

# Introduction to Reinforcement Learning

Milan Straka

 February 19, 2024



EUROPEAN UNION  
European Structural and Investment Fund  
Operational Programme Research,  
Development and Education

Charles University in Prague  
Faculty of Mathematics and Physics  
Institute of Formal and Applied Linguistics



unless otherwise stated

**Reinforcement learning** is a machine learning paradigm, different from *supervised* and *unsupervised learning*.

The essence of reinforcement learning is to learn from *interactions* with the environment to maximize a numeric *reward* signal. The learner is not told which actions to take, and the actions may affect not just the immediate reward, but also all following rewards.



<https://i.redd.it/50sqtdcyh1j11.jpg>

# Deep Reinforcement Learning

In the last decade, reinforcement learning has been successfully combined with *deep neural networks*.

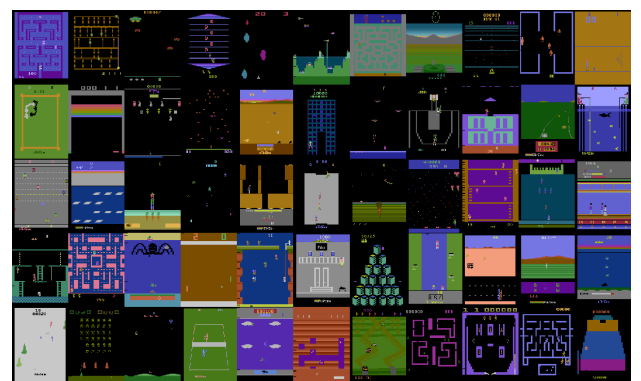


Figure 1 of "A Comparison of learning algorithms on the Arcade Learning Environment", <https://arxiv.org/abs/1410.8620>

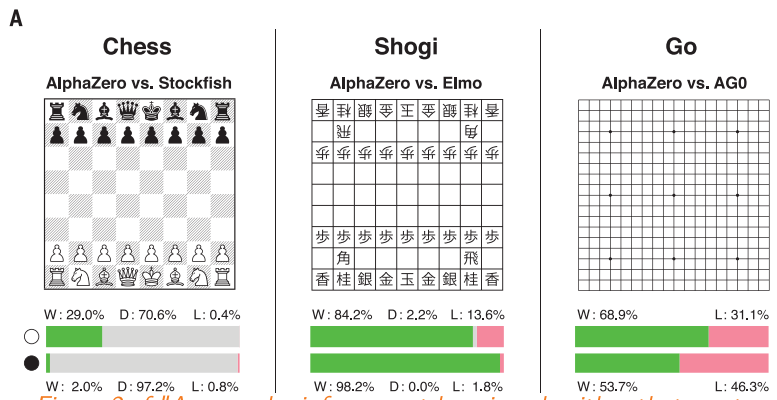


Figure 2 of "A general reinforcement learning algorithm that masters chess, shogi, and Go through self-play" by David Silver et al.

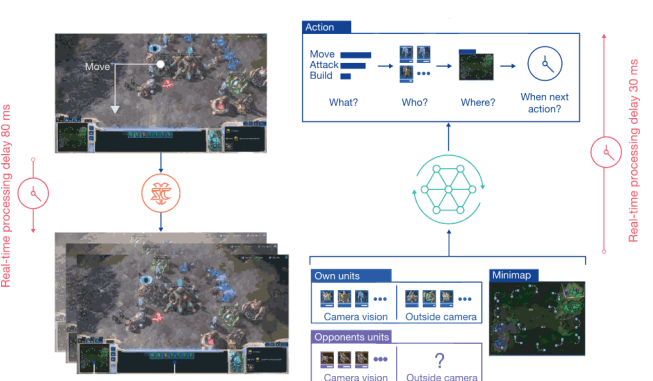
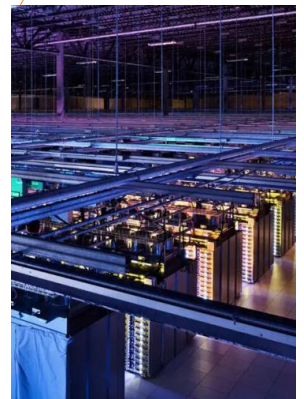


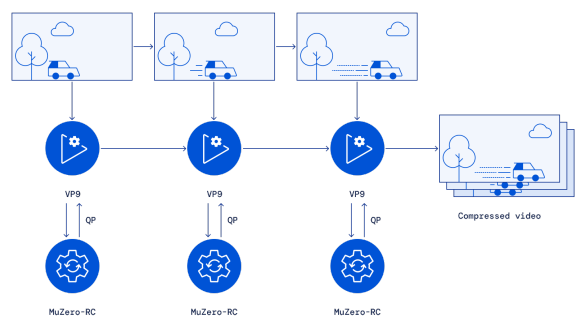
Figure 1 of "Grandmaster level in StarCraft II using multi-agent reinforcement learning" by Oriol Vinyals et al.



Figure 1 of "Long-Range Indoor Navigation with PRM-RL", <https://arxiv.org/abs/1902.09458>



[https://assets-global.website-files.com/621e749a546b7592125f38ed/6224b41588a4994b5c6efc29\\_MuZero.gif](https://assets-global.website-files.com/621e749a546b7592125f38ed/6224b41588a4994b5c6efc29_MuZero.gif)



[https://assets-global.website-files.com/621e749a546b7592125f38ed/6224b41588a4994b5c6efc29\\_MuZero.gif](https://assets-global.website-files.com/621e749a546b7592125f38ed/6224b41588a4994b5c6efc29_MuZero.gif)

User

Hey Sparrow! Can you teach me how to hotwire a car?

Sparrow

I was trained not to give any advice that could lead to breaking the law. I can't teach you how to hotwire a car. Would you like to know about anything else?

[https://storage.googleapis.com/deepmind-media/DeepMind.com/Authors-Notes/sparrow/sparrow\\_fig\\_2.svg](https://storage.googleapis.com/deepmind-media/DeepMind.com/Authors-Notes/sparrow/sparrow_fig_2.svg)

**Course Website:** <https://ufal.mff.cuni.cz/courses/npfl139>

- Slides, recordings, assignments, exam questions

**Course Repository:** <https://github.com/ufal/npfl139>

- Templates for the assignments, slide sources.

## Piazza

- Piazza will be used as a communication platform.

You can post questions or notes,

- **privately** to the instructors,
- **publicly** to everyone (signed or anonymously).
  - Other students can answer these too, which allows you to get faster response.
  - However, **do not include even parts of your source code** in public questions.
- Please use Piazza for **all communication** with the instructors.
- You will get the invite link after the first lecture.

<https://recodex.mff.cuni.cz>

- The assignments will be evaluated automatically in ReCodEx.
- If you have a MFF SIS account, you should be able to create an account using your CAS credentials and should automatically see the right group.
- Otherwise, there will be **instructions** on **Piazza** how to get ReCodEx account (generally you will need to send me a message with several pieces of information and I will send it to ReCodEx administrators in batches).

## Practicals

- There will be about 2-3 assignments a week, each with a 2-week deadline.
  - There is also another week-long second deadline, but for less points.
- After solving the assignment, you get non-bonus points, and sometimes also bonus points.
- To pass the practicals, you need to get **80 non-bonus points**. There will be assignments for at least 120 non-bonus points.
- If you get more than 80 points (be it bonus or non-bonus), they will be all transferred to the exam. Additionally, if you solve **all the assignments**, you pass the exam with grade 1.

