

# Deep Reinforcement Learning

## *2 - Multi-Armed Bandits to Sequences of Decisions*

Prof. Dr. Malte Schilling

Autonomous Intelligent Systems Group

# Overview Lecture

- Exploration-Exploitation Tradeoff
- Decision Making: Multi-Armed Bandit
  - Strategies
    - $\epsilon$ -greedy
    - UCB
    - Gradient-based Action Selection
- Sequences of States – towards Sequential Decision Making

# Admin - Exercises

- Exercise slot is Friday morning, 10:15 AM to 11:45, M5.
- Solutions (for programming exercises in python) should be submitted by small teams of two or three persons. Everybody has to be able to present the solution during meetings.
- First exercise sheet out today – due next Wednesday (upload in learnweb).

Register groups: Send an email to [malte.schilling@uni-muenster.de](mailto:malte.schilling@uni-muenster.de) with name and account information for each group member until Tuesday, 25.10.2022!

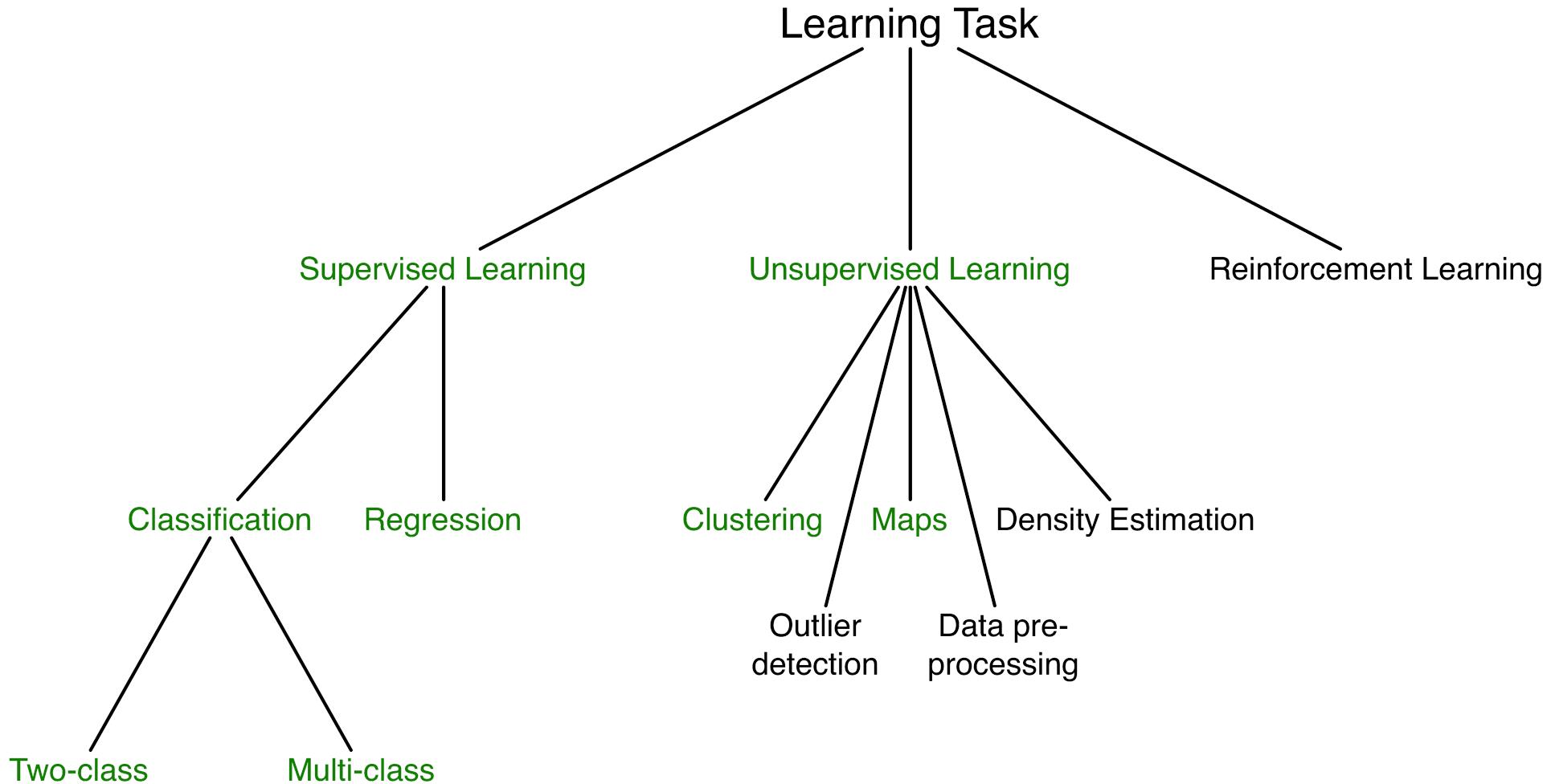
Join the learnweb course!

# Übungstermin Freitag, 21.10.2022

Inhalte:

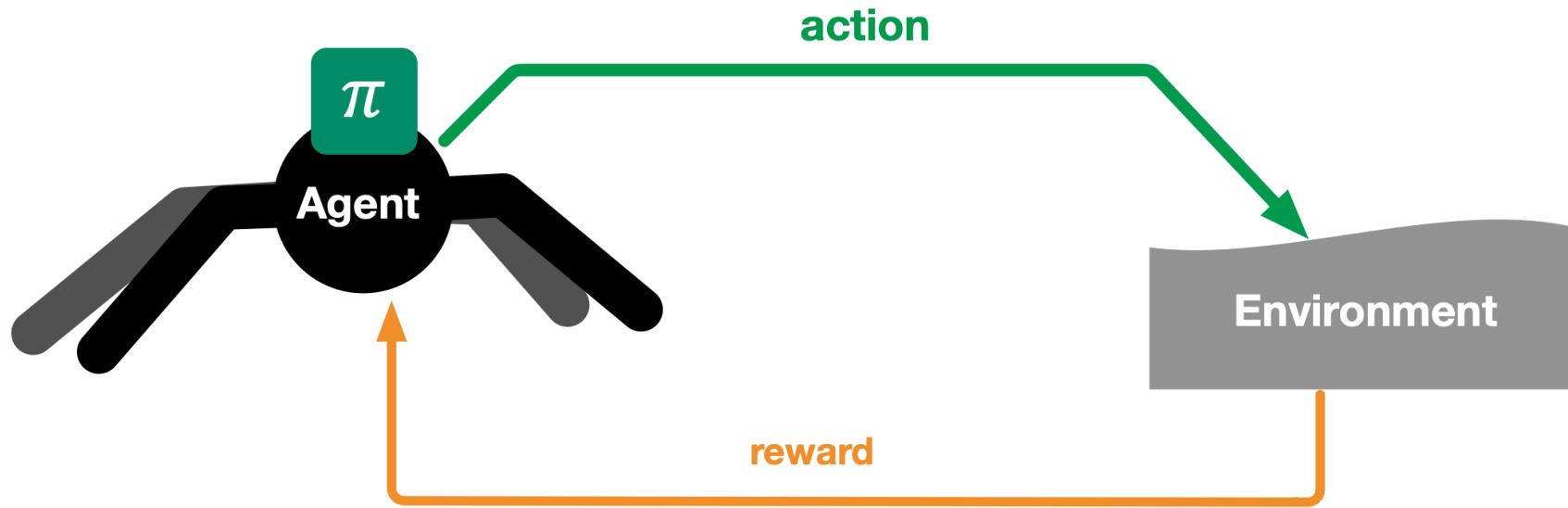
- Troubleshooting für Python Aufgaben
- Basierend auf dem Übungszettel: zusammenstellen der notwendigen numpy Funktionen (arrays, grundlegende statistische Funktionen, ...) und plot Funktionen (für eine random policy).

# Recap – Types of Learning



*Distinction of different types of learning*

# Parts of Decision Making



## Agent

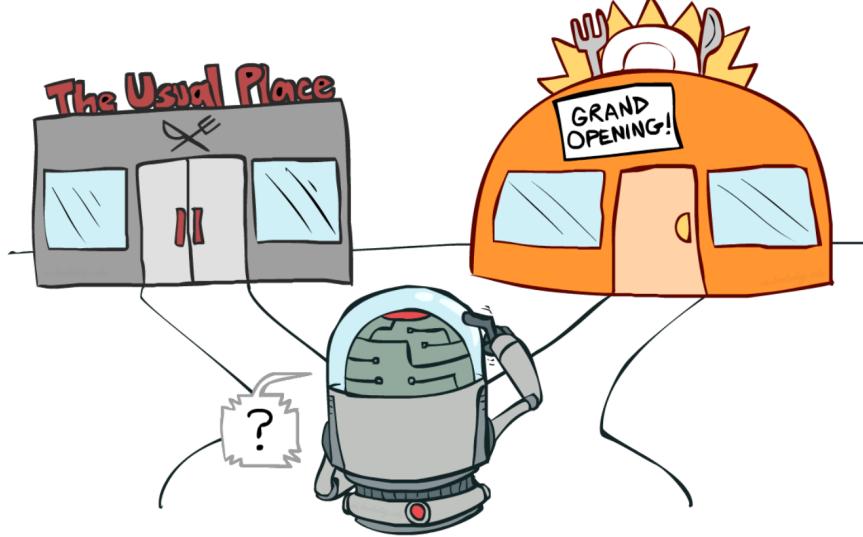
- A **policy** is the agent's behavior – chooses an action.

