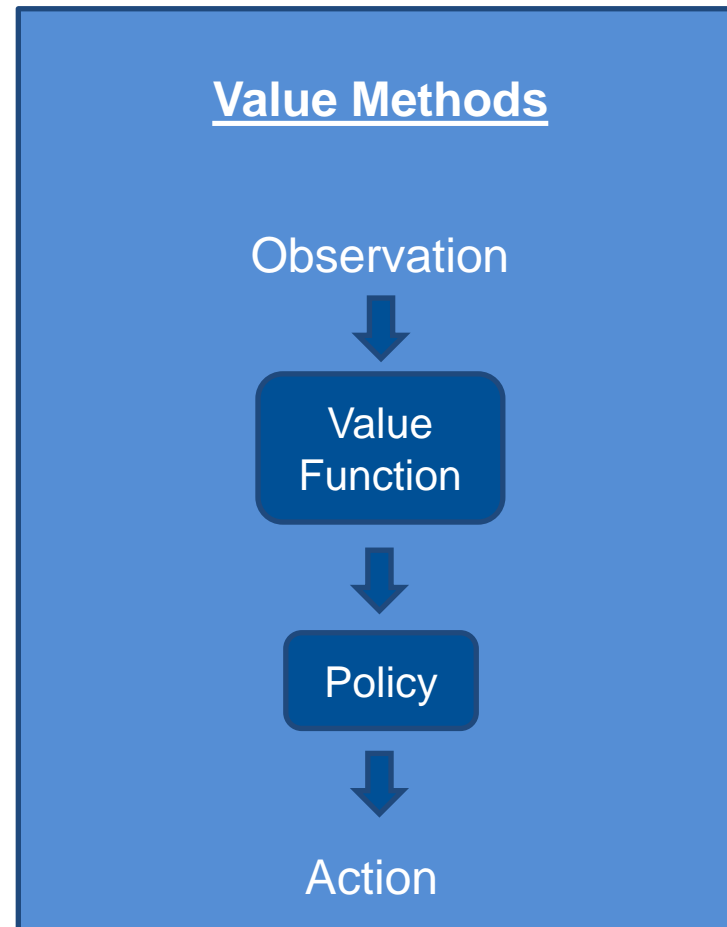
Reward Function: $\qquad R: \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$

$$R(s, a, s')$$

- Defines what the environment rewards

- Function is usually designed by a human and therefore known

- Tic-tac-toe: contains rules of the games, as well as the opponents behavior

# MDP – Markov Property

The probability of each possible value of $S_t$ and $R_t$ depends **only** on the **immediately preceding** state and action and **not on earlier** states and actions.

„The future only depends on the present and not on the past"

$$\mathbb{p}(S_{t+1} = s' \,|S_t = s, A_t = a) =$$
$$\mathbb{p}(S_{t+1} = s' \,|S_t = s, A_t = a, \dots, S_0 = s'', A_0 = a'', )$$

# VF

**Value Methods**

Observation

⬇

Value Function

⬇

Policy

⬇

Action

# VF – Policy

Policy: $\pi: \mathcal{S} \times \mathcal{A} \rightarrow [0,1]$

- defines how the <u>agent</u> behaves

- Rows are probability functions over the action space

- Example: random policy

| $\pi$ | $a_1$ | $a_2$ | $a_3$ | $a_4$ |
|-------|-------|-------|-------|-------|
| $s_1$ | 0.11 | 0.8 | 0.06 | 0.03 |
| $s_2$ | 0.01 | 0.1 | 0.7 | 0.19 |
| $s_3$ | 1.0 | 0.0 | 0.0 | 0.0 |
| $s_4$ | 0.65 | 0.32 | 0.01 | 0.02 |
| $s_5$ | 0.25 | 0.25 | 0.25 | 0.25 |
| $s_6$ | 0.4 | 0.1 | 0.2 | 0.4 |

# VF – Value Functions

Value Functions:
$$V : \mathcal{S} \to \mathbb{R}$$
$$Q : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$$