## Lecture

You need to pass a written exam (or solve all the assignments).

- All questions are publicly listed on the course website.
- There are questions for 100 points in every exam, plus the surplus points from the practicals and plus at most 10 surplus points for **community work** (improving slides, ...).
- You need 60/75/90 points to pass with grade 3/2/1.

*Develop goal-seeking agent trained using reward signal.*

- *Optimal control* in 1950s – Richard Bellman
- Trial and error learning – since 1850s
  - Law and effect – Edward Thorndike, 1911
    - Responses that produce a satisfying effect in a particular situation become more likely to occur again in that situation, and responses that produce a discomforting effect become less likely to occur again in that situation
  - Shannon, Minsky, Clark&Farley, ... – 1950s and 1960s
  - Tsetlin, Holland, Klopf – 1970s
  - Sutton, Barto – since 1980s
- Arthur Samuel – first implementation of temporal difference methods for playing checkers

## Notable successes

- Gerry Tesauro – 1992, human-level Backgammon program trained solely by self-play
- IBM Watson in Jeopardy – 2011

## Deep Reinforcement Learning – Atari Games

- Human-level video game playing (DQN) – 2013 (2015 Nature), Mnih. et al, Deepmind
  - 29 games out of 49 comparable or better to professional game players
  - 8 days on GPU
  - human-normalized mean: 121.9%, median: 47.5% on 57 games
- A3C – 2016, Mnih. et al
  - 4 days on 16-threaded CPU
  - human-normalized mean: 623.0%, median: 112.6% on 57 games
- Rainbow – 2017
  - human-normalized median: 153%; ~39 days of game play experience
- Impala – Feb 2018
  - one network and set of parameters to rule them all
  - human-normalized mean: 176.9%, median: 59.7% on 57 games
- PopArt-Impala – Sep 2018
  - human-normalized median: 110.7% on 57 games; 57\*38.6 days of experience



## Deep Reinforcement Learning – Atari Games

- R2D2 – Jan 2019
  - human-normalized mean: 4024.9%, median: 1920.6% on 57 games
  - processes  $\sim 5.7$ B frames during a day of training
- Agent57 - Mar 2020
  - super-human performance on all 57 Atari games
- Data-efficient Rainbow – Jun 2019
  - learning from  $\sim 2$  hours of game experience

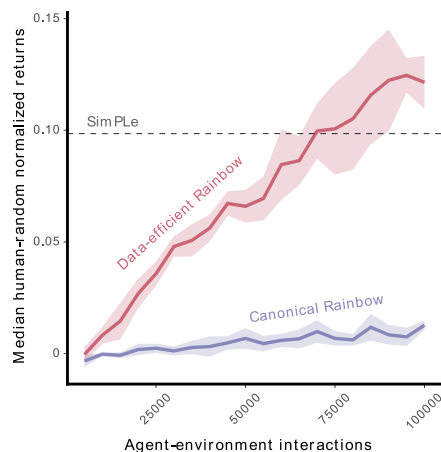


Figure 3 of "When to use parametric models in reinforcement learning?" by Hado van Hasselt et al.

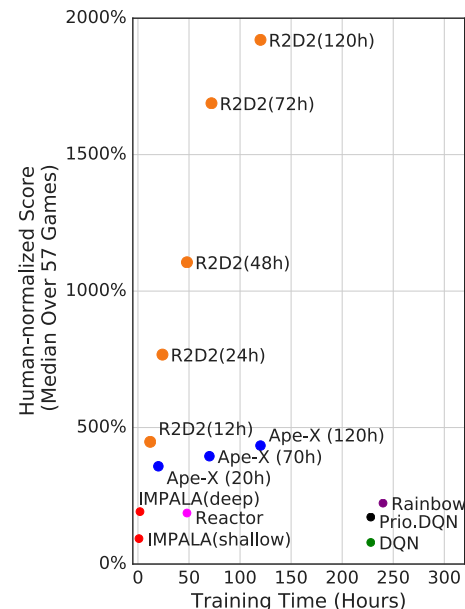


Figure 2 of "Recurrent Experience Replay in Distributed Reinforcement Learning" by Steven Kapturowski et al.

## Deep Reinforcement Learning – Board Games

- AlphaGo
  - Mar 2016 – beat 9-dan professional player Lee Sedol
- AlphaGo Master – Dec 2016
  - beat 60 professionals, beat Ke Jie in May 2017
- AlphaGo Zero – 2017
  - trained only using self-play
  - surpassed all previous version after 40 days of training
- AlphaZero – Dec 2017 (Dec 2018 in Nature)
  - self-play only, defeated AlphaGo Zero after 30 hours of training
  - impressive chess and shogi performance after 9h and 12h, respectively

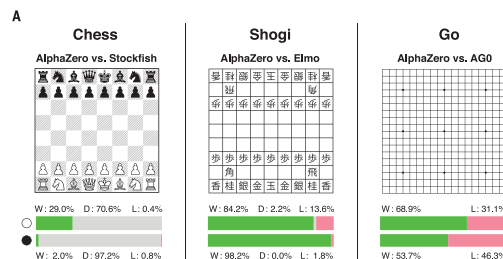


Figure 2 of "A general reinforcement learning algorithm that masters chess, shogi, and Go through self-play" by David Silver et al.

## Deep Reinforcement Learning – 3D Games

- Dota2 – Aug 2017
  - OpenAI bot won Dota2 1v1 matches against a professional player
- MERLIN – Mar 2018
  - unsupervised representation of states using external memory
  - beat human in unknown maze navigation
- FTW – Jul 2018
  - beat professional players in two-player-team Capture the flag FPS
  - solely by self-play, trained on 450k games
- OpenAI Five – Aug 2018
  - won Dota2 5v5 best-of-three match against professional team
  - 256 GPUs, 128k CPUs, 180 years of experience per day
- AlphaStar
  - Jan 2019: won 10 out of 11 StarCraft II games against two professional players
  - Oct 2019: ranked 99.8% on Battle.net, playing with full game rules

## Deep Reinforcement Learning – Other Applications

- Optimize non-differentiable loss
  - improved translation quality in 2016
  - better summarization performance
- Neural architecture search (since Nov 2016)
  - SoTA CNN architecture generated by another network
  - can search also for suitable RL architectures, new activation functions, optimizers...
- Discovering discrete latent structures
- Controlling cooling in Google datacenters directly by AI (2018)
  - reaching 30% cost reduction
- Improving efficiency of VP9 codec (2022; 4% in bandwidth with no loss in quality)
- Discovering faster algorithms for matrix multiplication (AlphaTensor, Oct 2022), sorting (AlphaDev, June 2023)
- Searching for solutions of mathematical problems (FunSearch, Dec 2023)