## Environment

- Provides reward

# Examples of Rewards

- Fly stunt manoeuvres in a helicopter
  - positive reward: if following a desired trajectory
  - negative reward: when crashing
- Manage an investment portfolio – reward is given as the money in that account
- playing computer games – reward directly given as score
- Locomotion of a robot
  - positive reward for movement in the correct direction
  - negative reward / cost: falling over (not maintaining height); energy consumed



Decision Making: sticking to a good past experience might make you miss out on even better options, but at least you can be confident to get something good.



- Homework answers: What is your example of an exploration-exploitation tradeoff?
- What would be the reward?

# Exploitation-Exploration Tradeoff

As information of a novel environment is incomplete, we need to gather information for good decisions and want to keep the risk under control.

## Exploitation

Taking advantage of the best known option.

An optimal long-term strategy may involve short-term sacrifices, e.g. learning from failure during exploration helps us avoid a certain action.

## Exploration

Take some risk to collect information about unknown options.

# Examples for Exploitation-Exploration Tradeoff

- Restaurant Selection
  - Exploitation: Go to your favourite restaurant
  - Exploration: Try a new restaurant
- Oil Drilling
  - Exploitation: Drill at the best known location
  - Exploration: Drill at a new location
- Game Playing
  - Exploitation: Play the move you believe is best
  - Exploration: Play an experimental move

# Rewards

- A **reward**  $R_t$  is a scalar feedback given to the agent from the environment
- that indicates how well the agent is doing (at time t).
- The agent aims to maximise the cumulative reward over time.

## Reward Hypothesis

All goals can be represented as maximization of a scalar reward (an expected cumulative reward).

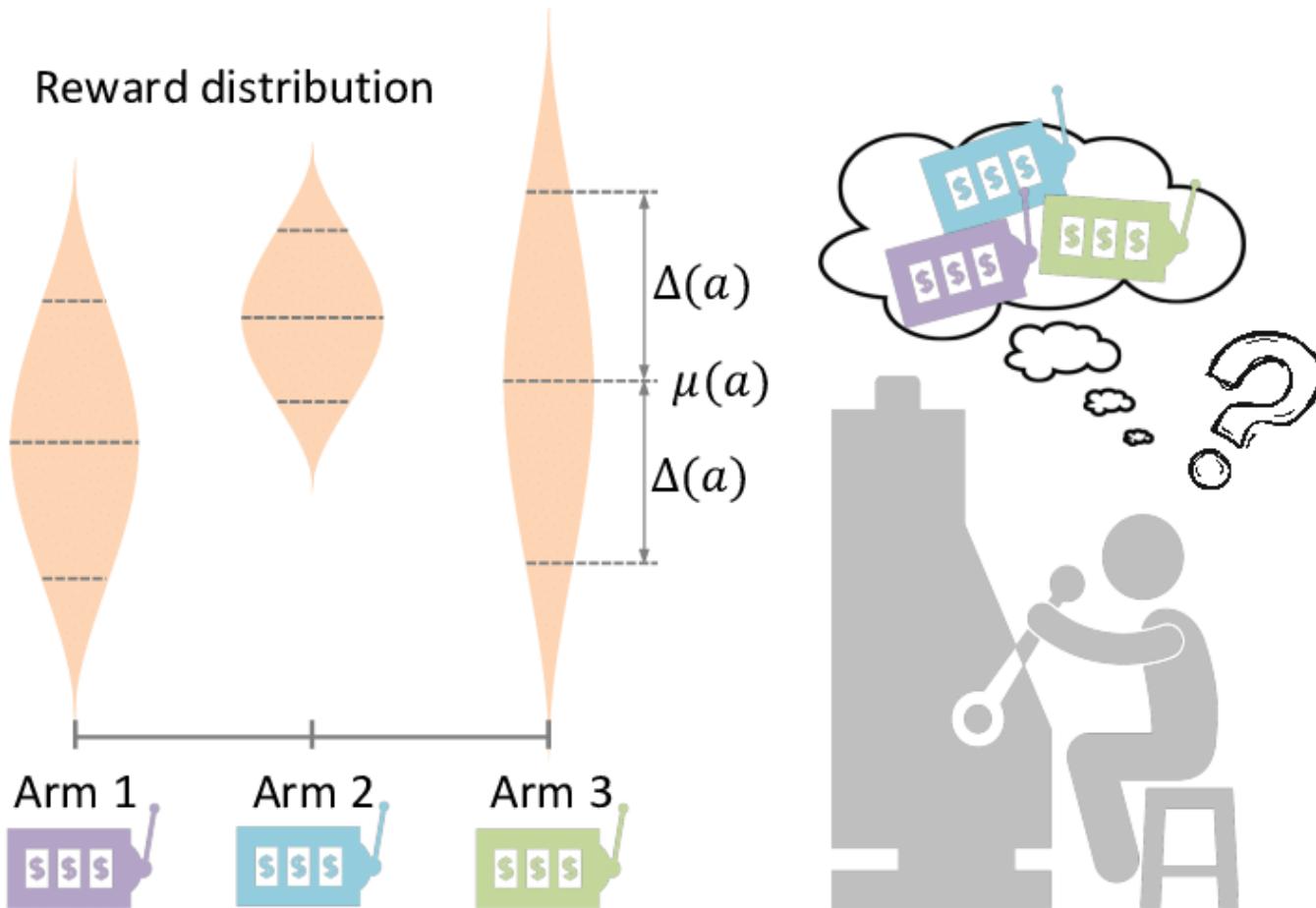
### Return (preliminary Def.)

The accumulated reward over time is called the return  $G_t$ .

$$G_t = R_{t+1} + R_{t+2} + \dots$$

# Multi-Armed Bandit

# Multi-Armed Bandit Overview

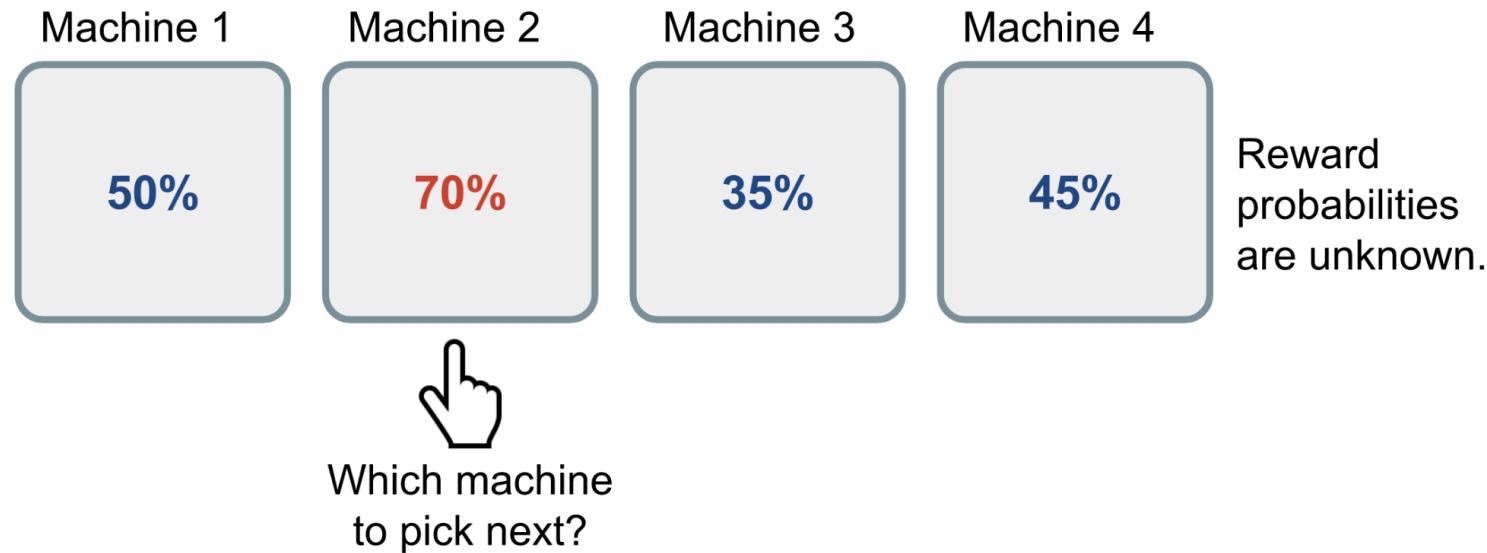


(Gao u. a. 2021)

# Multi-Armed Bandit

Idea: in a casino with multiple slot machines of unknown probabilities.

Which action (slot machine) should you choose for optimal reward?



Following a naive approach, one would gather information over a long time to get a true estimate of each of the probabilities. But as a consequence one will spend too much time on suboptimal actions.

# Multi-Armed Bandit

A Bernoulli multi-armed bandit can be described as a tuple of  $\langle \mathcal{A}, \mathcal{R} \rangle$ , where:

- $K$  machines with reward probabilities,  $\theta_1, \dots, \theta_K$
- for each time step  $t$ , take an action  $a$  on one slot machine and receive a reward  $r$ .
- $\mathcal{A}$  is a set of actions (one for each slot machine) – the value of action  $a$  is the expected reward  $Q(a) = \mathbb{E}[r_t|a] = \theta$  (when at time  $t$  action  $a_t$  is choosing the  $i$ -th machine  $Q(a_t) = \theta_i$ ).
- $\mathcal{R}$  is the reward function (a distribution on rewards). For a Bernoulli bandit, we observe a reward in a stochastic fashion ( $r_t = \mathcal{R}(a_t)$  may return 1 with probability  $Q(a_t)$ ).

This is a simplified version of a Markov decision process (there is no state  $\mathcal{S}$ ).

# Parts of Decision Making



## Agent

- A **policy** is the agent's behavior – chooses an action.
- **Value function:** Keep track of the value of an action.

