
The Mathematics of Reinforcement Learning

LEIF DÖRING

Contents

1	Stochastic bandits	1
1.1	A quick dive into two-stage stochastic experiments	1
1.2	Introduction to stochastic bandits	2
1.3	Algorithms: the exploration-exploitation trade-off	10
1.3.1	Basic explore-then-committ algorithm (ETC)	11
1.3.2	From greedy to UCB	15
1.3.3	Boltzmann exploration	25
1.3.4	Simple policy gradient for stochastic bandits	27
1.4	Lower bounds for stochastic bandits	31
1.4.1	A bit on relative entropy	32
1.4.2	Mini-Max lower bounds (model-dependent)	34
1.4.3	Asymptotic lower bound (model-dependent)	37
I	Tabular Reinforcement Learning	41
2	Basics: Stochastic Control and Dynamic Programming Methods	43
2.1	Markov decision process formalism	43
2.1.1	A quick dive into Markov (reward) chains	43
2.1.2	Markov decision processes	46
2.1.3	Elements of stochastic control theory	58
2.2	Basic tabular value iteration algorithm	68
2.3	Basic policy iteration (actor-critic) algorithm	71
2.3.1	Policy evaluation	71
2.3.2	Policy improvement	74
2.3.3	Policy iteration algorithms (tabular actor-critic)	77
2.3.4	Relation of policy iteration and value iteration	80
2.4	Stochastic control in finite time	80
2.4.1	Setting and dynamic programming	81
2.4.2	Dynamic programming algorithms	85
II	Non-Tabular Reinforcement Learning	87
3	A quick dive into gradient ascent/descent methods	91
3.1	Gradient descent for L -smooth functions	92
3.2	Gradient descent for L -smooth, convex functions	94
3.3	Gradient descent for L -smooth functions with PL inequality	97
3.4	Gradient descent with diminishing step-sizes	99
3.5	Stochastic gradient descent methods	101
3.6	Regularisation	104

4	Policy Gradient methods	105
4.1	Policy gradient theorems	107
4.1.1	Finite-time undiscounted MDPs, stationary policies	107
4.1.2	Infinite-time MDPs with discounted rewards	112
4.1.3	Convergence of REINFORCE	118
4.1.4	Baseline	123
4.2	PPO	125
4.2.1	Bias of the PPO surrogate gradient	125
5	Proximal policy optimization (PPO)	129
5.1	REINFORCE with cyclic off-policy updates	130
5.2	Generalized advantage estimation, actor-critic	131
5.2.1	Infinite horizon GAE	134
5.2.2	Finite time-horizon GAE (new!)	134
5.2.3	Termination-time GAE (new)	135
5.2.4	The backwards induction	135
5.2.5	Relationship of the finite and the stopping time GAE	135
5.2.6	Unbiasedness of the estimators	137
5.3	PPO as importance sampling variance reduction method	139
5.3.1	Bias of the PPO surrogate gradient	139
5.4	Effect of Large Mean KL on Gradient Magnitude	149
5.4.1	Setup and Definitions	149
5.4.2	Mean KL Divergence and Ratio Deviations	149
5.4.3	Expected Gradient Magnitude	150
5.4.4	High-Probability Small-Gradient Bound	150
5.4.5	Conclusion	150
5.4.6	Variance of the PPO surrogate gradient with GAE estimator	151
5.5	First rollout based variant of PPO	153
5.6	Theory	157
5.6.1	Improved GAE	157
5.6.2	Convergence	157
6	Reinforcement learning with function approximation	159
6.1	Policy evaluation with function approximation	159
6.1.1	Function approximation	160
6.1.2	Sample based policy evaluation with function approximation	161
6.2	Approximate policy improvement	165
6.3	Generic bounds for policy iteration with approximations	165
6.4	Q-learning with linear function approximation	169
6.5	Policy gradient with linear function approximation	169
6.6	A quick dive into neural networks	170
6.6.1	What are neural network functions?	170
6.6.2	Approximation properties	172
6.6.3	Differentiation properties	175
6.6.4	Using neural networks to approximate value functions	178
6.6.5	Using neural networks to parametrise policies	178
6.7	Deep Q-learning (DQN)	179
6.8	Deep policy gradient (actor critic methods)	179
6.8.1	Simple actor-critic (AC) and advantage actor-critic (A2C)	180
6.8.2	Soft actor-critic (SAC)	183
6.8.3	Proximal policy optimisation (PPO)	183
7	Monte Carlo Tree Search	185
8	Policy-Gradient Methods	187

Chapter 1

Stochastic bandits

.Multiarmed bandits can be considered to be the simplest situation in which optimal decision making can be learnt. Due to its simple structure many ideas get more visible that we will get to know much later for the more general setup of Markov decision processes. In fact, a vast literature of precise results exists for multiarmed bandits in contrast to many unproved methods for the general Markov decision situation. What we will try to achieve in this first chapter is to bridge the two worlds. We discuss ideas and methods that will be crucial for Markov decision processes mostly in the precise language of multiarmed bandits. The aim is to find the right compromise of the very imprecise Chapter 2 of Sutton and Barton and the very detailed book of Lattimore and Szepesvári.



The chapter focuses on algorithmic ideas towards the trade-off between exploration and exploitation in optimal decision making.

The language used in this chapter will thus be a combination of bandit and reinforcement learning (RL) notation. The chapter can also be seen as a motivation for later chapters to see how little we actually understand about general reinforcement learning compared to the well-understood case of multiarmed bandits.

1.1 A quick dive into two-stage stochastic experiments

From a probabilistic point of view reinforcement learning will always be about stochastic two-stage experiments. First choose an action (first experiment) and given that action observe a reward (second experiment). It's a bit of a mathematical overkill but let us quickly discuss the basics from two-stage experiments. If all experiments involved are discrete (i.e. take finitely or countably infinite values) then not much is needed and all computations work with the usual conditional probabilities defined by $\mathbb{P}(X = x|Y = y) = \frac{\mathbb{P}(X=x, Y=y)}{\mathbb{P}(Y=y)}$. The choice of actions indeed is discrete in most applications, unfortunately, the rewards often are not. Going towards non-discrete probability the concept of regular conditional expectation is needed. We are not going into detail and only summarize the most important facts. First, suppose (X, Y) is a pair of discrete random variables (or random vectors) on some probability space $(\Omega, \mathcal{A}, \mathbb{P})$, then rewriting the definition of conditional probabilities gives

$$\mathbb{P}(X = x, Y = y) = \underbrace{\mathbb{P}(X = x|Y = y)}_{\text{second stage, given result of first stage}} \cdot \underbrace{\mathbb{P}(Y = y)}_{\text{first stage}}$$

so that all quantities jointly involving X and Y can be computed by knowing the outcome of the first stage and the second stage given the outcome of the first stage. For computations one

can use the rules

$$\begin{aligned}\mathbb{E}[h(X, Y)] &= \sum_{x, y} h(x, y) \mathbb{P}(X = x | Y = y) \mathbb{P}(Y = y), \\ \mathbb{E}[h(X, Y) | Y = y] &= \sum_x h(x, y) \mathbb{P}(X = x | Y = y), \\ \mathbb{E}[h(X, Y) | Y] &= \sum_{x, y} h(x, y) \mathbb{P}(X = x | Y = y) \mathbf{1}_{Y=y}, \\ \mathbb{P}(X \in \cdot | Y) &= \sum_y \mathbb{P}(X \in \cdot | Y = y) \mathbf{1}_{Y=y}.\end{aligned}$$

The situation is more delicate if the distribution of the first stage (say Y) is not discrete. In that case $\mathbb{P}(X = x | Y = y)$ cannot be defined easily as the definition of (discrete) conditional probabilities would require division by 0 if $\mathbb{P}(Y = y) = 0$. Nonetheless, for computations there is an abstract way around it. One can prove the existence of a so-called Markov kernel¹ $k(\cdot, \cdot)$ such that the rules

$$\mathbb{E}[h(X, Y) | Y = y] = \int h(x, y) k(y, dx) \quad \text{and} \quad \mathbb{E}[h(X, Y) | Y] = \int h(x, Y) k(Y, dx)$$

hold. The kernel κ is called a conditional regular expectation of X given Y . In the discrete case it holds that $k(y, A) = \mathbb{P}(X \in A | Y = y)$ according to the usual definition. For absolutely continuous pairs (X, Y) the measure $k(y, \cdot)$ has density $k(y, x) = f_{(X, Y)}(x, y) / f_X(x) f_Y(y)$ but in general k is abstract. Nonetheless, to keep the analogy to the discrete setting one typically writes $\mathbb{P}(X \in A | Y = y)$ instead of $k(y, A)$ even if the conditional probability cannot be defined in the usual way. We also use the notation $\mathbb{P}(X \in A | Y) = \kappa(Y, A)$ from which it follows that $\mathbb{E}[h(X, Y) | Y] = \int h(x, Y) \mathbb{P}(X \in dx | Y)$. For later purposes it is crucial to know that for independent X and Y it holds that $k(y, A) = \mathbb{P}(X \in A)$ so that $\mathbb{E}[h(X, Y) | Y] = \int h(x, Y) \mathbb{P}(X \in dx)$.

1.2 Introduction to stochastic bandits

As there is no need to loose the reader in the mathematical formalisation of multiarmed bandits we will first gently introduce the ideas. For a multiarmed bandit there is a number of possible experiments among which a learner (in RL called agent) has the target to identify the best arm by observing random samples of the experiments. The goal is to learn efficiently which experiment yields the best outcome. There are different ways of defining what best means, in bandit theory best typically means the highest expectation. The simplest example to keep in mind is a daily visit to the university canteen. Let's say the canteen offers four choices: vegan, vegetarian, classic, oriental. These are the four possible experiments, the random outcomes are the quality of the dish (measured in some way, such as on a 1-10 scale). There are plenty of other situations that immediately come to mind such as

- medical treatment of patients with different doses, outcome measures the success, for instance 1 for healed and 0 otherwise,
- placing different advertisements on websites, outcome measures the click rates.

The wording multiarmed bandit is rather historic and comes from gambling. A one-armed bandit is a gambling machine in which every round of the game yields a random reward. Typically there are several one-armed bandits next to each other and a player aims to play the most favorable bandit. Since every round has a fixed cost, say one Euro, the player tries to figure out as quickly as possible which machine works best for him/her. In these notes we will only discuss so-called stochastic bandits. These are bandits where the random experiments are stationary, i.e. the distribution of the experiments do not change over time. More general bandit situations are adversarial bandits and contextual bandits which we will not touch in this section.

¹A Markov kernel is a measure in the second and a measurable mapping in the first entry



The basic mathematical model behind a stochastic bandits is as follows. There are distributions P_{a_1}, \dots, P_{a_K} (not necessarily discrete) from which an outcome is sampled if an actor decides to play a so-called arm a . The outcome is called X . If the choice of the actor is modeled using a random variable A taking values a_1, \dots, a_K then the two-stage experiment yields

$$\mathbb{P}(X \in \cdot | A = a) = P_a(\cdot) \quad \text{and} \quad \mathbb{P}(X \in \cdot | A) = \sum_{a \in \mathcal{A}} P_a(\cdot) \mathbf{1}_{Y=a} =: P_A(\cdot)$$

so that

$$\mathbb{P}(X \in \cdot, A = a) = \mathbb{P}(X \in \cdot | A = a) \mathbb{P}(A = a) = P_a(\cdot) \mathbb{P}(A = a).$$

In fact, a bandit model will depend not only on one action but on reward/action pairs of the past. This makes the mathematical framing a bit more tedious as rewards must not be discrete.

Before we go into more details let us discuss the optimization goal. Suppose we have a finite set $\mathcal{A} = \{a_1, \dots, a_K\}$ of experiments (arms to play) and denote by Q_a the expectation of the random outcome whose distribution we denote by P_a . Of course there could be other goals than finding the arm with the highest average outcome, but this is what we decide to optimize. Now fix a time-horizon n (which could be infinite), the number of rounds we can use for learning. A learning strategy is a sequence $(\pi_t)_{t=1, \dots, n}$ of probability distributions on \mathcal{A} that only depend on what has been played prior to time t . Here $\pi_t(\{a\})$ is the probability that the player chooses action a at time t and then receives the random outcome of the corresponding experiment.



Throughout these lecture notes we will be sloppy about measures on the power-set of discrete (finite or countably infinite) sets. Instead of writing $p(\{a\})$ we typically write $p(a)$ for probabilities of singleton sets if there is no danger of confusion.

The aim is to find a learning strategy that maximises the outcome of this procedure. To have an idea in mind think of the example of the university canteen. During your studies you might have $n = 500$ choices in total. If you are not vegan or vegetarian you probably started without preferences, i.e. $\pi_t(\text{vegan}) = \dots = \pi_t(\text{oriental}) = \frac{1}{4}$, and over time learnt from experience how to change the following distributions π_t in order to maximise the lunch experience.

Let's turn these basic ideas into mathematics.



Definition 1.2.1. Suppose \mathcal{A} is an index-set and $\nu = \{P_a\}_{a \in \mathcal{A}}$ is a family of real-valued distributions with finite expectations, called the reward distributions.

- The set ν is called a stochastic bandit model. In these lectures we will always assume $\mathcal{A} = \{a_1, \dots, a_K\}$ is finite, K is the number of arms. Often it will be useful to denote the arms by $1, \dots, K$ to simplify formulas.
- The action value (or Q -value) of an arm is defined by the expectation $Q_a := \int_{\mathbb{R}} x dP_a(x)$. A best arm, usually denoted by a_* , is an arm with highest Q -value, i.e.

$$a^* = \operatorname{argmax}_{a \in \mathcal{A}} Q_a.$$

Typically one abbreviates Q_* for the largest action value Q_{a^*} and if there are several optimal arms the argmax chooses any of them.

- A learning strategy for n rounds ($n = +\infty$ is allowed) consists of
 - an initial distribution π_1 on \mathcal{A} ,
 - a sequence $(\pi_t)_{t=2, \dots, n}$ of kernels on $\Omega_{t-1} \times \mathcal{A}$,



where Ω_t denotes all trajectories $(a_1, x_1, a_2, x_2, \dots, a_t, x_t) \in (\mathcal{A} \times \mathbb{R})^t$. We will write the kernels in reverse ordering of the arguments

$$\pi_t(\cdot; a_1, x_1, a_2, x_2, \dots, a_{t-1}, x_{t-1})$$

with the meaning that $\pi_t(a; a_1, x_1, a_2, x_2, \dots, a_t, x_t)$ is the probability arm a is chosen at time t if the past rounds resulted in actions/rewards $a_1, x_1, a_2, x_2, \dots, a_t, x_t$.

Recall that a probability distribution on a finite set is nothing but a probability vector (numbers in $[0, 1]$) that sum up to 1. An important special case occurs if the vector consists of one 1 (and 0s otherwise), i.e. the measure is a Dirac measure.



We will always assume that the action values Q_a are unknown but the bandit is a generative model, i.e. the random variables can be sampled. Everything that can be learnt about the model must be achieved by simulations (playing arms). In principle the learning strategy should be written down, most of the time we will construct the kernels π_t round by round using algorithms.

We will later see that learning strategies typically depend on different ingredients such as the maximal time-horizon if this is known in advance or certain quantities of the underlying bandit model. Of course it is desirable to have as little dependences as possible, but often additional information is very useful.



Definition 1.2.2. A learning strategy is called an index strategy if all appearing measures are Dirac measures, i.e. in all situations only a single arm is played with probability 1. The learning strategy is called soft if in all situations all arms have strictly positive probabilities.

Of course index strategies are just a small subset of all strategies but most algorithms we study are index strategies.

Next, we introduce stochastic bandit (learning) processes. This is a bit in analogy to Markov chains where first one introduces transition matrices and then defines a Markov chain and proves the existence. Or a random variable where first one defines the distribution function and then proves the existence of a random variable.



Definition 1.2.3. Suppose ν is a stochastic bandit model and $(\pi_t)_{t=1, \dots, n}$ a learning strategy for n rounds. Then a stochastic process $(A_t^\pi, X_t^\pi)_{t=1, \dots, n}$ on a probability space $(\Omega, \mathcal{F}, \mathbb{P})$ is called a stochastic bandit (learning) process with learning strategy π if $A_1^\pi \sim \pi_1$ and

- $\mathbb{P}(A_t^\pi = a | A_1^\pi, X_1^\pi, \dots, A_{t-1}^\pi, X_{t-1}^\pi) = \pi_t(a; A_1^\pi, X_1^\pi, \dots, A_{t-1}^\pi, X_{t-1}^\pi),$
- $\mathbb{P}(X_t^\pi \in B | A_1^\pi, X_1^\pi, \dots, A_t^\pi) = P_{A_t^\pi}(B) := \sum_{a \in \mathcal{A}} P_a(B) \mathbf{1}_{A_t^\pi = a}$

for all $t = 1, \dots, n$. We will call A_t^π the action (the chosen arm) at time t and X_t^π the outcome (or reward) when playing arm A_t^π at time t . For notational convenience the superscripts π will typically be dropped if the learning strategy is clear from the context.

In words, a stochastic bandit (learning) process is a stochastic process that works in a two-step fashion. Given a learning strategy π , in every round the strategy suggests probabilities for actions based on the past behavior. The sampled action A_t is then used to play the corresponding arm

and observe the outcome. The process (A_t, X_t) thus describes the sequence of action/rewards over time.

Just as for random variables or Markov chains it is not completely trivial that there is a probability space and a stochastic process (A_t, X_t) that satisfies the defining properties of a stochastic bandit process.



Theorem 1.2.4. For every stochastic bandit model and every learning strategy $(\pi_t)_{t=1, \dots, n}$ there is a corresponding stochastic bandit (learning) process (A^π, X^π) .

Proof. We give a construction that is known under the name random table model as the bandit (learning) process is constructed from a table of independent random variables.



Recall that sampling from a discrete distribution π on \mathcal{A} can be performed using $\mathcal{U}([0, 1])$ random variables. Suppose $\pi(\{a_k\}) = p_k$ and $[0, 1]$ is subdivided into disjoint intervals I_k of lengths p_k , then the discrete random variable defined by $\bar{U} = a_k$ if and only if $U \in I_k$ is distributed according to π .

Now suppose $(X_t^{(a)})_{a \in \mathcal{A}, t \in \mathbb{N}}$ is a table of independent random variables such that $X_t^{(a)} \sim P_a$ for all t and suppose $(U_t)_{t \in \mathbb{N}}$ is a sequence of independent $\mathcal{U}([0, 1])$. These random variables (all defined on some joint probability space $(\Omega, \mathcal{F}, \mathbb{P})$) exist due to the Kolmogorov extension theorem.

A_1	A_2	A_3	A_4	
$\uparrow \pi_1$	$\uparrow \pi_2$	$\uparrow \pi_3$	$\uparrow \pi_4$	
U_1	U_2	U_3	U_4	\dots
$X_1^{(\mathbf{a}_1)}$	$X_2^{(a_1)}$	$X_3^{(a_1)}$	$X_4^{(\mathbf{a}_1)}$	\dots
$X_1^{(a_2)}$	$X_2^{(a_2)}$	$X_3^{(\mathbf{a}_2)}$	$X_4^{(a_2)}$	\dots
\vdots	\vdots	\vdots	\vdots	\vdots
$X_1^{(a_K)}$	$X_2^{(\mathbf{a}_K)}$	$X_3^{(a_K)}$	$X_4^{(a_K)}$	\dots

Construction of bandit (learning) processes: bold entries were selected as $X_1, X_2, X_3, X_4, \dots$

If $(\pi_t)_{t \in \mathbb{N}}$ a learning strategy then the stochastic bandit (learning) process is constructed as follows:

- $t = 1$: Use U_1 to sample from the discrete measure $\pi_1(\cdot)$ an arm a , denote this arm by A_1 and set $X_1 := X_1^{(A_1)}$.
- $t \mapsto t + 1$: Use U_{t+1} to sample from the discrete measure $\pi(\cdot; A_1, X_1, \dots, A_t, X_t)$ (relying only on the table to the left of column $t + 1$) an arm a , denote this arm by A_t and set $X_{t+1} = X_{t+1}^{(A_{t+1})}$.

To have a picture in mind think of a large table of random variables. Only using the variables to the left of column t , the uniform variable U_t is used to produce the action A_t and the X_t is produced by choosing the reward from row A_t . To see that this process (A, X) is indeed a bandit (learning) process with learning strategy π we need to check the defining properties. Let us denote by $I_t^a(a_1, \dots, x_{t-1})$ a partition of $[0, 1]$ into K disjoint intervals with lengths $\pi_t(a; a_1, \dots, x_{t-1})$. Then, using $h(u, a_1, \dots, x_{t-1}) := \mathbf{1}_{u \in I_t^a(a_1, \dots, x_{t-1})}$ and

$$\kappa(a_1, \dots, x_{t-1}, \cdot) = \mathbb{P}(U_t \in \cdot | A_1 = a_1, \dots, X_{t-1} = x_{t-1}) \stackrel{\text{ind.}}{=} \mathbb{P}(U_t \in \cdot) \stackrel{\mathcal{U}([0,1])}{=} \lambda(\cdot)$$

yields

$$\begin{aligned}
\mathbb{P}(A_t = a | A_1, X_1, \dots, A_{t-1}, X_{t-1}) &= \mathbb{P}(U_t \in I_t^a(A_1, \dots, X_{t-1}) | A_1, X_1, \dots, A_{t-1}, X_{t-1}) \\
&= \mathbb{E}[h(U_t, A_1, \dots, X_{t-1}) | A_1, X_1, \dots, A_{t-1}, X_{t-1}] \\
&= \int h(u, A_1, \dots, X_{t-1}) \kappa(A_1, \dots, X_{t-1}, du) \\
&= \int \mathbf{1}_{u \in I_t^a(A_1, \dots, X_{t-1})} du \\
&= \pi_t(a; A_1, \dots, X_{t-1}).
\end{aligned}$$

The second property is derived as follows:

$$\begin{aligned}
\mathbb{P}(X_t \in B | A_1, X_1, \dots, A_t) &= \sum_{a \in \mathcal{A}} \mathbb{P}(X_t \in B, A_t = a | A_1, X_1, \dots, A_t) \\
&= \sum_{a \in \mathcal{A}} \mathbb{P}(X_t^{(a)} \in B, A_t = a | A_1, X_1, \dots, A_t) \\
&= \sum_{a \in \mathcal{A}} \mathbb{E}[\mathbf{1}_{X_t^{(a)} \in B} \mathbf{1}_{A_t = a} | A_1, X_1, \dots, A_t] \\
&\stackrel{\text{meas.}}{=} \sum_{a \in \mathcal{A}} \mathbf{1}_{A_t = a} \mathbb{E}[\mathbf{1}_{X_t^{(a)} \in B} | A_1, X_1, \dots, A_t] \\
&\stackrel{\text{ind.}}{=} \sum_{a \in \mathcal{A}} \mathbf{1}_{A_t = a} \mathbb{E}[\mathbf{1}_{X_t^{(a)} \in B}] \\
&= \sum_{a \in \mathcal{A}} \mathbb{P}(X_t^{(a)} \in B) \mathbf{1}_{A_t = a} \\
&\stackrel{\text{Def.}}{=} P_{A_t}(B).
\end{aligned}$$

□

There is an equivalent way of constructing the process that sometimes yields a more convenient notation.



The random table model can be slightly modified to what is known as the stack of rewards model. The only difference is that all random variables appearing in the random table model are used. If the a th arm is played for the n th time then the reward variable $X_n^{(a)}$ is used instead of the reward variable corresponding to the time at which the a th arm is played for the n th time. In formulas, one sets $X_t = X_{T_a(t)}^{(a)}$ instead of $X_t = X_t^{(a)}$, where $T_a(t) = \sum_{s \leq t} \mathbf{1}_{X_s = a}$ is the number of times the a th arm was played before time t .

In mathematics it is always good to have concrete examples in mind. Here are two examples to keep in mind. These examples will always be used in the practical exercises.

Example 1.2.5. • Gaussian bandit: all arms are Gaussians $\mathcal{N}(\mu_i, \sigma_i^2)$, for instance

$$(\mu, \sigma) \in \{(0, 1), (1.1, 1), (0.9, 1), (2, 1), (-5, 1), (-3, 2)\}.$$

- Bernoulli bandit: all arms take value 1 with probability p_i and value 0 with probability $1 - p_i$, for instance

$$p \in \{0.9, 0.85, 0.8, 0.5, 0.1, 0.88, 0.7, 0.3, 0.88, 0.75\}.$$

Now that bandit models are defined the next step is to discuss the questions of interest. In fact, with the targets of this lecture course in mind this is not completely clear as the goals

of stochastic bandit theory and reinforcement learning are different. The reason is that the research communities of statistics and AI traditionally have different examples in mind. While the stochastic bandit community originates from statistical question of optimal experimental design in medicine the AI community is more focused on artificial decision making (of computers). While both aim at developing and analysing algorithms that find the optimal arm, the different goals yield in different optimization goal. As an guiding example we go back to the two examples of medical treatment and advertisement. While in medical treatment every single round of learning refers to the medical treatment of an individual (which has the highest priority and the number of rounds is clearly limited say by $n = 200$) in online advertisement it might be much less problematic to play many rounds (say $n = 100k$) in the training procedure and to waste possible income. Here are two typical goals that we will formalise in due course:

- (A) For fixed $n \in \mathbb{N}$ find an algorithm that produces a learning strategy $(\pi_t)_{t=1, \dots, n}$ such that the expected reward $\mathbb{E}[\sum_{k=1}^n X_k^\pi]$ is maximised.
- (B) For fixed $n \in \mathbb{N}$ find an algorithm that produces a learning strategy $(\pi_t)_{t=1, \dots, n}$ such that the probability of choosing wrong arms is minimized.

Checking papers and text books you will realise that the first goal is typical in the stochastic bandit community (statistics), the second more typical in AI. The aim of these lecture notes is to introduce reinforcement learning, so why bother with questions from stochastic bandit theory? The reason is that a much better mathematical understanding is available from the stochastic bandit community, optimal choices of parameters have been derived theoretically for many bandit algorithms. In contrast, the AI community tends to deal with more realistic (more complicated) models in which choices of parameters are found by comparing simulations. It is the goal of these lecture to find the right compromise. To understand well enough the important mechanisms in simple models to improve the educated guesses in realistic models that might be untractable for a rigorous mathematical analysis.