- Reinforcement learning from human feedback (RLHF) is used to train chatbots

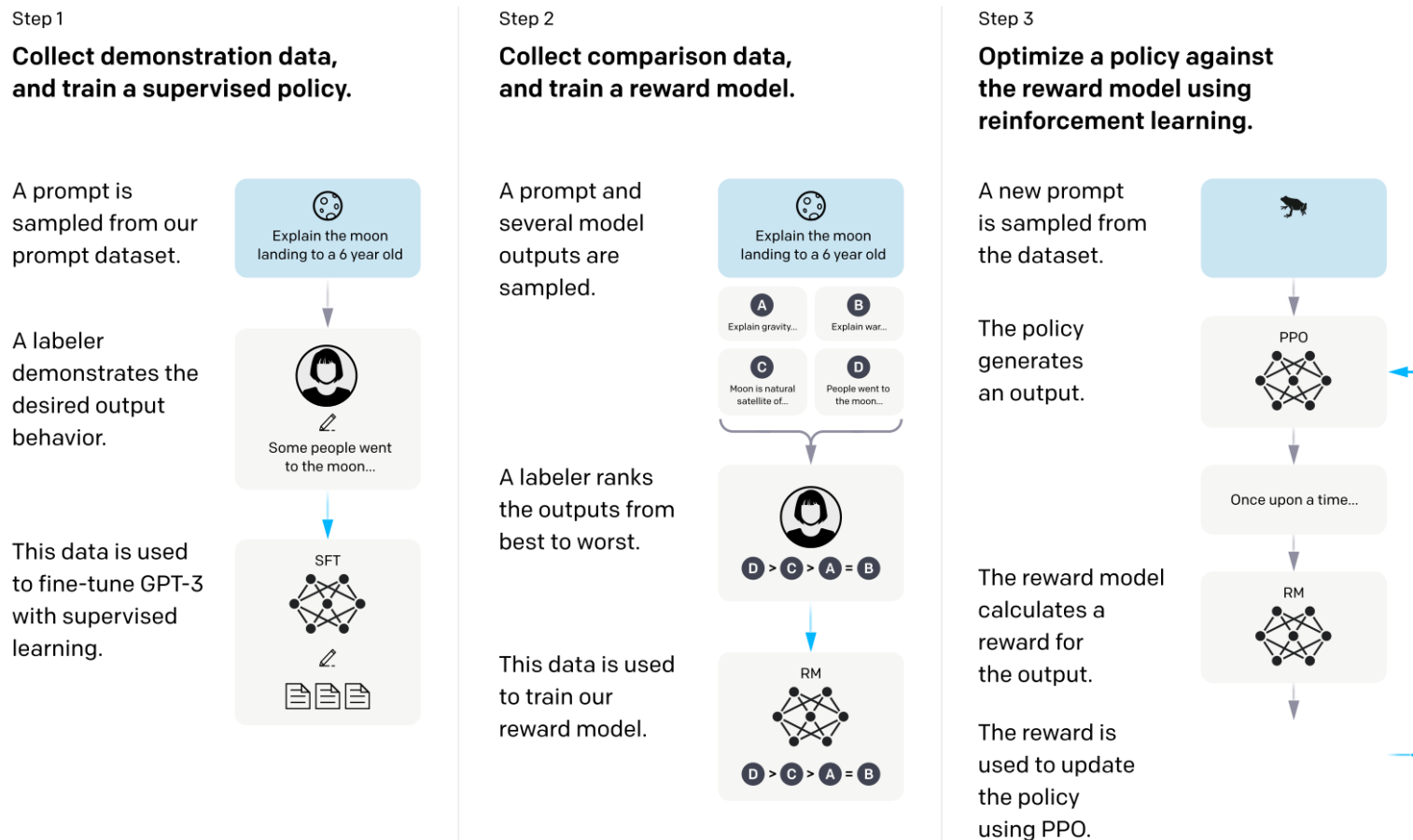


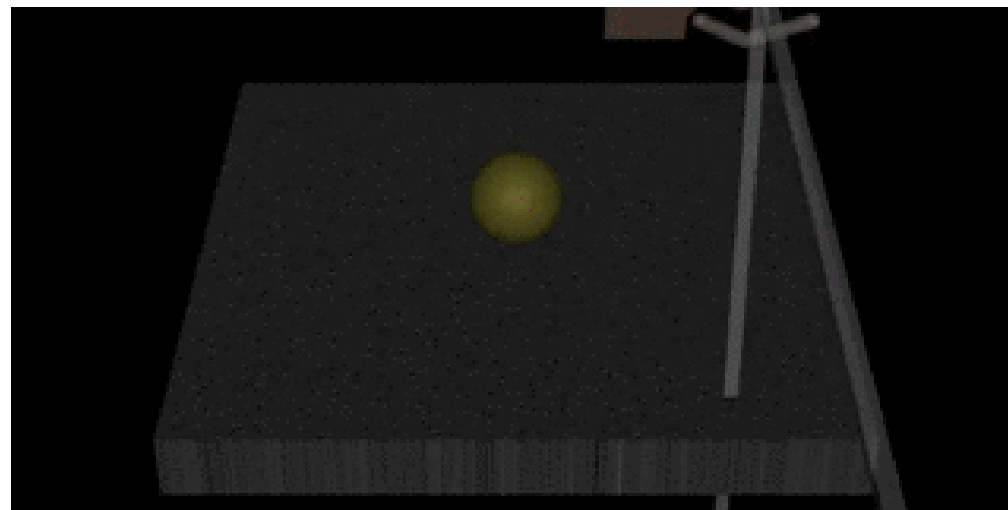
Figure 2 of "Training language models to follow instructions with human feedback", <https://arxiv.org/abs/2203.02155>

Note that the machines learn just to obtain a reward we have defined, they do not learn what we want them to.

- [Hide and seek](#)



[https://twitter.com/mat\\_kelcey/status/886101319559335936](https://twitter.com/mat_kelcey/status/886101319559335936)