## Environment

- Provides reward

# Action values

## Action value

The action value for action  $a$  is the expected reward

$$Q_t(a) = \mathbb{E}[R_t \mid A_t = a]$$

A simple estimate is the average of the sampled rewards:

$$\begin{aligned} Q_t(a) &= \frac{\text{sum of rewards when } a \text{ taken prior to } t}{\text{number of times } a \text{ taken prior to } t} \\ &= \frac{\sum_{i=1}^{t-1} R_i \cdot \mathbb{1}_{A_i=a}}{\sum_{i=1}^{t-1} \mathbb{1}_{A_i=a}} \end{aligned}$$

with  $\mathbb{1}$  being the indicator function, therefore we count choosing action  $a$  as

$$\sum_{i=1}^t \mathbb{1}(A_i = a)$$

# Maximization of cumulative reward

Goal is to maximize **cumulative reward**  $\sum_{t=1}^T r_t$  (the return  $G$ ).

The optimal action produces the maximal reward. Deviating from that action leads to a potential loss or **regret**.

The probability for the optimal reward  $\theta^*$  of the optimal action  $a^*$  is

$$\theta^* = Q(a^*) = \max_{a \in \mathcal{A}} Q(a) = \max_{1 \leq i \leq K} \theta_i$$

# Regret

The regret of an action  $a$  is given as

$$\Delta_a = v_* - q(a)$$

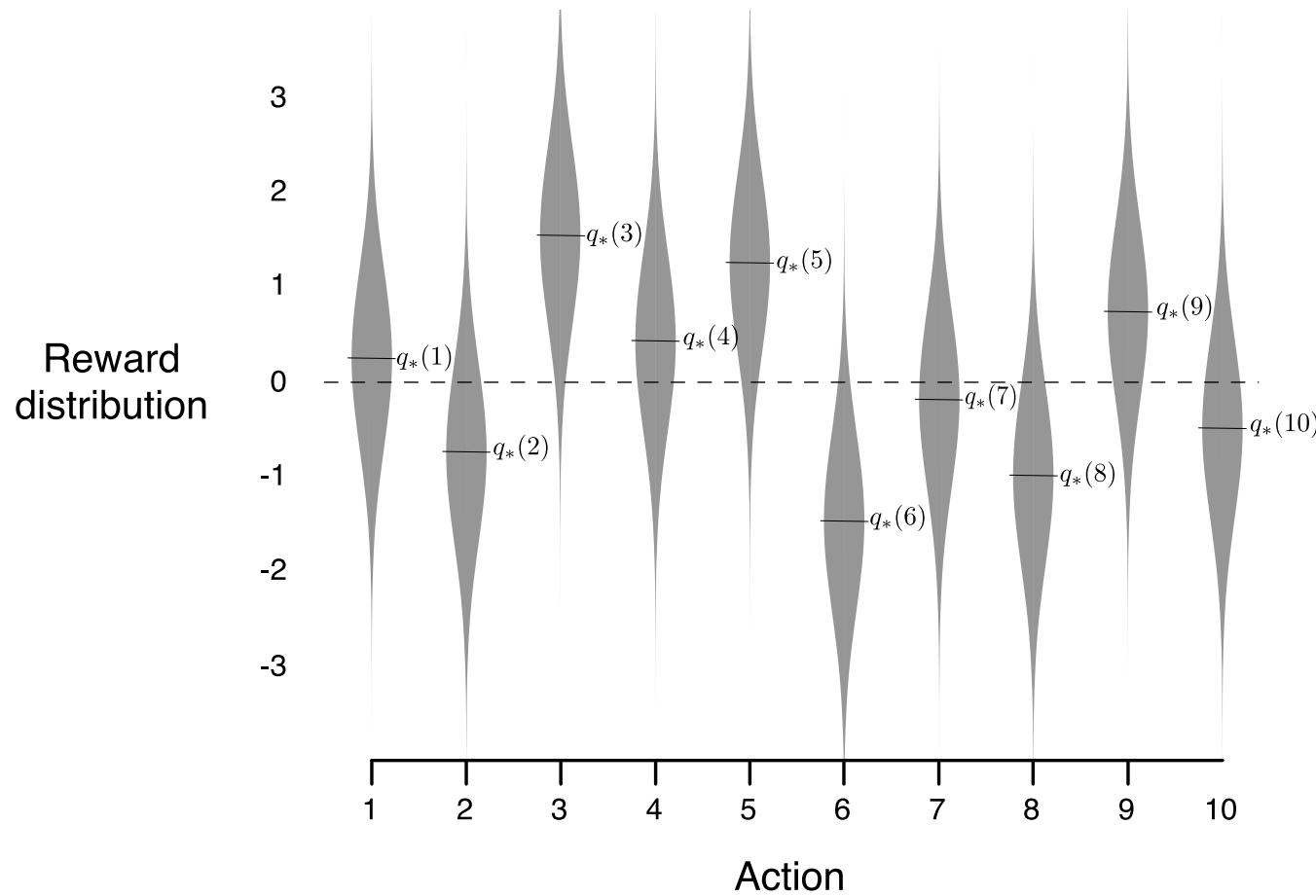
and the regret for choosing the optimal action is zero.

## Total regret as a loss

The loss function is the total regret for not selecting the optimal action up to the time step  $T$

$$L_T = \mathbb{E} \left[ \sum_{t=1}^T (\theta^* - Q(a_t)) \right] = \sum_{t=1}^T \Delta_{A_t}$$

# An example bandit problem



Random action values  $q_*(a)$ ,  $a = 1, \dots, 10$  (selected from a normal distribution, zero mean, unit std. dev.).

## Task – Maximize the return



Your task:

- \* Come up with an algorithm that maximizes the accumulated reward and improves over time.
- \* Consider advantages and disadvantages of your approach.
- \* Why is the experimental procedure given in the python example problematic?

# Beispiel in Python

Generating 10 bandits.

```
1 import time
2 import numpy as np
3 # Bandit class is taken from https://github.com/lilianweng/multi-armed-bandit
4 class BernoulliBandit(object):
5
6     def __init__(self, n, probas=None):
7         assert probas is None or len(probas) == n
8         self.n = n
9     if probas is None:
10         np.random.seed(int(time.time()))
11         self.probas = [np.random.random() for _ in range(self.n)]
12     else:
13         self.probas = probas
14
15         self.best_proba = max(self.probas)
16
17     def generate_reward(self, i):
18         # The player selected the i-th machine.
```

## The greedy policy

- Produce an estimate for action values
- Select action with highest value  $A_t = \arg \max_{a \in \mathcal{A}} \hat{Q}_t(a)$
- This is exploiting the currently available information.

# $\varepsilon$ -Greedy Strategy

- Take the best action most of the time:  $\hat{a}_t^* = \arg \max_{a \in \mathcal{A}} \hat{Q}_t(a)$
- But with  $p = \varepsilon$  do random exploration.

Best action is estimated from the collected action values from past experience (averaging the rewards for that action):

$$\hat{Q}_t(a) = \frac{1}{N_t(a)} \sum_{i=1}^t r_i \cdot \mathbb{1}(a_i = a)$$

$\mathbb{1}$  – binary indicator function for selecting an action

$N_t(a) = \sum_{i=1}^t \mathbb{1}(a_i = a)$  – counting how many times an action was selected

# $\varepsilon$ -greedy Algorithm

## A simple bandit algorithm

Initialize, for  $a = 1$  to  $k$ :

$$Q(a) \leftarrow 0$$

$$N(a) \leftarrow 0$$

Loop forever:

$$A \leftarrow \begin{cases} \arg \max_a Q(a) & \text{with probability } 1 - \varepsilon \quad (\text{breaking ties randomly}) \\ \text{a random action} & \text{with probability } \varepsilon \end{cases}$$