Let us first discuss the classical stochastic bandit approach (A). We already defined an optimal arm, an arm with maximal action value Q_a . Of course there might be several best arms, then a_* is chosen as any of them. Since the index set is not ordered there is no preference in which best arm to denote a_* . The goal is to maximise over all learning strategies the expectation $\mathbb{E}[\sum_{t=1}^n X_t] = \sum_{t=1}^n \mathbb{E}[X_t]$ for a fixed time-horizon n . There is a simple upper bound for the reward until time n , which is Q_* . Hence, if all Q_a would be known in advance then the stochastic bandit optimization problem would be trivial, just choose the best arm a_* in all rounds. Hence, we always assume the expectations are unknown but the outcomes of the arms (random variables) can be played (simulated). Since the expected rewards are upper bounded by nQ_* it is common practice not to maximise the expected reward but instead the difference to the best case as this gives an objective criterion that does not depend on the bandit model itself.



Definition 1.2.6. Suppose ν is a bandit model and $(\pi_t)_{t=1, \dots, n}$ a learning strategy. Then the (cumulated) regret is defined by

$$R_n(\pi) := nQ_* - \mathbb{E}\left[\sum_{t=1}^n X_t^\pi\right].$$

The stochastic bandit problem consists in finding learning strategies that minimise the regret. Algorithms are only allowed to use samples of the reward distribution. If π is clear from the context then we will shorten to R_n .

The regret is called regret because (in expectation) this is how much is lost by not playing the best arm from the beginning. To get acquainted with the definition please check the following facts:



- (i) Consider a two-armed bandit with $Q_1 = 1$ and $Q_2 = -1$ and a learning strategy π given by

$$\pi_t = \begin{cases} \delta_1 & : t \text{ even,} \\ \delta_2 & : t \text{ odd.} \end{cases}$$

Calculate the regret $R_n(\pi)$.

- (ii) Define a stochastic bandit and a learning strategy such that the regret is 5 for all $n \geq 5$.
- (iii) Show for all learning strategies π that $R_n(\pi) \geq 0$ and $\limsup_{n \rightarrow \infty} \frac{R_n(\pi)}{n} < \infty$.
- (iv) Let $R_n(\pi) = 0$. Prove that π only chooses best arms. If there is only one best arm, then π is deterministic, i.e. $\Pi_t(a^*) = 1$ for all t .

What is considered to be a good learning strategy? Linear regret (as a function in n) is always possible by just uniformly choosing arms as this (stupid) learning strategy yields

$$R_n(\pi) = nQ_* - n\mathbb{E}[X_1^\pi] = n \left(\underbrace{Q_* - \sum_{k=1}^K \frac{1}{K} Q_a}_{\geq 0} \right).$$

Thus, a linearly increasing regret can always be achieved when learning nothing. As a consequence only learning strategies with sublinearly increasing regret are considered reasonable.



In stochastic bandit theory any algorithm that produces a learning strategies with linearly growing regret is considered useless. The aim is to significantly improve on linear regret.

There are different kind of bounds that one can aim for. First of all, one can aim for upper bounds and lower bounds for regret. In these lectures notes we mainly focus on upper bounds. Nonetheless, there are celebrated lower bounds due to Lai and Robbins that are not too hard to prove, see Section 1.4. These theoretical lower bounds are important as they tell us if there is any hope to search for better algorithms as the one we discuss (the actual answer is that one cannot do much better than the UCB algorithm presented below). Furthermore, the kind of estimates differ:

- bounds that depend on the bandit model ν are called model-based, such bounds typically involve the differences between the action values Q_a ,
- bounds that only depend on n are called model independent, they are typically proved for entire classes of bandit models for which certain moment conditions are assumed.

For the committ-then-explore and UCB algorithms below model-based upper bounds will be derived from which also model independent upper bounds can be deduced. We will see that it is not too hard to obtain algorithms that achieve model-based upper bounds that are logarithmic in n regret bounds that also involve differences of action values Q_a that can make the estimates as terrible as possible by choosing models where the best and second best arms are hard to distinguish. In fact, the model independent Lai-Robbins lower bounds shows that the best algorithm on all subgaussian bandits can only have a regret as good as $C\sqrt{Kn}$ for some constant C .



From the practical perspective one might wonder why to deal with regret upper bounds. If the bounds are reasonably good, then they can be used in order to tune appearing parameters to optimize the algorithms with guaranteed performance bounds. As an example, we will use the bounds for the explore-then-commit algorithm to tune the



exploration lengths. Even though the resulting parameters might involve unrealistic quantities the understanding can still help us to understand how to work with the algorithms.

In principle, algorithms could depend on a time-horizon n if n is specified in advance. In that case asymptotic regret bounds are non-sense and we aim for finite n bounds only. Sometimes algorithms also depend on the unknown expectations Q_a through the so-called reward gaps.



Definition 1.2.7. The differences $\Delta_a := Q_* - Q_a$ are called reward gaps.

Of course it is not desirable to have algorithms that depend on the reward gaps as a priori knowledge of the expectations Q_a would turn the stochastic bandit problem into a trivial one (just chose the arm with the largest expectation). Nonetheless, the analysis of such algorithms can be of theoretial interest to better understand the mechanism of learning strategies. Also we will see some examples below, the explore-then-commit algorithm depends on the time-horizon n , so does the simple UCB algorithm, whereas the ε_n -algorithm is independent of n but mildly depends on the Δ_a through a constant. Bounds on the regret typically depend on n and the expectations μ_a often in the form $\Delta_a := Q_* - Q_a$.

In order to analyse the regret of a given algorithm in many instances (such as explore-then-commit and UCB) one always uses the regret decomposition lemma:



Lemma 1.2.8. (Regret decomposition lemma)

Defining $T_a(n) := \sum_{t=1}^n \mathbf{1}_{A_t=a}$ the following decomposition holds:

$$R_n(\pi) = \sum_{a \in \mathcal{A}} \Delta_a \mathbb{E}[T_a(n)].$$

Proof. If you know a bit of probability theory it is clear what we do, we insert a clever 1 that distinguishes the appearing events:

$$R_n(\pi) = nQ_* - \mathbb{E}\left[\sum_{t \leq n} X_t\right] = \sum_{t \leq n} \mathbb{E}[Q_* - X_t] = \sum_{t \leq n} \sum_{a \in \mathcal{A}} \mathbb{E}[(Q_* - X_t) \mathbf{1}_{A_t=a}].$$

To compute the right hand side note that

$$\begin{aligned} \mathbb{E}[(Q_* - X_t) \mathbf{1}_{A_t=a} \mid A_1, X_1, \dots, A_t] &= \mathbf{1}_{A_t=a} \mathbb{E}[Q_* - X_t \mid A_1, X_1, \dots, A_t] \\ &= \mathbf{1}_{A_t=a} (Q_* - Q_{A_t}) \\ &= \mathbf{1}_{A_t=a} (Q_* - Q_a) \\ &= \mathbf{1}_{A_t=a} \Delta_a. \end{aligned}$$

Here we used the general fact $\mathbb{E}[X|Y] = \int x \mathbb{P}(X \in dx|Y)$ and that $\mathbb{P}(X_t \in \cdot | A_1, X_1, \dots, A_t) \sim P_{A_t} = \sum_a P_a \mathbf{1}_{A_t=a}$. Using the tower property a combination of both computations yields

$$R_n = \sum_{t \leq n} \sum_{a \in \mathcal{A}} \mathbb{E}[\mathbb{E}[(Q_* - X_t) \mathbf{1}_{A_t=a} \mid A_1, X_1, \dots, A_t]] = \sum_{a \in \mathcal{A}} \Delta_a \mathbb{E}\left[\underbrace{\sum_{t \leq n} \mathbf{1}_{A_t=a}}_{T_a(n)}\right].$$

□

The statistical regret analysis for many bandit algorithm follows the same approach, using the regret decomposition lemma to reduce regret estimates to so-called concentration inequalities. Under suitable assumptions on the distributions of the arms one can plug-in different concentration inequalities from probability theory to derive regret bounds.

We continue our discussion with the second perspective on how to analyse bandit algorithms. Approach (C) is more refined than (A), as an analogie to function (C) is similar to studying the asymptotic behavior of a function through that of its derivative. To get this idea clear let us introduce a new notation:



Definition 1.2.9. Suppose ν is a bandit model and π a learning strategy. Then the probability the learner choses a suboptimal arm in round t , i.e.

$$\tau_t(\pi) := \mathbb{P}(Q_{A_t} \neq Q_*)$$

is called the failure probability in round t .

It is clearly desirable to have $\tau_t(\pi)$ decay to zero as fast as possible. Note that the failure probability is typically not the target for stochastic bandits but connects well to ideas in reinforcement learning. Since $\mathbb{E}[T_n(a)] = \sum_{t=1}^n \mathbb{E}[\mathbf{1}_{A_t=a}] = \sum_{t=1}^n \mathbb{P}(A_t = a)$ the regret decomposition lemma can be reformulated as follows:



Lemma 1.2.10.

$$R_n(\pi) = \sum_{t=1}^n \sum_{a \in \mathcal{A}} \Delta_a \mathbb{P}(A_t = a),$$

and, in particular,

$$R_n(\pi) \leq \max_{a \in \mathcal{A}} \Delta_a \sum_{t=1}^n \tau_t(\pi) \quad \text{and} \quad R_n(\pi) \geq \min_{a \neq a^*} \Delta_a \sum_{t=1}^n \tau_t(\pi).$$

As a consequence we see that the study of regret and failure probability is ultimately connected. If we interpret the sum as an integral, then understanding the failure probability instead of the regret is just as studying the asymptotics of a function by its derivate (which is typically harder). Here are two observations that we will use later for the examples:



- (i) If the failure probabilities do not decay to zero (no learning of the optimal arm), then the regret grows linearly.
- (ii) If the failure probabilities behave (make this precise) like $\frac{1}{n}$, then the regret behaves like $\sum_{a \in \mathcal{A}} \Delta_a \log(n)$ with constants that depend on the concrete bandit model. Hint: Recall from basic analysis that $\int_1^t \frac{1}{x} dx = \log(t)$ and how to relate sums and integrals for monotone integrands.

1.3 Algorithms: the exploration-exploitation trade-off

We will now go into a few of the most popular algorithms. There is not much choice on how to design an algorithm. Essentially, all that can be done is to learn about arms that one is not too sure about (called exploration) and play arms that one expects to be good (called exploitation). We will not only present algorithms but also discuss theoretical regret bounds. Even though those won't be directly useful for our later understanding of reinforcement learning there is one important learning: theoretical results allow to understand how to chose optimally the parameters involved, in contrast to learn by experimenting which is always restricted to particular examples. In spirit we follow the exposition of Chapter 2 in Sutton and Barto, but we try to mix in more mathematics to push the understanding further than just simulations.

1.3.1 Basic explore-then-committ algorithm (ETC)

Without any prior understanding of stochastic bandits here is a simple algorithm that everyone would come up with himself/herself. Recalling the law of large numbers $\frac{1}{n} \sum_{t=1}^n Y_t \rightarrow \mathbb{E}[Y_1]$ we first produce estimates \hat{Q}_a for the expectations Q_a and then play the arm with the largest estimated expectation. How do we use the law of large numbers? By just playing every arm m times and for the remaining $n - mK$ rounds play the best estimated arm. That's it, that is the commit-then-exploit algorithm. Before turning the basic idea into an algorithm let us fix a notation that will occur again and again.



Definition 1.3.1. If (A_t, X_t) is a stochastic process on $\mathcal{A} \times \mathbb{R}$, $a \in \mathcal{A}$, and $T_a(t) := \sum_{k=1}^t \mathbf{1}_{A_k=a}$, then

$$\hat{Q}_a(t) := \frac{1}{T_a(t)} \sum_{k=1}^t X_k \mathbf{1}_{A_k=a}$$

is called an estimated action value.

The Q -estimator for a stochastic bandit model estimates Q_a by averaging over all past rewards seen when playing arm a . The ETC algorithm can be written in different ways; either try all arms m -times in a row or to alternate between the arms. For the pseudocode of Algorithm 1 we chose the latter. The learning strategy π_t is an index strategy and can be written as

$$\pi_t(a; a_1, x_1, \dots, x_t) = \begin{cases} 1 & : a = a_{t \bmod K+1} \text{ and } t \leq mK \\ 1 & : \operatorname{argmax}_a \hat{Q}_a(mK) \text{ and } t > mK \\ 0 & : \text{otherwise} \end{cases}$$

for arms denoted by $1, \dots, K$. As a first example on how to estimate the regret of bandit algorithms

Data: m, n , bandit model ν

Result: actions A_1, \dots, A_n and rewards X_1, \dots, X_n

while $t \leq mK$ **do**

$$A_t := \begin{cases} a_{t \bmod K+1} & : t \leq mK; \\ \operatorname{argmax}_a \hat{Q}_a(mK) & : t > mK; \end{cases}$$

Obtain reward X_t by playing arm A_t ;

Compute $\hat{Q}_{A_t}(t)$;

end

Algorithm 1: Basic m -rounds ETC algorithm

we prove the following upper bound. The notion 1-subgaussian will be introduced in the course of the proof, keep in mind Bernoulli bandits or Gaussian bandits with variance at most σ^2 .



Theorem 1.3.2. (Regret bound for simple explore-then-commit)

Suppose ν is a σ -subgaussian bandit model, i.e. all P_a are σ -subgaussian (see below), with K arms and $Km \leq n$ for some $n \in \mathbb{N}$, then

$$R_n(\pi) \leq \underbrace{m \sum_{a \in \mathcal{A}} \Delta_a}_{\text{exploration}} + \underbrace{(n - mK) \sum_{a \in \mathcal{A}} \Delta_a \exp\left(-\frac{m\Delta_a^2}{4\sigma^2}\right)}_{\text{exploitation}}.$$

Since the regret bound looks a bit frightening on first view let us discuss the ingredients first. What do we believe a bound should depend on? Certainly on the total number of rounds n , the number of exploration rounds m , and the number of arms. Probably also on the distributions of

the arms. Why is this? If all arms have the same law, i.e. $\Delta_a = 0$ for all arms, then the regret would be 0 as we always play an optimal arm. If the best arm is much better than other arms, then the exploration phase forces a larger regret as we decided to also play the worst arm m times. Looking into the regret bound, the summand $m \sum_{a \in \mathcal{A}} \Delta_a$ is clearly the regret from the exploration phase.



Summands of the form $\sum_{a \in \mathcal{A}} \Delta_a$ must appear in all reasonable regret bounds as every reasonable algorithm must try every arm at least once.

The interesting question is the regret obtained during the exploitation phase if a suboptimal arm is played. It seems clear that the best arm is the most likely to be exploited as $\hat{Q}_a(n) \approx Q_a(n)$ for large n by the law of large numbers. The regret thus comes from deviations of this large n principle, if the rewards of an arm exceed what they would yield on average. Since the simple explore-then-commit algorithm involves a lot of independence probabilities overestimation of arms can easily be estimated by inequalities which are known as concentration inequalities in probability theory.

Proof, decomposing exploration and exploitation: Without loss of generality (the order or the arms does not matter) we may assume that $Q_* = Q_{a_1}$. Using the regret decomposition and the algorithm yields the following decomposition into exploitation and exploration:

$$R_n(\pi) = \sum_{a \in \mathcal{A}} \Delta_a \mathbb{E}[T_a(n)] = m \sum_{a \in \mathcal{A}} \Delta_a + \sum_{a \in \mathcal{A}} \Delta_a (n - mK) \mathbb{P}(\hat{Q}_a(mK) \geq \max_{b \in \mathcal{A}} \hat{Q}_b(mK)).$$

Where does this come from? Each arm is explored m times and, if the arm was the best, then another $n - mK$ times. Hence, $T_a(n) = m + (n - mK) \mathbf{1}_{\{a \text{ was best up to } mK\}}$. Computing the expectation the probability appears due to the construction of the learning strategy π . The probability can be estimated from above by replacing the maximal arm by some other arm (we chose the first). This leads to

$$\begin{aligned} R_n(\pi) &\leq m \sum_{a \in \mathcal{A}} \Delta_a + (n - mK) \sum_{a \in \mathcal{A}} \Delta_a \mathbb{P}(\hat{Q}_a(mK) \geq \hat{Q}_{a_1}(mK)) \\ &= m \sum_{a \in \mathcal{A}} \Delta_a + (n - mK) \sum_{a \in \mathcal{A}} \Delta_a \mathbb{P}((\hat{Q}_a(mK) - \hat{Q}_{a_1}(mK)) - (Q_a - Q_{a_1}) \geq \Delta_a). \end{aligned}$$

The appearing probability has the particularly nice form

$$\mathbb{P}(Z_a - \mathbb{E}[Z_a] \geq \Delta_a), \quad \text{with} \quad Z_a = \frac{1}{m} \sum_{j=1}^m (X_j^{(a)} - X_j^{(1)}),$$

where $X_1^{(a)}, \dots, X_m^{(a)}$ are distributed according to arm a and all of them are independent. If we can estimate the probabilities by $\exp(-m\Delta_a^2/(4\sigma^2))$ the proof is complete. In order to do so we first need an excursion to probability theory. \square

Let's have a short excursion into the topic of concentration. A concentration inequality is a bound on the deviation of a random variable from its mean, either in a two- or one-sided fashion:

$$c(a) \leq \mathbb{P}(|X - \mathbb{E}[X]| > a) \leq C(a) \quad \text{or} \quad c(a) \leq \mathbb{P}(X - \mathbb{E}[X] > a) \leq C(a).$$

The faster the function $C(a)$ decreases in a the more the random variable is concentrated (takes values close to its expectation with larger probability, the randomness is less significant). The idea is that a random variable is stronger concentrated (has less mass away from its expectation) if larger moments are finite. Where this idea comes from is easily seen from the expectation formula

$$\mathbb{E}[g(X)] = \begin{cases} \int_{\mathbb{R}} g(y) f_X(y) dy & : X \text{ has the probability density function } f_X \\ \sum_{k=1}^N g(a_k) p_k & : X \text{ takes the values } a_k \text{ with probabilities } p_k \end{cases},$$

so that finite expectation for a (strongly) increasing function g such as an exponential function forces the density (or probability weights) to decrease (strongly) at infinity and thus to be more concentrated. Here is an example: Markov's inequality states that

$$\mathbb{P}(|X - \mathbb{E}[X]| > a) \leq \frac{\mathbb{V}[X]}{a^2}$$

for every random variable for which $\mathbb{E}[X^2] < \infty$. Markov's (concentration) inequality is useful as it holds for many random variables but the concentration inequality is very bad, the upper bound only decreases like a polynomial as we move away from the mean. A more useful inequality holds for so-called subgaussian random variables.



Definition 1.3.3. A random variable X on a probability space $(\Omega, \mathcal{A}, \mathbb{P})$ is called σ -subgaussian for $\sigma > 0$, if $\mathcal{M}_{X - \mathbb{E}[X]}(\lambda) = \mathbb{E}[e^{\lambda(X - \mathbb{E}[X])}] \leq e^{\frac{\lambda^2 \sigma^2}{2}}$ for all $\lambda \in \mathbb{R}$.

The wording subgaussian of course comes from the fact that $\mathbb{E}[e^{\lambda X}] = e^{\lambda^2 \sigma^2 / 2}$ if $X \sim \mathcal{N}(0, \sigma^2)$. Here is a little exercise to get acquainted to the definition:



- Show that every σ -subgaussian random variable satisfies $\mathbb{V}[X] \leq \sigma^2$.
- If X is σ -subgaussian, then cX is $|c|\sigma$ -subgaussian.
- Show that $X_1 + X_2$ is $\sqrt{\sigma_1^2 + \sigma_2^2}$ -subgaussian if X_1 and X_2 are independent σ_1 -subgaussian and σ_2 -subgaussian random variables.
- Show that a Bernoulli-variable is $\frac{1}{4}$ -subgaussian by explicitly computing $\log \mathcal{M}_{X-p}(\lambda)$, checking for which p the formula for $\log \mathcal{M}_{X-p}(\lambda)$ is maximal and then estimating the remaining function by $\frac{\lambda^2}{8}$.
- Every centered bounded random variable, say bounded below by a and above by b is $\frac{(b-a)}{2}$ -subgaussian (this is called Hoeffding's lemma).

It is important to note that every σ -subgaussian random variable is also σ' -subgaussian for every $\sigma' > \sigma$ but this is not interesting as we will use σ to bound the variability (σ somehow measures the variance) as good as possible. This becomes clear in the next proposition, using σ larger than necessary only weakens the concentration inequality:



Proposition 1.3.4. If X is σ -subgaussian, then

$$\mathbb{P}(X \geq a) \leq e^{-\frac{a^2}{2\sigma^2}} \quad \text{and} \quad \mathbb{P}(|X| \geq a) \leq 2e^{-\frac{a^2}{2\sigma^2}}$$

for all $a > 0$.

Proof. The proof is based on a trick called Cramér-Chernoff method. The trick uses the Markov inequality for a parametrized family of functions and then optimising over the parameter to find the best estimate:

$$\mathbb{P}(X \geq a) = \mathbb{P}(e^{\lambda X} \geq e^{\lambda a}) \leq \frac{\mathbb{E}[e^{\lambda X}]}{e^{\lambda a}} \leq \frac{e^{\frac{\lambda^2 \sigma^2}{2}}}{e^{\lambda a}} = e^{\frac{\lambda^2 \sigma^2}{2} - \lambda a}.$$

Minimizing the right hand side for λ (differentiation!) shows that $\lambda = \frac{a}{\sigma^2}$ yields the smallest bound and this is the first claim of the proposition. Since the same holds for $-X$ we obtain the second claim by writing $\mathbb{P}(|X| \geq a) = \mathbb{P}(X \geq a \text{ or } X \leq -a) \leq \mathbb{P}(X \geq a) + \mathbb{P}(-X \geq a)$. \square

As an application of the above we get a first simple concentration inequality for sums of random variables:


Corollary 1.3.5. (Hoeffding's inequality)

Suppose X_1, \dots, X_n are iid random variables on a probability space $(\Omega, \mathcal{A}, \mathbb{P})$ with expectation $\mu = \mathbb{E}[X_1]$ such that X_1 is σ -subgaussian. Then

$$\mathbb{P}\left(\frac{1}{n} \sum_{k=1}^n X_k - \mu \geq a\right) \leq e^{-\frac{na^2}{2\sigma^2}} \quad \text{and} \quad \mathbb{P}\left(\left|\frac{1}{n} \sum_{k=1}^n X_k - \mu\right| \geq a\right) \leq 2e^{-\frac{na^2}{2\sigma^2}}, \quad \forall a > 0.$$

Proof. This follows from the exercise and the proposition above because $\frac{1}{n} \sum_{k=1}^n X_k - \mu = \frac{1}{n} \sum_{k=1}^n (X_k - \mathbb{E}[X_k])$ is a centered $\frac{\sigma}{\sqrt{n}}$ -subgaussian random variable. \square

We can now combine the exploration exploration decomposition with the concentration inequality to derive the upper bound of the regret in the explore-then-commit algorithm:

Completing the proof of Theorem 1.3.2. Since we assumed that the exploitation phase consists of independent runs of the same arm we are exactly in the situation of Corollary 1.3.5. Hence, with the notation from the first part of the proof we get the concentration bound

$$\mathbb{P}(Z_a - \mathbb{E}[Z_a] \geq \Delta_a) \leq \exp\left(-\frac{\Delta_a^2}{2\frac{2\sigma^2}{m}}\right) = \exp\left(-\frac{m\Delta_a^2}{4\sigma^2}\right).$$

Plugging-in yields the upper bound from the theorem. \square

Also without the estimates the following logic is clear: If m is large (a lot of exploration) then the exploration regret is large (this is the first summand) and the exploitation regret is small (second summand). In the exercises you will explore numerically how to properly chose m in different examples. Let's explore the regret upper bound to find a reasonable choice of m . Of course, this approach is only reasonable if the upper bound is reasonably good. Indeed, the first summand is an equality, the second summand only uses Hoeffding's inequality which is a reasonably good estimate. To show how to find a good exploration length let us consider the simple case $K = 2$ of two arms (again, the first arm is assumed to be the optimal one). Since $\Delta_{a_1} = 0$, we abbreviate $\Delta = \Delta_{a_2}$, the estimate simplifies to

$$R_n(\pi) \leq m\Delta + (n - 2m)\Delta \exp\left(-\frac{m\Delta^2}{4\sigma^2}\right) \leq \Delta\left(m + n \exp\left(-\frac{m\Delta^2}{4\sigma^2}\right)\right).$$

Pretending that m is a continuous parameter the righthand side can be minimized (in m) by differentiation. Do this to solve the following exercise:



The regret upper bound is minimized by

$$m = \max\left\{1, \left\lceil \frac{4\sigma^2}{\Delta^2} \log\left(\frac{n\Delta^2}{4\sigma^2}\right) \right\rceil\right\}. \quad (1.1)$$

Thus, in terms of regret optimisation, our best guess is the explore-then-committ algorithm with this particular m . For this choice of m the regret is upper bounded by

$$R_n(\pi) \leq \min\left\{n\Delta, \Delta + \frac{4\sigma^2}{\Delta} \left(1 + \max\left\{0, \log\left(\frac{n\Delta^2}{4\sigma^2}\right)\right\}\right)\right\} \quad (1.2)$$

which for $n \geq \frac{4\sigma^2}{\Delta^2}$ gives a model-dependent logarithmic regret bound $R_n(\pi) \leq C_\Delta + \frac{4\sigma^2}{\Delta} \log(n)$.

In the programming exercises you will be asked to compare this theoretical m with the best m that can be „seen“ from simulations in a special example. No doubt, tuning parameters by simulating examples is not very appealing as the parameter might be useless for other examples.



Do you think the explore-then-committ strategy with m from (1.1) is reasonable? No, it's cheating. The algorithm relies on n , σ , and Δ through the choice of m .

- The number of rounds n might be fixed in advance for some examples it might be not for other examples. For infinite time-horizon there is a trick called the "doubling-trick" that allows to combine fixed-time algorithms into an infinite time-horizon algorithm so we do not accept dependence on n .
- The variance parameter σ might be known in some situations (for instance Gaussian rewards with fixed variance but unknown mean) but will typically be unknown as well. There are algorithms that estimate σ on the run.
- The dependence on $\Delta = Q_{a_1} - Q_{a_2}$ is much more severe as the action values are never known in advance (otherwise we could just choose a best arm to obtain zero regret). Hence, the only non-trivial situation is that of unknown action values but known difference $\Delta = Q_{a_1} - Q_{a_2}$, but this is extremely special.

It will turn out below in the Lai-Robbins Theorem 1.4.8 that the regret upper bound from (1.2) is close to optimal for large n because of the learning strategy independent lower bound $\liminf_{n \rightarrow \infty} \frac{R_n(\pi)}{\log(n)} \geq \frac{\sigma^2}{\Delta}$ at least for Gaussian bandit models. In the next section we will show how to construct a similarly good algorithm without cheating in the choice of the algorithm parameters.