- indicate how „good" a state or an action given a state is

$$V(\ \ ) = 0.64$$

$$Q(\ \ , \ \ ) = 1.0$$

| | $V$ |
|---|---|
| $s_1$ | 2.3 |
| $s_2$ | -6.8 |
| $s_3$ | 10.0 |
| $s_4$ | -5.4 |
| $s_5$ | 4.1 |
| $s_6$ | -0.4 |

| $Q$ | $a_1$ | $a_2$ | $a_3$ | $a_4$ |
|---|---|---|---|---|
| $s_1$ | -0.1 | 4.3 | 0.3 | 1.8 |
| $s_2$ | 7.5 | 3 | 8.2 | -5.7 |
| $s_3$ | 0.0 | 10.0 | 0.0 | 0.0 |
| $s_4$ | 6.5 | 0.32 | 0.01 | -0.2 |
| $s_5$ | 3.0 | 0.21 | -7.2 | 0.25 |
| $s_6$ | 4.5 | 0.1 | -2.0 | 0.4 |

# Value Function

# Value Function
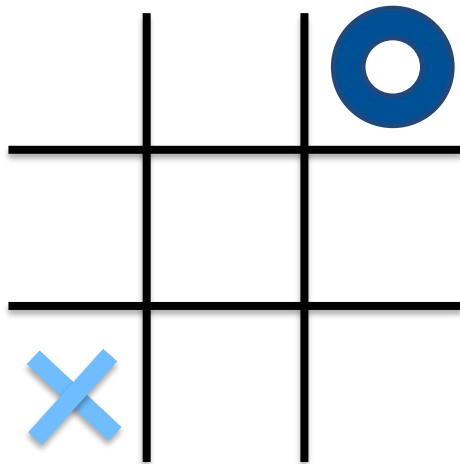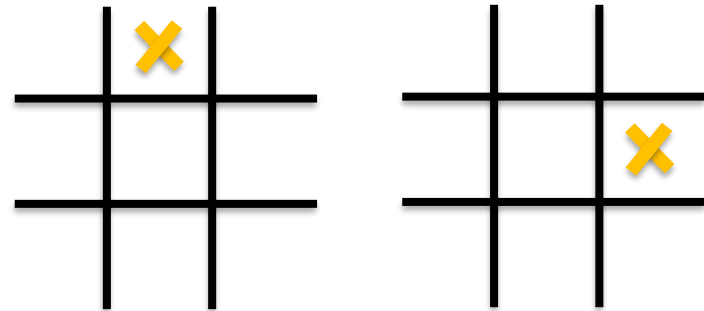
# Value Function

# Value Function

# Value Function

# Value Function

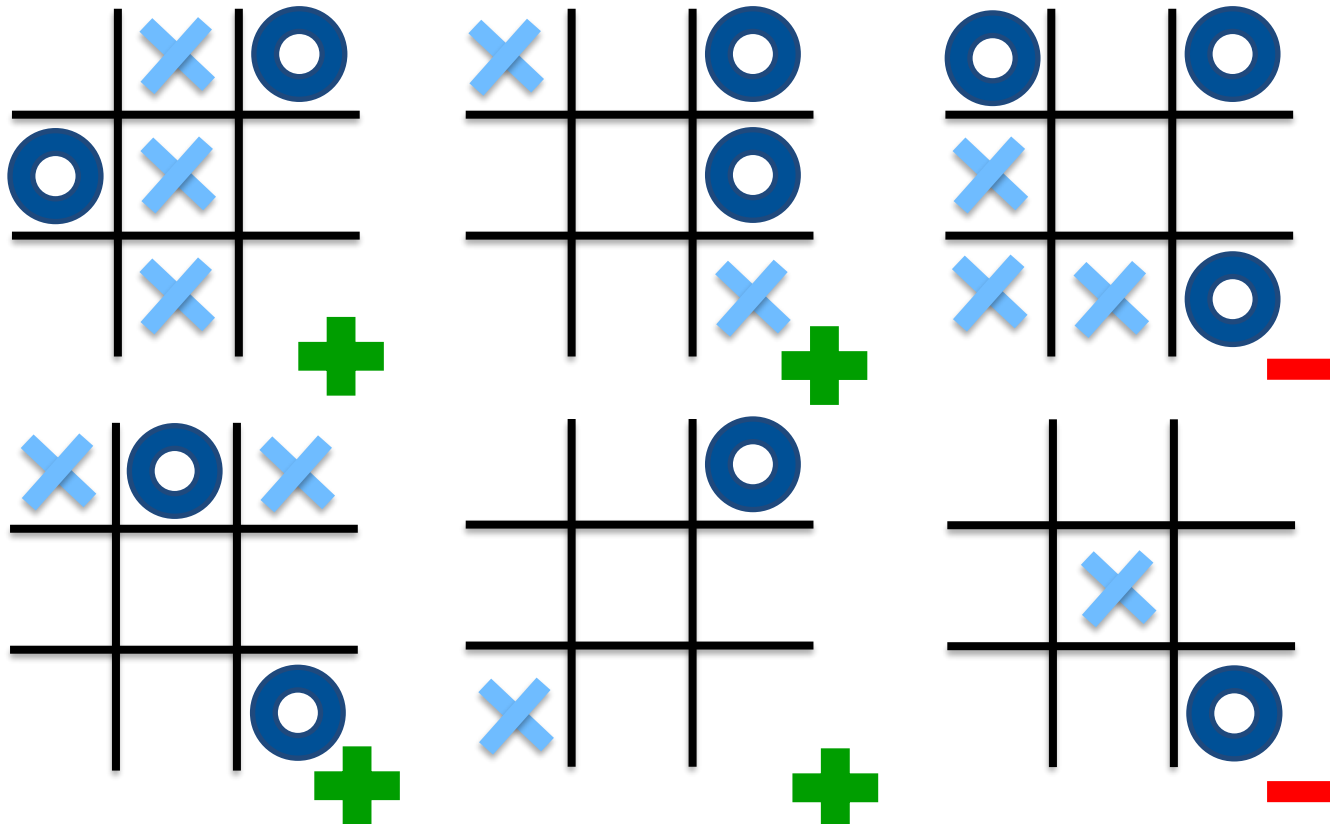State:



Greedy Actions:

# VF – Motivation

You are: ✖         Your opponent is: ⬤

It is your turn!