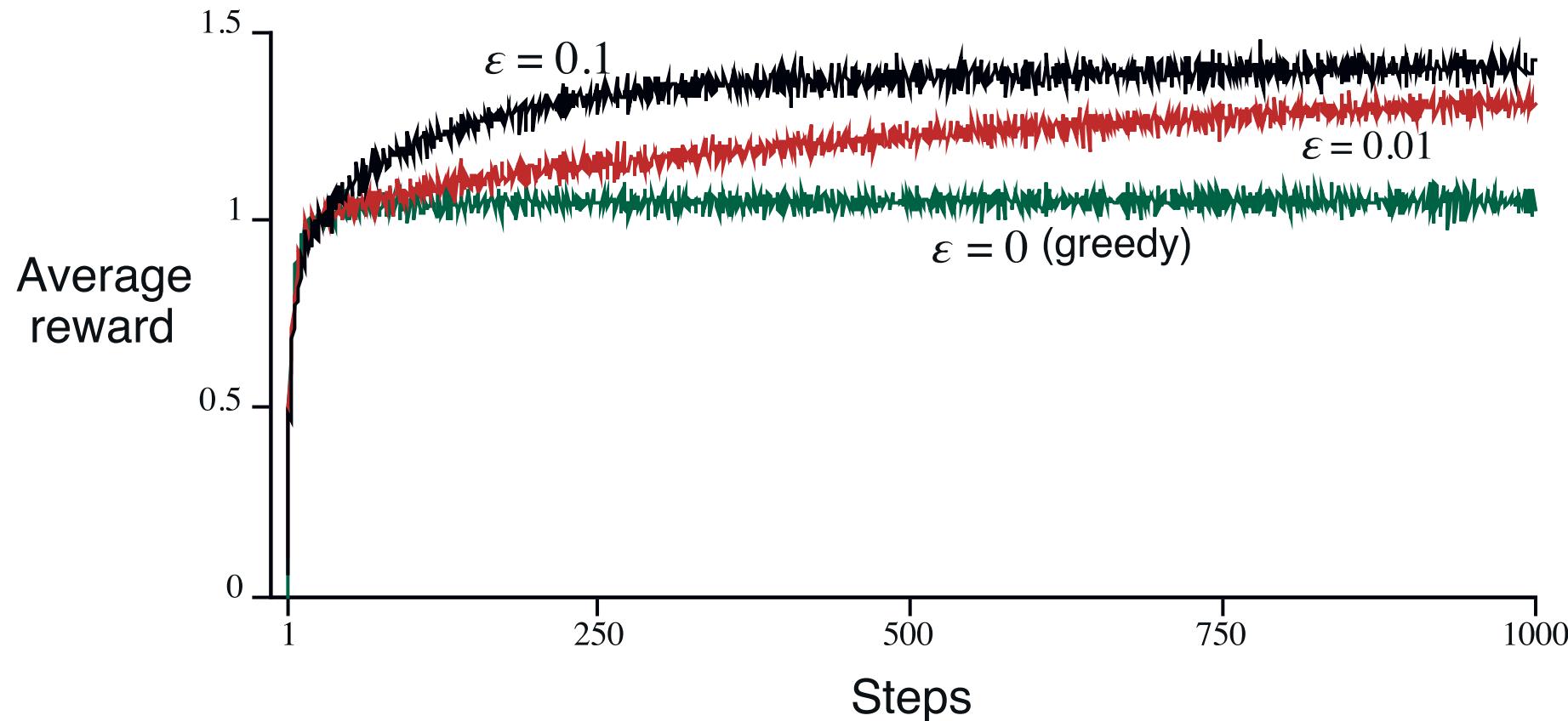
$$R \leftarrow \text{bandit}(A)$$

$$N(A) \leftarrow N(A) + 1$$

$$Q(A) \leftarrow Q(A) + \frac{1}{N(A)} [R - Q(A)]$$

Note: Incremental implementation to update estimates of value function is efficient.

# Performance of $\varepsilon$ -Greedy Strategy



Average performance of  $\varepsilon$ -greedy method (averaged over 2000 runs) for 10-arm bandit problem.

# Adaptation of $\epsilon$ -Greedy Strategy

## Drawback of $\epsilon$ -Greedy Strategy:

- during exploration we are randomly selecting actions – even though we might already have established bad actions

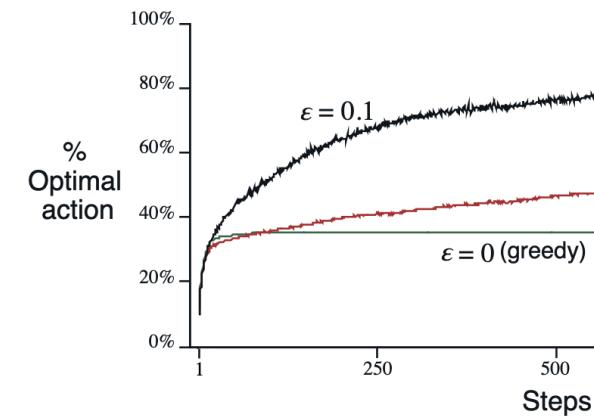
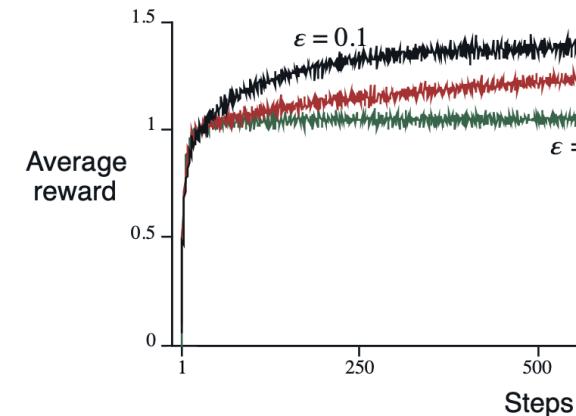
## Possible Solutions

- decrease  $\epsilon$  over time
- keep track of actions – of an estimate of how uncertain we are about this action (addressing the exploitation-exploration tradeoff)

# Übungszettel Hinweise

# Aufgabe 1.1: Implementation Multi-Armed Bandit

- Implementierung in Python eines  $k = 4$ -Bandit-Problems
- Abbildungen zu return über die Zeit und Anteil an Auswahl der optimalen Option (mehrere runs durchführen und auswerten)
- Einfluss des Parameters  $\epsilon$
- und der Initialbedingungen.
- $\epsilon$  über die Zeit verringern und Beobachtungen erklären.



(Sutton und Barto 2018)

## Aufgabe 2: Stationarität von Multi-Armed Bandits

Voraussetzung der untersuchten Verfahren: **Stationäre Probleme** – die Wahrscheinlichkeiten für das Vergeben von rewards (und später Zustandsübergängen) bleibt gleich.

Aufgabe:

- Warum ist dies ein Problem für die vorgestellten Verfahren?
- Stellen sie dies dar durch eine Anpassung des Multi-Armed-Bandit-Problems.

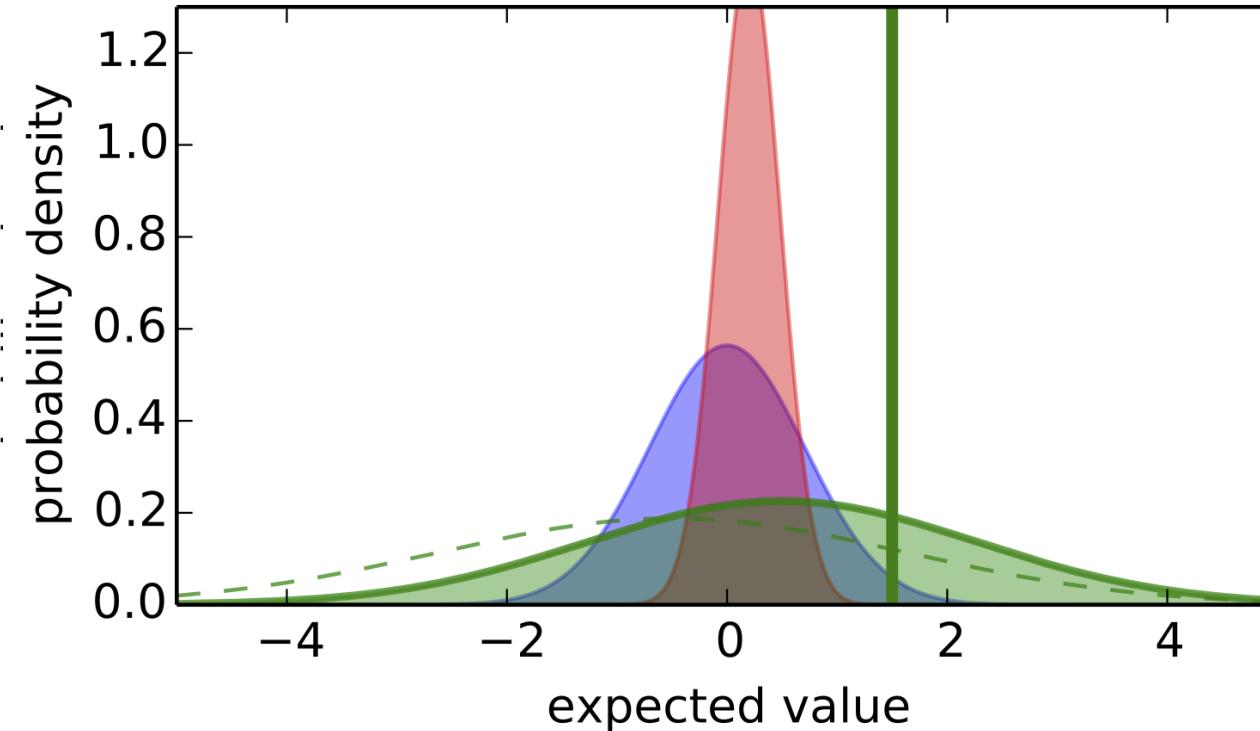
(als Zusatzpunkte!)

# Übungstermin Freitag, 21.10.2022

Inhalte:

- Troubleshooting für Python Aufgaben
- Basierend auf dem Übungszettel: zusammenstellen der notwendigen numpy Funktionen (arrays, grundlegende statistische Funktionen, ...) und plot Funktionen
- für eine random policy

# Optimism in the Face of Uncertainty



The more uncertain we are about an action-value, the more important it is to explore that action as it could turn out to be the best action.

One approach: Keep an (over-)optimistic estimate for each action.

# Upper Confidence Bounds (UCB)

Idea: favor exploration of actions that still have a strong potential to have an optimal value.

This potential is measured as an upper confidence bound of the reward value  $\hat{U}_t(a)$ . It depends on how often we have tried an action ( $N_t(a)$ ).