1.3.2 From greedy to UCB

A greedy learning strategy is a strategy that always chooses the option the algorithm currently believes to be the best. For bandits this is the arm with the highest believed reward, typically measured in terms of the empirical mean, that is the mean of the already observed rewards $\hat{Q}_a(t-1) = \frac{1}{T_a(t-1)} \sum_{k=1}^{t-1} X_k \mathbf{1}_{A_k=a}$. As similar concepts will reappear in reinforcement learning we will spend some time on this topic. It will turn out that pure greedy algorithms do not work at all, but simple modifications that force additional exploration work very well.

Purely greedy bandit algorithm

The pure greedy algorithm turns out to be complete non-sense. Nonetheless, it is a good start into the discussion to get acquainted with the notation and simple modifications give useful algorithms. The plain idea is as follows: Take the past observations of each arm to define estimates $\hat{Q}_a(t-1)$ of the action values Q_a at time t and then choose the maximal estimated arm, i.e. $A_t := \operatorname{argmax}_a \hat{Q}_a(t-1)$. As always, since there is no preference order for the arms, if several estimated action values are equal we randomly choose one of them. In the algorithm we only use

Data: bandit model ν , vector \hat{Q} , n
Result: actions A_1, \dots, A_n and rewards X_1, \dots, X_n
 Initialise $T_a = 0$ for all a ;
while $t \leq n$ **do**
 Set $A_t = \operatorname{argmax}_a \hat{Q}_a$;
 Obtain reward X_t by playing arm A_t ;
 Set $T_{A_t} = T_{A_t} + 1$;
 Set $\hat{Q}_{A_t} = \hat{Q}_{A_t} + \frac{1}{T_{A_t}}(X_t - \hat{Q}_{A_t})$;
end

Algorithm 2: Purely greedy bandit algorithm

one vector \hat{Q} to store the estimates $\hat{Q}(t)$ as we only need the current estimated action values. The computation of \hat{Q} looks a bit strange and relies on a simple memory trick that allows to

compute successive averages without storing all numbers. Suppose R_1, \dots are real numbers and all averages $Q_n = \frac{1}{n} \sum_{k=1}^n R_k$ should be computed without storing the rewards R_n forever. To compute Q_n it is actually only needed to know Q_{n-1} and R_n :

$$\begin{aligned}
 Q_n &= \frac{1}{n} \sum_{k=1}^n R_k \\
 &= \frac{1}{n} \left(R_n + \sum_{k=1}^{n-1} R_k \right) \\
 &= \frac{1}{n} \left(R_n + (n-1) \frac{1}{n-1} \sum_{k=1}^{n-1} R_k \right) \\
 &= \frac{1}{n} \left(R_n + (n-1) Q_{n-1} \right) \\
 &= Q_{n-1} + \frac{1}{n} \left(R_n - Q_{n-1} \right)
 \end{aligned} \tag{1.3}$$

Under the initialisation $\hat{Q} \equiv 0$ the \hat{Q} are nothing but the sample means and the greedy algorithm plays greedily the arms with the largest empirical mean. We can also think differently about the algorithm. The actor suggests a vector with some initial estimates \hat{Q}_a of the action values. If nothing is known the actor will always choose $\hat{Q} \equiv 0$. Then the actor plays greedily by always playing the arm that is believed to be best. After playing the Q -values are updated by adding $\frac{1}{T_{A_t}} (X_t - \hat{Q}_{A_t})$. They are increased (resp. decreased) if the reward was higher (resp. lower) than what was believed before.



An algorithm is called a tabular algorithm if there is a table of real numbers (here a vector) that is constantly updated and used to make the decisions. We will get to know a similar approach as Q -learning in Chapter ??.

The pure greedy algorithm depends extremely on the initialisation of the vector \hat{Q} and the distribution of the arms. Suppose $\hat{Q} \equiv 0$ and suppose one arm only returns positive values and this arm is chosen at the beginning. Then the estimated action value is increased to a positive number and no other arm will be played in the future. Similarly, if an arm takes positive values with high probability then future exploration is very unlikely to occur. Also for a Gaussian bandit model a similar phenomenon arises. Suppose (at least) one arm has positive expectation and suppose $3\sigma < \mu$. If initially that arm is played then with probability at least 0.997 (three- σ -rule) the result will be positive so that the arm will be played again. Continuing like that it will take a long time until eventually a different arm will be explored. Since that phenomenal will occur again we give it a name:



Definition 1.3.6. The phenomenon that a bandit (later: reinforcement learning) algorithm focuses too early on a suboptimal decision is called committal behavior.

From a practical point of view there are workarounds for instance using different initialisations \hat{Q} . As an example one might start with large \hat{Q} . In that case the algorithms would start with many rounds of exploration before starting the greedy update. A problem remains: How large would be good without knowing details of the bandit model? If \hat{Q} is chosen too large there would be too much exploration before playing greedy the arms with the largest estimated action values. If \hat{Q} would not be large enough then there would be too little exploration. A second workaround would be trying to center the distributions of all arms by subtracting the same constant from all arms to reduce the committal behavior effect described above. Again, it is unclear what constant should be subtracted without a priori knowledge on the action values. We will get back to this idea when we discuss the policy gradient method where this idea will return as the base-line trick.

ε -greedy bandit algorithms

The simplest variant to make the pure greedy algorithm more reasonable is to force completely random additional exploration. The idea originates from the reinforcement learning community as from the point of view of minimizing regret the algorithm is completely useless. To get a

Data: bandit model ν , exploration rate $\varepsilon \in (0, 1)$, vector \hat{Q} , n
Result: actions A_1, \dots, A_n and rewards X_1, \dots, X_n
 Initialise $T_a = 0$ for all a ;
while $t \leq n$ **do**
 Sample $U \sim \mathcal{U}([0, 1])$;
 if $U < \varepsilon$ **then**
 [exploration part];
 Uniformly chose an arm A_t ;
 Obtain reward X_t by playing arm A_t ;
 Set $T_{A_t} = T_{A_t} + 1$;
 Set $\hat{Q}_{A_t} = \hat{Q}_{A_t} + \frac{1}{T_{A_t}}(X_t - \hat{Q}_{A_t})$;
 end
 if $U \geq \varepsilon$ **then**
 [greedy part];
 Set $A_t = \operatorname{argmax}_a \hat{Q}_a$;
 Obtain reward X_t by playing arm A_t ;
 Set $T_{A_t} = T_{A_t} + 1$;
 Set $\hat{Q}_{A_t} = \hat{Q}_{A_t} + \frac{1}{T_{A_t}}(X_t - \hat{Q}_{A_t})$;
 end
end

Algorithm 3: ε -greedy bandit algorithm

Lecture 3

feeling for the algorithm you can run different simulations in the exercises. Simulations of this kind can be very useful to understand basic mechanisms, but not much more. Of course, for this particular example we could repeat the simulations again and again to fine-tune the choice of ε . But this does not result in further understanding of the basic mechanisms for an unknown problem. The following exercise shows that there is no use of the simple greedy learning strategy, it has linear regret even if each arm is explored once.



Let π the learning strategy that first explores each arm once and then continues according to ε -greedy for some $\varepsilon \in (0, 1)$ fixed. Show that the regret grows linearly:

$$\lim_{n \rightarrow \infty} \frac{R_n(\pi)}{n} = \frac{\varepsilon}{K} \sum_{a \in \mathcal{A}} \Delta_a.$$

Hint: Lower bound $\mathbb{E}[T_a(n)]$ by the random exploration and upper bound $\mathbb{E}[T_a(n)]$ using arguments inspired by the proof of Theorem 1.3.7 below.

There are many versions of ε -greedy with reasonable regret bounds. We will only touch upon an example (introduced and studied in Auer et al.²) that one could call „explore-then- ε -greedy with decreasing exploitation rate“. The algorithm replaces in the simple ε -greedy algorithm ε by the time-dependent exploration rate $\varepsilon_t = \min\{1, \frac{CK}{d^2 t}\}$. To justify the name note that for the first $\frac{CK}{d^2}$ rounds the exploration rate is 1. Thus, the algorithm first explores randomly and then plays ε -greedy with decreasing exploration rate ε .

²P. Auer, N. Cesa-Bianchi, P. Fischer: „Finite-time Analysis of the Multiarmed Bandit Problem“, Machine Learning, 47:235-256, (2002)


Theorem 1.3.7. (Explore-then- ε -greedy with decreasing ε)

Suppose that all arms take values in $[0, 1]$, $d < \min_{a: Q_a \neq Q_{a^*}} \Delta_a$, and $C > \max\{5d^2, 2\}$. Then the ε -greedy algorithm with decreasing exploration rate $\varepsilon_t = \min\{1, \frac{CK}{d^2 t}\}$ satisfies

$$\limsup_{n \rightarrow \infty} \tau_n(\pi) \cdot n \leq \frac{(K-1)C}{d^2}.$$

Using Lemma 1.2.10 (note that $\Delta_a \leq 1$ if the arms only take values in $[0, 1]$) a consequence of the theorem is logarithmic upper bound

$$\limsup_{n \rightarrow \infty} \frac{R_n(\pi)}{\log(n)} \leq \frac{(K-1)C}{d^2}. \quad (1.4)$$

While logarithmically growing regret is very good (compare the UCB algorithm in Theorem 1.3.8 below) there are two disadvantages. First, the constants are pretty big (the possibly very small reward gap appears squared in the numerator) and will dominate the logarithm for reasonably sized n ($\log(100.000) \approx 11.5$ so that even an additional factor 5 matters quite a lot). Secondly, with the appearing constant d the algorithm assumes prior knowledge on the bandit model reward gaps, the algorithm cheats!



Set $K = 2$ and Δ small and then compare the algorithm with the explore-then-committ algorithm for two arms.

Proof (not part of the course). Auer et al. proved a much more precise estimate. Suppose j is a suboptimal arm and $n > \frac{CK}{d^2}$. Then we prove for all $C > 0$ that

$$\begin{aligned} \mathbb{P}(A_t = j) &\leq \frac{C}{d^2 n} + 2 \left(\frac{C}{d^2} \log \left(\frac{(n-1)d^2 e^{1/2}}{CK} \right) \right) \left(\frac{CK}{(n-1)d^2 e^{1/2}} \right)^{C/(5d^2)} \\ &\quad + \frac{4e}{d^2} \left(\frac{CK}{(n-1)d^2 e^{1/2}} \right)^{C/2}. \end{aligned}$$

Since there are at most $K-1$ suboptimal arms an upper bound for $\tau_t(\pi)$ is obtained by multiplying by $(K-1)$. The choice $C > 5$ ($C > 5d$ actually suffices) implies that the first summand dominates in the limit.



The failure probability for the first $\frac{CK}{d^2}$ rounds (uniform exploration) is easily seen to be $\tau_t(\pi) = 1 - \frac{K^*}{K}$, the probability to chose one of the suboptimal arms.

The proof is mostly a brute force estimation of the probabilities plus Bernstein's inequality a concentration inequality that is a bit more general than Hoeffding's inequality:



Suppose X_1, \dots, X_n are independent (not necessarily identically distributed!) random variables with variances σ_k^2 and expectations $\mu_k = \mathbb{E}[X_k]$. If all X_i are bounded by M , i.e. $|X_i| \leq M$ for all i , and $\sigma^2 := \sum_{k=1}^n \sigma_k^2$, then

$$\mathbb{P}\left(\frac{1}{n} \sum_{k=1}^n X_k - \frac{1}{n} \sum_{k=1}^n \mu_k \geq a\right) \leq e^{-\frac{n^2 a^2}{2(\sigma^2 + \frac{1}{3} n a M)}}.$$

In particular, if X_1, \dots, X_n are iid with variance σ^2 and expectation μ then Bernstein's inequality becomes

$$\mathbb{P}\left(\frac{1}{n} \sum_{k=1}^n X_k - \mu \geq a\right) \leq e^{-\frac{n a^2}{2(\sigma^2 + \frac{1}{3} a M)}}.$$

Let $x_0 := \frac{1}{2K} \sum_{s=1}^t \varepsilon_s$ with $\varepsilon_n = \frac{CK}{d^2 n}$ the exploration rate from the algorithm. Suppose j is a suboptimal arm, we are going to estimate $\mathbb{P}(A_t = j)$. From the algorithm we obtain

$$\begin{aligned} \mathbb{P}(A_t = j) &= \frac{\varepsilon_n}{K} + \left(1 - \frac{\varepsilon_n}{K}\right) \mathbb{P}(\hat{Q}_j(t-1) \geq \max_a \hat{Q}_a(t-1)) \\ &\leq \frac{\varepsilon_n}{K} + \left(1 - \frac{\varepsilon_n}{K}\right) \mathbb{P}(\hat{Q}_j(t-1) \geq \hat{Q}_*(t-1)) \\ &\leq \frac{\varepsilon_n}{K} + \left(1 - \frac{\varepsilon_n}{K}\right) \left(\mathbb{P}(\hat{Q}_j(t-1) \geq Q_j + \Delta_j/2) + \mathbb{P}(\hat{Q}_*(t-1) < Q_{a_*} - \Delta_j/2) \right), \end{aligned}$$

where Q_* denotes the estimated action value for a fixed optimal arm, say the first. Here we used that, by definition of Δ_j , $Q_* - \frac{\Delta_j}{2} = Q_j + \frac{\Delta_j}{2}$ and the elementary estimate

$$\mathbb{P}(X \geq Y) = \mathbb{P}(X \geq Y, Y \geq a) + \mathbb{P}(X \geq Y, Y < a) \leq \mathbb{P}(X \geq a) + \mathbb{P}(Y < a).$$

From the above we estimate both probabilities separately (the argument is the same). Denote by $T_j^R(t)$ the numbers of random explorations of the arm j before time t . Then

$$\begin{aligned} \mathbb{P}(\hat{Q}_j(t) \geq Q_j + \Delta_j/2) &= \sum_{s=1}^t \mathbb{P}(\hat{Q}_j(t) \geq Q_j + \Delta_j/2, T_j(t) = s) \\ &= \sum_{s=1}^t \mathbb{P}(T_j(t) = s \mid \hat{Q}_j(t) \geq Q_j + \Delta_j/2) \mathbb{P}(\hat{Q}_j(t) \geq Q_j + \Delta_j/2) \\ &\stackrel{1.3.5}{\leq} \sum_{s=1}^t \mathbb{P}(T_j(t) = s \mid \hat{Q}_j(t) \geq Q_j + \Delta_j/2) e^{-\Delta_j^2 s/2}, \end{aligned}$$

using that random variables with values in $[0, 1]$ are $\frac{1}{2}$ -subgaussian. Splitting the sum into the sum up to $\lfloor x_0 \rfloor$ and the rest, using the estimate $\sum_{t=x_0+1}^{\infty} e^{-\kappa t} \leq \frac{1}{\kappa} e^{-\kappa x}$ (think of the integral!) yields the upper bounded

$$\begin{aligned} &\sum_{s=1}^{\lfloor x_0 \rfloor} \mathbb{P}(T_j(t) = s \mid \hat{Q}_j(t) \geq Q_j + \Delta_j/2) + \frac{2}{\Delta_j^2} e^{-\Delta_j^2 \lfloor x_0 \rfloor / 2} \\ &\leq \sum_{s=1}^{\lfloor x_0 \rfloor} \mathbb{P}(T_j^R(t) \leq \lfloor x_0 \rfloor \mid \hat{Q}_j(T) \geq Q_j + \Delta_j/2) + \frac{2}{\Delta_j^2} e^{-\Delta_j^2 \lfloor x_0 \rfloor / 2} \\ &= \lfloor x_0 \rfloor \mathbb{P}(T_j^R(n) \leq \lfloor x_0 \rfloor) + \frac{2}{\Delta_j^2} e^{-\Delta_j^2 \lfloor x_0 \rfloor / 2}, \end{aligned}$$

where the conditioning could be dropped because the exploration is independent. Using Bienaymé and mean, variance of Bernoulli random variables yields

$$\mathbb{E}[T_j^R(t)] = \frac{1}{K} \sum_{s=1}^t \varepsilon_s = 2x_0 \quad \text{and} \quad \mathbb{V}[T_j^R(t)] = \sum_{s=1}^t \frac{\varepsilon_s}{K} \left(1 - \frac{\varepsilon_s}{K}\right) \leq \frac{1}{K} \sum_{s=1}^t \varepsilon_s$$

so that Bernstein's inequality gives $\mathbb{P}(T_j^R(t) \leq \lfloor x_0 \rfloor) \leq e^{-x_0/5}$. In total we derived the bound

$$\mathbb{P}(\hat{Q}_j(t) \geq Q_j + \Delta_j/2) \leq \lfloor x_0 \rfloor e^{-x_0/5} + \frac{2}{\Delta_j^2} e^{-\Delta_j^2 \lfloor x_0 \rfloor / 2}$$

From the probabilistic point the proof is complete but we have not taking into account the choice of ε . It only remains to play around with $x_0 = \frac{1}{2K} \sum_{t=1}^n \varepsilon_t$ to simplify the presentation. Recall the choice of the exploitation rate ε_t and set $n' = \frac{CK}{d^2}$. The exploitation rate is constant 1 up to n' and then equal to $\varepsilon_t = \frac{CK}{d^2 t}$, hence,

$$x_0 = \frac{1}{2K} \sum_{t=1}^{n'} \varepsilon_t + \frac{1}{2K} \sum_{t=n'+1}^n \varepsilon_t \geq \frac{n'}{2K} + \frac{C}{d^2} \log\left(\frac{n}{n'}\right) \geq \frac{C}{d^2} \log\left(\frac{nd^2 e^{1/2}}{CK}\right).$$

Putting everything together yields, $x \mapsto xe^{-x/5}$ is decreasing for $x > 5$, for a suboptimal arm j and $n \geq n'$,

$$\begin{aligned} \mathbb{P}(A_t = j) &\leq \frac{\varepsilon_n}{K} + 2\lfloor x_0 \rfloor e^{-\lfloor x_0 \rfloor / 5} + \frac{4}{\Delta_j^2} e^{-\Delta_j^2 \lfloor x_0 \rfloor / 2} \\ &\leq \frac{C}{d^2 n} + 2 \left(\frac{C}{d^2} \log \left(\frac{(n-1)d^2 e^{1/2}}{CK} \right) \right) \left(\frac{CK}{(n-1)d^2 e^{1/2}} \right)^{C/(5d^2)} \\ &\quad + \frac{4e}{d^2} \left(\frac{CK}{(n-1)d^2 e^{1/2}} \right)^{C/2}. \end{aligned}$$

□

We are not going to discuss further the ε -greedy algorithm. It turns out that that the dependence on the parameters K, Δ_a, C is horribly bad compared to the UCB algorithm that is discussed next. Still, it is important to keep the previous algorithm in mind as ε -greedy algorithms are used frequently in reinforcement learning.

UCB algorithm - optimism in the face of uncertainty

The UCB algorithm (upper confidence bound algorithm) is usually not considered as a version of the greedy algorithm but more a further development of explore-then-commit. UCB follows similar ideas that involve concentration inequalities but with more thoughts. Still, since UCB can also be seen as greedy with an additional exploration bonus we prefer to interpret UCB as an extension of greedy. Here is the main idea, optimism in the face of uncertainty. The principle suggests to be more optimistic, more curious, in situations that are less certain. For instance, trying more new dishes when travelling unknown regions of the world. In the context of bandits the principle states to add an exploration bonus for less good estimated action values. In the context of estimations involving the approximate action values \hat{Q} one should thus add an exploration bonus that decreases with T_a because the uncertainty in the estimate $\hat{Q}_a \approx Q_a$ decreases as T_a increases. In fact, the estimated empirical variance should also play a role but for the simple UCB algorithm this is ignored. Following the principle in the greedy algorithm we should add an exploration bonus to the estimated action values that decreases as T_a increases. The bonus should be something of the kind

$$\hat{Q}_a(t) + \frac{\dots}{T_a(t)}.$$

The choice of the exact exploration bonus is critical. For later purposes one typically uses the

$$\text{UCB}_a(t, \delta) := \begin{cases} \infty & : T_a(t) = 0 \\ \underbrace{\hat{Q}_a(t)}_{\text{greedy}} + \underbrace{\sqrt{\frac{2 \log(1/\delta)}{T_a(t)}}}_{\text{exploration bonus}} & : T_a(t) \neq 0. \end{cases}$$

The exploration bonus is motivated by concentration inequalities, this is the upper bound of confidence intervals for sums of iid random variables, justifying the name UCB. More precisely, for 1-subgaussian random variables with mean Q , plugging-into Hoeffdings inequality 1.3.5 yields

$$\mathbb{P}\left(\frac{1}{n} \sum_{k=1}^n X_k \geq Q + \sqrt{\frac{2 \log(1/\delta)}{n}}\right) \leq \exp\left(-\frac{n \left(\sqrt{\frac{2 \log(1/\delta)}{n}}\right)^2}{2}\right) = \delta. \quad (1.5)$$

Thus, the UCB exploration gives a simple control of the overestimation probabilities and with it the exploration caused by the optimism principle. The parameter δ can be chosen and it turns out later that $\delta = \frac{1}{n^2}$ is a good choice of the time-horizon is fixed and known. Instead of writing down the algorithm to run the bandit we could also write down the learning strategy explicitly:

$$\pi_t(a; a_1, x_1, \dots, x_t) = \operatorname{argmax}_a \text{UCB}_a(t-1, \delta)$$

Data: bandit model ν , $\delta \in (0, 1)$, n , vector \hat{Q}
Result: actions A_1, \dots, A_n and rewards X_1, \dots, X_n
 Initialise $T_a = 0$ for all a ;
while $t \leq n$ **do**
 $A_t = \operatorname{argmax}_a \operatorname{UCB}_a(t - 1, \delta)$;
 Obtain reward X_t by playing arm A_t ;
 Set $T_{A_t} = T_{A_t} + 1$;
 Set $\hat{Q}_{A_t}(t) = \hat{Q}_{A_t}(t) + \frac{1}{T_{A_t}}(X_t - \hat{Q}_{A_t}(t))$;
end

Algorithm 4: UCB algorithms with parameter δ

with the dependence $\hat{Q}_a(t) = \frac{1}{T_a(t)} \sum_{k=1}^t x_k \mathbf{1}_{a_k=a}$ and $T_a(t) = \sum_{k=1}^t \mathbf{1}_{a_k=a}$ of the past. Note that the initialisation $\operatorname{UCB}(0, \delta) \equiv +\infty$ forces the algorithm to explore every arm at least once, a condition that reasonable algorithms should fulfill.



Theorem 1.3.8. Suppose ν is a bandit model with 1-subgaussian arms, $n \in \mathbb{N}$, $\delta = \frac{1}{n^2}$. Then the UCB algorithms initialised with $\hat{Q} \equiv 0$ has the following regret upper bound:

$$R_n(\pi) \leq 3 \sum_{a \in \mathcal{A}} \Delta_a + 16 \log(n) \sum_{a: Q_a \neq Q_*} \frac{1}{\Delta_a}.$$

Simulations show that the simple UCB algorithms performs very well on most examples. In contrast to the optimal version of ETC (cheating by using reward gaps) there is no need to use reward gaps.

Proof. We will assume without loss of generality that $a_* = a_1$ and estimate the expected times a suboptimal arm a will be played. Since δ is fixed as $\frac{1}{n^2}$ we skip the δ from $\operatorname{UCB}(t, \delta)$. Combined with the regret decomposition Lemma 1.2.8 this will give the upper bound on the regret.



Here is the idea of the analysis. Using the regret decomposition it suffices to estimate $\mathbb{E}[T_a(n)]$. A bit similarly to the analysis of ETC we decompose

$$\mathbb{E}[T_a(n)] = \mathbb{E}[T_a(n) \mathbf{1}_{H_m}] + \mathbb{E}[T_a(n) \mathbf{1}_{H_m^c}],$$

where the events H_m should be {arm a is played at most m times}. In that case we can estimate

$$\mathbb{E}[T_a(n)] \leq m\mathbb{P}(H_m) + n\mathbb{P}(H_m^c) \quad (1.6)$$

since $T_a(n)$ can be at most n . If now the probability can be estimated then one can optimise over m . The proof is a bit more delicate, we cannot compute H_m . Instead we use smaller events $G_m \subseteq H_m$ for which the probabilities can be computed. The decomposition (1.6) works equally but the way the G_m are defined their probabilities can be estimated using the independence of all rewards.

The estimates are based on Hoeffding's inequality for sums of iid random variables. This is why we use the random stack construction of the bandit process (A, X) . For that sake recall the table $\{X_t^{(a)}\}_{t \leq n, a \in \mathcal{A}}$ of independent random variables with $X_t^{(a)} \sim P_a$. Using

$$\bar{Q}_s^{(a)} = \frac{1}{s} \sum_{k=s}^s X_k^{(a)}$$

this means that $\hat{Q}_a(t) = \bar{Q}_s^{(a)}$ if $T_a(t) = s$. From now on we fix arm a and estimate $\mathbb{E}[T_a(n)]$. Define $G_m = G_1 \cap G_{2,m}$ with

$$G_1 = \left\{ \omega : Q_{a_1} < \min_{t \leq n} \text{UCB}_{a_1}(t)(\omega) \right\},$$

$$G_{2,m} = \left\{ \omega : \bar{Q}_m^{(a)}(\omega) + \sqrt{\frac{2 \log(1/\delta)}{m}} < Q_{a_1} \right\},$$

with a natural number m to be specified later. Thus, G_m is the event when Q_{a_1} is never underestimated by the upper confidence bound of the first arm, while at the same time the upper confidence bound for the mean of arm a after m observations are taken from this arm is below the action value of the optimal arm. In what follows we will estimate the probability of G_m and show that $G_m \subseteq H_m$ with H_m from above. Let us start with the latter and prove that

$$\text{if } \omega \in G_m, \text{ then } T_a(n)(\omega) \leq m. \quad (1.7)$$

First in words: If $\omega \in G_{2,m}$ and $\omega \in G_1$ then the UCB value for the m th round of playing arm a is smaller than that for playing arm a_1 (because one is smaller, the other bigger than Q_{a_1}). Thus, arm a is not played more than m times. Now more formally. Suppose $\omega \in G_m$ but $T_a(n)(\omega) > m$ holds. Then there must be some time index $t \leq n$ with $T_a(t-1)(\omega) = m$ and $A_t(\omega) = a$. Using the definitions of $G_1, G_{a,2}$, and UCB yields

$$\begin{aligned} \text{UCB}_a(t-1)(\omega) &\stackrel{\text{Def.}}{=} \hat{Q}_a(t-1)(\omega) + \sqrt{\frac{2 \log(1/\delta)}{T_a(t-1)(\omega)}} \\ &= \bar{Q}_m^{(a)}(\omega) + \sqrt{\frac{2 \log(1/t\delta)}{m}} \\ &\stackrel{G_{a,2}}{<} Q_{a_1} \\ &\stackrel{G_1}{<} \text{UCB}_{a_1}(t-1)(\omega), \end{aligned}$$

a contradiction to $a = A_t(\omega) = \text{argmax}_b \text{UCB}_b(t-1)$. Hence, (1.7) holds.