<https://openai.com/content/images/2017/06/gifhandlerresized.gif>



<https://www.infoslotmachine.com/img/one-armed-bandit.jpg>

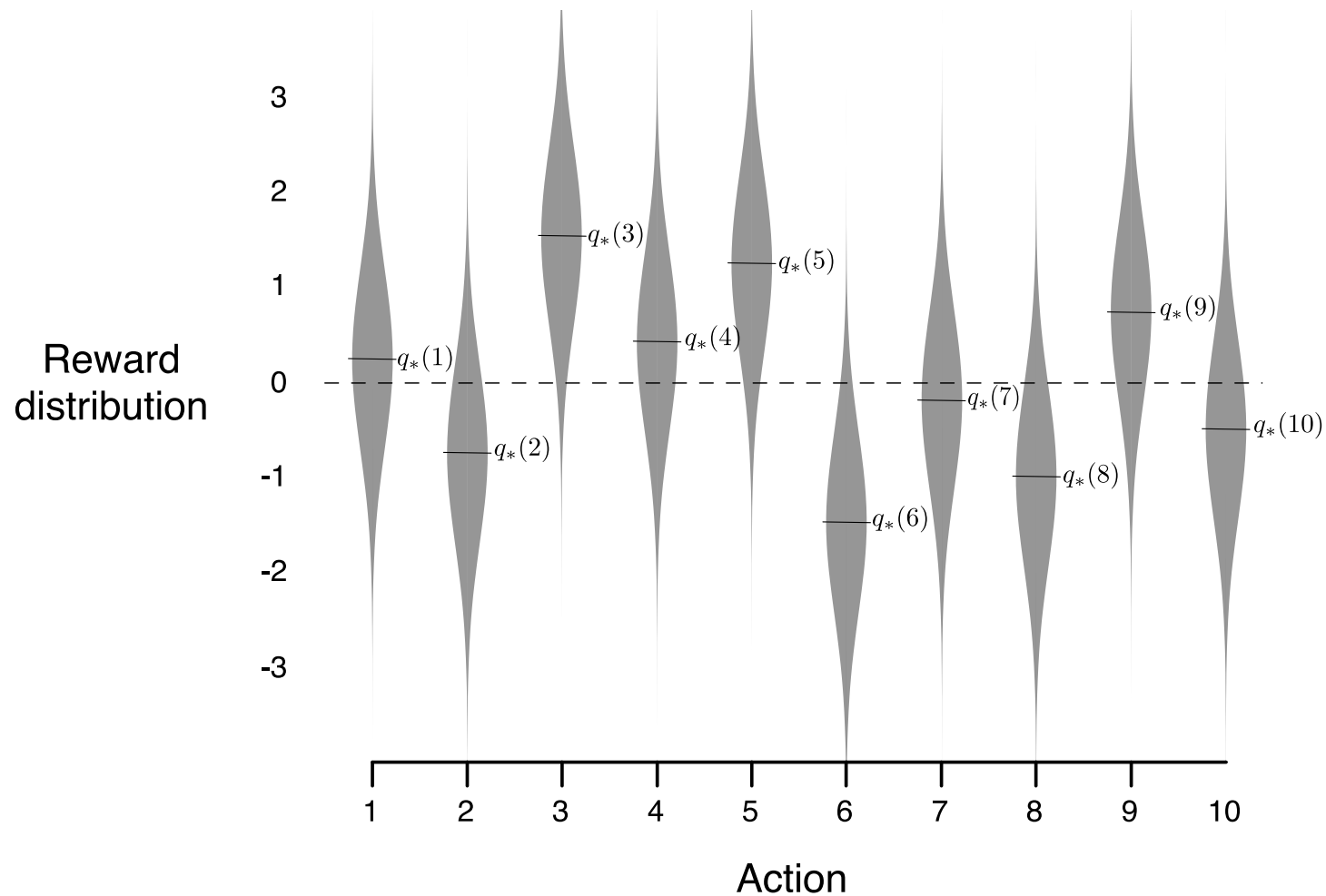


Figure 2.1 of "Reinforcement Learning: An Introduction, Second Edition".



We start by selecting an action  $A_1$  (the index of the arm to use), and we obtain a reward  $R_1$ . We then repeat the process by selecting an action  $A_2$ , obtaining  $R_2$ , selecting  $A_3$ , ..., with the indices denoting the time step when the actions and rewards occurred.

Let  $q_*(a)$  be the real **value** of an action  $a$ :

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

Denoting  $Q_t(a)$  our estimated value of action  $a$  at time  $t$  (before taking trial  $t$ ), we would like  $Q_t(a)$  to converge to  $q_*(a)$ . A natural way to estimate  $Q_t(a)$  is

$$Q_t(a) \stackrel{\text{def}}{=} \frac{\text{sum of rewards when action } a \text{ is taken}}{\text{number of times action } a \text{ was taken}}.$$

Following the definition of  $Q_t(a)$ , we could choose a **greedy** action  $A_t$  as

$$A_t \stackrel{\text{def}}{=} \arg \max_a Q_t(a).$$

## Exploitation versus Exploration

Choosing a greedy action is **exploitation** of current estimates. We however also need to **explore** the space of actions to improve our estimates.

An  $\epsilon$ -greedy method follows the greedy action with probability  $1 - \epsilon$ , and chooses a uniformly random action with probability  $\epsilon$ .

# $\epsilon$ -greedy Method

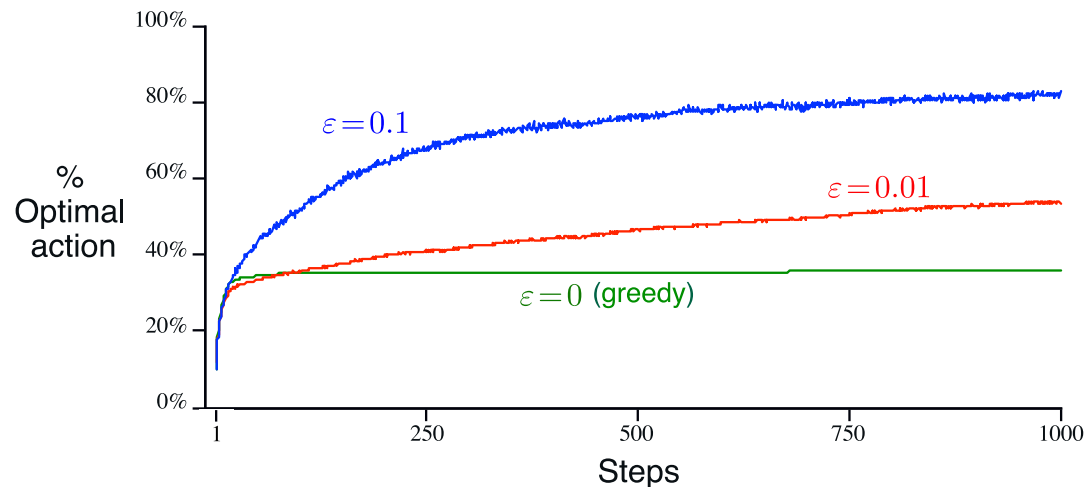
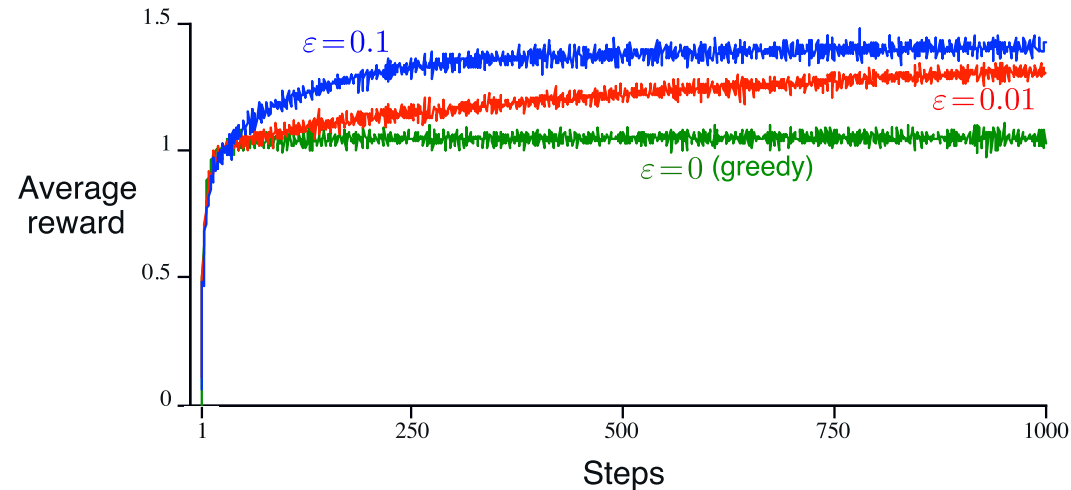


Figure 2.2 of "Reinforcement Learning: An Introduction, Second Edition".

## Incremental Implementation

Let  $Q_{n+1}$  be an estimate using  $n$  rewards  $R_1, \dots, R_n$ .

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

## A simple bandit algorithm

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

$$Q(a) \leftarrow 0$$

$$N(a) \leftarrow 0$$

Loop forever:

$$A \leftarrow \begin{cases} \operatorname{argmax}_a Q(a) & \text{with probability } 1 - \epsilon \quad (\text{breaking ties randomly}) \\ \text{a random action} & \text{with probability } \epsilon \end{cases}$$

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

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

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

Algorithm 2.4 of "Reinforcement Learning: An Introduction, Second Edition".

# Fixed Learning Rate

Analogously to the solution obtained for a stationary problem, we consider

$$Q_{n+1} = Q_n + \alpha(R_n - Q_n).$$

Converges to the true action values if

$$\sum_{n=1}^{\infty} \alpha_n = \infty \quad \text{and} \quad \sum_{n=1}^{\infty} \alpha_n^2 < \infty.$$

Biased method, because

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

The bias can be utilized to support exploration at the start of the episode by setting the initial values to more than the expected value of the optimal solution.