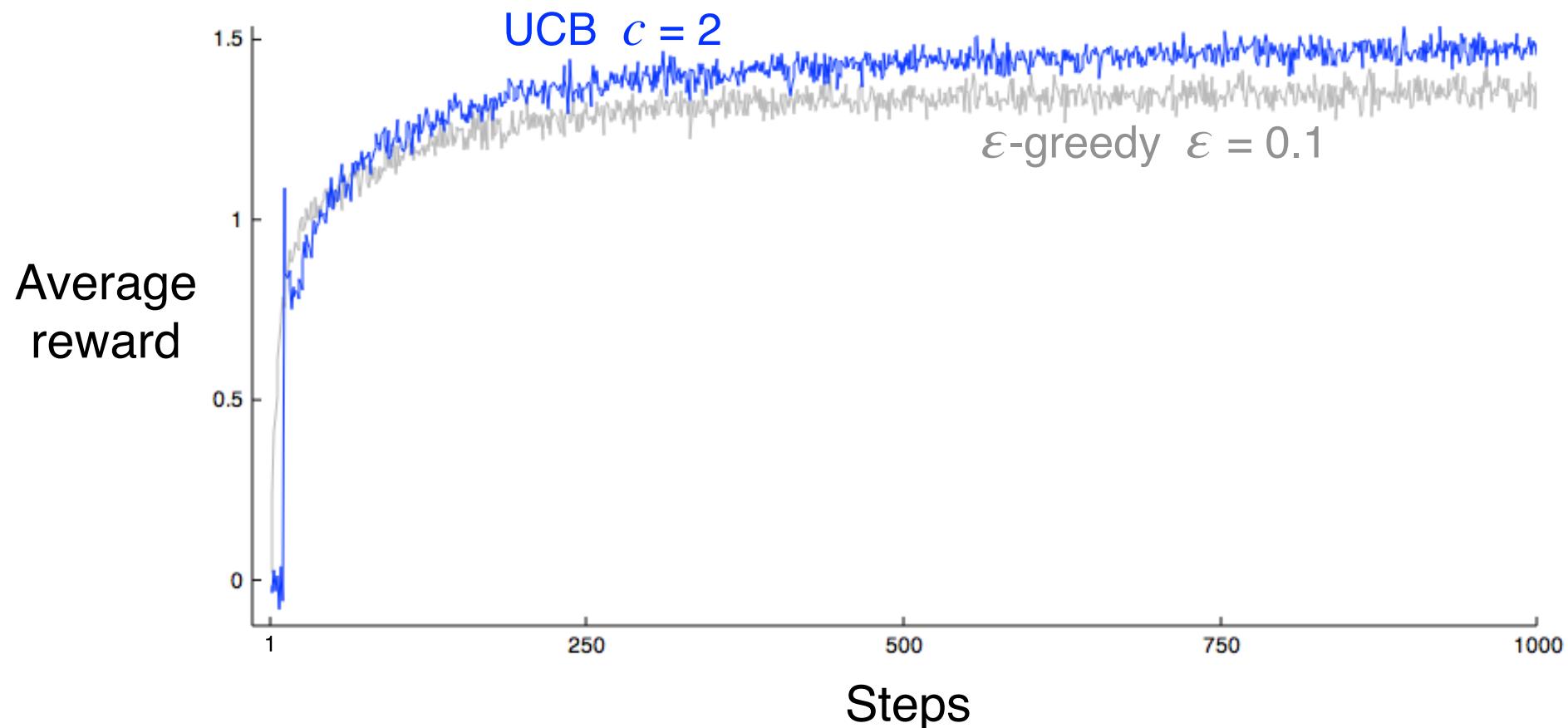
Therefore, the true reward value is bound to:

$$Q(a) \leq \hat{Q}_t(a) + \hat{U}_t(a)$$

In UCB algorithm, actions are selected greedily in order to maximize the upper confidence bound:

$$a_t^{UCB} = \arg \max_{a \in \mathcal{A}} \hat{Q}_t(a) + \hat{U}_t(a) = \arg \max_{a \in \mathcal{A}} \hat{Q}_t(a) + c \sqrt{\frac{\text{XXX } v}{N_t(a)}}$$

# Average Performance of UCB



Average performance of UCB (averaged over 2000 runs) for 10-arm bandit problem.

UCB outperforms  $\epsilon$ -greedy.

# Considering regret

Regret was defined as  $\Delta_a = v_* - q(a)$ .

It depends on

- the action count – how often each action was selected
- and on the overall collected regret  $L_T = \mathbb{E} \left[ \sum_{t=1}^T (\theta^* - Q(a_t)) \right]$ .

Therefore, we should aim for algorithm that avoid large regrets.

# Regret in the case of UCB

Selection of  $a_t$ :

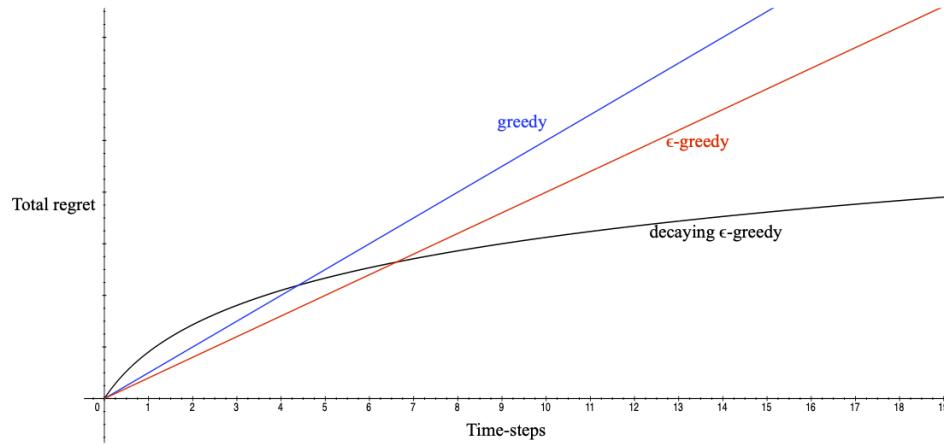
$$a_t = \arg \max_{a \in \mathcal{A}} Q_t(a) + c \sqrt{\frac{\log t}{N_t(a)}}$$

Intuitively,

- When  $\Delta_a$  is large, the action will not be selected
- unless  $N_t(a)$  is (comparatively) small.

It can be shown, that  $\Delta_a \cdot N_t(a) \leq O(\log t)$  for all  $a$ .

# Comparison of regret



## Asymptotic total regret (Lai und Robbins 1985)

Asymptotic total regret is at least logarithmic in number of steps

$$\lim_{t \rightarrow \infty} L_t \geq \log t \sum_{a | \Delta_a > 0} \frac{\Delta_a}{KL(\mathcal{R}_a || \mathcal{R}_{a^*})}$$

With  $KL(\mathcal{R}_a || \mathcal{R}_{a^*}) \sim \Delta_a^2$  – note: Hard problems have similar-looking arms with (Silver 2015)

$\forall x \exists y (\forall a \mid \forall a^*) \wedge \Delta_a$  Note. Many predicates have similar sounding names with different means.

# Policy - Directly learning how to act

Consider action selection as a probability distribution:

- For each action: Consider an (estimated) preference  $H_t(a)$  of that action which
- can be directly used to express a probability for selecting that action (as a soft-max distribution)

$$p(A_t = a) = \frac{e^{H_t(a)}}{\sum_{b=1}^k e^{H_t(b)}} = \pi_t(a)$$

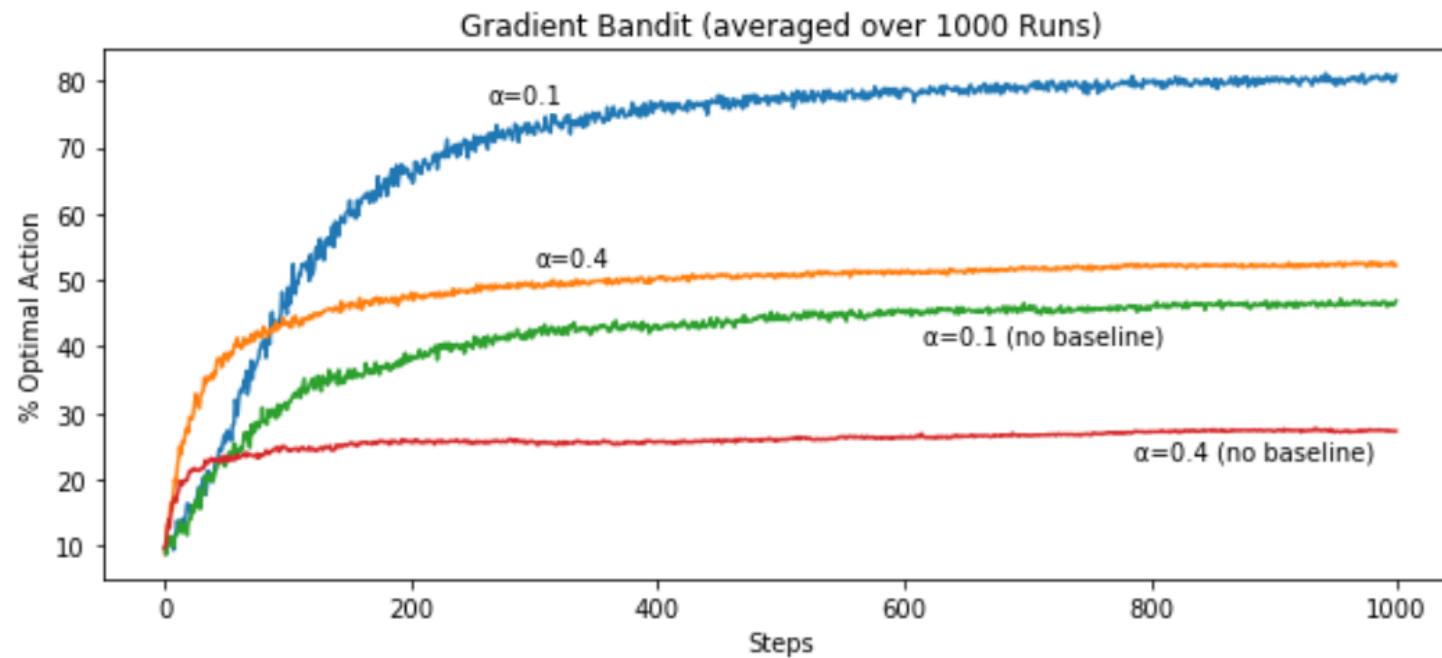
## Gradient Bandit Algorithm

Learn / adapt the action preference function directly using stochastic gradient ascent:

$$H_{t+1}(A_t) = H_t(A_t) + \alpha(R_t - \bar{R}_t)(1 - \pi_t(A_t)), \quad \text{and}$$

$$H_{t+1}(a) = H_t(A_t) - \alpha(R_t - \bar{R}_t)\pi_t(a) \quad \text{for all } a \neq A_t$$

# Average Performance of Gradient Bandit Algorithm



Shown are variations of the  $\alpha$  parameter and the importance of including the mean reward in the update for unbiasing (for brown curve this baseline is removed and the added bias of +4 affects learning).