# VF – Motivation
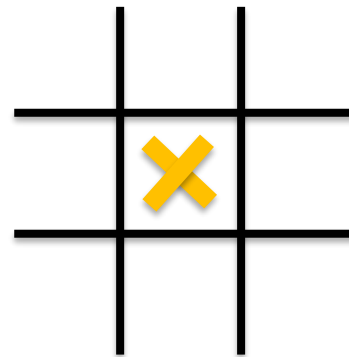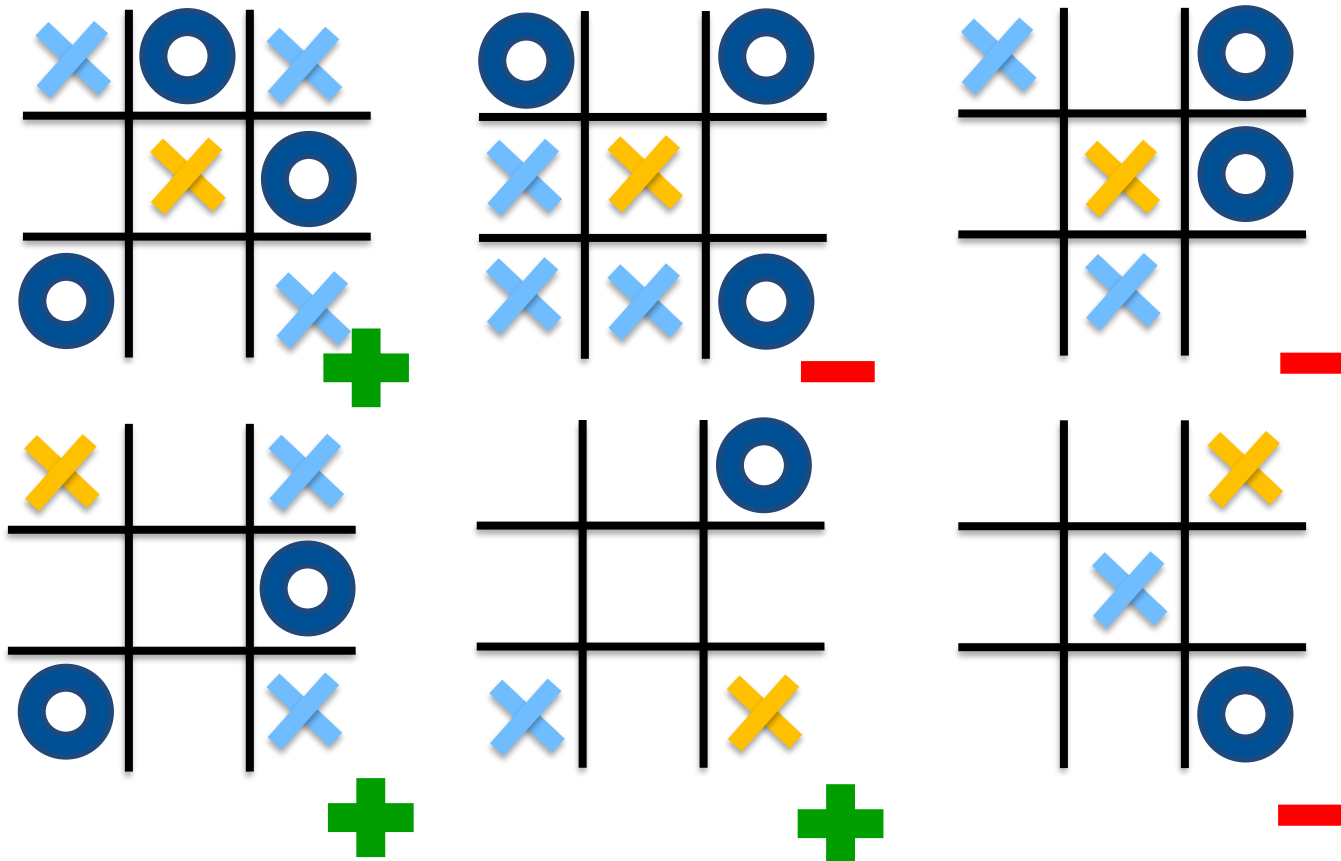
You are: ✕

Your opponent is: ○

It is your turn!

indicates the next move

Which Value Function to choose?