We now follow the idea sketched above. Since $T_a(n)$ is trivially bounded by n the following key estimate holds:

$$\mathbb{E}[T_a(n)] = \mathbb{E}[T_a(n)\mathbf{1}_{G_m}] + \mathbb{E}[T_a(n)\mathbf{1}_{G_m^c}] \leq m + n(\mathbb{P}(G_1^c) + \mathbb{P}(G_{2,m}^c)). \quad (1.8)$$

It now suffices to estimate separately both summands on the right hand side of (1.8). First, it holds that

$$\begin{aligned} \mathbb{P}(G_1^c) &= \mathbb{P}(Q_{a_1} \geq \text{UCB}_{a_1}(t) \text{ for some } t \leq n) \\ &\leq \sum_{s \leq n} \mathbb{P}\left(Q_{a_1} - \sqrt{\frac{2 \log(1/\delta)}{s}} \geq \bar{Q}_s(a_1)\right) \\ &\stackrel{\text{Hoeffding}}{\leq} \sum_{t \leq n} \delta = n\delta. \end{aligned}$$

Next, we chose m which so far was left unspecified. Let us chose m large enough so that

$$\Delta_a - \sqrt{\frac{2 \log(1/\delta)}{m}} \geq \frac{1}{2} \Delta_a \quad (1.9)$$

holds, for instance $m = \left\lceil \frac{2 \log(1/\delta)}{\frac{1}{4} \Delta_a^2} \right\rceil$. Then

$$\begin{aligned} \mathbb{P}(G_{2,m}^c) &= \mathbb{P}\left(\bar{Q}_m^{(a)} + \sqrt{\frac{2 \log(1/\delta)}{m}} \geq Q_{a_1}\right) \\ &= \mathbb{P}\left(\bar{Q}_m^{(a)} - Q_a \geq \Delta_a - \sqrt{\frac{2 \log(1/\delta)}{m}}\right) \\ &\stackrel{(1.9)}{\leq} \mathbb{P}\left(\bar{Q}_m^{(a)} \geq Q_a + \frac{1}{2} \Delta_a\right) \\ &\stackrel{\text{Hoeffding}}{\leq} \exp\left(-\frac{m \Delta_a^2}{8}\right). \end{aligned}$$

Combining the above yields

$$\mathbb{E}[T_a(n)] \leq m + n \left(n \delta + \exp\left(-\frac{m \Delta_a^2}{8}\right) \right). \quad (1.10)$$

It remains to plug-in. Using $\delta = \frac{1}{n^2}$ yields

$$m = \left\lceil \frac{2 \log(n^2)}{\frac{1}{4} \Delta_a^2} \right\rceil \leq 1 + \frac{16 \log(n)}{\Delta_a^2}$$

so that

$$\mathbb{E}[T_a(n)] \leq m + 1 + n n^{-2} \leq m + 2 \leq 3 + \frac{16 \log(n)}{\Delta_a^2}.$$

Combined with the regret decomposition the claim follows. \square

Lecture 4

The previous bound is logarithmic in n with inverse reward gaps. While the dependency in n is very good the inverse reward gaps can be arbitrarily large if the second best arm is close to the best arm. Estimating slightly differently in the final step of the proof we can derive a different upper bound which is less favorable in the dependency of n but more favorable in term of the reward gap (which of course is a priori unknown).



Theorem 1.3.9. Suppose ν is a bandit model with 1-subgaussian arms, $n \in \mathbb{N}$, and $\delta = \frac{1}{n^2}$. Then the UCB algorithms also has the alternative regret upper bound

$$R_n(\pi) \leq 8\sqrt{Kn \log(n)} + 3 \sum_{a \in \mathcal{A}} \Delta_a.$$

The term $\sum_{a \in \mathcal{A}} \Delta_a$ appearing in both bounds is very natural as every reasonable algorithm plays each arm at least once and thus, by the regret decomposition lemma, gives $\sum_{a \in \mathcal{A}} \Delta_a$. In many examples the sum of the reward gaps can be bounded. If for instance all arms are Bernoulli distributed, then the sum can be replaced by k and be neglected as it will be dominated by the first summand. More interesting is the first summand for which one can decide to emphasize more the dependence on n or the regret gap.

Proof. The proof of the regret bound given above revealed that $\mathbb{E}[T_a(n)] \leq 3 + \frac{16 \log(n)}{\Delta_a^2}$. The idea of the proof is as follows. From the regret decomposition we know that small reward gaps should not pose problems as they appear multiplicatively. Separating the arms by a threshold into those with small and those with large reward gaps we only use the estimate for the ones

with large reward gaps and then minimise over the threshold. Pursuing this way leads to

$$\begin{aligned}
 R_n(\pi) &= \sum_{a \in \mathcal{A}} \Delta_a \mathbb{E}[T_a(n)] \\
 &= \sum_{a \in \mathcal{A}: \Delta_a < \Delta} \Delta_a \mathbb{E}[T_a(n)] + \sum_{a \in \mathcal{A}: \Delta_a \geq \Delta} \Delta_a \mathbb{E}[T_a(n)] \\
 &\leq n\Delta + \sum_{a \in \mathcal{A}: \Delta_a \geq \Delta} \left(\Delta_a 3 + \frac{16 \log(n)}{\Delta_a} \right) \\
 &\leq n\Delta + \frac{16K \log(n)}{\Delta} + 3 \sum_{a \in \mathcal{A}} \Delta_a.
 \end{aligned}$$

Since Δ can be chosen arbitrarily it suffices to minimise the righthand side as a function in Δ . The minimum can be found easily by differentiation in Δ to be located at $\Delta = \sqrt{16K \log(n)/n}$. Plugging-in yields

$$R_n(\pi) \leq 8\sqrt{Kn \log(n)} + 3 \sum_{a \in \mathcal{A}} \Delta_a.$$

□



For σ -subgaussian bandit models the UCB exploration bonus is modified as

$$\text{UCB}_a(t) := \begin{cases} \infty & : T_a(t) = 0 \\ \underbrace{\hat{Q}_a(t)}_{\text{greedy}} + \underbrace{\sqrt{\frac{4\sigma^2 \log(n)}{T_a(t)}}}_{\text{exploration bonus}} & : T_a(t) \neq 0. \end{cases}$$

Check that the regret bound in Theorem 1.3.8 changes to

$$R_n(\pi) \leq 3 \sum_{a \in \mathcal{A}} \Delta_a + 16\sigma^2 \log(n) \sum_{a: Q_a \neq Q_*} \frac{1}{\Delta_a},$$

and this leads to

$$R_n(\pi) \leq 8\sigma\sqrt{Kn \log(n)} + 3 \sum_{a \in \mathcal{A}} \Delta_a$$

in Theorem 1.3.9. Thus, for Bernoulli bandits the exploration bonus should be $\sqrt{\frac{1}{4} \frac{\log(n)}{T_a(t)}}$ and, as you should check in simulations (!) the constant $\frac{1}{4}$ is crucial for a good performance.

The idea of the last proof is very important. The model based regret bounds are often completely useless as they emphasise too strongly small reward gaps which by means of the regret decomposition should not be important at all. To get a better feeling please think a moment about the following exercise:



Recall (1.2), the upper bound for ETC in the case of two arms. Use the idea from the proof of Theorem 1.3.9 to derive the following upper bound of the ETC regret:

$$R_n(\pi) \leq \Delta + C\sqrt{n},$$

for some model-free constant C (for instance $C = 8 + \frac{2}{\epsilon}$) so that, in particular, $R_n(\pi) \leq 1 + C\sqrt{n}$ for all bandit models with regret bound $\Delta \leq 1$ (for instance for Bernoulli bandits).

It is very crucial to compare the constants appearing in the regret bounds. As mathematicians we tend to overestimate the importance of n and always think of $n \rightarrow \infty$ (as we did in the formulation of Theorem 1.3.7). Keeping in mind logarithmic growth one quickly realises the importance of constants. As an example $\log(1.000.000) \approx 13$ so that a constant \sqrt{K} or even worse $1/\Delta$ can be much more interesting. As an example the constants (5 and $1/\Delta$ appears even squared) in the regret bound of explore-then- ε greedy are extremely bad even though the logarithm in n is right order.



UCB is typically used as benchmark algorithm to compare other bandit algorithms. The reason is that UCB is simple to state, simple to analyse, and also pretty close to optimal both in the model-based and model-independent regret bounds. Authors typically compare their bounds with the growth in n (typically logarithmic model-based and square-root model-independent), the appearing generic constants (such as 8 in UCB), and how the reward bounds enter the regret upper bounds.

Comparing with the Lai-Robbins lower bounds for subgaussian bandit algorithms shows that the UCB algorithm is pretty close to optimal. There are more refined versions of the simple UCB algorithm presented above, such as the MOOS algorithm. What changes are different choices for additional exploitation. For further reading we refer to the wonderful book „Bandit Algorithms“ of Tor Lattimore and Csaba Szepesvári which is available online for free.

1.3.3 Boltzmann exploration

In this section we present a method that connects greedy exploration and the UCB algorithm in a surprising way. Before explaining the algorithm the concept of softmax distributions is needed.



Definition 1.3.10. Suppose $x, \theta \in \mathbb{R}^d$, then x defines a discrete distribution $\text{SM}(\theta, x)$ on finite sets with d elements with probability weights

$$p_k := \frac{e^{\theta_k x_k}}{\sum_{i=1}^d e^{\theta_i x_i}}, \quad k = 1, \dots, d.$$

The weights are called Boltzmann weights of the vector x .

Typically, all θ_i are equal, in which case θ is called inverse temperature. The origin of such distributions lies in statistical physics. The softmax distributions is a so-called categorical distribution, written $\text{Categorical}(p_1, \dots, p_d)$ which is also called a multinoulli or generalised Bernoulli distribution with probabilities p_1, \dots, p_d . For us, the following idea is much more important. If all θ_k are non-negative then sampling from $\text{SM}(\theta, x)$ is somewhat similar to the deterministic distribution M_x that only charges mass on the index $\text{argmax}_k x_k$ because $\text{SM}(\theta, x)$ assigns the highest probabilities to the coordinate with largest value x_k . But unlike the deterministic maximum distribution the softmax distribution only weights stronger the maximal element than the smaller elements. The role of θ is important: for large θ the softmax resembles the argmax distribution while for small θ the distribution resembles the uniform distribution on $\{1, \dots, d\}$.

Replacing sampling from the (deterministic) argmax distribution in the greedy algorithm yields the Boltzmann exploration algorithm. In contrast to the greedy algorithm there is continuous exploration as all arms have positive probabilities. It is quite obvious that the choice of θ is crucial as the algorithm resembles two unfavorable algorithms with linear regret, both for small and large θ (uniform and greedy exploration). In fact³, both constant and decreasing choices of θ are unfavorable. To guess a better choice for θ we need the following surprising lemma:

³N. Cesa-Bianchi, C. Gentile, G. Lugosi, G. Neu: Boltzmann exploration done right, NeuRIPS, (2017)

Data: bandit model ν , vector \hat{Q} , parameter θ
Result: actions A_1, \dots, A_n and rewards X_1, \dots, X_n
 Initialise $T_a = 0$ for all a ;
while $t \leq n$ **do**
 Sample A_t from $\text{SM}(\theta, \hat{Q})$;
 Obtain reward X_t by playing arm A_t ;
 Set $T_{A_t} = T_{A_t} + 1$;
 Set $\hat{Q}_{A_t} = \hat{Q}_{A_t} + \frac{1}{T_{A_t}}(X_t - \hat{Q}_{A_t})$;
end

Algorithm 5: Simple Boltzmann exploration



Lemma 1.3.11. Let $x, \theta \in \mathbb{R}^d$, then

$$\text{SM}(\theta, x) \stackrel{(d)}{=} \operatorname{argmax}_k \{\theta_k x_k + g_k\},$$

where g_1, \dots, g_d are iid Gumbel random variables with scale 1 and mode 0, i.e. have probability density function $f(x) = e^{-(x+e^{-x})}$ or CDF $F(t) = e^{-e^{-t}}$.

Proof. Reformulated right the proof is quite simple. We first prove very well-known statement from elementary probability theory on exponential random variables. Suppose E_1, E_2 are independent exponential random variables with parameters λ_1, λ_2 . Then $\mathbb{P}(E_1 \leq E_2) = \frac{\lambda_1}{\lambda_1 + \lambda_2}$. This is a simple integration exercise: with $A = \{0 \leq x_1 \leq x_2\} \in \mathcal{B}(\mathbb{R}^2)$ the computation rules for vectors of independent random variables yields

$$\begin{aligned} \mathbb{P}(E_1 \leq E_2) &= \mathbb{E}[\mathbf{1}_A(E_1, E_2)] \\ &= \int_A f_{(E_1, E_2)}(x_1, x_2) \, d(x_1, x_2) \\ &= \int_0^\infty \int_0^{x_2} \lambda_1 e^{-\lambda_1 x_1} \lambda_2 e^{-\lambda_2 x_2} \, dx_1 \, dx_2 \\ &= \lambda_1 \lambda_2 \int_0^\infty e^{-\lambda_2 x_2} \frac{1}{\lambda_1} (1 - e^{-\lambda_1 x_2}) \, dx_2 \\ &= \lambda_2 \int_0^\infty (e^{-\lambda_2 x_2} - e^{-(\lambda_1 + \lambda_2)x_2}) \, dx_2 \\ &= 1 - \frac{\lambda_2}{\lambda_1 + \lambda_2} = \frac{\lambda_1}{\lambda_1 + \lambda_2}. \end{aligned}$$

Next, we can compute the argmin distribution of independent exponential variables. Suppose E_1, \dots, E_n are independent exponential random variables with parameters $\lambda_1, \dots, \lambda_n$, then $\mathbb{P}(\operatorname{argmin}_{j \leq n} E_j = i) = \frac{\lambda_i}{\lambda_i + \sum_{k \neq i} \lambda_k}$. The identity is a direct consequence from the above:

$$\mathbb{P}(\operatorname{argmin}_j E_j = i) = \mathbb{P}(E_i \leq E_k, \forall k) = \mathbb{P}(E_i \leq E) = \frac{\lambda_i}{\lambda_i + \sum_{k \neq i} \lambda_k},$$

where $E := \min_{k \neq i} E_k$ is independent of E_i and exponentially distributed with parameter $\sum_{k \neq i} \lambda_k$. Here recall from elementary probability that the minimum of independent exponentials is again exponential and the parameter is the sum of the parameters. The second ingredient is a connection between exponential and Gumbel variables. If $X \sim \text{Exp}(1)$, then $Z := -\log(X)$ is Gumbel with scale parameter $\beta = 1$ and mode $\mu = 0$. The claim is checked by computing the cumulative distribution function:

$$\mathbb{P}(Z \leq t) = \mathbb{P}(X \geq e^{-t}) = e^{-e^{-t}}, \quad t \in \mathbb{R}.$$

As a consequence, with g Gumbel with scale $\beta = 1$ and mode $\mu = 0$ and $E_\lambda \sim \text{Exp}(\lambda)$, we obtain the identity in law

$$c + g \sim -\log(e^{-c}) - \log(E_1) = -\log(e^{-c}E_1) \sim -\log(E_{e^c}).$$

For the last equality in law we have also used the scaling property of the exponential distribution. If g_1, \dots, g_k are iid Gumbel with scale $\beta = 1$ and mode $\mu = 0$ we finally obtain

$$\begin{aligned} \mathbb{P}(\operatorname{argmax}_k(\theta_k x_k + g_k) = i) &= \mathbb{P}(\operatorname{argmax}_k -\log(E_{e^{\theta_k x_k}}) = i) \\ &= \mathbb{P}(\operatorname{argmin}_k E_{e^{\theta_k x_k}} = i) \\ &= \frac{e^{\theta_i x_i}}{\sum_k e^{\theta_k x_k}}. \end{aligned}$$

In other words, the Gumbel argmax is distributed according to $\text{SM}(\theta, x)$. \square

Using the lemma the Boltzmann algorithm can now be interpreted differently: arms are chosen according to the distribution

$$A_t \sim \operatorname{argmax}_a \{ \theta \hat{Q}_a(t-1) + g_a \} = \operatorname{argmax}_a \{ \hat{Q}_a(t-1) + \frac{g_a}{\theta} \},$$

where the g_a are iid Gumbel and independent of \hat{Q} . Does this look familiar? Of course, this is the greedy strategy with an additional random exploration bonus as we have seen in the UCB algorithm. Since typical values of g are around 1, motivated by the UCB algorithm a first idea should immediately come to mind: θ^{-1} should be arm-dependent and $\sqrt{\frac{C}{T_a}}$ might be a good idea. In fact, that's true as was shown in the article of Cesa-Bianchi et al. We will not go into detail here. The aim of this section was to link the Boltzmann exploration with greedy-type algorithms. Here is an interesting research question to think about. Let us forget that UCB-type exploration with Gumbel random exploration bonus is linked to the rather natural exploitation strategy given by Boltzmann softmax weights. It is then very natural to replace the Gumbel distribution by some other distribution. It seems most plausible that non-negative distributions should be more reasonable as a negative value turns the exploration bonus into an exploration malus which was never intended.



Use the simulation code provided for the exercises to play with different distributions and see if replacing the Gumbel distribution can be favorable. Choosing distributions that concentrate around 1 (Dirac measure at 1 gives UCB) might be reasonable, such as a Gamma distribution with shape parameter larger than 1 and scale parameter that forces expectation 1.

Lecture 5

1.3.4 Simple policy gradient for stochastic bandits

We now come to a completely different approach, the so-called policy gradient approach. The approach presented here is not really part of the research field of multiarmed bandits but is a very simple special case of which will later be called policy gradient method for reinforcement learning. We use the softmax policy gradient method to connect the topic of multiarmed bandits with reinforcement learning that will be started in the next lecture. For that sake, here is a new way of thinking of bandits. The bandit game is a one-step game, the game of choosing one of the K arms and obtaining the corresponding outcome.



Definition 1.3.12. A probability distribution π on \mathcal{A} is called a policy for the bandit game, the expected reward $V(\pi) := Q_\pi := \sum_{a \in \mathcal{A}} Q_a \pi(a)$ when playing the policy is called the value of the policy.

Keep in mind that since \mathcal{A} is assumed to be finite a distribution on \mathcal{A} is nothing but a probability vector, a vector of non-negative numbers that sums up to 1. A vector $(0, \dots, 1, \dots, 0)$ with 1 for the a th arm (position) corresponds to playing only arm a and in this case $Q_\pi = Q_a$. A policy is optimal if the expected outcome Q_π is maximal, i.e. equals Q_* . If there are several optimal arms, i.e. arms satisfying $Q_a = Q_*$, then any policy is optimal that has mass only on optimal arms.



The goal in reinforcement learning is similar to that of multiarmed bandits, but different. In reinforcement learning we are typically concerned with finding the best (at least very good) policy (here: in the bandit game the best arm) as quickly as possible. In reinforcement learning we are not aiming to minimise regret. Essentially, we do not care if very bad arms are played as long as the best arm is found quickly. A learning strategy is used to find the optimal policy.

One of the reasons for this different point of view is the field of applications. While one of the main motivations for the multiarmed bandit stems from medicine where every single life counts a major field of applications that pushed the development of reinforcement is automated learning of optimal strategies in gaming or online advertisement where obviously the impact of suboptimal attempts is much less severe.

In this section we start with a simple approach to policy search. A set of possible policies is defined and then a learning strategy tries to find the best among these policies.



Definition 1.3.13. If $\Theta \subseteq \mathbb{R}^d$, then a set $\{\pi_\theta : \theta \in \Theta\}$ of probability distributions on \mathcal{A} is called a parametrised family of policies.

The policy gradient idea is as follows: given a parametrised policy over a continuous index set we aim to maximise the value function J over the parameter set:

$$\theta \mapsto J(\theta) := Q_{\pi_\theta} = \sum_{a \in \mathcal{A}} \pi_\theta(a) Q_a.$$

The value function J is nothing but a multidimensional function, hence, maximisation algorithms can be used. If the parametric family can approximate Dirac-measures δ_a then one could hope to identify the optimal arm using an optimisation procedure. So how can we find the optimal arm with policy gradient? By typical optimization methods from lectures on numerical analysis. One example is the classical gradient ascent method:

$$\theta_{n+1} := \theta_n + \alpha \nabla J(\theta_n), \quad n \in \mathbb{N}.$$

Under suitable assumptions on J (such as convexity) that algorithm converges to a maximum θ_* of J . Unfortunately, that approach has diverse difficulties which are topics of ongoing research:

1. Given a bandit model, what is a good parametric family of policies?
2. If J is unknown (because the expectations Q_a are unknown) how can gradient descent be carried out (most efficiently)?
3. Even if the function J would be known explicitly, will the gradient ascent algorithm converge to an optimal policy? How fast?

In this first section on the policy gradient method we will address the first two issues by a discussion of one particular parametrised family and a sketch of what later will be discussed in details in the chapter on the policy gradient method. The third question is widely open with some recent progress⁴. We will start with the second issue on how to deal with the gradient of J if J is unknown but samples can be generated. The trick that will occur a few times in this lecture course is the so-called log-trick, here in its simplest form:

⁴see for instance A. Agarwal, S. Kakade, J. Lee, G. Mahajan: „On the Theory of Policy Gradient Methods: Optimality, Approximation, and Distribution Shift“, JMLR, 1-76, (2021)



The following short computation is typically called log-trick (or score-function trick):

$$\begin{aligned}
 \nabla J(\theta) &= \nabla \sum_{a \in \mathcal{A}} Q_a \pi_\theta(a) \\
 &= \sum_{a \in \mathcal{A}} Q_a \nabla \pi_\theta(a) \frac{\pi_\theta(a)}{\pi_\theta(a)} \\
 &= \sum_{a \in \mathcal{A}} Q_a \nabla \log(\pi_\theta(a)) \pi_\theta(a) \\
 &= \mathbb{E}_{\pi_\theta} [Q_A \nabla \log(\pi_\theta(A))] \\
 &= \mathbb{E}_{\pi_\theta} [X_A \nabla \log(\pi_\theta(A))],
 \end{aligned}$$

where the last equality follows from the tower-property and $\mathbb{E}[X_A|A] = Q_A$ for instance by using $X_t = \sum_{a \in \mathcal{A}} X_t^{(a)} \mathbf{1}_{A_t=a}$ the random table model from the proof of Theorem 1.2.4. Note that whenever we take the expectation of a vector of random variables (the gradient is a vector) the expectation is taken coordinate wise.

Now that the gradient is expressed as the expectation of a random vector the idea is to replace in the gradient ascent algorithm the true gradient by samples of the underlying random variable. Either by one sample

$$\tilde{\theta}_{n+1} := \tilde{\theta}_n + \alpha X_{A_n} \nabla \log(\pi_{\tilde{\theta}_n}(A_n)), \quad n \in \mathbb{N}, \quad (1.11)$$

or by a so-called batch of independent samples:

$$\tilde{\theta}_{n+1} := \tilde{\theta}_n + \alpha \frac{1}{N} \sum_{i=1}^N X_{A_n^i} \nabla \log(\pi_{\tilde{\theta}_n}(A_n^i)), \quad n \in \mathbb{N}, \quad (1.12)$$

where the A_n^i are independently sampled according to $\tilde{\pi}_{\tilde{\theta}_n}$ and the $X_{A_n^i}$ are independent samples from the arms A_n^i . The sequence of $\pi_{\tilde{\theta}_n}$ of policies obtained from gradient is a learning strategy that hopefully approximates the optimal policy.



Definition 1.3.14. The appearing function $(a, \theta) \mapsto \nabla \log \pi_\theta(a)$ is called the score-function of π_θ .

Of course the algorithm is most useful if the score-function is nice so that the gradient disappears. Here is the most prominent example:

Example 1.3.15. Suppose $d = |\mathcal{A}|$, so that there is a parameter θ_a for each arm a . Furthermore, define

$$\pi_\theta(a) = \frac{e^{\theta_a}}{\sum_{k \in \mathcal{A}} e^{\theta_k}}.$$

Then the family $\{\pi_\theta : \theta \in \mathbb{R}^d\}$ is called the tabular softmax family. The softmax family is huge (way too big to be used in practice) and rich enough to approximate all optimal policies. Indeed, if an arm a is optimal, then the parameter θ_a needs to be sent to $+\infty$. The score-function is easily computed to be

$$(\nabla \log \pi_\theta(a))_i = \mathbf{1}_{a=i} - \left(\nabla \log \left(\sum_{a \in \mathcal{A}} e^{\theta_a} \right) \right)_i = \mathbf{1}_{a=i} - \frac{e^{\theta_i}}{\sum_{a \in \mathcal{A}} e^{\theta_a}} = \mathbf{1}_{a=i} - \pi_\theta(i).$$

Plugging-into (1.11) yields the updates (now written component wise)

$$\begin{aligned}\theta_{1,n+1} &= \theta_{1,n} - \alpha X_i \pi_\theta(1), \\ &\dots = \dots \\ \theta_{i,n+1} &= \theta_{i,n} + \alpha X_i (1 - \pi_\theta(i)), \\ &\dots = \dots \\ \theta_{d,n+1} &= \theta_{d,n} - \alpha X_i \pi_\theta(d)\end{aligned}$$

if $N = 1$, $A_n = i$, and X_i is a sample of arm i . Here is the first explanation why reinforcement learning is called reinforcement learning. Remember that the θ_a are indirect parametrisations of the probabilities to play arm a . If at time n the policy current policy decides to play arm i then the probabilities are reinforced as follows. If the reward obtained by playing arm i is positive, let's interpret the positive reward as positive feedback, then the probability to play arm i is increased and all other probabilities are decreased. Similarly, if the feedback from playing arm i is negative, then the probability to do so in the next round is decreased and all other probabilities are increased. But how about the situation in which all rewards are positive? How do we learn to distinguish good and bad arms? It would be much easier if bad arms have negative rewards leading to a reduction of likelihood in the softmax update. In that situation the actor only learns indirectly. The actor only learns an arm is bad by playing good arms which reduces the likelihood to play the bad arms. It would be much more effective to learn directly that bad arms are bad. One might think about the mensa food. If a dish is bad we can learn it is bad by trying it. If all dishes are good and a few are pretty good we can only learn a dish is relatively bad by eating better food.

```

Input: bandit model  $\nu$ , step-size parameter  $\alpha$ 
for  $t \geq 0$  do
  Sample an action  $a_t \sim \pi_{\theta_t}(\cdot)$  and observe reward
   $X_t \sim P_{a_t}$ 
  for each action  $a$  do
    if  $a = a_t$  then
       $\theta_{t+1}(a) \leftarrow \theta_t(a) + \alpha X_t (1 - \pi_{\theta_t}(a))$ 
    else
       $\theta_{t+1}(a) \leftarrow \theta_t(a) - \alpha X_t \pi_{\theta_t}(a)$ 
    end if
  end for
end for
Algorithm 6: Softmax policy gradient for bandits (without baselines)

```

Another way of seeing the same problem is to think about committal behavior. Suppose for instance that the starting vector is $\theta_0 \equiv 0$ and all arms return positive values (for instance Bernoulli bandits). Then the first arm is chosen uniformly as π_{θ_0} is uniform on \mathcal{A} . If the first chosen arm, say a , yields a positive return then the the updated parameter vector θ_1 only has one positive entry, namely the a th entry. Hence, for the second update the a th arm is most likely to be chosen again. A way to reduce committal behavior is to subtract a baseline that is supposed to help distinguish good and bad arms⁵.