# Further variations for Solving Bandit Problems

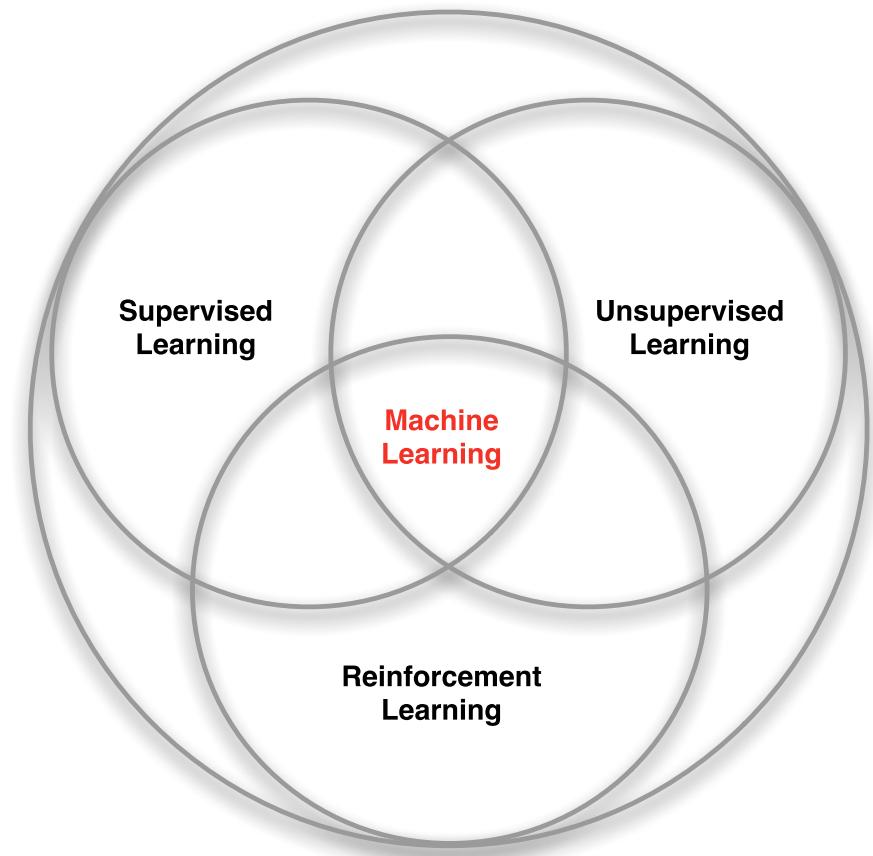
- employing an optimistic estimate of the action value instead of a default initialization
- Bayesian UCB: introduce a prior assumption for the reward distribution – as a Gaussian – and use confidence intervals
- Thompson Sampling – formulate action selection as a probabilistic process itself (selecting an action with a probability that estimates it is optimal)

# Sequences of States

# Recap – Forms of Learning

## What is special in Reinforcement Learning?

- There is no supervisor, only a reward signal.
- **Sequential:** Time really matters (non i.i.d data).
- Feedback can be delayed, not instantaneous.
- Agent's actions *affect the subsequent data* it receives



# *Sequential Decision Making*

The interaction between agent and environment is a **sequence** of actions and returned observations plus rewards.

Goal of the agent: select actions to maximise total future reward

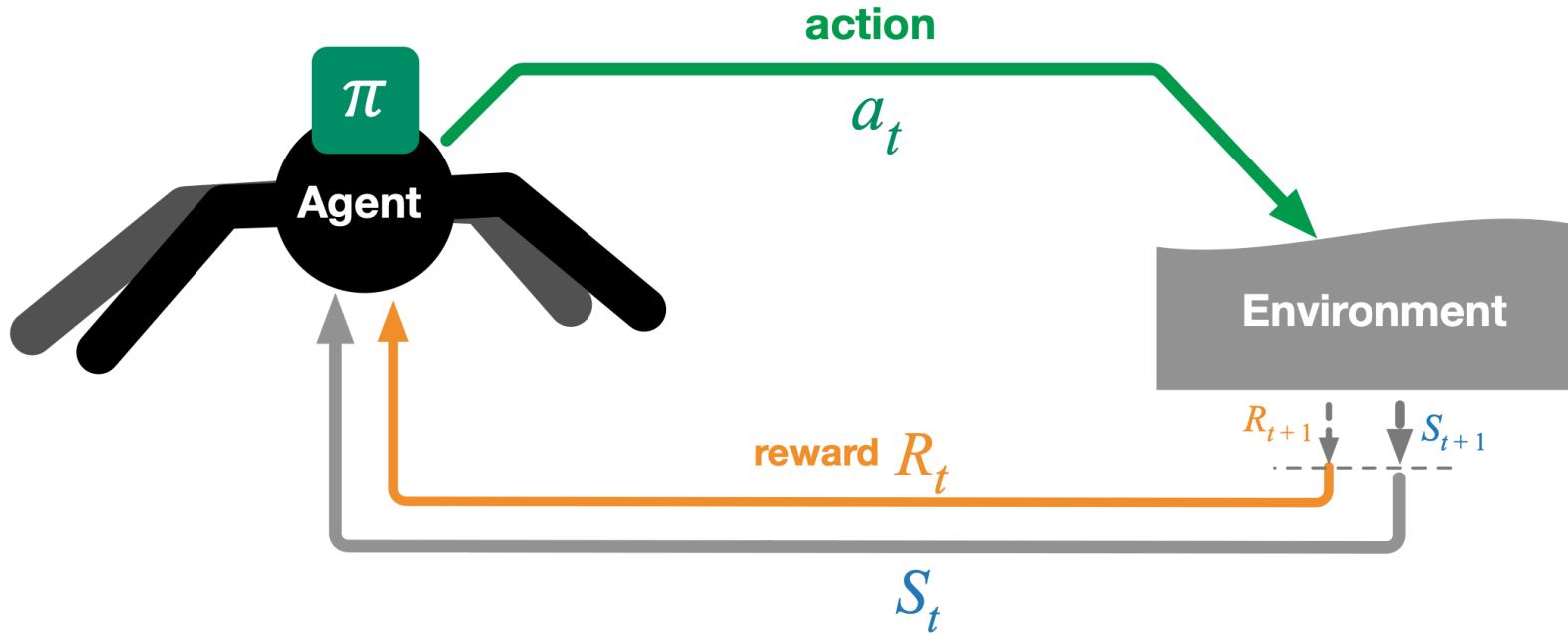
- But actions may have long term consequences and reward may be delayed.
- It may be better to sacrifice immediate reward to gain more long-term reward

An agent's **policy**  $\pi(s) = a$  describes which action  $a$  an agent selects depending on the current state.

For the stochastic state, a policy is a probability distribution over actions:

$$\pi(a|s) = \mathbb{P}_\pi[A = a|S = s].$$

# RL Cycle – Sequential Decision Making



## Agent

- Policy: choose an action.
- Value-Function: Estimate of achievable return from a state.

## Environment

- Reward: Describing goal
- State: observation for agent

# Markov Process

A Markov process is a memoryless random process – a sequence of random states  $S_1, S_2, \dots$  with the Markov property.

## Markov Property

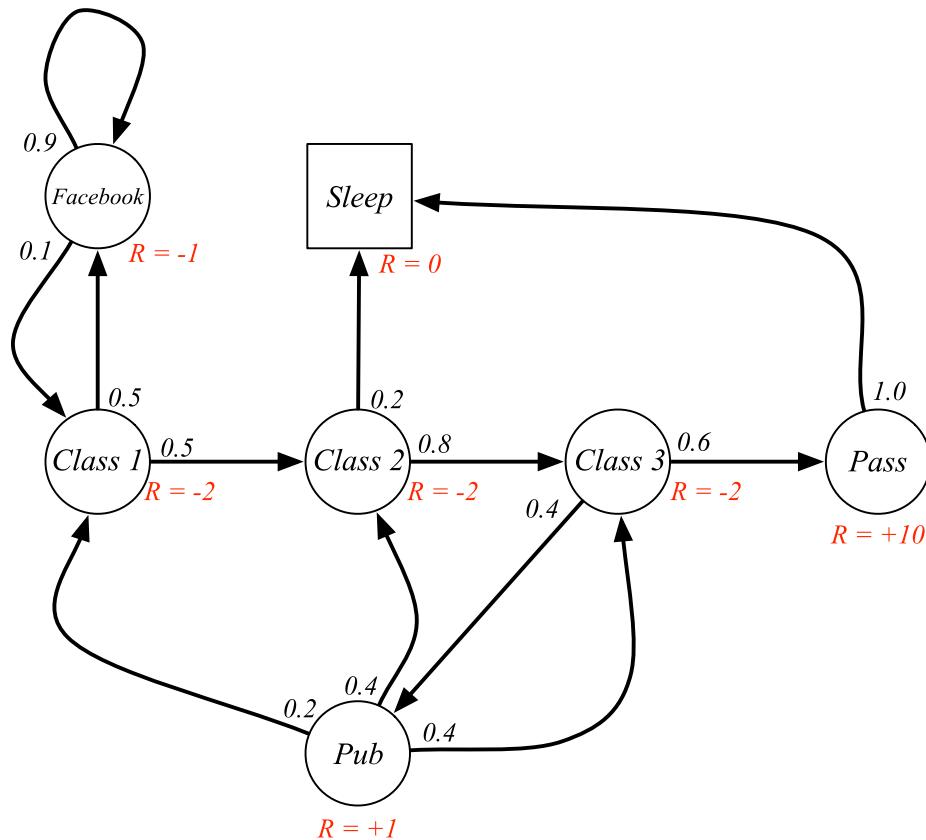
A state  $S_t$  is Markov iff  $P(S_{t+1}|S_t) = P(S_{t+1}|S_1, \dots, S_t)$

States captures all relevant information from the history (“The future is independent of the past given the present”).