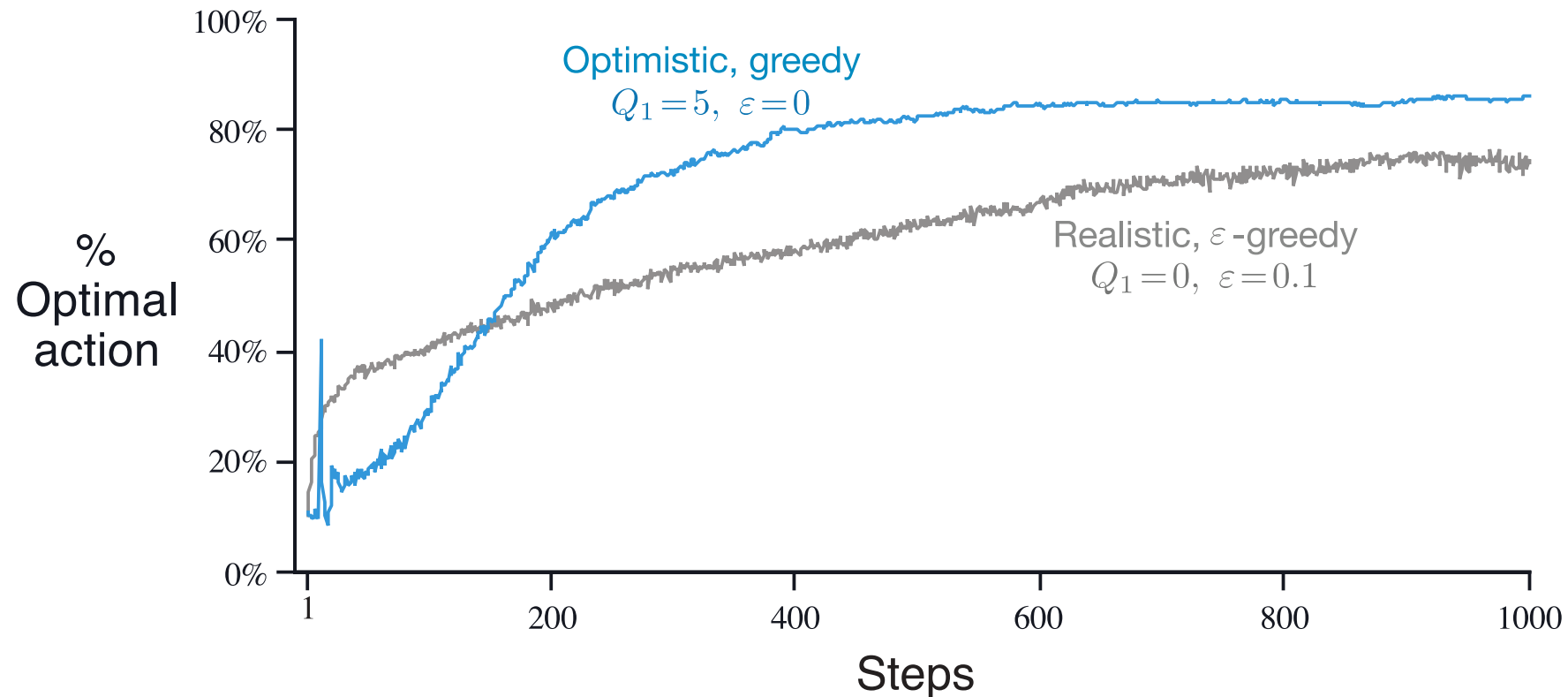


Figure 2.3 of "Reinforcement Learning: An Introduction, Second Edition".

# Method Comparison

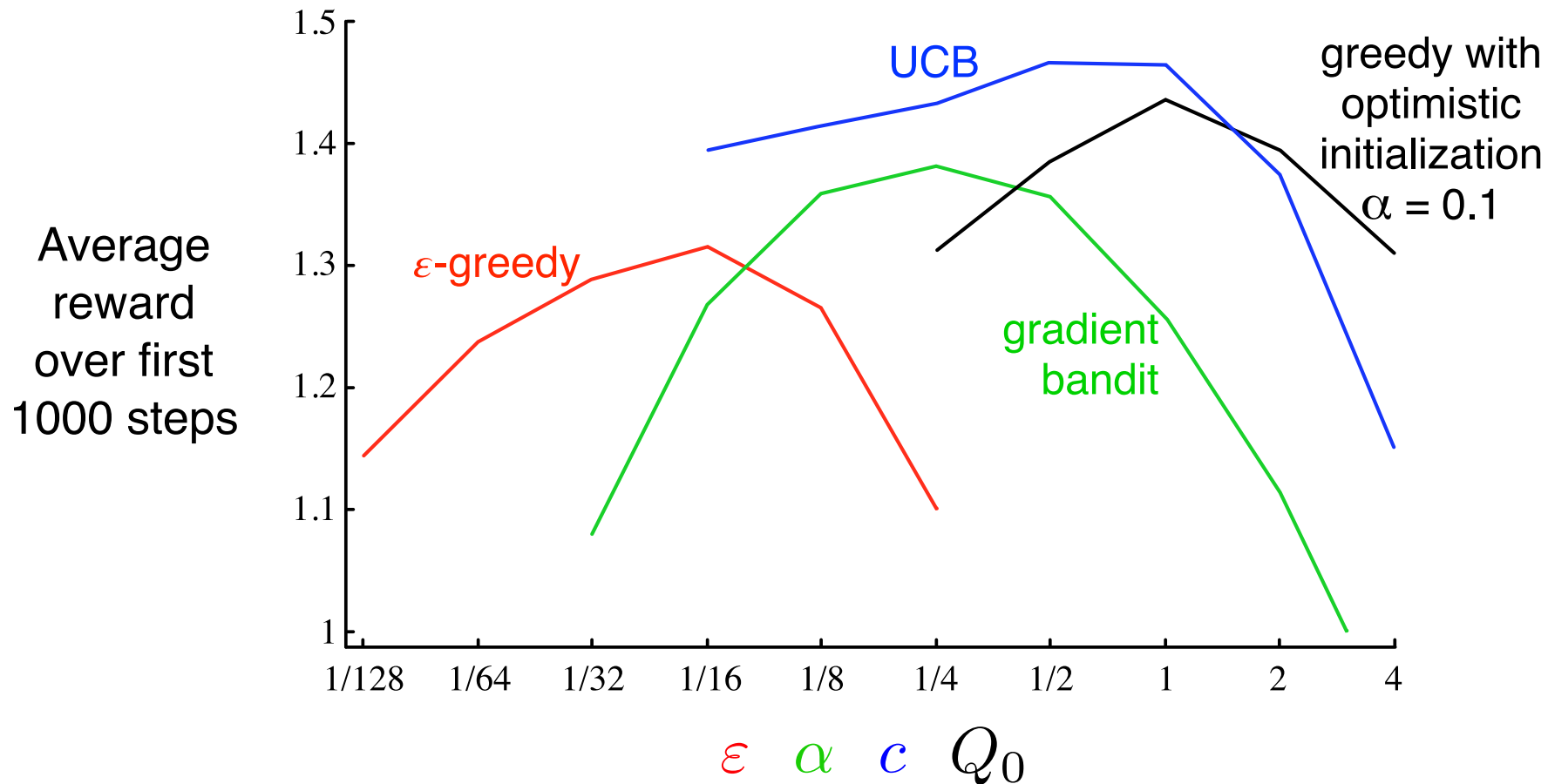
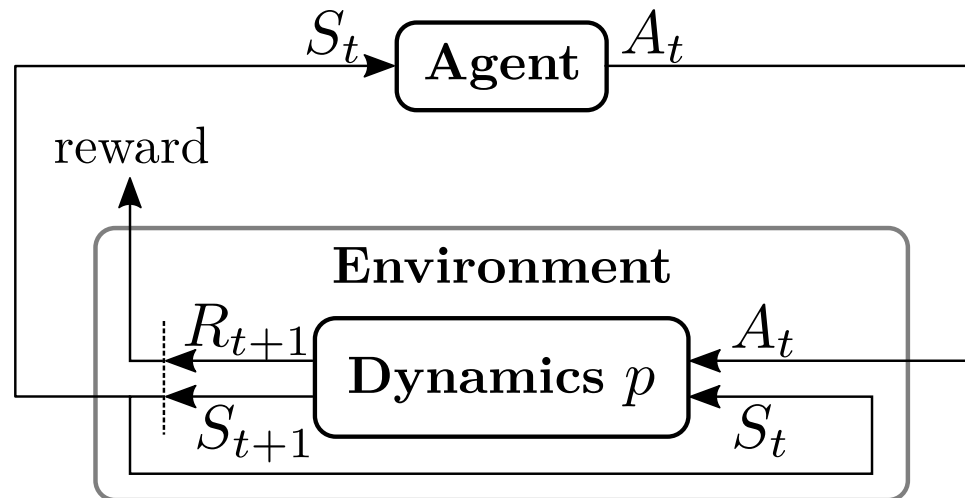