Here is another representation for the gradient:

$$\nabla J(\theta) = \mathbb{E}_{\pi_\theta} [(X_A - b) \nabla \log(\pi_\theta(A))],$$

where b is any constant. The reason is simply that, using the same computation as

⁵W. Chung, V. Thomas, M. Machado, N. Le Roux: „Beyond Variance Reduction: Understanding the True Impact of Baselines on Policy Optimization“, ICML, (2021)



above,

$$\begin{aligned}\mathbb{E}_{\pi_\theta}[b\nabla \log(\pi_\theta(A))] &= b \sum_{a \in \mathcal{A}} \nabla \log(\pi_\theta(a)) \pi_\theta(a) \\ &= b \sum_{a \in \mathcal{A}} \nabla \pi_\theta(a) \frac{\pi_\theta(a)}{\pi_\theta(a)} \\ &= b \nabla \sum_{a \in \mathcal{A}} \pi_\theta(a) \\ &= b \nabla 1 = 0.\end{aligned}$$

It is important to realise that b can be chosen arbitrary without changing the gradient and as such not changing the gradient ascent algorithm

$$\theta_{n+1} := \theta_n + \alpha \nabla J(\theta), \quad n \in \mathbb{N}.$$

Still, it drastically changes the stochastic version (here for the softmax parametrisation)

$$\tilde{\theta}_{i,n+1} := \tilde{\theta}_{i,n} + \alpha(X_{A_n} - b)(\mathbf{1}_{A_n=i} - \pi_{\tilde{\theta}_n}(i)), \quad n \in \mathbb{N}, i \in \mathcal{A},$$

as now the effect of committal behavior can be reduced drastically. An important question is what baseline to choose so that the algorithm is speed up the most. Usually people choose b to minimise the variance of the estimator of ∇J , but on the pathwise level this choice is not optimal. Thinking about the pathwise behavior it would be more reasonable to use $b = V(\pi_\theta)$ as it turns rewards negatives if the arm returns less than expectation. Unfortunately, this b is unknown and should be estimated again. Later we will call such algorithms actor-critic.

Even though not optimal on a pathwise level it can be instructive to compute the minimum variance baseline for every step of the policy gradient scheme. Intuitively that makes sense as variance in the gradient approximations leads to wrong update directions and thus slows down the convergence.



The variance of a random vector X is defined by to be $\mathbb{V}[X] = \mathbb{E}[|X^2|] - \mathbb{E}[X]^T \mathbb{E}[X]$. Show by differentiation that

$$b_* = \frac{\mathbb{E}_{\pi_\theta}[X_A |\nabla \log \pi_\theta(A)|_2^2]}{\mathbb{E}_{\pi_\theta}[|\nabla \log \pi_\theta(A)|_2^2]}$$

is the baseline that minimises the variance of the unbiased estimators

$$(X_A - b) \nabla \log(\pi_\theta(A)), \quad A \sim \pi_\theta,$$

of $J(\theta)$.

To finish the discussion we might wish to estimate the regret $R_n(\pi)$ of a learning strategy $(\pi_{\theta_n})_{n \geq 1}$. In fact, almost nothing is known, only some rough estimates for the softmax parametrisation can be found in the literature. Give it a try!

1.4 Lower bounds for stochastic bandits

In the previous sections we have discussed a number of explicit algorithms to learn the optimal arm of a bandit model. For reward distributions with very small tails (such as Gaussian or even bounded) we derived $\log(n)$ and \sqrt{n} type estimates for the regret. In this section we derive lower

bounds using estimates from classical statistics (essentially hypothesis testing).

1.4.1 A bit on relative entropy

Let us first recall some notion from probability theory. Suppose P and Q are probability measures on some probability space (Ω, \mathcal{A}) . One says $P \ll Q$, P is dominated by Q or P is absolutely continuous with respect to Q , if $P(A)$ implies $Q(A)$ for all $A \in \mathcal{A}$. The Radon-Nykodym theorem states that in that situation there is a non-negative measurable mapping f such that $Q(A) = \int_A f dP$. While it is not too hard to show the domination property computing densities might be a complicated matter.



Definition 1.4.1. (Relative Entropie or Kullback-Leibler divergence)

Suppose P and Q are probability measures on a probability space (Ω, \mathcal{F}) . Then

$$D(P, Q) = \begin{cases} \int \log\left(\frac{dP}{dQ}\right) dP & : \text{if } P \ll Q \\ \infty & : \text{otherwise} \end{cases}$$

is called **relative entropie** (or KL-divergence) of P with respect to Q .

The expression divergence refers to a non-negative mapping on the cartesian product with properties similar (but weaker) than a metric. For instance, the KL divergence is not symmetric and does not satisfy the triangle inequality. Nonetheless, the KL-divergence is meant to measure in a way the difference of the measures. It is non-negative and zero if and only if $P = Q$. There is a lot of information theoretic meaning hidden in the definition, we will use the divergence for hypotheses testing. The assumption of absolute continuity is nothing but saying the integral is well-defined. Most importantly, if both measures are absolutely continuous, then they are clearly absolutely continuous and the Radon-Nykodym derivative is the ratio of the densities.



Lemma 1.4.2. Suppose $P \ll Q \ll \nu$, then $\frac{dP}{dQ} \frac{dQ}{d\nu} = \frac{dP}{d\nu}$.

Solving the equation yields $\frac{dP}{dQ} = \frac{\frac{dP}{d\nu}}{\frac{dQ}{d\nu}}$ but one should be a bit careful with division by zero. There is no big problem as $\{\frac{dQ}{d\nu} = 0\}$ is a Q -zero set and thus also a P zero set. Since densities are only unique up to nullsets one typically just modifies $\frac{dP}{dQ} = \mathbf{1}_{\frac{dQ}{d\nu} > 0} \frac{\frac{dP}{d\nu}}{\frac{dQ}{d\nu}}$

Proof.

$$P(B) = \int_B \frac{dP}{dQ} dQ = \int_B \frac{dP}{dQ} \frac{dQ}{d\nu} d\nu$$

Note that there is no problem with the denomination as $P \ll Q$ implies that zeros of q are also zeros of p (except on a Q -zero set). \square

The lemma shows how to compute densities $\frac{dP}{dQ}$ if one can compute densities with respect to the same dominating measure.



There are two important examples to which the lemma is usually applied. Two discrete (dominating measures the counting measure) or two absolutely continuous measures (dominating measure the Lebesgue measure). For two absolutely continuous measures P and Q with positive densities p and q the formula yields $\frac{dP}{dQ}(x) = \frac{p(x)}{q(x)}$. For discrete measures on a_1, \dots, a_N the formula yields $\frac{dP}{dQ}(a_k) = \frac{p_k}{q_k}$.

The formulas can be used to compute relative entropies for two discrete (resp. absolutely continuous) measures. Let's check two particularly important examples, discrete Bernoulli

variables and Gaussian random variables. If $P \sim \text{Ber}(p)$ and $Q \sim \text{Ber}(q)$ with $p, q \in (0, 1)$, then

$$D(P, Q) = p \log \left(\frac{p}{q} \right) + (1-p) \log \left(\frac{1-p}{1-q} \right)$$

because both are absolutely continuous with respect to the counting measures on $\{0, 1\}$. If $P \sim \mathcal{N}(\mu_1, \sigma^2)$ and $Q \sim \mathcal{N}(\mu_2, \sigma^2)$, then, using that both are absolutely continuous with respect to the Lebesgues measure,

$$\begin{aligned} d(P, Q) &= \frac{1}{\sqrt{2\pi\sigma^2}} \int_{\mathbb{R}} \left(-\frac{(x-\mu_1)^2}{2\sigma^2} - \frac{(x-\mu_2)^2}{2\sigma^2} \right) e^{-\frac{(x-\mu_1)^2}{2\sigma^2}} dx \\ &= \frac{1}{\sqrt{2\pi\sigma^2}} \int_{\mathbb{R}} e^{-\frac{(x-\mu_1)^2}{2\sigma^2}} \frac{1}{\sigma^2} x(\mu_1 - \mu_2) dx + \frac{1}{\sqrt{2\pi\sigma^2}} \int_{\mathbb{R}} e^{-\frac{(x-\mu_2)^2}{2\sigma^2}} \frac{\mu_2^2 - \mu_1^2}{2\sigma^2} dx \\ &= \frac{1}{\sigma^2} \left(\mu_1(\mu_1 - \mu_2) + \frac{\mu_2^2 - \mu_1^2}{2} \right) \\ &= \frac{(\mu_2 - \mu_1)^2}{2\sigma^2}. \end{aligned}$$

There is also a sufficiently handy formula for the relative entropy of two general Gaussian vectors. We will later see how to compute relative entropies for bandit processes using absolute continuity with respect to the uniform path measure.

Here is the key application of KL-divergence needed for our purposes:



Theorem 1.4.3. (Bretagnolle-Huber)
 Suppose P and Q are probability measures on a probability space (Ω, \mathcal{F}) . Then

$$P(A) + Q(A^C) \geq \frac{1}{2} \exp(-D(P, Q)), \quad \forall A \in \mathcal{F}.$$

Usually the Bretagnolle-Huber inequality is stated in terms of bounding the total-variation distance of two probability measures. The version we state here is a direct consequence of the total-variation version.

Proof. Let us denote $a \wedge b = \min\{a, b\}$ and $a \vee b = \max\{a, b\}$. We set $v = P + Q$ and denote by p and q the densities of P and Q with respect to v . Those exist as clearly $P \ll v$ and $Q \ll v$. According to the remark above it follows that

$$D(P, Q) = \int_{\Omega} \log \left(\frac{p(\omega)}{q(\omega)} \right) P(d\omega). \quad (1.13)$$

Next, we show that

$$\int (p \wedge q) dv \geq \frac{1}{2} \exp(-D(P, Q)). \quad (1.14)$$

Since $\{\omega: q(\omega) = 0\}$ is a Q -zero set it follows that

$$\begin{aligned} \exp(-D(P, Q)) &= \exp \left(- \int_{\{q>0\}} \log \left(\frac{p}{q} \right) dP \right) \\ &= \exp \left(\int_{\{q>0\}} \log \left(\frac{q}{p} \right) dP \right) \\ &\stackrel{\text{Jensen}}{\leq} \exp \left(2 \log \left(\int \sqrt{\frac{q}{p}} dP \right) \right) \\ &= \exp \left(2 \log \left(\int \sqrt{\frac{q}{p}} p dv \right) \right) \end{aligned}$$

$$\begin{aligned}
&= \left(\int \sqrt{pq} \, dv \right)^2 \\
&= \left(\int \sqrt{(p \wedge q)(q \vee p)} \, dv \right)^2 \\
&\stackrel{\text{Cauchy-Schwarz}}{\leq} \left(\int (p \wedge q) \, dv \right) \left(\int (p \vee q) \, dv \right) \\
&\leq 2 \left(\int (p \wedge q) \, dv \right).
\end{aligned}$$

For the final step we used that $p \wedge q + p \vee q = p + q$ so that $\int (p \vee q) \, dv = 2 - \int (p \wedge q) \, dv \leq 2$. This proves (1.14). Finally, using

$$\int (p \wedge q) \, dv = \int_A \underbrace{(p \wedge q)}_{\leq p} \, dv + \int_{A^c} \underbrace{(p \wedge q)}_{\leq q} \, dv \leq P(A) + Q(A^c).$$

the claim follows. \square

Here is an interesting example that explains the importance for hypothesis testing. Suppose we are given an observation and know it is Gaussian with a given variance σ^2 . What is unknown is the mean, it might be 0 or Δ . Now we want to give a rule on how to determine if 0 or Δ was the mean underlying the observation. How good can such a rule be? A rule consists of a measurable set $A \subseteq \mathbb{R}$, if the observation falls in A the rule predicts $\mu = 0$, otherwise $\mu = \Delta$. Denote $P \sim \mathcal{N}(\Delta, \sigma^2)$ and $Q \sim \mathcal{N}(0, \sigma^2)$. Then the error made if the underlying distribution is P is $P(A)$, $Q(A^c)$ if the true distribution is Q . Now the Bretagnolle-Huber inequality yields

$$P(A) + Q(A^c) \geq \frac{1}{2} \exp(-D(P, Q)) = \frac{1}{2} \exp\left(-\frac{\Delta^2}{2\sigma^2}\right).$$

If for instance $\Delta^2 < \sigma^2$ (signal (expectation) to noise (variance) ratio is small), then the righthand side is larger than $\frac{3}{10}$ implying that $\max\{P(A), Q(A^c)\} \geq \frac{3}{10}$. What does it mean? Either false-positive or false-negative must occur with at least probability $\frac{3}{10}$ if only one sample is used.

1.4.2 Mini-Max lower bounds (model-dependent)



Definition 1.4.4. If $\xi = \{v^{(i)}\}_{i \in I}$ is a family of stochastic bandit models, then

$$R_n^*(\xi) = \inf_{\pi} \sup_{v \in \xi} R_n(\pi, v)$$

is called the **minimax-regret** of ξ .

Minimising the minimax-regrets means finding a learning strategy that performs on all models, not only a fixed one. The worst possible regret should be small. In order to give a lower bound on the minimax-regret we will use the Bretagnolle-Huber inequality. To get the inequality working in the context of bandits the following lemma is needed:



Lemma 1.4.5. (Entropy-decomposition)

Let k the number of arms and $v = \{P_1, \dots, P_k\}$, $v' = \{P'_1, \dots, P'_k\}$ the reward distributions of two k -armed bandits. Furthermore, let π a learning strategy and \mathbb{P} and \mathbb{P}' the probability distributions on $(\mathbb{R} \times \{1, \dots, k\})^n$, that describe the n -step bandit process for the two models under the learning strategy π . In this notation \mathbb{R} describes the reward and $\mathcal{A} = \{1, \dots, k\}$ the arms. Then the following



decomposition holds:

$$D(\mathbb{P}_\pi, \mathbb{P}'_\pi) = \sum_{a \in \mathcal{A}} \mathbb{E}_\pi[T_a(n)] D(P_a, P'_a).$$

Proof. For simplicity let us first assume the bandit model is discrete, i.e. the reward distributions P_a are discrete (absolutely continuous with respect to a counting measure λ). The main observation is that the law of the bandit process for a given policy π is absolutely continuous with respect to counting measure on the finite set of action-reward paths taking values in the support of P_a . The counting measure on paths is $(\rho \otimes \lambda)^{\otimes n}$, where ρ is the counting measure on $\{1, \dots, K\}$. Keep in mind, we are only talking about discrete measures which are absolutely continuous with respect to the counting measure of the underlying set and the density is giving by the singleton probabilities. Hence,

$$\mathbb{P}_\pi(A) = \int_A p(a_1, x_1, \dots, a_n, x_n) (\rho \otimes \lambda)^{\otimes n} (d(a_1, x_1, \dots, a_n, x_n))$$

with the density being the singleton probabilities $A = \{(a_1, x_1, \dots, a_n, x_n)\}$ of paths:

$$p(a_1, x_1, \dots, a_n, x_n) = \prod_{t=1}^n \pi_t(\{a_t\}; a_1, x_1, \dots, a_{t-1}, x_{t-1}) P_{a_t}(\{x_t\})$$

Similarly, \mathbb{P}'_π has density

$$p'(a_1, x_1, \dots, a_n, x_n) = \prod_{t=1}^n \pi_t(\{a_t\}; a_1, x_1, \dots, a_{t-1}, x_{t-1}) P'_{a_t}(\{x_t\})$$

with respect to $(\rho \otimes \lambda)^{\otimes n}$. Now we are in the situation that both measures \mathbb{P} and \mathbb{P}' are absolutely continuous to the same domination measure. Hence, the density of \mathbb{P} with respect to \mathbb{P}' is

$$\frac{d\mathbb{P}_\pi}{d\mathbb{P}'_\pi}(a_1, x_1, \dots, a_n, x_n) = \frac{\prod_{t=1}^n \pi_t(\{a_t\}; a_1, x_1, \dots, a_{t-1}, x_{t-1}) P_{a_t}(x_t)}{\prod_{t=1}^n \pi_t(\{a_t\}; a_1, x_1, \dots, a_{t-1}, x_{t-1}) P'_{a_t}(\{x_t\})} = \frac{\prod_{t=1}^n P_{a_t}(\{x_t\})}{\prod_{t=1}^n P'_{a_t}(\{x_t\})}$$

Plugging-into the definition of relative entropy yields

$$\begin{aligned} D(\mathbb{P}_\pi, \mathbb{P}'_\pi) &= \mathbb{E}_\pi \left[\log \left(\prod_{t \leq n} \frac{P_{A_t}(\{X_t\})}{P'_{A_t}(\{X_t\})} \right) \right] \\ &= \sum_{t \leq n} \mathbb{E}_\pi \left[\log \left(\frac{P_{A_t}(\{X_t\})}{P'_{A_t}(\{X_t\})} \right) \right] \\ &= \sum_{t \leq n} \sum_{a=1}^k \mathbb{E}_\pi \left[\mathbf{1}_{\{A_t=a\}} \log \left(\frac{P_a(\{X_t\})}{P'_a(\{X_t\})} \right) \right] \\ &= \sum_{a=1}^k \sum_{t \leq n} \mathbb{P}_\pi(A_t = a) D(P_a, P'_a) \\ &= \sum_{a=1}^k \mathbb{E}_\pi[T_a(n)] D(P_a, P'_a), \end{aligned}$$

using that $\mathbb{P}(X_t \in B \mid A_t = a) = P_a(B)$.

If all arms are absolutely continuous the same argument works choosing λ to be Lebesgues measure replacing P_a by the densities p_a of P_a . \square

Here is the main theorem, a minimax-regret bound for Gaussian rewards.



Theorem 1.4.6. Let ξ^k the family of all k -armed bandits with Gaussian rewards unit variance and expectations $\mu_1, \dots, \mu_k \in [0, 1]$. Then

$$R_n^*(\xi^k) \geq \frac{1}{27} \sqrt{(k-1)n}$$

holds for all $n \geq k$.

Proof of Theorem 1.4.6. Every model is described by an expectation vector $\mu = (\mu_1, \dots, \mu_k) \in [0, 1]^k$. We show that for all $n \geq k-1$ and every learning strategy π there is an expectation vector μ with

$$R_n(\pi, v_\mu) \geq \frac{1}{27} \sqrt{(k-1)n}.$$

Let $\Delta \in [0, \frac{1}{2}]$ arbitrary and set $\mu = \mu(\Delta) = (\Delta, 0, \dots, 0)$. Next, let $a = \operatorname{argmin}_{i \leq k} \mathbb{E}_{v_\mu, \pi}[T_i(n)]$ the arm which is played the least under learning strategy π on the bandit model v_μ . Furthermore, define

$$\mu' = (\Delta, 0, \dots, 0, \underbrace{2\Delta}_{a\text{th position}}, 0, \dots, 0).$$

To save some notation we write \mathbb{P} and \mathbb{P}' for the laws of the bandit processes playing both policy π but on the different bandit models μ and μ' . In what follows we show that

$$\min\{R_n(\pi, v_\mu), R_n(\pi, v_{\mu'})\} \geq \frac{1}{27} \sqrt{(k-1)n}$$

for a suitably chosen Δ . Thus, there is (at least) one bandit model with regret lower bound $\frac{1}{27} \sqrt{(k-1)n}$. To do so, we start with a first computation based on the regret decomposition:

$$\begin{aligned} R_n(\pi, v_\mu) &= \sum_i \mathbb{E}[T_i(n)] \Delta_i \\ &= \sum_{i \neq 1} \mathbb{E}[T_i(n)] \Delta \\ &= \Delta(n - \mathbb{E}[T_1(n)]) \\ &= \Delta(n - \mathbb{E}[T_1(n)(\mathbf{1}_{T_1(n) > n/2} + \mathbf{1}_{T_1(n) \leq n/2})]) \\ &\stackrel{T_1(n) \leq n}{\geq} \Delta\left(n - n\mathbb{P}\left(T_1(n) > \frac{n}{2}\right) - \frac{n}{2}\left(\mathbb{P}\left(T_1(n) \leq \frac{n}{2}\right)\right)\right) \\ &= \frac{n\Delta}{2} \mathbb{P}\left(T_1(n) \leq \frac{n}{2}\right) \end{aligned}$$

Estimating $R_n(\pi, v_{\mu'})$ with the complementary probability is simpler as the bandit model was chosen in a way to produce regret gaps Δ when playing arm 1:

$$R_n(\pi, v_{\mu'}) = \sum_i \mathbb{E}'[T_i(n)] \Delta'_i \geq \mathbb{E}'[T_1(n)] \Delta \geq \mathbb{E}'[T_1(n) \mathbf{1}_{T_1(n) > \frac{n}{2}}] \Delta \geq \frac{n\Delta}{2} \mathbb{P}'\left(T_1(n) > \frac{n}{2}\right).$$

Combining both and applying Bretagnolle-Huber with $A = \{T_1(n) \leq \frac{n}{2}\}$ yields

$$\begin{aligned}
R_n(\pi, v_\mu) + R_n(\pi, v_{\mu'}) &\geq \frac{n\Delta}{2} \left(\mathbb{P} \left(T_1(n) \leq \frac{n}{2} \right) + \mathbb{P}' \left(T_1(n) > \frac{n}{2} \right) \right) \\
&\stackrel{\text{Thm 1.4.3}}{\geq} \frac{n\Delta}{4} \exp(-D(\mathbb{P}, \mathbb{P}')) \\
&\stackrel{\text{Thm 1.4.5}}{=} \frac{n\Delta}{4} \exp \left(- \sum_i \mathbb{E}_\mu [T_i(n)] \underbrace{D(\mathcal{N}(\mu_i, 1), \mathcal{N}(\mu'_i, 1))}_{=0, \text{ except } i=a} \right) \\
&= \frac{n\Delta}{4} \exp \left(- \mathbb{E}_\mu [T_a(n)] \frac{(2\Delta)^2}{2} \right) \\
&\geq \frac{n\Delta}{4} \exp \left(- \frac{2n\Delta^2}{k-1} \right).
\end{aligned}$$

⁶ The final inequality holds because a is the arm with smallest expected number of play, hence, $(k-1)\mathbb{E}[T_a(n)] \leq \sum_i \mathbb{E}[T_i(n)] = n$. Minimising the righthand side over Δ (or just choosing $\Delta = \sqrt{\frac{k-1}{4n}}$ yields

$$R_n(\pi, v_\mu) + R_n(\pi, v_{\mu'}) \geq \sqrt{(k-1)n} \frac{1}{8} e^{-\frac{1}{2}} \geq \frac{2}{27} \sqrt{(k-1)n}.$$

□

1.4.3 Asymptotic lower bound (model-dependent)

We next turn towards model-dependent asymptotic bounds. What we try to do is to find lower bounds (asymptotically in n) for a policy for a fixed model from a certain class of models. For instance a lower bound of the UCB algorithm for a given sub-Gaussian bandit model. A general solution is problematic for the following reason. If a policy only plays a fixed arm then the regret is either 0 or grows linearly depending on that arm to be optimal or not. To rule out such extremal cases we only consider policies that certainly have sublinear growth in the following sense:



Definition 1.4.7. A policy π is called consistent over a class of bandits \mathcal{E} if for all $\mu \in \mathcal{E}$ and all $p > 0$ it holds that

$$\lim_{n \rightarrow \infty} \frac{R_n(\pi, \nu)}{n^p} = 0.$$

An example of a consistent policy over all sub-Gaussian bandits is the policy obtained from the UCB algorithm. Here is the famous Lai-Robbins theorem on asymptotic lower bounds⁷:



Theorem 1.4.8. (Lai-Robbins lower bound)

If π is a consistent policy for a set $\mathcal{E} = \mathcal{M}_1 \times \dots \times \mathcal{M}_k$ of bandit models, then the regret for all bandits $\nu = (P_1, \dots, P_k)$ from \mathcal{E} satisfies

$$\liminf_{n \rightarrow \infty} \frac{R_n(\pi)}{\log(n)} \geq c^*(\nu, \mathcal{E}) := \sum_{a \in \mathcal{A}} \frac{\Delta_a}{d_a},$$

where $d_a = \inf_{P' \in \mathcal{M}_a} \{D(P_a, P') : Q_{P'} > Q_*\}$ and Q_P is the expectation of P .

⁶warum nicht k statt $k-1$?

⁷T.L Lai, H. Robbins: "Asymptotically efficient adaptive allocation rules", Advances in Applied Mathematics, 1985, pp. 4-22

Lai Robbins actually assumed all \mathcal{M}_a to be equal and an additional continuity property on the relative entropies of that class (continuity of the KL-distance as a function of the expectations). In that case the lower bound looks simpler because $d_a = D(P_a, P_*)$, where P_* is the law of the optimal arm:

$$\liminf_{n \rightarrow \infty} \frac{R_n(\pi)}{\log(n)} \geq \sum_{a \in \mathcal{A}} \frac{\Delta_a}{D(P_a, P_*)}$$

Here is an example in which a directly computation shows the same. If for instance $\mathcal{M}_a = \{\mathcal{N}(\mu, \sigma^2) : \mu \in \mathbb{R}\}$ for all a then $d_a = \frac{(Q_a - Q_*)^2}{\sigma^2} = \frac{\Delta_a^2}{\sigma^2} = D(P_a, P_*)$. Thus, if all arms are Gaussian with variance σ^2 , then the lower bound for UCB is

$$\liminf_{n \rightarrow \infty} \frac{R_n}{\log(n)} \geq \sum_{a \neq a_*} \frac{\sigma^2}{\Delta_a}.$$

Comparing with the UCB upper bound of Theorem 1.3.8, $\sigma = 1$ this is pretty close. In fact, it is not too difficult to modify the UCB bonus so that the upper bound matches the lower bound, thus, the algorithm is asymptotically optimal. The term d_a in a way also measures the reward gap but rather the KL-distance from the optimal arm, not the mean distance.