## Markov Process (Markov Chain)

- $\mathcal{S}$  is a (finite) set of states
- $\mathcal{P}$  is a state transition probability matrix:  $P_{ss'} = P(S_{t+1} = s'|S_t = s)$

# A Markov Reward Process



Before turning to action, we focus on a simpler class of Markov Chains: the Markov Reward Process. We could sample trajectories from it.

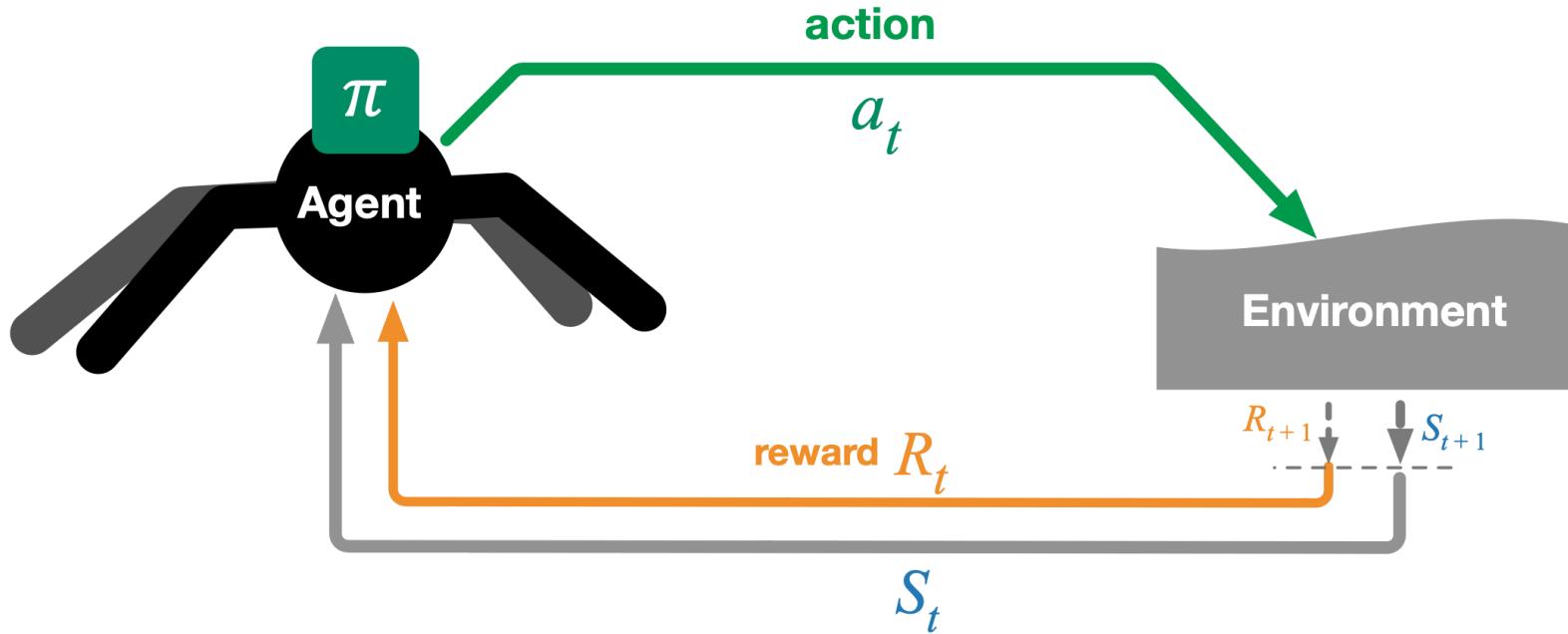
# Markov Reward Process

A Markov reward process is a Markov Chain with associated reward values.

**Markov reward process consists of (def.)**

- $\mathcal{S}$  is a (finite) set of states
- $\mathcal{P}$  is a state transition probability matrix:  $P_{ss'} = P(S_{t+1} = s' | S_t = s)$
- $\mathcal{R}$  is a reward function,  $\mathcal{R}_s = \mathbb{E}[R_{t+1} | S_t = s]$
- $\gamma$  is a discount factor,  $\gamma \in [0, 1]$  (used to give a higher value to rewards that are closer in time)

# RL Cycle – Sequential Decision Making



## Agent

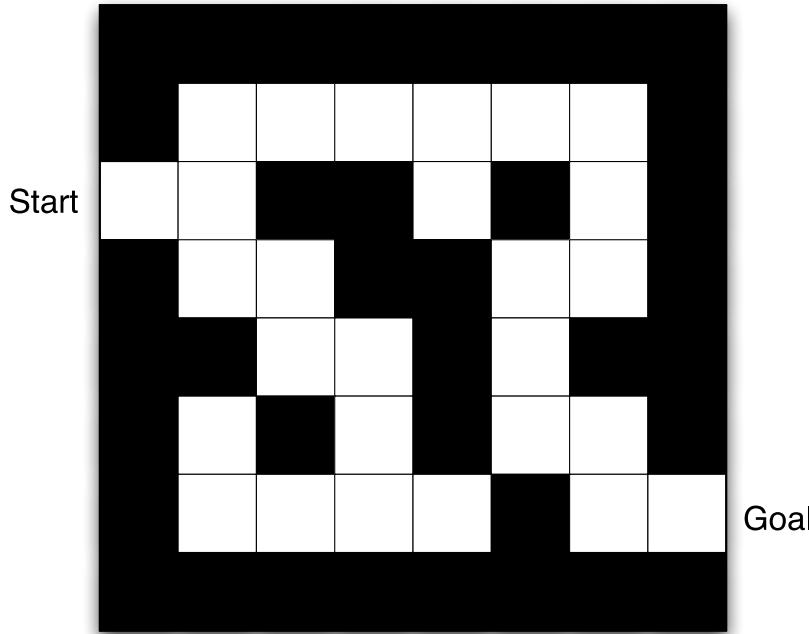
- Policy: choose an action **depending on current state**.
- Value-Function: Estimate of achievable

## Environment

- Reward: Describing goal
- State: observation for agent

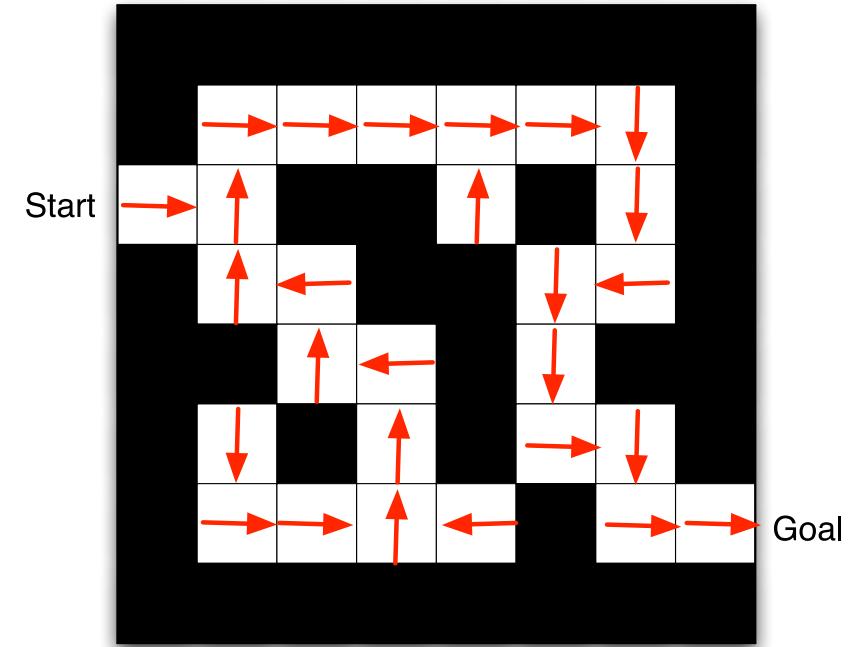
# Maze Example: Policy

## Task



- Reward of  $-1$  per time step in maze
- Actions are move N, S, W, E
- State is location