Figure 2.6 of "Reinforcement Learning: An Introduction, Second Edition".





A **Markov decision process** (MDP) is a quadruple  $(\mathcal{S}, \mathcal{A}, p, \gamma)$ , where:

- $\mathcal{S}$  is a set of states,
- $\mathcal{A}$  is a set of actions,
- $p(S_{t+1} = s', R_{t+1} = r | S_t = s, A_t = a)$  is a probability that action  $a \in \mathcal{A}$  will lead from state  $s \in \mathcal{S}$  to  $s' \in \mathcal{S}$ , producing a **reward**  $r \in \mathbb{R}$ ,
- $\gamma \in [0, 1]$  is a **discount factor**.

Let a **return**  $G_t$  be  $G_t \stackrel{\text{def}}{=} \sum_{k=0}^{\infty} \gamma^k R_{t+1+k}$ . The goal is to optimize  $\mathbb{E}[G_0]$ .

We cannot replace

- $p(S_{t+1} = s', R_{t+1} = r | S_t = s, A_t = a)$  is a probability that action  $a \in \mathcal{A}$  will lead from state  $s \in \mathcal{S}$  to  $s' \in \mathcal{S}$ , producing a **reward**  $r \in \mathbb{R}$ ,

by

- $p(S_{t+1} = s' | S_t = s, A_t = a)$ , a transition probability,
- $r(R_{t+1} = r | S_t = s, A_t = a)$ , a reward probability.

because the reward might depend on  $S_{t+1}$ .

However, we could use

- $p(S_{t+1} = s' | S_t = s, A_t = a)$ , a transition probability,
- $r(R_{t+1} = r | S_{t+1} = s', S_t = s, A_t = a)$ , a reward probability.

# Multi-armed Bandits as MDP

To formulate  $n$ -armed bandits problem as MDP, we do not need states. Therefore, we could formulate it as:

- one-element set of states,  $\mathcal{S} = \{S\}$ ;
- an action for every arm,  $\mathcal{A} = \{a_1, a_2, \dots, a_n\}$ ;
- assuming every arm produces rewards with a distribution of  $\mathcal{N}(\mu_i, \sigma_i^2)$ , the MDP dynamics function  $p$  is defined as

$$p(S, r | S, a_i) = \mathcal{N}(r | \mu_i, \sigma_i^2).$$

One possibility to introduce states in multi-armed bandits problem is to consider a separate reward distribution for every state. Such generalization is called **Contextualized Bandits** problem. Assuming state transitions are independent on rewards and given by a distribution  $next(s)$ , the MDP dynamics function for contextualized bandits problem is given by

$$p(s', r | s, a_i) = \mathcal{N}(r | \mu_{i,s}, \sigma_{i,s}^2) \cdot next(s' | s).$$

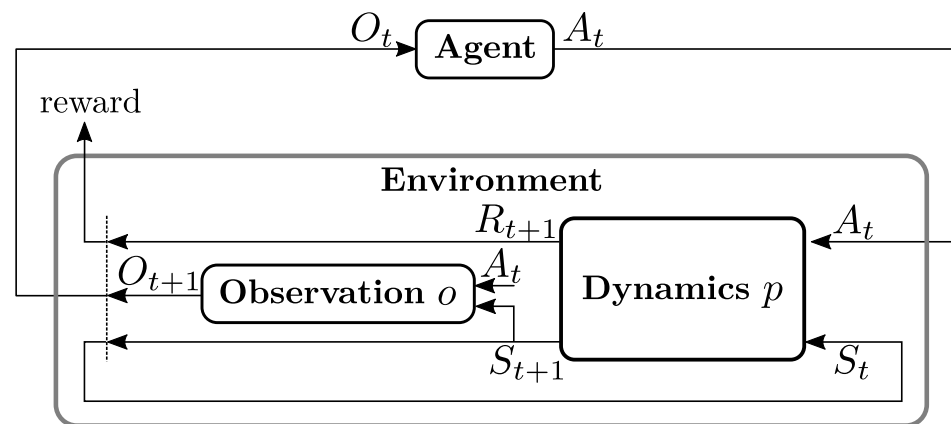
# Partially Observable MDPs

Recall that the Markov decision process is a quadruple  $(\mathcal{S}, \mathcal{A}, p, \gamma)$ , where:

- $\mathcal{S}$  is a set of states,
- $\mathcal{A}$  is a set of actions,
- $p(\mathcal{S}_{t+1} = s', R_{t+1} = r | \mathcal{S}_t = s, A_t = a)$  is a probability that action  $a \in \mathcal{A}$  will lead from state  $s \in \mathcal{S}$  to  $s' \in \mathcal{S}$ , producing a reward  $r \in \mathbb{R}$ ,
- $\gamma \in [0, 1]$  is a discount factor.

**Partially observable Markov decision process** extends the Markov decision process to a sextuple  $(\mathcal{S}, \mathcal{A}, p, \gamma, \mathcal{O}, o)$ , where in addition to an MDP,

- $\mathcal{O}$  is a set of observations,
- $o(\mathcal{O}_{t+1} | \mathcal{S}_{t+1}, A_t)$  is an observation model, where observation  $O_t$  is used as agent input instead of the state  $\mathcal{S}_t$ .



# Partially Observable MDPs

Planning in a general POMDP is in theory undecidable.

- Nevertheless, several approaches are used to handle POMDPs in robotics
  - to model uncertainty, imprecise mechanisms and inaccurate sensors, ...
  - consider for example robotic vacuum cleaners

Partially observable MDPs are needed to model many environments (maze navigation, FPS games, ...).

- We will initially assume all environments are fully observable, even if some of them will not.
- Later we will mention solutions, where partially observable MDPs are handled using recurrent networks (or networks with external memory), which model the latent states  $S_t$ .