Proof. Fix a consistent policy and suppose $\nu = (P_1, \dots, P_k) \in \mathcal{E}$. Recalling the regret decomposition it is enough to prove that

$$\liminf_{n \rightarrow \infty} \frac{\mathbb{E}_\pi[T_a(n)]}{\log(n)} \geq \frac{1}{d_a}$$

holds for all suboptimal arm. Fix such an arm a and fix an $\varepsilon > 0$ and define ν' as the bandit model in which the arm distribution P_a is replaced by a distribution $P'_a \in \mathcal{M}_a$ such that $D(P_a, P'_a) \leq d_a + \varepsilon$ and $Q_{P'_a} > Q_*$. By the definition of the infimum such a bandit model exists in the class \mathcal{E} . Since only one arm differs the entropic decomposition for bandits yields $D(\mathbb{P}, \mathbb{P}') \leq \mathbb{E}[T_a(n)](d_a + \varepsilon)$. We next follow an argument that is inspired by the proof of Theorem 1.4.6. First note that

$$R_n(\pi) + R'_n(\pi) \geq \frac{n}{2} \left(\mathbb{P} \left(T_a(n) > \frac{n}{2} \right) \Delta_a + \mathbb{P} \left(T_a(n) \leq \frac{n}{2} \right) (Q'_a - Q_*) \right).$$

First, from the regret decomposition it clearly holds that

$$R_n(\pi) \geq \Delta_a \mathbb{E}[T_a(n)] \geq \Delta_a \mathbb{E}[T_a(n) \mathbf{1}_{T_a(n) > n/2}] \geq \Delta_a \frac{n}{2} \mathbb{P} \left(T_a(n) > \frac{n}{2} \right).$$

Secondly, note that for ν' the best arm is a as the new arm distribution has a mean strictly larger than Q_* . Hence, if arm a is played less than $\frac{n}{2}$ times, then all other arms in are played at least $\frac{n}{2}$ times in total. Thus,

$$\begin{aligned} R'_n(\pi) &= \sum_{i \neq a} \Delta'_i \mathbb{E}'[T_i(n)] \\ &\geq \sum_{i \neq a} (Q'_a - Q_*) \mathbb{E}'[T_i(n)] \\ &\geq \sum_{i \neq a} (Q'_a - Q_*) \mathbb{E}'[T_i(n) \mathbf{1}_{T_a(n) < n/2}] \\ &\geq \frac{n}{2} (Q'_a - Q_*) \mathbb{P}' \left(T_a(n) \leq \frac{n}{2} \right). \end{aligned}$$

Chosing $A = \{T_a(n) > n/2\}$ we continue with Bretagnolle-Huber:

$$\begin{aligned} R_n(\pi) + R'_n(\pi) &\geq \frac{n}{2} \left(\mathbb{P} \left(T_a(n) > \frac{n}{2} \right) \Delta_a + \mathbb{P}' \left(T_a(n) \leq \frac{n}{2} \right) (Q'_a - Q_*) \right) \\ &\geq \frac{n}{2} \min\{\Delta_a, Q'_a - Q_*\} (\mathbb{P}(A) + \mathbb{P}'(A^c)) \\ &\geq \frac{n}{4} \min\{\Delta_a, Q'_a - Q_*\} \exp(-\mathbb{E}[T_a(n)](d_a + \varepsilon)). \end{aligned}$$

Rearranging and using the consistency property of the policy yields the claim:

$$\begin{aligned} \liminf_{n \rightarrow \infty} \frac{\mathbb{E}[T_a(n)]}{\log(n)} &\geq \frac{1}{d_a + \varepsilon} \liminf_{n \rightarrow \infty} \frac{\log\left(\frac{n \min\{\Delta_a, Q'_a - Q_*\}}{4(R_n(\pi) + R'_n(\pi))}\right)}{\log(n)} \\ &\geq \frac{1}{d_a + \varepsilon} \left(1 - \limsup_{n \rightarrow \infty} \frac{\log(R_n(\pi) + R'_n(\pi))}{\log(n)}\right) = \frac{1}{d_a + \varepsilon}. \end{aligned}$$

For the last equality we used that, for n large enough, $R_n(\pi) \leq Cn^p$ and $R'_n(\pi) \leq C'n^p$ so that

$$0 \leq \limsup_{n \rightarrow \infty} \frac{\log(R_n(\pi) + R'_n(\pi))}{\log(n)} \leq \limsup_{n \rightarrow \infty} \frac{\log(Cn^p + C'n^p)}{\log(n)} = \limsup_{n \rightarrow \infty} \frac{\log(C) + p \log(n)}{\log(n)} = p.$$

Since p can be chosen arbitrarily close to 0 the limit superior exists and is equal to 0. \square

Part I

Tabular Reinforcement Learning

Chapter 2

Basics: Stochastic Control and Dynamic Programming Methods

2.1 Markov decision process formalism

After the first introductory chapter on bandits we will now turn towards the main topic of this book: optimal decision making in which decisions also influence the system (in contrast to stochastic bandits). This basic chapter covers the following topics:

- Introduce the basic setup for decision making in complex systems under uncertainty, we will use so-called Markov decision processes (MDP).
- Understand optimal decision policies and their relations to Bellman optimality and expectation equations.
- Understand how to turn the theoretical results into algorithms, so-called value and policy iteration algorithms.

All topics covered in this chapter are standard, they can e.g. also be found in the comprehensive overview of Putterman¹. It should be noted that the Q -functions appearing in this chapter are not commonly used in stochastic optimal control. They are, however, far more relevant to the reinforcement learning approaches developed in the next chapters.

2.1.1 A quick dive into Markov (reward) chains

Before starting with Markov decision processes let us briefly recall some facts on Markov chains. A finite-state Markov chain is a discrete-time stochastic process $(S_t)_{t \in \mathbb{N}_0}$ with values in some finite set \mathcal{S} on a probability space $(\Omega, \mathcal{F}, \mathbb{P})$ that satisfies the Markov property:

$$\mathbb{P}(S_{t+1} = s_{t+1} \mid S_0 = s_0, \dots, S_t = s_t) = \mathbb{P}(S_{t+1} = s_{t+1} \mid S_t = s_t)$$

for all $t \in \mathbb{N}_0$ and $s_0, \dots, s_{t+1} \in \mathcal{S}$. In words, the transition probability from a state to another does not depend on the past. The most important special case is that of a time-homogeneous Markov chain for which transition probabilities do not change over time, i.e.

$$\mathbb{P}(S_{t+1} = s' \mid S_t = s) = \mathbb{P}(S_1 = s' \mid S_0 = s)$$

for all $t \in \mathbb{N}_0$ and $s, s' \in \mathcal{S}$. Time-homogeneous Markov chains are closely related to stochastic matrices

¹M. Putterman: Markov decision processes: discrete stochastic dynamic programming, Wiley

**Definition 2.1.1. (Stochastic Matrix)**

A matrix $P = (p_{s,s'})_{s,s' \in \mathcal{S}} \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$ is called a stochastic matrix if all its entries are non-negative and each row sums to one, i.e.,

$$p_{s,s'} \geq 0 \quad \text{for all } s, s' \in \mathcal{S}, \quad \text{and} \quad \sum_{s' \in \mathcal{S}} p_{s,s'} = 1 \quad \text{for all } s \in \mathcal{S}.$$

Given an initial distribution μ on \mathcal{S} and a stochastic matrix P , a time-homogeneous Markov chain on some probability space $(\Omega, \mathcal{F}, \mathbb{P})$ with initial distribution μ and transition matrix P is a stochastic process that satisfies the basic path probabilities

$$\mathbb{P}(S_0 = s_0, \dots, S_t = s_t) = \mu(s_0) p_{s_0, s_1} \cdots p_{s_{t-1}, s_t}. \quad (2.1)$$

Using (2.1), one can easily show that the distribution of a time-homogeneous Markov chain $(S_t)_{t \in \mathbb{N}_0}$ is uniquely determined by the initial distribution μ and the transition matrix P .

Alternatively, and this is how Markov chains are generalised to more than finitely many states, the transition matrix can be interpreted as a Markov kernel on $\mathcal{S} \times \mathcal{S}$, that is, a family of probability measures $P(\cdot, s)$ on \mathcal{S} for all $s \in \mathcal{S}$. Then, the path probabilities can be written as

$$\mathbb{P}(S_0 = s_0, \dots, S_t = s_t) = \mu(\{s_0\}) P(\{s_1\}, s_0) \cdots P(\{s_t\}, s_{t-1}).$$

Many probabilities can be computed from the basic path formula (2.1), for instance:



Lemma 2.1.2. Let $(S_t)_{t \in \mathbb{N}_0}$ be a time-homogeneous Markov chain with initial distribution μ and transition matrix P . Then,

$$\mathbb{P}(S_{t+1} = s_{t+1}, \dots, S_{t+k} = s_{t+k} \mid S_t = s_t) = p_{s_t, s_{t+1}} \cdots p_{s_{t+k-1}, s_{t+k}}$$

for every $t \in \mathbb{N}_0$, $k \in \mathbb{N}$ and $s_t, \dots, s_{t+k} \in \mathcal{S}$ with $\mathbb{P}(S_t = s_t) > 0$.

Proof. The computation tricks are always the same: spell out the conditional probability explicitly, write the events of interest as disjoint union of simple paths, plug-in the path probability formula, and finally cancel terms in the resulting products.

In this case, we can write

$$\mathbb{P}(S_{t+1} = s_{t+1}, \dots, S_{t+k} = s_{t+k} \mid S_t = s_t) = \frac{\mathbb{P}(S_t = s_t, S_{t+1} = s_{t+1}, \dots, S_{t+k} = s_{t+k})}{\mathbb{P}(S_t = s_t)}.$$

Both numerator and denominator can be computed by taking disjoint unions of paths which probabilities can be expressed via (2.1). Thus,

$$\begin{aligned} \mathbb{P}(S_t = s_t) &= \sum_{s_0 \in \mathcal{S}} \cdots \sum_{s_{t-1} \in \mathcal{S}} \mathbb{P}(S_0 = s_0, \dots, S_t = s_t) \\ &= \sum_{s_0 \in \mathcal{S}} \cdots \sum_{s_{t-1} \in \mathcal{S}} \mu(s_0) p_{s_0, s_1} \cdots p_{s_{t-1}, s_t}. \end{aligned} \quad (2.2)$$

and

$$\begin{aligned} \mathbb{P}(S_t = s_t, S_{t+1} = s_{t+1}, \dots, S_{t+k} = s_{t+k}) &= \sum_{s_0 \in \mathcal{S}} \cdots \sum_{s_{t-1} \in \mathcal{S}} \mathbb{P}(S_0 = s_0, \dots, S_{t+k} = s_{t+k}) \\ &= \sum_{s_0 \in \mathcal{S}} \cdots \sum_{s_{t-1} \in \mathcal{S}} \mu(s_0) p_{s_0, s_1} \cdots p_{s_{t+k-1}, s_{t+k}}. \end{aligned} \quad (2.3)$$

Dividing (2.3) by (2.2) yields the statement. \square

Another way of expressing the Markov property is as follows: If $\mathbb{P}(S_n = s) > 0$ and $\tilde{\mathbb{P}} := \mathbb{P}(\cdot \mid S_n = s)$, then on $(\Omega, \mathcal{F}, \tilde{\mathbb{P}})$ the shifted process $(\tilde{S}_t)_{t \in \mathbb{N}} := (S_{t+n})_{t \in \mathbb{N}}$ is again a Markov chain with transition matrix P but started from the dirac distribution δ_s . To train yourself in using the path probabilities please check the next claim:



Check that \tilde{S} is indeed a Markov chain on $(\Omega, \mathcal{F}, \tilde{\mathbb{P}})$ with the same transition probabilities as S . Hint: Compute the path probabilities.

A Markov chain can be visualised as a graph whose nodes represent the possible states, and whose directed edges indicate the allowed transitions, each labeled with the corresponding transition probability. This idea is illustrated in Figure 2.1.

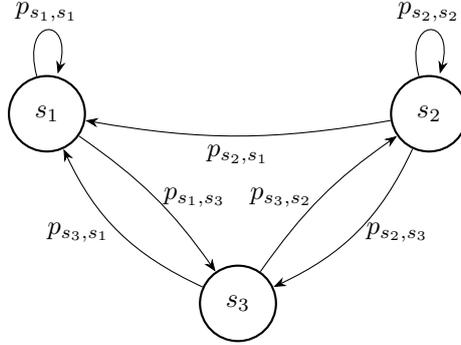


Figure 2.1: Illustration of a Markov chain with 3 states, i.e. $\mathcal{S} = \{s_1, s_2, s_3\}$. We include a directed edge from s to s' iff $p_{s,s'} \neq 0$. In this example, we have set $p_{s_1,s_2} = p_{s_3,s_3} = 0$.

I've set R_1 to be the reward that is sampled together with S_1 . I think, for MRP it would be more natural, although for MDPs it makes sense to call it R_0 , since it is the reward that corresponds to taking action A_0 in state S_0 .

A Markov chain can be seen as a random process that can be observed from the outside, for instance a game that changes states over time. The concept of a Markov reward process is less well-known. A Markov reward process is a Markov chain with an additional \mathcal{R} -valued stochastic process $(R_t)_{t \in \mathbb{N}}$ that is sampled together with the next state. We assume that \mathcal{R} is a finite set equipped with the power set σ -algebra $\bar{\mathcal{R}} = \{A : A \subset \mathcal{R}\}$. Typically, \mathcal{R} is a subset of the real numbers. The next state and reward are sampled according to a transition/reward-kernel p on $(\mathcal{R} \times \mathcal{S}) \times \mathcal{S}$, i.e. $p(r, s'; s)$ denotes the probability to transition the chain from s to s' and obtain reward r . More formally, the defining property of a Markov reward process is

$$\mathbb{P}_\mu(S_{t+1} = s_{t+1}, R_{t+1} = r_{t+1} | S_0 = s_0, R_1 = r_1, \dots, R_t = r_t, S_t = s_t) = p(r_{t+1}, s_{t+1}; s_t), \quad (2.4)$$

for all $t \in \mathbb{N}_0$, $s_0, \dots, s_{t+1} \in \mathcal{S}$ and $r_1, \dots, r_{t+1} \in \mathcal{R}$. Here, the subscript μ refers to the initial distribution of S_0 .

If we think of a Markov chain as describing a game then the rewards might be direct consequences of its rules. As an example, we set $R_t = 1$ if and only if the player that we observe has scored a goal and $R_t = -1$ if the opponent scored a goal. The path probabilities are given by

$$\mathbb{P}_\mu(S_0 = s_0, R_1 = r_1, \dots, R_n = r_n, S_n = s_n) = \mu(s_0) \cdot p(r_1, s_1; s_0) \cdot \dots \cdot p(r_n, s_n; s_{n-1}). \quad (2.5)$$

Another way of expressing the Markov reward property is as follows: If $\mathbb{P}(S_n = s) > 0$ and $\tilde{\mathbb{P}} := \mathbb{P}(\cdot | S_n = s)$ for an $n \in \mathbb{N}$ and $s \in \mathcal{S}$, then on $(\Omega, \mathcal{F}, \tilde{\mathbb{P}})$ the shifted process $(\tilde{S}_t, \tilde{R}_{t+1})_{t \in \mathbb{N}} := (S_{t+n}, R_{t+n+1})_{t \in \mathbb{N}}$ is again a Markov reward chain started from s' . To train yourself in using the path probabilities please check the next claim:



Compute the path probabilities in (2.5) and check that $(\tilde{S}_t, \tilde{R}_{t+1})_{t \in \mathbb{N}_0}$ is indeed a Markov reward process on $(\Omega, \mathcal{F}, \tilde{\mathbb{P}})$ that also satisfies (2.4). Moreover, show that $(S_{t+1}, R_{t+1})_{t \in \mathbb{N}_0}$ is a time-homogeneous Markov chain.

Lastly, we show a very important property that readily follows from the formula for the path probabilities (2.5). If $\mu = \delta_s$ for an $s \in \mathcal{S}$ we write $\mathbb{E}_\mu = \mathbb{E}_s$ and $\mathbb{P}_\mu = \mathbb{P}_s$.

**Lemma 2.1.3. (Functional Markov property)**

Let $t \in \mathbb{N}$ and $s, s' \in \mathcal{S}$ such that $\mathbb{P}_s(S_n = s') > 0$. Then,

$$\mathbb{P}_s((R_{n+t})_{t \in \mathbb{N}} \in A | S_t = s') = \mathbb{P}_{s'}((R_t)_{t \in \mathbb{N}} \in A) \quad (2.6)$$

for all $A \in (\bar{\mathcal{R}})^{\otimes \mathbb{N}}$ and

$$\mathbb{E}_s[f(R_{t+1}, R_{t+2}, \dots) | S_t = s'] = \mathbb{E}_{s'}[f(R_1, R_2, \dots)] \quad (2.7)$$

for all measurable functions $f : \mathcal{R}^{\mathbb{N}} \rightarrow [0, \infty)$.

Proof. One only has to show that

$$\mathbb{P}_s(R_{t+1} = r_1, \dots, R_N = r_N | S_t = s') = \mathbb{P}_{s'}(R_1 = r_1, \dots, R_N = r_N)$$

for all $t, N \in \mathbb{N}$ and $r_1, \dots, r_N \in \mathcal{R}$. Then, (2.6) follows from a Dynkin argument and (2.7) follows from a monotone class argument since

$$\{\{r_1\} \times \dots \times \{r_N\} \times \mathcal{R}^{\mathbb{N}} : N \in \mathbb{N}, r_1, \dots, r_N \in \mathcal{R}\}$$

is a \cap -stable generator of $(\bar{\mathcal{R}})^{\mathbb{N}}$ and

$$\mathbb{P}_{s'}(R_1 = r_1, \dots, R_N = r_N) = \mathbb{E}_{s'}[\mathbf{1}_{R_1=r_1} \cdot \dots \cdot \mathbf{1}_{R_N=r_N}].$$

This, however, follows immediately from computing the basic path probabilities. \square

2.1.2 Markov decision processes

The situation of Markov decision processes (MDPs) is slightly more complicated. For MDPs we do not simply observe a game but take the role of a participant. We observe the Markov chain describing the match but can influence the transitions by taking actions. As an example, by substituting as a coach an additional attack player we increase the probability of players scoring goals (positive reward) but decrease the probabilities of players winning duels which leads to larger probabilities to receive a goal (negative reward). The target will be to find a strategy (called policy) that maximises the expected reward of the game. Finding such a strategy, this is what reinforcement learning is all about!

Is this paragraph still correct? It seems that we immediately consider the finite situation here.

Unfortunately, the formal notion in Markov decision processes is more sophisticated. In the beginning, we keep the definitions quite general to emphasise that everything we will later prove in the discrete setting could be generalised by replacing vectors and matrices by Markov kernels. For this book the intention is to keep the level of sophistication high enough to observe most difficulties but stay simple enough to not disturb too much the understanding of concepts with technicalities. We will write down definitions and the existence theorem in a more general setting but then work only with finite Markov decision processes which makes our lives much easier without losing most features of interest.

any reason we have to in-

**Definition 2.1.4. (Markov decision model)**

A (discrete) Markov decision model is a tuple $(\mathcal{S}, \mathcal{A}, \mathcal{R}, p)$ consisting of the following ingredients:

- (i) A finite state-space \mathcal{S} .
- (ii) Finite action-spaces \mathcal{A}_s for all states $s \in \mathcal{S}$. The entire action space is defined as $\mathcal{A} = \bigcup_{s \in \mathcal{S}} \mathcal{A}_s$, i.e. \mathcal{A} contains all possible actions.
- (iii) A finite reward set $\mathcal{R} \subseteq \mathbb{R}$ including 0.



(iv) A transition-function

$$p : (\mathcal{S} \times \mathcal{R}) \times (\mathcal{S} \times \mathcal{A}) \rightarrow [0, 1]$$

written as $p(s', r; s, a)$ and to be read "probability to end up in state s' with reward r when action a was taken in state s . For state-action pairs $(s, a) \in \mathcal{S} \times \mathcal{A}$ the values $p(s', r; s, a)$ are probability weights indexed by $\mathcal{S} \times \mathcal{R}$, i.e.

$$\sum_{s' \in \mathcal{S}} \sum_{r \in \mathcal{R}} p(s', r; s, a) = 1 \quad \text{for all } s \in \mathcal{S} \text{ and } a \in \mathcal{A}.$$

We use the abbreviations $p(r; s, a) = \sum_{s' \in \mathcal{S}} p(s', r; s, a)$ and $p(s'; s, a) = \sum_{r \in \mathcal{R}} p(s', r; s, a)$.

In RL literature you will almost exclusively find the notion $p(\cdot | s, a)$. In this course the notion " $|$ " will be used exclusively for conditional probabilities, and for conditional probability it must be ensured that there is no conditioning on zero-probability events.



For notational simplicity, we set the domain of the transition-function to be $(\mathcal{S} \times \mathcal{R}) \times (\mathcal{S} \times \mathcal{A})$. Since, in a state $s \in \mathcal{S}$ we only allow actions from the action space \mathcal{A}_s it suffices to define $p(s', r; s, a)$ for all $s, s' \in \mathcal{S}$, $r \in \mathcal{R}$ and $a \in \mathcal{A}_s$. All derived properties and statements will be independent on the values of the transition-function for actions $a \notin \mathcal{A}_s$.

The definition is simplified as far as possible to make it accessible for a lecture course but still be general enough to allow for good examples. Assuming finite sets makes notation and proofs easier, there is a number of difficulties appearing Markov decision theory for countable state/action-spaces. Next, a general notion for non-discrete state- and/or action-spaces and/or reward distributions uses Markov kernels and a tedious notion involving the σ -algebras for state-, action-, and reward-spaces. The interested reader may always replace the probability weights by kernels and appearing sums by integrals. The main ideas of reinforcement learning become fully visible in the setting of finite decision problems, so for didactic reasons we stick to the simple setting for these lecture notes.



The key ingredient of the definition is the transition-function p that describes the transition of the environment when an action is played. According to the definition the reward can depend on the current state-action pair **and** the next state, but in most examples the reward and the next state will be independent (see Example 2.1.9 below).

I don't know whether this is very helpful. It makes no sense when talking about variables.

A Markov decision model does not fully describe a process that runs from state to state and returns rewards. Only the transitions $(s, a) \mapsto (s', r)$ is described but not how the next action a' is chosen. For that purpose an additional ingredient is needed, the policy that can be chosen by the actor.



Definition 2.1.5. (Time-dependent policy)

For a Markov decision model $(\mathcal{S}, \mathcal{A}, \mathcal{R}, p)$ a policy $(\pi_t)_{t \in \mathbb{N}_0}$ consists of

- initial probability vectors π_0 on \mathcal{A} indexed by \mathcal{S} , i.e. $\pi_0 : \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ satisfying

$$\sum_{a \in \mathcal{A}_s} \pi_0(a; s) = 1, \quad (2.8)$$



for all $s \in \mathcal{S}$,

- a sequence $\pi = (\pi_t)_{t \in \mathbb{N}}$ of probability vectors on \mathcal{A} indexed by $(\mathcal{S} \times \mathcal{A})^{t-1} \times \mathcal{S}$ such that

$$\sum_{a \in \mathcal{A}_{s_t}} \pi_t(a; s_0, a_0, \dots, s_{t-1}, a_{t-1}, s_t) = 1 \quad (2.9)$$

for all $(s_0, a_0, \dots, s_{t-1}, a_{t-1}, s_t) \in (\mathcal{S} \times \mathcal{A})^{t-1} \times \mathcal{S}$.

The set of all policies is denoted by Π .

A policy governs the choices of actions given the current state and all previous states-action pairs (not the rewards) seen before. The condition (2.8) and (2.9) mean that only allowed actions from \mathcal{A}_s can be played in state s . Sometimes all \mathcal{A}_s are identical but in most examples they are not. As an example, playing a game like chess the allowed actions clearly depend on the current state of the game.

We define a subclass of the



Definition 2.1.6. (Stationary policy)

A policy is called stationary if there exists a probability vector φ on \mathcal{A} such that

$$\pi_t(\cdot; s_0, a_0, \dots, s_t) = \varphi(\cdot; s_t), \quad \forall (s_0, a_0, \dots, s_t) \in (\mathcal{S} \times \mathcal{A})^{t-1} \times \mathcal{S}.$$

With the definition of a policy we can now define a stochastic process on $\mathcal{S} \times \mathcal{A} \times \mathcal{R}$ whose dynamics are entirely defined by the function p and policy π (plus an initial probability distribution over which state the stochastic process will start in). Let us formalize this stochastic process and prove its existence.



Theorem 2.1.7. (Existence of (discrete) MDPs)

Let $(\mathcal{S}, \mathcal{A}, \mathcal{R}, p)$ be a (discrete) Markov decision model, π a policy, μ a probability measure on \mathcal{S} . Then there exists a probability space $(\Omega, \mathcal{F}, \mathbb{P}_\mu^\pi)$ carrying a stochastic process $(S_t, A_t, R_t)_{t \in \mathbb{N}_0}$ with values in $\mathcal{S} \times \mathcal{A} \times \mathcal{R}$ on $(\Omega, \mathcal{F}, \mathbb{P}_\mu^\pi)$ such that

- (i) **Initialisation:**

$$\mathbb{P}_\mu^\pi(S_0 = s_0, A_0 = a_0) = \mu(s_0)\pi_0(a_0; s_0)$$

- (ii) **Choice of actions:**

$$\mathbb{P}_\mu^\pi(A_t = a_t | S_0 = s_0, A_0 = a_0, \dots, S_t = s_t) = \pi_t(a_t; s_0, a_0, \dots, s_t)$$

- (iii) **Environment transition:**

$$\mathbb{P}_\mu^\pi(S_{t+1} = s_{t+1}, R_t = r_t | S_0 = s_0, A_0 = a_0, \dots, S_t = s_t, A_t = a_t) = p(s_{t+1}, r_t; s_t, a_t)$$

To increase readability we will oftne skipp π and μ from the notation.

² In words the mechanism goes as follows. In a state s an action is sampled according to the policy which is allowed to refer to the entire past of states and actions (not to the rewards!). The current and past state-action pairs (s, a) are used to sample the next state s' and the reward. Note: The reward is allowed to depend both on the current state-action pair (s, a) and the future state s' ! Typically, the reward will only depend on (s, a) , not on s' , but some examples require this greater generality.

²todo: Schreibe lieber allgemeinen Beweis fuer Markov reward chains auf und zitiere hier. Vermutlich besser zugaenglich

Proof. We only give the proof in the discrete setting to save quite a bit of time. The proof is essentially identical to the existence proof for Markov chains.