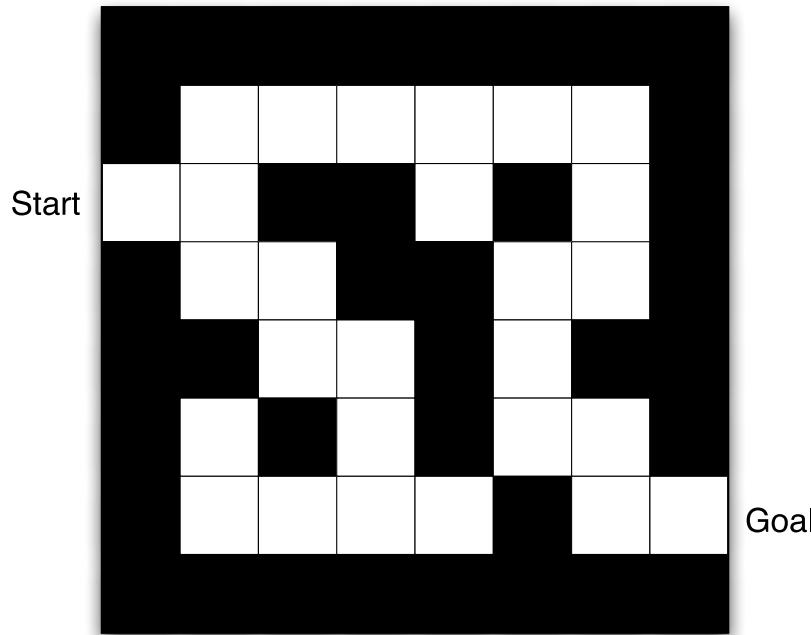
## Policy Representation



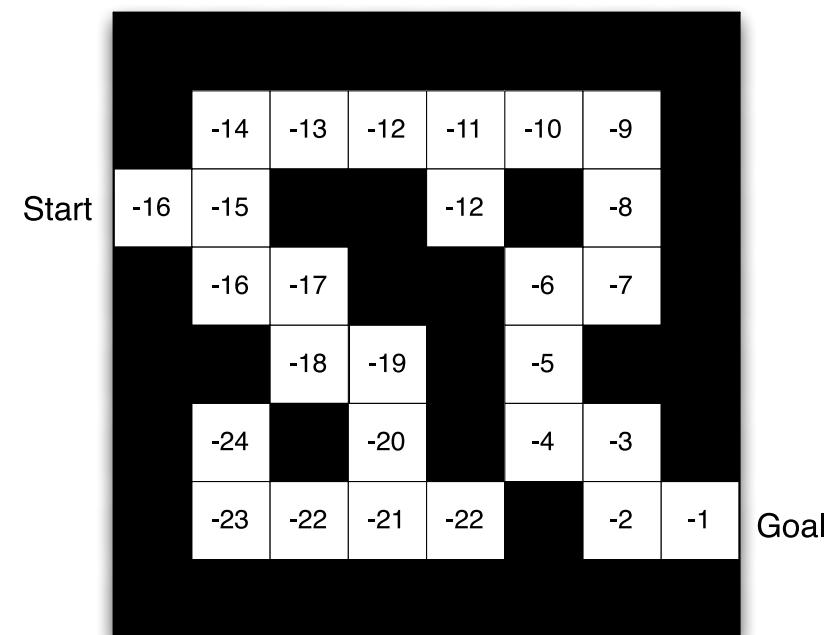
Arrows represent policy  $\pi(s)$  for all the states.

# Maze Example: Value Function

Task



State Value



- Reward of  $-1$  per time step in maze
- Actions are move N, S, W, E
- State is location

Shown are values  $v_{\pi}(s)$  for the different states.

# Partially Observable Environments

We will focus mostly on fully observable environments: The agent state contains all necessary information required for making an informed decision.

## Partial Observability

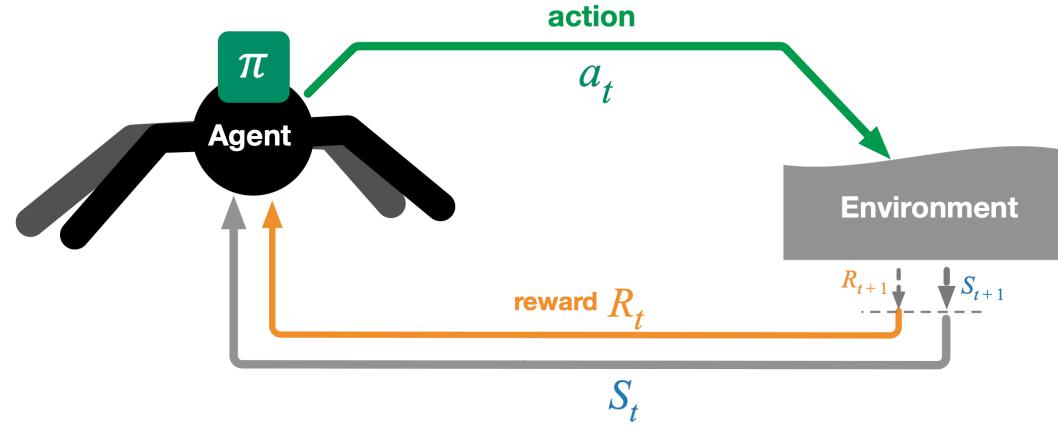
The agent only indirectly experiences the environment, e.g. no exact position information, but only relying on a camera.

Importantly, he might not be able to distinguish states.

The agent, therefore, must construct its own internal state representation (which could, e.g., include information on history).



# RL Cycle – Sequential Decision Making



## Agent

- Policy: choose an action *depending on current state*.
- Value-Function: Estimate of achievable return from a state (following  $\pi$ ).
- Model: A *predictor* of the environment.

## Environment

- Reward: Describing goal
- State: observation for agent

# Model

A model allows to predict how the environment will react (as a probability distribution).

- $\mathcal{P}$  predicts the subsequent state:

$$\mathcal{P}_{ss'}^a \approx p(S_{t+1} = s' | S_t = s, A_t = a)$$

- $\mathcal{R}$  predicts the next reward:

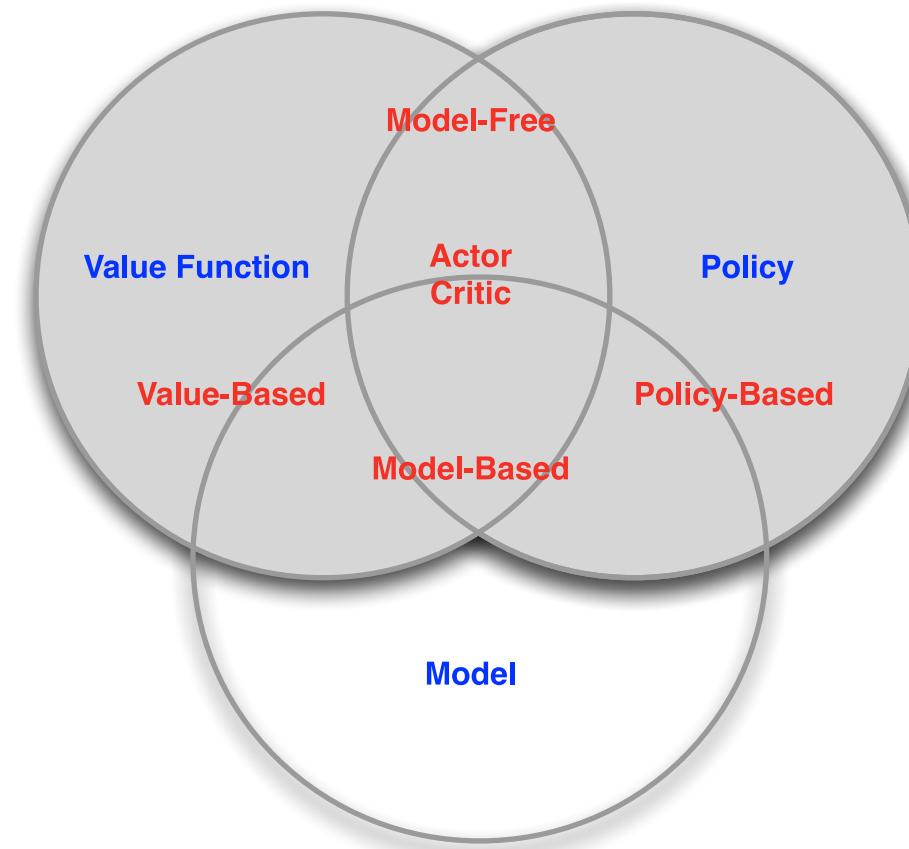
$$\mathcal{R}_s^a = \mathbb{E}(R_{t+1} | S_t = s, A_t = a)$$

A model does not give us immediately a good policy.

But it allows us to plan – test possible alternative actions.

# Categorization of Reinforcement Learning Agents

- Value Based
  - No Policy (Implicit)
  - Value Function
- Policy Based
  - Policy
  - No Value Function
- Actor Critic
  - Policy
  - Value Function



# References

- Gao, Chongming, Wenqiang Lei, Xiangnan He, Maarten Rijke, und Tat-Seng Chua. 2021. „Advances and Challenges in Conversational Recommender Systems: A Survey“.
- Hasselt, Hado van, und Diana Borsa. 2021. „Reinforcement Learning Lecture Series 2021“. <https://www.deepmind.com/learning-resources/reinforcement-learning-lecture-series-2021>.
- Klein, Dan, und Pieter Abbeel. 2014. „UC Berkeley CS188 Intro to AI“. <http://www0.cs.ucl.ac.uk/staff/d.silver/web/Teaching.html>.
- Lai, T. L., und H. Robbins. 1985. „Asymptotically Efficient Adaptive Allocation Rules“. *Advances in Applied Mathematics* 6: 4–22.
- Silver, David. 2015. „UCL Course on RL UCL Course on RL UCL Course on Reinforcement Learning“. <http://www0.cs.ucl.ac.uk/staff/d.silver/web/Teaching.html>.
- Sutton, Richard S., und Andrew G. Barto. 2018. *Reinforcement Learning: An Introduction*. Second. The MIT Press.
- Weng, Lilian. 2018. „The Multi-Armed Bandit Problem and Its Solutions“. [The Multi-Armed Bandit Problem and Its Solutions](#).