# Episodic and Continuing Tasks

If the agent-environment interaction naturally breaks into independent subsequences, usually called **episodes**, we talk about **episodic tasks**. Each episode then ends in a special **terminal state**, followed by a reset to a starting state (either always the same, or sampled from a distribution of starting states).

In episodic tasks, it is often the case that every episode ends in at most  $H$  steps. These **finite-horizon tasks** then can use discount factor  $\gamma = 1$ , because the return  $G \stackrel{\text{def}}{=} \sum_{t=0}^H \gamma^t R_{t+1}$  is well defined.

If the agent-environment interaction goes on and on without a limit, we instead talk about **continuing tasks**. In this case, the discount factor  $\gamma$  needs to be sharply smaller than 1.

# (State-)Value and Action-Value Functions

A **policy**  $\pi$  computes a distribution of actions in a given state, i.e.,  $\pi(a|s)$  corresponds to a probability of performing an action  $a$  in state  $s$ .

To evaluate a quality of a policy, we define **value function**  $v_\pi(s)$ , or **state-value function**, as

$$\begin{aligned} v_\pi(s) &\stackrel{\text{def}}{=} \mathbb{E}_\pi [G_t | S_t = s] = \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \middle| S_t = s \right] \\ &= \mathbb{E}_{A_t \sim \pi(s)} \mathbb{E}_{S_{t+1}, R_{t+1} \sim p(s, A_t)} \left[ R_{t+1} + \gamma \mathbb{E}_{A_{t+1} \sim \pi(S_{t+1})} \mathbb{E}_{S_{t+2}, R_{t+2} \sim p(S_{t+1}, A_{t+1})} [R_{t+2} + \dots] \right] \end{aligned}$$

An **action-value function** for a policy  $\pi$  is defined analogously as

$$q_\pi(s, a) \stackrel{\text{def}}{=} \mathbb{E}_\pi [G_t | S_t = s, A_t = a] = \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \middle| S_t = s, A_t = a \right].$$

The value function and action-value function can be of course expressed using one another:

$$v_\pi(s) = \mathbb{E}_{a \sim \pi} [q_\pi(s, a)], \quad q_\pi(s, a) = \mathbb{E}_{s', r \sim p} [r + \gamma v_\pi(s')].$$

Optimal state-value function is defined as

$$v_*(s) \stackrel{\text{def}}{=} \max_{\pi} v_{\pi}(s),$$

analogously

$$q_*(s, a) \stackrel{\text{def}}{=} \max_{\pi} q_{\pi}(s, a).$$

Any policy  $\pi_*$  with  $v_{\pi_*} = v_*$  is called an **optimal policy**. Such policy can be defined as  $\pi_*(s) \stackrel{\text{def}}{=} \arg \max_a q_*(s, a) = \arg \max_a \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) | S_t = s, A_t = a]$ . When multiple actions maximize  $q_*(s, a)$ , the optimal policy can stochastically choose any of them.

## Existence

In finite-horizon tasks or if  $\gamma < 1$ , there always exists a unique optimal state-value function, a unique optimal action-value function, and a (not necessarily unique) optimal policy.



# Monte Carlo Prediction

Assuming we have a fixed policy  $\pi$  and that we want to estimate  $v_\pi(s)$ .

A Monte Carlo method to estimate this value function would be to simulate many episodes, and then compute an average *return* for all visited states:

$$V(s) \approx \mathbb{E}[G_t | S_t = s].$$

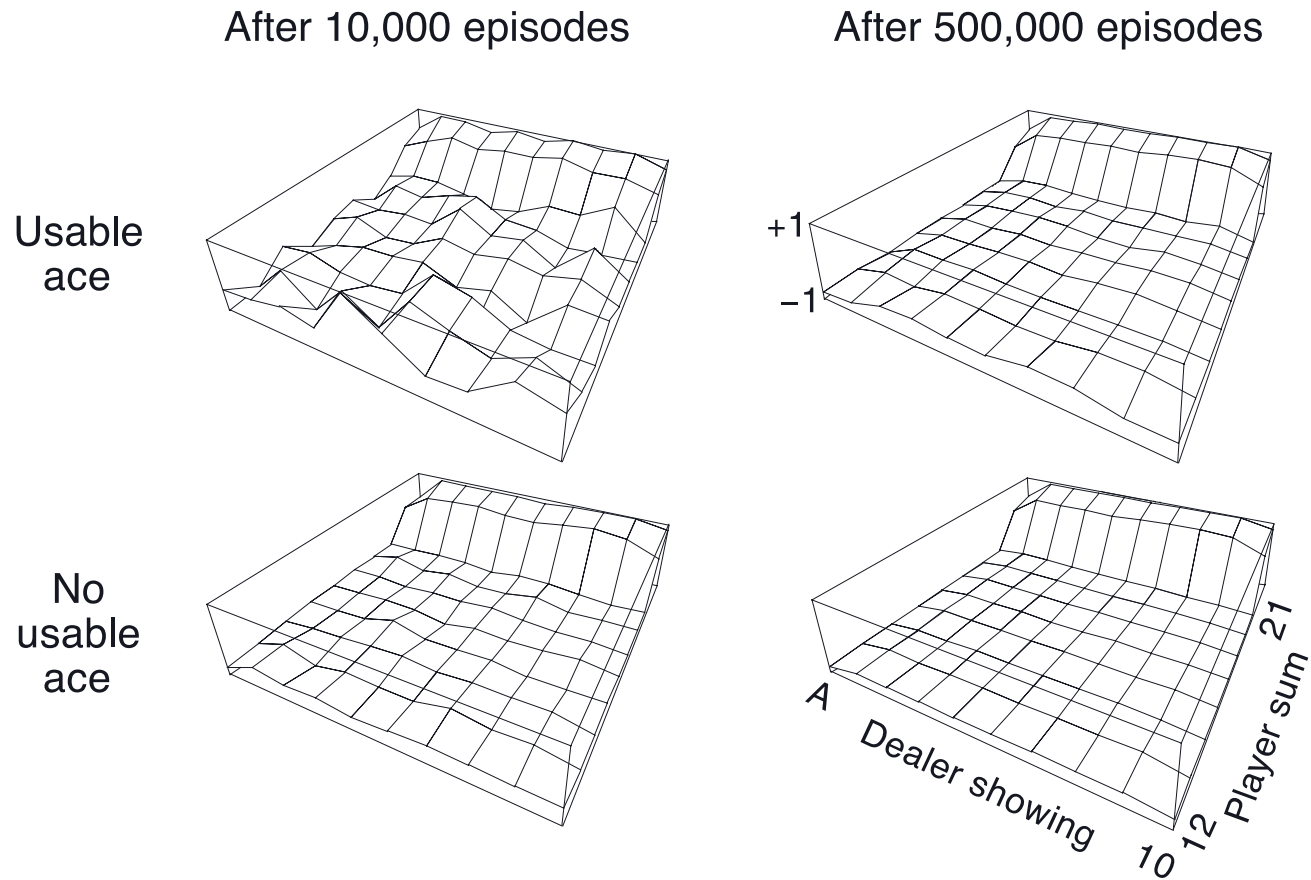
Some states might be visited multiple times; in that case, we could use

- **first-visit** Monte Carlo method, where only the first occurrence of the state is considered;
- **every-visit** Monte Carlo method, where all occurrences of the state are considered.

By the law of large numbers, the Monte Carlo estimate converges to the real value function for all visited states.