Measure for finite time-horizon $T < \infty$.

The probability space is written down explicitly as the set of trajectories, the σ -algebra is chosen to be the power set, a canonical measure is defined for every element of the probability space by an explicit definition, and the process using the identity mapping. Define the measurable space $(\Omega_T, \mathcal{F}_T)$ by $\Omega := (\mathcal{S} \times \mathcal{A} \times \mathcal{R})^T$, the trajectories of length T , i.e.

$$\omega = (s_0, a_0, r_0, \dots, s_T, a_T, r_T),$$

and \mathcal{F}_T as the powerset. As for Markov chains the probabilities for paths are written down explicitly (skipping brackets to increase readability):

$$\begin{aligned} & \mathbb{P}_T((s_0, a_0, r_0, \dots, s_T, a_T, r_T)) \\ & := \mu(s_0) \cdot \pi_0(a_0; s_0) \cdot \prod_{i=1}^T p(s_i, r_{i-1}; s_{i-1}, a_{i-1}) \cdot \pi_i(a_i; s_0, a_0, \dots, a_{i-1}, s_i) \cdot p(\mathcal{S} \times \{r_T\}; s_T, a_T). \end{aligned}$$

In words: An initial state is chosen by μ , the first action is sampled using π_0 , and the first reward is set to 0. For every further time-step the new state s_i and the reward r_i are determined by p , an action a_i is chosen according to π_i .



Show that defining \mathbb{P}_T on the singletons as above yields a probability measure.

The exercise is important to get a feeling for the path probabilities. You will have to sum over all possible paths, i.e. over $\sum_{a_0, r_0, s_0} \dots \sum_{a_T, r_T, s_T}$, then pull out the factors that are independent of the summands to simplify backwards using the kernel property to obtain (here for the summands corresponding to T)

$$\begin{aligned} & \sum_{a_T, s_T} p(s_T, r_T; s_{T-1}, a_{T-1}) \cdot \pi_T(a_T; s_0, a_0, \dots, a_{T-1}, s_T) \\ & = \sum_{s_T} p(s_T, r_T; s_{T-1}, a_{T-1}) \sum_{a_T} \pi_T(a_T; s_0, a_0, \dots, a_{T-1}, s_T) = 1. \end{aligned}$$

Summing over all trajectories shows easily that \mathbb{P}_T is a probability measure on the paths.



Measure for infinite time-horizon.

The same construction cannot work for $T = \infty$ as the infinite products would be zero. Instead, Kolmogorov's extension theorem is employed. Define the measurable space (Ω, \mathcal{F}) by $\Omega := (\mathcal{S} \times \mathcal{R} \times \mathcal{A})^\infty$, the trajectories of infinite length, and the corresponding σ -algebra $\mathcal{F} := (\bar{\mathcal{S}} \otimes \bar{\mathcal{A}} \otimes \bar{\mathcal{R}})^{\otimes \infty}$, the cylinder σ -algebra. The elements of Ω can be written as infinite sequences of the form

$$\omega = (s_0, r_0, a_0, \dots).$$

Consider the projections $\pi_T^{T+1} : \Omega^{T+1} \rightarrow \Omega^T, \omega \mapsto \omega|_T$ and $\pi_T : \Omega \rightarrow \Omega^T, \omega \mapsto \omega|_T$. The projections simply remove the triple corresponding to time $T + 1$ from the sequence. We now show the consistency property $\mathbb{P}_T = \mathbb{P}_{T+1} \circ (\pi_T^{T+1})^{-1}$ for any $T \in \mathbb{N}$. If the consistency can be proved, then Kolmogorov's extension theorem gives a unique probability measure \mathbb{P} on (Ω, \mathcal{F}) such that $\mathbb{P}_T = \mathbb{P} \circ (\pi_T)^{-1}$. The consistency property simply follows from the product structure in the measures \mathbb{P}_T defined above. Since the measures are discrete the equality can be checked

on all elements separately:

$$\begin{aligned}
& \mathbb{P}_{T+1}((\pi_T^{T+1})^{-1}(s_0, a_0, r_0, \dots, s_T, a_T, r_T)) \\
&= \mathbb{P}_{T+1}(\cup_{s \in \mathcal{S}} \cup_{a \in \mathcal{A}} \cup_{r \in \mathcal{R}} \{(s_0, a_0, r_0, \dots, s_T, a_T, r_T, s, a, r)\}) \\
&= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \sum_{r \in \mathcal{R}} \mathbb{P}_{T+1}(s_0, a_0, r_0, \dots, s_T, a_T, r_T, s, a, r) \\
&\stackrel{\text{linearity}}{=} \mu(s_0) \cdot \pi_0(a_0; s_0) \cdot \prod_{i=1}^T p(s_i, r_{i-1}; s_{i-1}, a_{i-1}) \cdot \pi_i(a_i; s_0, a_0, \dots, a_{i-1}, s_i) \\
&\quad \times \underbrace{\sum_{s \in \mathcal{S}} p(s, r_T; s_T, a_T)}_{p(\mathcal{S} \times \{r_T\}; S_T, a_T)} \cdot \underbrace{\sum_{a \in \mathcal{A}} \pi_{T+1}(a; s_0, a_0, \dots, a_T, s_{T+1})}_{=1} \underbrace{\sum_{r \in \mathcal{R}} p(\mathcal{S} \times \{r\}; s_T, a_T)}_{=1} \\
&= \mu(s_0) \cdot \pi_0(a_0; s_0) \cdot \prod_{i=1}^T p(s_i, r_{i-1}; s_{i-1}, a_{i-1}) \cdot \pi_i(a_i; s_0, a_0, \dots, a_T, s_{T+1}) \\
&\quad \cdot p(\mathcal{S} \times \{r_T\}; S_T, a_T) \\
&= \mathbb{P}_T(s_0, a_0, r_0, \dots, s_T, a_T, r_T).
\end{aligned}$$

Kolmogorov's extension theorem now yields a unique measure \mathbb{P} extending all \mathbb{P}_T . We set $\mathbb{P}_\mu^\pi := \mathbb{P}$ to indicate the dependency of μ and π .



Canonical construction $(\Omega, \mathcal{F}, \mathbb{P}_\mu^\pi, (S, A, R)_{t \in \mathbb{N}_0})$.

The measurable space (Ω, \mathcal{F}) has been defined above. On Ω we define $(S_t, A_t, R_t)(\omega) = (s_t, a_t, r_t)$ if ω takes the form (s_0, a_0, r_0, \dots) .



The properties claimed in the theorem hold.

The first two claims follow directly from the definition of the measure \mathbb{P}_μ^π and (S_0, A_0, R_0) . We check the claim for A and leave the other claims as an exercise. First note that putting no restrictions to the rewards in the path probabilities leads to setting $p(\{s'\} \times \mathcal{R}; s, a)$ in the state-reward transitions (summing over all possibilities). Using the extension property of \mathbb{P}_μ^π (there is no dependence later than t) yields

$$\begin{aligned}
& \mathbb{P}_\mu^\pi(S_0 = s_0, A_0 = a_0, \dots, A_{t-1} = a_{t-1}, S_t = s_t, A_t = a_t) \\
&= \mu(s_0) \cdot \pi_0(a_0; s_0) \cdot \prod_{i=1}^t p(\{s_i\} \times \mathcal{R}; s_{i-1}, a_{i-1}) \cdot \pi_i(a_i; s_0, a_0, \dots, a_{i-1}, s_i)
\end{aligned}$$

and

$$\begin{aligned}
& \mathbb{P}_\mu^\pi(S_0 = s_0, A_0 = a_0, \dots, A_{t-1} = a_{t-1}, S_t = s_t) \\
&= \mathbb{P}_\mu^\pi(S_0 = s_0, A_0 = a_0, \dots, A_{t-1} = a_{t-1}, S_t = s_t, A_t \in \mathcal{A}) \\
&= \mu(s_0) \cdot \pi_0(a_0; s_0) \cdot \prod_{i=1}^{t-1} p(\{s_i\} \times \mathcal{R}; s_{i-1}, a_{i-1}) \cdot \pi_i(a_i; s_0, a_0, \dots, a_{i-1}, s_i) \\
&\quad \times \underbrace{p(\{s_t\} \times \mathcal{R}; s_{t-1}, a_{t-1}) \cdot \sum_{a \in \mathcal{A}} \pi_t(a; s_0, a_0, \dots, s_{t-1}, a_{t-1}, s_t)}_{=1}.
\end{aligned}$$

Hence, using the definition of conditional probability the products cancel in the fraction to give

$$\mathbb{P}_\mu^\pi(A_t = a_t | S_0 = s_0, A_0 = a_0, \dots, A_{t-1} = a_{t-1}, S_t = s_t) = \pi_t(a_t; s_0, a_0, \dots, s_{t-1}, a_{t-1}, s_t)$$



Check the other two claimed conditional probability identities.

□



Definition 2.1.8. (Markov decision process (MDP))

Given a discrete Markov decision model $(\mathcal{S}, \mathcal{A}, \mathcal{R}, p)$, a policy π , and an initial distribution μ on S , the stochastic process (S_t, A_t) from the previous theorem is called discrete-time Markov decision process (MDP), (R_t) the corresponding reward process. To abbreviate, we write \mathbb{P}_s^π instead of $\mathbb{P}_{\delta_s}^\pi$. The super- and subscript are sometimes removed from \mathbb{E}_μ^π and \mathbb{P}_μ^π if there is no possible confusion about the initial distribution of S and the policy.

In the literature, a Markov decision model is mostly called an MDP directly and the terms model and process are not distinguished. We keep in mind that the Markov decision model just provides the prevailing instructions and together with any policy and initial distribution induces an MDP. However, for reasons of convenience, we will mostly use the term MDP interchangeable, if we do not want to emphasize the differences.



It is very important to remember the formula

$$\begin{aligned} & \mathbb{P}_\mu^\pi(S_0 = s_0, A_0 = a_0, R_0 = r_0, \dots, S_T = s_T, A_T = a_T, R_T = r_T) \\ &= \mu(s_0) \cdot \pi_0(a_0; s_0) \cdot \prod_{i=1}^T p(s_i, r_{i-1}; s_{i-1}, a_{i-1}) \cdot \pi_i(a_i; s_0, a_0, \dots, a_{i-1}, s_i) \quad (2.10) \\ & \cdot p(r_T; s_T, a_T), \end{aligned}$$

which gives the probability an MDP follows a given path $(s_0, a_0, r_0, \dots, s_T, a_T, r_T)$ of states, actions, and rewards. The formula will later explain the practical importance of the policy gradient theorem.

It is always very instructive to compare with the situation of an ordinary Markov chain that appears as a special case if $\mathcal{A} = \{a\}$ and $\mathcal{R} = \{0\}$. In that case then the process (S_t) is a Markov chain with transition matrix $p_{s_i, s_j} = p(s_j, 0; s_i, a)$ and the path probabilities immediately simplify to the Markov chain formula (2.1). No doubt, for MDPs the formula is more complicated but many probabilities can be computed similarly to the situation of a Markov chain. Playing with conditional probabilities it is instructive to check the following formulas (using a cancellation of the form $\sum_i (a \cdot b_i) / \sum_i b_i = a$):

$$\begin{aligned} & \mathbb{P}(S_t = s_t, A_t = a_t, R_t = r_t, \dots, S_T = s_T, A_T = a_T, R_T = r_T \mid S_{t-1} = s_{t-1}, A_{t-1} = a_{t-1}) \\ &= p(\{s_t\} \times \mathcal{R}; s_{t-1}, a_{t-1}) \cdot \pi_i(a_t; s_0, a_0, \dots, a_{t-1}, s_t) \\ & \times \prod_{i=t+1}^T p(s_i, r_{i-1}; s_{i-1}, a_{i-1}) \cdot \pi_i(a_i; s_0, a_0, \dots, a_{i-1}, s_i) \cdot p(\mathcal{S} \times \{r_T\}; s_T, a_T) \end{aligned}$$

and

$$\begin{aligned} & \mathbb{P}(S_t = s_t, A_t = a_t, R_t = r_t, \dots, S_T = s_T, A_T = a_T, R_T = r_T \mid S_{t-1} = s) \\ &= \sum_{a \in \mathcal{A}_s} \pi_{t-1}(a; s_0, a_0, \dots, a_{t-2}, s) \prod_{i=t}^T p(s_i, r_i; s_{i-1}, a_{i-1}) \cdot \pi_i(a_i; s_0, a_0, \dots, a_{i-1}, s_i), \end{aligned}$$

which are the analogues to the Markov chain formulas recalled above.

Here are two more special situations in which more simple Markov decision models can be abstracted in our general definition:

Example 2.1.9. There are many situations in which the reward and the transition both only depend on (s, a) but do not depend on each other. Suppose there are

- a kernel h on $\bar{\mathcal{S}} \times (\mathcal{S} \times \mathcal{A})$ for the transitions from (s, a) to s' ,
- a kernel q on $\bar{\mathcal{R}} \times (\mathcal{S} \times \mathcal{A})$ for the rewards r obtained from (s, a) .

Note that in the discrete setting measurability is always satisfied, kernel only means that $h(\cdot; s, a)$ and $q(\cdot; s, a)$ are discrete measures for all state-action pairs (s, a) . If both transitions are assumed to be independent then we can define the product kernel through

$$p(s', r; s, a) := h(s'; s, a) \cdot q(r; s, a) \quad (2.11)$$

and the corresponding Markov decision process samples next states/rewards independently. Now the Markov decision model in the sense of Definition 2.1.4 with Markov decision process from Theorem 2.1.7 has exactly the claimed transitions. Plugging-in p immediately gives

$$\begin{aligned} & \mathbb{P}(R_t = r_t \mid S_t = s_t, A_t = a_t, S_{t+1} = s_{t+1}) \\ \stackrel{\text{cond. prob.}}{=} & \frac{\mathbb{P}(S_{t+1} = s_{t+1}, R_t = r_t \mid S_t = s_t, A_t = a_t)}{\mathbb{P}(S_{t+1} = s_{t+1} \mid S_t = s_t, A_t = a_t)} \\ \stackrel{2.1.7}{=} & \frac{p(s_{t+1}, r_t; s_t, a_t)}{p(\{s_{t+1}\} \times \mathcal{R}; s_t, a_t)} \\ \stackrel{(2.11)}{=} & p(\mathcal{S} \times \{r_t\}; s_t, a_t) \\ \stackrel{2.1.7}{=} & \mathbb{P}(R_t = r_t \mid S_t = s_t, A_t = a_t), \end{aligned}$$

so that the proclaimed independence of the future reward from the future state indeed holds.

Example 2.1.10. In many examples the rewards are deterministic functions of state, actions, and next state but themselves do not influence the next state (say $r = R(s, a, s')$). An example appears in the ice vendor example below. In that situation the kernel p is simply

$$p(s', r; s, a) = h(s'; s, a) \cdot \mathbf{1}_{r=R(s, a, s')},$$

where h is the transition function from (s, a) to the next state s' . An even simpler situation that is very common is a relation $r = R(s, a)$ where the rewards are determined deterministically by the prior state-action pair.

The discussion is sometimes reversed as follows, with slight abuse of notation:



Definition 2.1.11. For a Markov decision model we define state-action-state probabilities and expected rewards as well as the state-action expected rewards as follows:

$$\begin{aligned} r(s, a) &:= \sum_{r \in \mathcal{R}} r p(r; s, a), \\ r(s, a, s') &:= \sum_{r \in \mathcal{R}} r \frac{p(s', r; s, a)}{p(s'; s, a)} \end{aligned}$$

for $s, s' \in \mathcal{S}$ and $a \in \mathcal{A}$ if the denominator is non-zero.

The notation will be used frequently, in reinforcement algorithms of the upcoming chapters. If for instance we are only interested in the state-action process, then it is more convenient to use $p(s'; s, a)$ instead of $p(\{s'\} \times \mathcal{R}; s, a)$ in computations with the path probabilities.



Use the formal definition of the stochastic process (S, A, R) to check that

$$\begin{aligned}
 p(s' ; s, a) &= \mathbb{P}(S_{t+1} = s' \mid S_t = s, A_t = a), \\
 p(r ; s, a) &= \mathbb{P}(R_t = r \mid S_t = s, A_t = a), \\
 r(s ; a) &= \mathbb{E}[R_t \mid S_t = s, A_t = a], \\
 r(s, a, s') &= \mathbb{E}[R_t \mid S_t = s, A_t = a, S_{t+1} = s']
 \end{aligned}$$

holds if the events in the condition have positive probability.

To get an idea about Markov decision processes we will discuss two examples of simple Markov decision models, grid world and ice vendor. Many games, such as grid world, stop once a particular state is reached. Such states are called terminating states.



Definition 2.1.12. A state $s \in S$ satisfying $p(s ; s, a) = 1$ for all $a \in \mathcal{A}_s$ is called a terminating state. The set of terminating states will be denoted by Δ and we will always assume $p(s, 0 ; s, a) = 1$ for all $s \in \Delta, a \in \mathcal{A}_s$.

In words: Once a terminal state is hit the MDP will stop moving and all future rewards are 0. To get a feeling for the definitions we will spell out a representative example in detail:

Example 2.1.13. (Grid world)

Suppose we have a robot that we want to teach to walk through our garden. The garden is set up as a small grid and at each time the robot can go up, down, left or right, however the robot cannot move through a wall (or additional obstacles). The aim is to move from the starting position S to the target position G . To make the problem a bit more difficult, there is a bomb B just before the goal G . If the robot lands in the bomb or in the target, the game is over. In

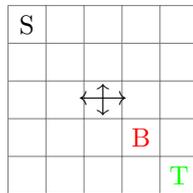


Figure 2.2: Grid world of size 5. The start is marked by S , the bomb by B , and the goal by G .

the following we will formulate this example as a Markov Decision Model. For this purpose, the individual objects from Definition 2.1.4 must be determined. The state space should contain all relevant information about the game. The garden can be represented by a two dimensional grid similar to a chess board. Thus, for a grid with side length $Z \in \mathbb{N}, Z > 2$, the state space is given by $S := \{(i, j), i, j \in \{0, 1, \dots, Z - 1\}\}$. Since the state space S is finite, we will use as σ -algebra the power set, i.e. $\mathcal{S} := \mathcal{P}(S)$.

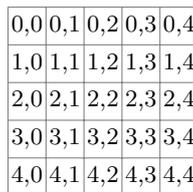


Figure 2.3: Representation of the state space of Example 2.1.13 for a grid of length 5.

The second component of the Markov decision model is the entire action space \mathcal{A} and the action spaces \mathcal{A}_s for states $s \in S$. In many cases it is easier to first define the set of all possible actions

\mathcal{A} and then exclude actions that are not possible for a certain state $s \in S$ to determine \mathcal{A}_s . In this example we want to use this method to define the action spaces of the individual states as well as the whole action space. An action describes the movement of the robot down, up, right or left. Often actions in a Markov decision problem are simply ‘counted through’. For our example, the action space could be given by the set $\{0, 1, 2, 3\}$, where action 0 describes the movement of the robot to the right, action 1 describes the movement of the robot downwards and so on:

$$\{\text{right, up, left, down}\} = \mathcal{A} = \{(1, 0), (0, 1), (-1, 0), (0, -1)\} \subseteq \mathbb{R}^2$$

The advantage of this definition is that by simply adding states and actions we get the new state of the robot. In the following we want to determine the action space for a state s using what is called masking (excluding actions from the set \mathcal{A} of all possible actions). Let $s = (i, j) \in S$ then the valid actions for the state are given by

$$\mathcal{A}_s = \mathcal{A}_{(i,j)} = \mathcal{A} \setminus \begin{cases} \{(-1, 0)\} & : \text{if } i = 0 \\ \{(1, 0)\} & : \text{if } i = Z - 1 \\ \emptyset & : \text{otherwise} \end{cases} \setminus \begin{cases} \{(0, -1)\} & : \text{if } j = 0 \\ \{(0, 1)\} & : \text{if } j = Z - 1 \\ \emptyset & : \text{otherwise} \end{cases}$$

The agent may play all actions, except the agent is on the edge of the playing field. The third components of the Markov decision model are the rewards. We introduce the rewards and motivate them afterwards. Let \mathcal{R} be given by $\{-10, -1, 0, 10\}$ and $\bar{\mathcal{R}} := \mathcal{P}(\mathcal{R})$. If the agent lands in the target, the reward is 10. If the agent runs into the bomb, the reward is -10 . We want to determine the rewards in such a way that desirable scenarios are rewarded and undesirable scenarios are punished. If the agent lands on another field, the agent should receive the reward -1 . In this way we want to avoid that the agent is not penalized for not choosing the fastest possible way to the destination. This will become clear when we deal with the optimization problem in Markov decision problems in the following chapters. Next, we define the fourth and most complicated component of the Markov decision process, the transition/reward-function p . To define this function, we will first define an auxiliary function $g : S \rightarrow \mathcal{R}$ which should return the associated reward for a state $s \in S$. The function is given by

$$g : S \rightarrow \mathcal{R}, s \mapsto \begin{cases} 10 & : s = T \\ -10 & : s = B \\ -1 & : \text{otherwise} \end{cases} .$$

In addition, we still need to define the set of terminal states Δ . In our example, the set of terminal states consists of the target and the bomb, $\Delta = \{(Z - 2, Z - 2), (Z - 1, Z - 1)\}$. The transition function is defined by

$$p(s', r; s, a) = \mathbf{1}_{s \in \Delta} \cdot \mathbf{1}_{s'=s} \cdot \mathbf{1}_{r=0} + \mathbf{1}_{s \notin \Delta} \cdot \mathbf{1}_{s'=s+a} \cdot \mathbf{1}_{r=g(s+a)}$$

for $(s, a) \in S \times \mathcal{A}_s$. The transition function looks a bit confusing and complicated at first sight, it is a compact way to write down all possible cases. Thus all components of a Markov decision model are defined. Since those will be used in algorithms discussed in future chapters let us compute the functions from Definition 2.1.11. For all state-action pairs the transition probability are either 0 or 1, the next states and the rewards are deterministic for a chosen action:

$$p(s'; s, a) = \mathbf{1}_{s \in \Delta} \cdot \mathbf{1}_{s'=s} + \mathbf{1}_{s \notin \Delta} \cdot \mathbf{1}_{s'=s+a} .$$

We next formulate a simple supply optimisation problem in the setting of Markov decision processes, an ice vendor who has to decide every day about the amount of produced/bought ice cream. In the course of these lecture notes it will be discussed how to optimise the production decision in order to maximise the profit.

Example 2.1.14. (Ice vendor)

The ice vendor owns an ice cream truck that can store up to $M \in \mathbb{N}$ ice creams. Every morning

the ice vendor can produce/buy a discrete amount $0, 1, 2, \dots, M$ of ice cream. Throughout the day, a random amount of ice cream is consumed, the demand can be higher than the stored amount of ice cream. In this simplified model, seasonality in demand is not considered. If less ice cream is consumed than the ice vendor has in stock, the ice cream must be stored overnight and can then be sold the next day. There are costs for storing ice cream overnight, thus, the ice vendor should carefully decide about the amount of ice cream produced/bought in the morning. In addition, there are costs for the production/purchase of ice cream. For the sale of ice cream the vendor receives a fixed price per ice cream scoop. For simplicity, we also assume that revenues and costs do not change over time (which is realistic for ice cream, the price is typically adjusted once per year). The trade-off in this problem is simple. How much ice should be ordered in the morning so that the vendor can maximize the revenue (serve demand) but keep the costs low. The ice vendor situation can be formulated as Markov decision process. The state space is given by the stored amount of ice: $\mathcal{S} = \{0, 1, \dots, M\}$. An action models the amount of ice cream produced/ordered in a morning. If there are already $s \in \mathcal{S}$ scoops in stock, the vendor can produce/order at most $M - s$ many scoops. Thus, $\mathcal{A}_s = \{0, 1, \dots, M - s\}$ and $\mathcal{A} = \{0, \dots, M\}$. For demand, we assume that it is independently identically distributed regardless of timing. Thus $\mathbb{P}(D_t = d) = p_d$ for $d \in \mathbb{N}$ holds for all time $t \in \mathbb{N}$. For simplicity, let us assume the demand can only take values in $\{0, 1, \dots, M\}$. In order to define the transition reward function and the reward, auxiliary functions are used to determine the revenues and costs. The selling price of a scoop of ice cream is given by a function $f : \mathcal{S} \rightarrow \mathbb{R}_+$, for instance $f(s) = c \cdot s$, where $c > 0$ is the selling price for a scoop of ice cream. Similarly, we define a mapping $o : \mathcal{A} \rightarrow \mathbb{R}_+$ for the production/purchase cost of the products and $k : \mathcal{S} \rightarrow \mathbb{R}_+$ for the storage cost. Thus, for a pair $(s, a) \in \mathcal{S} \times \mathcal{A}_s$ and another state $s' \in \mathcal{S}$, the gain is given by the mapping $R : (\mathcal{S} \times \mathcal{A}_s) \times \mathcal{S} \rightarrow \mathbb{R}$ with

$$R(s, a, s') := \underbrace{f(s + a - s')}_{\text{sold ice cream}} - \underbrace{o(a)}_{\text{production costs}} - \underbrace{k(s + a)}_{\text{storage costs}}.$$

Moreover, let us define a mapping $h : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$

$$h(s'; s, a) := \begin{cases} p_{s+a-s'} & : 1 \leq s' < s + a \\ \sum_{i \geq s+a} p_i & : s' = 0 \\ 0 & : \text{otherwise} \end{cases}$$

and from this the transition probabilities

$$p(s', r; s, a) := h(s'; s, a) \cdot \mathbf{1}_{r=R(s, a, s')}$$

As for grid world we get the state-action-state transition probabilities as

$$p(s'; s, a) = p(\{s'\} \times \mathcal{R}; s, a) = h(s'; s, a)$$

and the reward expectations as

$$\begin{aligned} r(s, a) &= \sum_{r \in \mathcal{R}} r p(\mathcal{S} \times \{r\}; s, a) \\ &= \sum_{s' \in \mathcal{S}} \sum_{r \in \mathcal{R}} r p(s', r; s, a) \\ &= \sum_{s' \in \mathcal{S}} \sum_{r \in \mathcal{R}} r \cdot h(s'; s, a) \cdot \mathbf{1}_{r=R(s, a, s')} \\ &= \sum_{s' \in \mathcal{S}} R(s, a, s') \cdot h(s'; s, a). \end{aligned}$$

So far the definition of Markov decision models and policies is very general. It will later turn out that much smaller classes of policies are needed to solve optimal control problems in the MDP setting.