# VF

Naive approach:
5 outcomes
3 positive

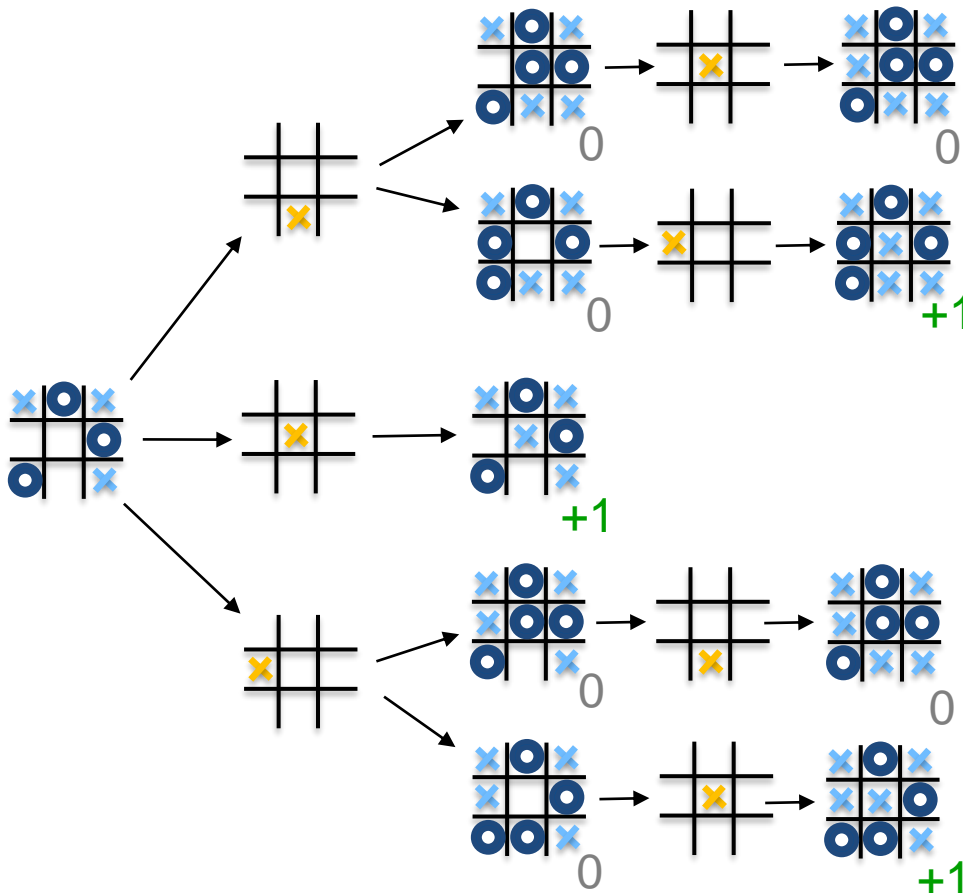$$\Rightarrow \quad V(\text{board}) = \frac{3}{5} = 0.6$$

What if intermediate reward is given?

$\Rightarrow$ Need to add all reward until the end

$$G_t := \sum_{n=t+1}^{T} \gamma^n R_n$$

is also a random variable

# VF

Naive approach:
5 outcomes
3 positive

$$\Rightarrow \quad V(\text{⊞}) = \frac{3}{5} = 0.6$$

Implies equal probability for all outcomes.
What if agent only chooses actions from the second column if possible?

$$\Rightarrow \quad V(\text{⊞}) = \frac{2}{3} = 0.667$$

Value function depends on current policy (which we know)

$$\Rightarrow V_\pi(s)$$

# VF

Value function also depends on the environment (which we don't know)

=> Sampling

Sampling targets:

$$V_\pi(s) \approx \mathbb{E}_\pi[G_t | S_t = s]$$

$$Q_\pi(s, a) \approx \mathbb{E}_\pi[G_t | S_t = s, A_t = a]$$

# VF – Monte Carlo Method

1) Initialize with random policy and 0 for value funciton $Q(s,a)$

| $Q$ | $a_1$ | $a_2$ | $a_3$ | $a_4$ |
|-----|-------|-------|-------|-------|
| $s_1$ | 0.0 | 0.0 | 0.0 | 0.0 |
| $s_2$ | 0.0 | 0.0 | 0.0 | 0.0 |
| $s_3$ | 0.0 | 0.0 | 0.0 | 0.0 |
| $s_4$ | 0.0 | 0.0 | 0.0 | 0.0 |
| $s_5$ | 0.0 | 0.0 | 0.0 | 0.0 |
| $s_6$ | 0.0 | 0.0 | 0.0 | 0.0 |

| $\pi$ | $a_1$ | $a_2$ | $a_3$ | $a_4$ |
|-------|-------|-------|-------|-------|
| $s_1$ | 0.25 | 0.25 | 0.25 | 0.25 |
| $s_2$ | 0.25 | 0.25 | 0.25 | 0.25 |
| $s_3$ | 0.25 | 0.25 | 0.25 | 0.25 |
| $s_4$ | 0.25 | 0.25 | 0.25 | 0.25 |
| $s_5$ | 0.25 | 0.25 | 0.25 | 0.25 |
| $s_6$ | 0.25 | 0.25 | 0.25 | 0.25 |

2) Play episode with current policy $\pi$

# VF – Monte Carlo Method

2) Play episode with current policy $\pi$



3) Iterate episode from back to front, calculate the return $G_t$ and update the value function

4) Make policy $\epsilon$-greedy with regard to value function

5) continue with step 2)