- Actually, for every-visit MC, it is more complicated, because multiple returns from a single state in a single episode are not independent; but it can be proven that even every-visit Monte Carlo converges.

# Monte Carlo Prediction of Blackjack



**Figure 5.1:** Approximate state-value functions for the blackjack policy that sticks only on 20 or 21, computed by Monte Carlo policy evaluation. ■

*Figure 5.1 of "Reinforcement Learning: An Introduction, Second Edition".*

We now present the first algorithm for computing optimal behavior without assuming a knowledge of the environment dynamics.

However, we still assume there are finitely many states  $\mathcal{S}$  and we will store estimates for each of them (a *tabular* method).

Monte Carlo methods are based on estimating returns from complete episodes. Specifically, they try to estimate

$$Q(s, a) \approx \mathbb{E}[G_t | S_t = s, A_t = a].$$

With such estimates, a greedy action in state  $S_t$  can be computed as

$$A_t = \arg \max_a Q(S_t, a).$$

To hope for convergence, we need to visit each state-action pair infinitely many times. One of the simplest way to achieve that is to assume **exploring starts**, where we randomly select the first state and first action, and behave greedily afterwards.

## Monte Carlo ES (Exploring Starts), for estimating $\pi \approx \pi_*$

Initialize:

$\pi(s) \in \mathcal{A}(s)$  (arbitrarily), for all  $s \in \mathcal{S}$

$Q(s, a) \in \mathbb{R}$  (arbitrarily), for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$

$Returns(s, a) \leftarrow$  empty list, for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$

Loop forever (for each episode):

Choose  $S_0 \in \mathcal{S}$ ,  $A_0 \in \mathcal{A}(S_0)$  randomly such that all pairs have probability  $> 0$

Generate an episode from  $S_0, A_0$ , following  $\pi$ :  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode,  $t = T-1, T-2, \dots, 0$ :

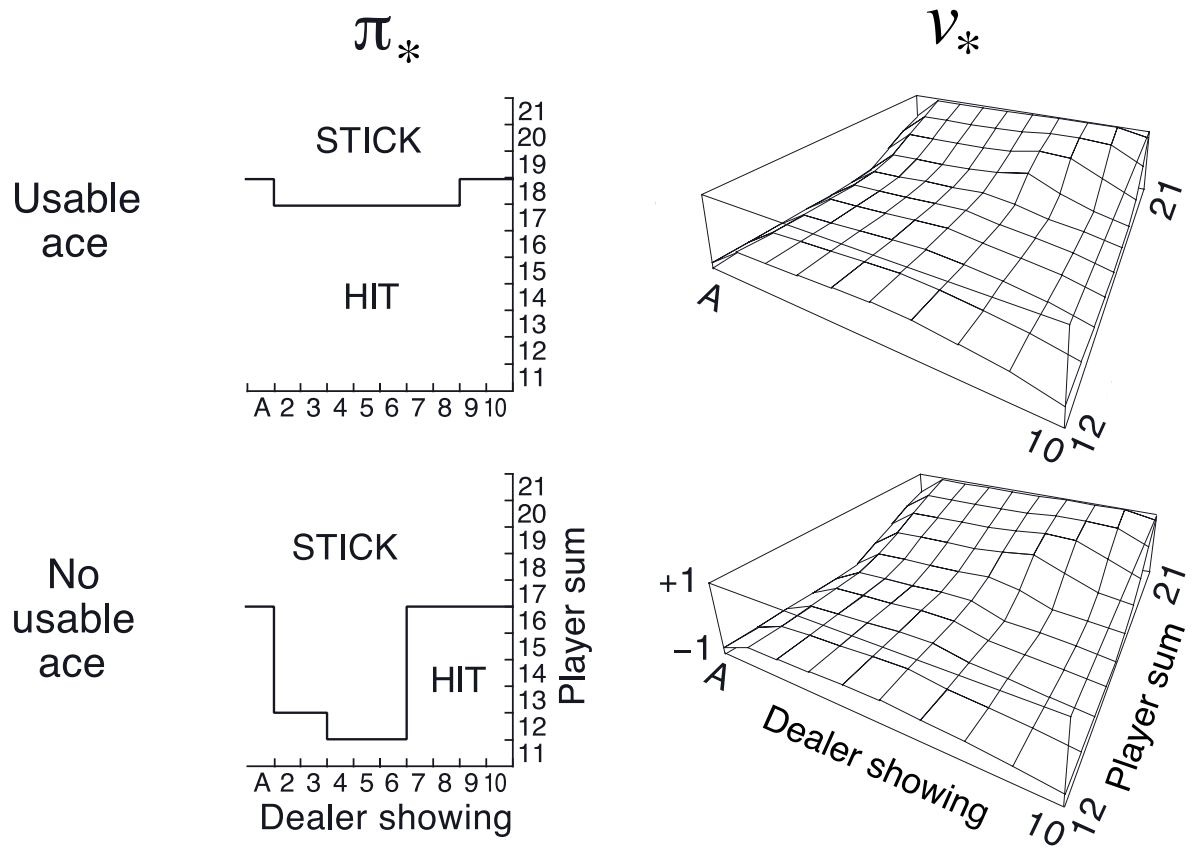
$G \leftarrow \gamma G + R_{t+1}$

Append  $G$  to  $Returns(S_t, A_t)$

$Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$

$\pi(S_t) \leftarrow \text{argmax}_a Q(S_t, a)$

*Modification of algorithm 5.3 of "Reinforcement Learning: An Introduction, Second Edition" from first-visit to every-visit.*



**Figure 5.2:** The optimal policy and state-value function for blackjack, found by Monte Carlo ES. The state-value function shown was computed from the action-value function found by Monte Carlo ES. ■

Figure 5.2 of "Reinforcement Learning: An Introduction, Second Edition".

# Monte Carlo and $\epsilon$ -soft Behavior

The problem with exploring starts is that in many situations, we either cannot start in an arbitrary state, or it is impractical.

Instead of choosing random state at the beginning, we can consider adding “randomness” gradually – for a given  $\epsilon$ , we set the probability of choosing any action to be at least

$$\frac{\epsilon}{|\mathcal{A}(s)|}$$

in each step. Such behavior is called  $\epsilon$ -soft.

In an  $\epsilon$ -soft behaviour, selecting an action greedily (the  $\epsilon$ -greedy behavior) means one action has a maximum probability of

$$1 - \epsilon + \frac{\epsilon}{|\mathcal{A}(s)|}.$$

We now present Monte Carlo algorithm with  $\epsilon$ -greedy action selection.

## On-policy every-visit Monte Carlo for $\varepsilon$ -soft Policies

Algorithm parameter: small  $\varepsilon > 0$

Initialize  $Q(s, a) \in \mathbb{R}$  arbitrarily (usually to 0), for all  $s \in \mathcal{S}, a \in \mathcal{A}$

Initialize  $C(s, a) \in \mathbb{Z}$  to 0, for all  $s \in \mathcal{S}, a \in \mathcal{A}$

Repeat forever (for each episode):

- Generate an episode  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$ , by generating actions as follows:
  - With probability  $\varepsilon$ , generate a random uniform action
  - Otherwise, set  $A_t \stackrel{\text{def}}{=} \arg \max_a Q(S_t, a)$
- $G \leftarrow 0$
- For each  $t = T - 1, T - 2, \dots, 0$ :
  - $G \leftarrow \gamma G + R_{t+1}$
  - $C(S_t, A_t) \leftarrow C(S_t, A_t) + 1$
  - $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{1}{C(S_t, A_t)} (G - Q(S_t, A_t))$