Definition 2.1.15. (Markov and stationary policies)

A policy $\pi = (\pi_t)_{t \in \mathbb{N}_0} \in \Pi$ is called

- (i) a Markov policy if there exists a sequence of kernels $(\varphi_t)_{t \in \mathbb{N}_0}$ on $\bar{\mathcal{A}} \times \mathcal{S}$ such that

$$\pi_t(\cdot; s_0, a_0, \dots, s_t) = \varphi_t(\cdot; s_t), \quad \forall (s_0, a_0, \dots, s_t) \in (\mathcal{S} \times \mathcal{A})^{t-1} \times \mathcal{S}.$$

The set of all Markov policies is denoted by Π_M .

- (ii) a stationary policy if there exists a kernel φ on $\bar{\mathcal{A}} \times \mathcal{S}$ such that

$$\pi_t(\cdot; s_0, a_0, \dots, s_t) = \varphi(\cdot; s_t), \quad \forall (s_0, a_0, \dots, s_t) \in (\mathcal{S} \times \mathcal{A})^{t-1} \times \mathcal{S}.$$

The set of all stationary policies is denoted by Π_S .

- (iii) a deterministic stationary policy if there exists a kernel φ on $\bar{\mathcal{A}} \times \mathcal{S}$ taking only values in $\{0, 1\}$ such that

$$\pi_t(\cdot; s_0, a_0, \dots, s_t) = \varphi(\cdot; s_t), \quad \forall (s_0, a_0, \dots, s_t) \in (\mathcal{S} \times \mathcal{A})^{t-1} \times \mathcal{S}.$$

The set of all deterministic stationary policies is denoted by Π_S^D .

In the stationary cases we typically write π instead of φ .

From the definition it holds that

$$\Pi_S^D \subseteq \Pi_S \subseteq \Pi_M.$$

In words: A Markov policy only uses the actual state (not the past) to choose the action, a stationary policy does not depend on time (and the past), a deterministic stationary policy only chooses one action (like an index strategy for bandits). In fact, it will turn out that only deterministic stationary policies are needed to solve the central optimisation problem that will be defined below.

As noted earlier, stochastic bandit models can be seen as special case of MDPs. If $|\mathcal{S}| = 1$, then a one-step Markov decision model is nothing but a bandit model that is played once. There are only actions (arms) that are played once according to a one-step policy, R_1 is the reward obtained by playing the arm. In fact, there is another way of linking stochastic bandits to MDPs. The learning process using a learning strategy can also be seen as an MDP played with a non-Markovian policy as follows. If $|\mathcal{S}| = 1$ and $T = n$, then a Markov decision model is a bandit model and a policy is a learning strategy. The rewards R_1, \dots, R_n are the outcomes of playing the arms chosen according to π . The way a learning strategy was defined for bandits it is neither Markovian nor stationary policy.

There is a good reason that MDPs carry the word Markov. For Markov policies they are indeed Markov (reward) processes. The state-action process (S, A) is Markov so that together with the rewards the process (S, A, R) is a Markov reward process.


Proposition 2.1.16. (Markov property)

If $\pi \in \Pi_S$, then $(S_t, A_t)_{t \in \mathbb{N}}$ is a time-homogeneous Markov chain on $\mathcal{S} \times \mathcal{A}$ with two-step transitions

$$p_{(a,s),(a',s')} = p(s'; s, a) \cdot \pi(a'; s').$$

Proof. The proof is a computation with the path probabilities from (2.10). Since we only care for state-actions (S_t, A_t) we use the short-hand notation $p(s'; s, a)$ instead of $p(\{s'\} \times \mathcal{R}; s, a)$.

Plugging-in yields

$$\begin{aligned}
& \mathbb{P}(S_{t+1} = s_{t+1}, A_{t+1} = a_{t+1} | S_t = s_t, A_t = a_t) \\
&= \frac{\mathbb{P}(S_{t+1} = s_{t+1}, A_{t+1} = a_{t+1}, S_t = s_t, A_t = a_t)}{\mathbb{P}(S_t = s_t, A_t = a_t)} \\
&= \frac{\sum_{s_0, a_0} \dots \sum_{s_{t-1}, a_{t-1}} \mathbb{P}(S_0 = s_0, A_0 = a_0, \dots, S_{t+1} = s_{t+1}, A_{t+1} = a_{t+1})}{\sum_{s_0, a_0} \dots \sum_{s_{t-1}, a_{t-1}} \mathbb{P}(S_0 = s_0, A_0 = a_0, \dots, S_t = s_t, A_t = a_t)} \\
&= \frac{\sum_{s_0, a_0} \dots \sum_{s_{t-1}, a_{t-1}} \mu(s_0) \pi_0(a_0; s_0)}{\sum_{s_0, a_0} \dots \sum_{s_{t-1}, a_{t-1}} \mu(s_0) \pi_0(a_0; s_0)} \\
&\quad \times \frac{\prod_{i=1}^{t+1} p(s_i; s_{i-1}, a_{i-1}) \cdot \pi_i(a_i; s_0, a_0, \dots, a_{i-1}, s_i)}{\prod_{i=1}^t p(s_i; s_{i-1}, a_{i-1}) \cdot \pi_i(a_i; s_0, a_0, \dots, a_{i-1}, s_i)} \\
&\stackrel{\pi \in \Pi_M}{=} \frac{\sum_{s_0, a_0} \dots \sum_{s_{t-1}, a_{t-1}} \mu(s_0) \pi_0(a_0; s_0) \prod_{i=1}^{t+1} p(s_i; s_{i-1}, a_{i-1}) \cdot \varphi_i(a_i; s_i)}{\sum_{s_0, a_0} \dots \sum_{s_{t-1}, a_{t-1}} \mu(s_0) \pi_0(a_0; s_0) \prod_{i=1}^t p(s_i; s_{i-1}, a_{i-1}) \cdot \varphi_i(a_i; s_i)} \\
&= p(s_{t+1}; s_t, a_t) \cdot \varphi_{t+1}(a_{t+1}; s_{t+1}) \\
&\quad \times \frac{\sum_{s_0, a_0} \dots \sum_{s_{t-1}, a_{t-1}} \mu(s_0) \cdot \pi_0(a_0; s_0) \cdot \prod_{i=1}^t p(s_i; s_{i-1}, a_{i-1}) \cdot \varphi_i(a_i; s_i)}{\sum_{s_0, a_0} \dots \sum_{s_{t-1}, a_{t-1}} \mu(s_0) \cdot \pi_0(a_0; s_0) \cdot \prod_{i=1}^t p(s_i; s_{i-1}, a_{i-1}) \cdot \varphi_i(a_i; s_i)} \\
&= p(s_{t+1}; s_t, a_t) \cdot \varphi_{t+1}(a_{t+1}; s_{t+1}) \\
&= \frac{\mathbb{P}(S_0 = s_0, A_0 = a_0, \dots, S_t = s_t, A_t = a_t, S_{t+1} = s_{t+1}, A_{t+1} = a_{t+1})}{\mathbb{P}(S_0 = s_0, A_0 = a_0, \dots, S_t = s_t, A_t = a_t)} \\
&= \mathbb{P}(S_{t+1} = s_{t+1}, A_{t+1} = a_{t+1} | S_0 = s_0, A_0 = a_0, \dots, S_t = s_t, A_t = a_t).
\end{aligned}$$

Finally, if $\varphi_t = \varphi$ the computation (compare the third line from below) shows that $\mathbb{P}(S_{t+1} = s', A_{t+1} = s' | S_t = s, A_t = a)$ is independent of t . \square



Proposition 2.1.17. (Markov reward property)

If $\pi \in \Pi_S$ then $(S_t, A_t, R_t)_{t \in \mathbb{N}}$ satisfies the Markov reward process property

$$\begin{aligned}
& \mathbb{P}((S_{t+1}, A_{t+1}, R_{t+1}) = (s_{t+1}, a_{t+1}, r_{t+1}) | (S_t, A_t) = (s_t, a_t), \dots, (S_0, A_0) = (s_0, a_0)) \\
&= \mathbb{P}((S_{t+1}, A_{t+1}, R_{t+1}) = (s_{t+1}, a_{t+1}, r_{t+1}) | (S_t, A_t) = (s_t, a_t))
\end{aligned}$$

with time-homogeneous state/reward transition probabilities

$$p_{(s,a),(s',a',r')} = p(s', r; s, a) \pi(a'; s').$$

Proof.



The computation is almost identical to that of the Markov property for (S, A) , do it! \square



In what follows we will use the notation $\mathbb{P}_s^\pi, \mathbb{P}_{s,a}^\pi$ in different initialisations.

- Under $\mathbb{P}_{s,a}^\pi$ the Markov reward process (S, A, R) is started in (s, a) . In state s we start with action a and then continue the Markov chain transitions and reward payoffs.
- Under \mathbb{P}_s^π the Markov reward process (S, A, R) is started in the random initialisation $\delta_s \otimes \pi_0(\cdot; s)$. That is, in state s the first action is chosen by π_0 and then continue the Markov chain transitions and reward payoffs.



In particular, by definition it holds that

$$\mathbb{P}_s^\pi(\cdot) = \sum_{a \in \mathcal{A}_s} \mathbb{P}_{s,a}^\pi(\cdot) \pi_0(a; s).$$

The notation is a bit annoying but the only way to avoid non-justified conditioning that one can see everywhere in the reinforcement learning literature.

2.1.3 Elements of stochastic control theory

So far we have introduced the concept of a Markov decision processes for a given transition function p and a policy π . But what is the actual question that we are after? The question is to find a policy that maximise what we will call the expected discounted reward. Since a policy can be used to control a system this problem is also known under stochastic optimal control. In order to do so there will be two main steps. How to compute the expected discounted reward (this will be called prediction) and then how to find the optimal policy (this will be called control). Here is the optimization target that stochastic optimal control is about and reinforcement learning tries to solve:



Given a Markov reward model and some terminal time T to defined later find a policy π that maximizes the expected sum of discounted future rewards

$$\mathbb{E}^\pi \left[\sum_{t=0}^T R_t \right].$$

There are three typical choices for T :

- $T = 0$ is called contextual bandit, for $|\mathcal{S}| = 1$ this is a stochastic bandit problem.
- $T = \min\{t : S_t = s'\}$ for some fixed state s' . The choice $R_t \equiv 1$ then counts the number of steps needed to reach s' , justifying the name stochastic shortest path problems for this choice of T . We will not discuss stochastic shortest path problems in these notes.
- $T \in \mathbb{N}$ fixed. This is called finite-time stochastic control problem, we briefly touch upon finite-time problems in Section 2.4.
- The main focus of these notes is the case $T \sim \text{Geo}(1 - \gamma)$, $\gamma \in (0, 1)$, for some geometric random variable T independent of (S, A, R) . It turns out that this situation is much simpler as there is additional Markovian structure (geometric random variables are memoryless). Given a fixed time the memoryless property says that the end is as far away as it was at the beginning. Integrating out the independent time-horizon gives

$$\begin{aligned} \mathbb{E}^\pi \left[\sum_{t=0}^T R_t \right] &= \mathbb{E}^\pi \left[\sum_{k=0}^{\infty} (1 - \gamma) \gamma^k \sum_{t=0}^k R_t \right] \\ &= (1 - \gamma) \mathbb{E}^\pi \left[\sum_{t=0}^{\infty} \sum_{k=t}^{\infty} \gamma^k R_t \right] \\ &= (1 - \gamma) \mathbb{E}^\pi \left[\sum_{t=0}^{\infty} \frac{\gamma^t}{1 - \gamma} R_t \right] \\ &= \mathbb{E}^\pi \left[\sum_{t=0}^{\infty} \gamma^t R_t \right] \end{aligned} \tag{2.12}$$

Thus, from the optimisation point of view it is equivalent to optimise the expected accumulated rewards up to a finite independent geometric time or to optimise the expected

accumulated discounted rewards. Since it is more standard we will mainly focus on the second point of view but keep in mind that the geometric interpretation might be more reasonable from a sampling perspective.

Before going into the details let us fix some conditions that we want to assume, mostly to make our lives easier:

Assumption: The state space \mathcal{S} , the action space \mathcal{A} , and the set of rewards \mathcal{R} are assumed to be finite. All appearing σ -algebras are chosen to be the power sets.

The assumption of bounded rewards is not necessary but makes our lives much easier for the presentation as everything will be finite and exchanging expectations and sums will always follow from dominated convergence.



Definition 2.1.18. (Value functions)

For $\pi \in \Pi$ and $\gamma \in (0, 1)$, the function $Q^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ defined by

$$Q^\pi(s, a) = \mathbb{E}_{s,a}^\pi \left[\sum_{t=0}^{\infty} \gamma^t R_t \right]$$

is called state-action value function; or Q-function. The value function is defined by

$$V^\pi(s) = \mathbb{E}_s^\pi \left[\sum_{t=0}^{\infty} \gamma^t R_t \right] = \sum_{a \in \mathcal{A}_s} \pi_0(a; s) Q^\pi(s, a).$$

In words: The state value function is the expected discounted total reward when $S_0 = s$ is fixed and A_0 is played according to π_0 . In contrast, Q also fixes the first action $a_0 = a$. Think about chess. The state-value function shall describe the expected success following a policy while state-action-value function is the expected success when fixing the first move.

For stationary policies the value function and state-value functions satisfy simple systems of equations. That will help us later to get our hands on V^π and Q^π .



Proposition 2.1.19. (Bellman equations for Q^π and V^π)

Suppose π is a stationary policy, then Q^π and V^π satisfy the following systems of linear equations:

$$Q^\pi(s, a) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \sum_{a' \in \mathcal{A}_{s'}} p(s'; s, a) \pi(a'; s') Q^\pi(s', a'), \quad s \in \mathcal{S}, a \in \mathcal{A},$$

and

$$V^\pi(s) = \sum_{a \in \mathcal{A}_s} \pi(a; s) \left(r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'; s, a) V^\pi(s') \right), \quad s \in \mathcal{S}.$$

Proof. Recall from Proposition 2.1.17 that (S, A, R) is a Markov reward process so that by (2.7)

$$\mathbb{E}_{s,a}^\pi \left[R_1 + \gamma R_2 + \gamma^2 R_3 + \dots \mid S_1 = s', A_1 = a' \right] = \mathbb{E}_{s',a'}^\pi \left[R_0 + \gamma R_1 + \gamma^2 R_2 + \dots \right] = Q^\pi(s', a').$$

Thus, using the formula of total probability,

$$\begin{aligned}
& Q^\pi(s, a) \\
&= \mathbb{E}_{s,a}^\pi \left[\sum_{t=0}^{\infty} \gamma^t R_t \right] \\
&= \mathbb{E}_{s,a}^\pi [R_0] + \mathbb{E}_{s,a}^\pi [\gamma R_1 + \gamma^2 R_2 + \gamma^3 R_3 + \dots] \\
&= r(s, a) + \gamma \sum_{s' \in \mathcal{S}, a' \in \mathcal{A}_s} \mathbb{E}_{s,a}^\pi [R_1 + \gamma R_2 + \gamma^2 R_3 + \dots \mid S_1 = s', A_1 = a'] \mathbb{P}_{s,a}^\pi(S_1 = s', A_1 = a') \\
&= r(s, a) + \gamma \sum_{s' \in \mathcal{S}, a' \in \mathcal{A}_s} \mathbb{P}_{s,a}^\pi(S_1 = s', A_1 = a') Q^\pi(s', a') \\
&= r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \sum_{a' \in \mathcal{A}_{s'}} p(s'; s, a) \pi(a'; s') Q^\pi(s', a').
\end{aligned}$$

The equation for V follows directly by plugging-in $V^\pi(s) = \sum_{a \in \mathcal{A}_s} \pi(a; s) Q^\pi(s, a)$ twice:

$$\begin{aligned}
V^\pi(s) &= \sum_{a \in \mathcal{A}_s} \pi(a; s) Q^\pi(s, a) \\
&= \sum_{a \in \mathcal{A}_s} \pi(a; s) \left(r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \sum_{a' \in \mathcal{A}_{s'}} p(s'; s, a) \pi(a'; s') Q^\pi(s', a') \right) \\
&= \sum_{a \in \mathcal{A}_s} \pi(a; s) \left(r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'; s, a) V^\pi(s') \right).
\end{aligned}$$

□

The definition of V^π immediately allows to compute V^π from Q^π . For stationary policies the opposite also holds true:



Corollary 2.1.20. Given a stationary policy $\pi \in \Pi_S$, the following relations between the state- and action-value functions hold:

$$Q^\pi(s, a) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'; s, a) V^\pi(s'), \quad s \in \mathcal{S}, a \in \mathcal{A}.$$

Proof. This follows directly from the Bellman equation for Q plugging-in the definition $V^\pi(s) = \sum_{a \in \mathcal{A}_s} \pi(a; s) Q^\pi(s, a)$ for V . □

We will later see that linking Q and V is crucial, in particular computing Q from V via Corollary (2.1.20).

The classical theory of optimising Markov decision processes is very much about functional analysis. One of the most important theorem from basic functional analysis is Banach's fixed point theorem:



Theorem 2.1.21. (Banach fixed-point theorem)

Let $(U, \|\cdot\|)$ be a Banach space, i.e. a complete normed vector space, and $T : U \rightarrow U$ a contraction, i.e. there exists $\lambda \in [0, 1)$ such that

$$\|Tu_1 - Tu_2\| \leq \lambda \|u_1 - u_2\|$$

for all $u_1, u_2 \in U$. Then

- (i) there exists a unique fixed point, i.e. $u^* \in U$ such that $Tu^* = u^*$; and



(ii) for arbitrary $u_0 \in U$, the sequence $(u_n)_{n \in \mathbb{N}}$ defined by

$$u_{n+1} = Tu_n = T^{n+1}u_0$$

converges in U to u^* .

The Banach fixed point theorem is useful because the condition can be checked in many cases and also the algorithm yields a fast converging algorithm. One of the major insights of optimal control theory is the observation that contractions appear very naturally and lead to solution algorithms. For that sake we will always use the Banach space $(V, \|\cdot\|_\infty)$ consisting of $V := \{v : \mathcal{S} \rightarrow \mathbb{R}\}$ equipped with the maximum-norm. Since we assume that \mathcal{S} is finite, V is nothing but $\mathbb{R}^{|\mathcal{S}|}$ where a vector is written as a function (mapping index to coordinate value). Similarly, we define $U := \{q : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}\}$ equipped with the supremum norm which is then nothing but the Banach space $\mathbb{R}^{|\mathcal{S}| \cdot |\mathcal{A}|}$.



Definition 2.1.22. (Bellman expectation operator)

Given a Markov decision model and a stationary policy π the operators

$$(T^\pi v)(s) := \sum_{a \in \mathcal{A}_s} \pi(a; s) \left(r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'; s, a) v(s') \right)$$

mapping X into X and

$$(T^\pi q)(s, a) := r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \sum_{a' \in \mathcal{A}_s} p(s'; s, a) \pi(a'; s) q(s', a')$$

mapping Y into Y are called Bellman expectation operators.

There is a bit of ambiguity by denoting both operators by T . Since the number of arguments differs it will always be clear from the context which operator is meant. The operators might look familiar. We have actually already proved in Proposition 2.1.19 that value functions are fixed points for Bellman expectation operators. Combined with the simple contractivity property with respect to the maximum norm the first major theorem follows:



Theorem 2.1.23. Both Bellman expectation operators are contractions with contraction constant γ . Their unique fixed points are the value functions V^π and Q^π , i.e. $T^\pi V^\pi = V^\pi$ and $T^\pi Q^\pi = Q^\pi$.

Proof. We already proved the fixed point claims. Let us check the contraction property. Plugging-in yields

$$\begin{aligned} \|T^\pi v_1 - T^\pi v_2\|_\infty &= \gamma \max_s \left| \sum_{a \in \mathcal{A}_s} \sum_{s' \in \mathcal{S}} \pi(a; s) p(s'; s, a) (v_1(s') - v_2(s')) \right| \\ &\leq \gamma \|v_1 - v_2\|_\infty \max_s \underbrace{\left| \sum_{a \in \mathcal{A}_s} \pi(a; s) \underbrace{\sum_{s' \in \mathcal{S}} p(s'; s, a)}_{=1} \right|}_{=1} \\ &= \gamma \|v_1 - v_2\|_\infty. \end{aligned}$$

Please check yourself the same computation for the second Bellman operator. □

Since Bellman's expectation equations are linear systems they can in principle be solved by linear algebra methods. Since the linear operators are also contractions the linear systems $T^\pi v = v$ (resp. $T^\pi q = q$) have unique solutions. Still, it might be more reasonable to solve using Banach's

fixed point iteration by iterating $v_{n+1} = T^\pi v_n$ for some starting vector v_0 (resp. $q_{n+1} = T^\pi q_n$ for some starting matrix q).

With the value functions we can define the concept of an optimal policies. The search for optimal policies in MDPs will be the main content of the subsequent chapters.



Definition 2.1.24. (Optimal policy & optimal value function)

For a given Markov decision model the following quantities will be of central importance:

- (i) The function $V^* : \mathcal{S} \rightarrow \mathbb{R}$ that takes values

$$V^*(s) := \sup_{\pi \in \Pi} V^\pi(s), \quad s \in \mathcal{S},$$

is called optimal (state) value function.

- (ii) The function $Q^* : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ that takes values

$$Q^*(s, a) := \sup_{\pi \in \Pi} Q^\pi(s, a), \quad s \in \mathcal{S}, a \in \mathcal{A},$$

is called optimal state-action value function.

- (iii) A policy $\pi^* \in \Pi$ that satisfies

$$V^{\pi^*}(s) = V^*(s), \quad s \in \mathcal{S},$$

is called optimal policy.

It should be emphasised that we make our lives much simpler by assuming a finite Markov decision model. In that case the maximum is always attained and plenty of technical problems do not appear. In fact, the results that we are going to prove next do not necessarily hold for infinite Markov decision models. We will see later that optimality of policies could have been defined alternatively by $Q^{\pi^*}(s, a) = Q^*(s, a)$ for all state-action pairs. We stick to the notion in terms of V as this is more common in the literature.



It is important to keep in mind that a priori V^* and Q^* are not the value functions for the best policy but instead pointwise the best possible expected rewards. A priori it is absolutely not clear if there is one policy that is optimal for all starting states. It is also not clear if there is any way of characterising such a policy that leads to an algorithm. It is also not clear if that policy is very complicated. In fact, Bellman's results show that there is actually a stationary, even deterministic, policy that can be characterised by solving (non-linear) equations.

We now turn towards the second operator of Bellman, the optimality operator. As for the expectation operator we work on $U = \{v : \mathcal{S} \rightarrow \mathbb{R}\}$ and $V := \{q : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}\}$. As long as we assume the Markov decision model to be finite X is nothing but $\mathbb{R}^{|\mathcal{S}|}$ and Y is $\mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|}$. To make X and Y Banach spaces we will fix the respective maximum-norms.



Definition 2.1.25. (Bellman optimality equations/operators)

For a given Markov decision model we define the following operators:

- (i) The non-linear system of equations

$$v(s) = \max_{a \in \mathcal{A}_s} \left\{ r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s' ; s, a) v(s') \right\}, \quad s \in \mathcal{S},$$



is called Bellman optimality equation. The operator $T^* : U \rightarrow U$ defined by

$$(T^*v)(s) := \max_{a \in \mathcal{A}_s} \left\{ r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'; s, a)v(s') \right\}, \quad s \in \mathcal{S},$$

is called the Bellman optimality operator (for state-value functions).

(ii) The non-linear system of equations

$$q(s, a) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'; s, a) \max_{a' \in \mathcal{A}_{s'}} q(s', a'), \quad (s, a) \in \mathcal{S} \times \mathcal{A},$$

is called Bellman state-action optimality equation. The state-action Bellman optimality operator $T^* : V \rightarrow V$ is defined as

$$(T^*q)(s, a) := r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'; s, a) \max_{a' \in \mathcal{A}_{s'}} q(s', a'), \quad (s, a) \in \mathcal{S} \times \mathcal{A}.$$

Warning: both optimality operators are denoted by T^* but are safely distinguished by the number of arguments.

It will turn out that the Bellman optimality operators are directly linked to optimal value functions and solving the corresponding fixed point equations (i.e. the Bellman optimality equations) is equivalent to finding optimal policies. Unfortunately, the equations are non-linear and as such not easy to solve. To get a first little hand on the operator please do the following exercise.



All Bellman operators are monotone, i.e. if $u_1 \leq u_2$ it also holds that $T^\pi u_1 \leq T^\pi u_2$ and $T^* u_1 \leq T^* u_2$. The same holds for the matrix-valued operators T^π and T^* .

Similarly to the expectation operators also the optimality operators are contractions:



Theorem 2.1.26. Bellman's optimality operators are contractions and, thus, have unique fixedpoints.

Proof. To deal with the operators a simple inequality from analysis will be used:

$$\left| \max_{a \in \mathcal{A}} f(a) - \max_{a \in \mathcal{A}} g(a) \right| \leq \max_{a \in \mathcal{A}} |f(a) - g(a)|.$$

If the number in the absolute value of the left hand side is positive, then

$$\left| \max_{a \in \mathcal{A}} f(a) - \max_{a \in \mathcal{A}} g(a) \right| \leq \max_{a \in \mathcal{A}} (f(a) - g(a)) \leq \max_{a \in \mathcal{A}} |f(a) - g(a)|.$$

Otherwise, the role of f and g is reversed. With $v_1, v_2 \in X$, the Bellman optimality operator and the triangle inequality yield

$$\begin{aligned} \|T^*v_1 - T^*v_2\|_\infty &= \max_{s \in \mathcal{S}} |T^*v_1(s) - T^*v_2(s)| \\ &\leq \max_{s \in \mathcal{S}} \max_{a \in \mathcal{A}} \gamma \sum_{s' \in \mathcal{S}} p(s'; s, a) |v_1(s') - v_2(s')| \\ &\leq \max_{s \in \mathcal{S}} \max_{a \in \mathcal{A}} \gamma \sum_{s' \in \mathcal{S}} p(s'; s, a) \|v_1 - v_2\|_\infty \\ &= \gamma \|v_1 - v_2\|_\infty. \end{aligned}$$

Hence, T^* is a contraction on $V = \{v : \mathcal{S} \rightarrow \mathbb{R}\}$ equipped with the maximum norm and thus has a unique fixed point. Next, we show the same for T^* on $U = \{q : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}\}$: With $q_1, q_2 \in U$, the second Bellman operator and the triangle inequality yield

$$\begin{aligned} \|T^*q_1 - T^*q_2\|_\infty &= \max_{s \in \mathcal{S}} \max_{a \in \mathcal{A}_s} |T^*q_1(s, a) - T^*q_2(s, a)| \\ &\leq \max_{s \in \mathcal{S}} \max_{a \in \mathcal{A}} \left| \gamma \sum_{s' \in \mathcal{S}} p(s'; s, a) \left(\max_{a' \in \mathcal{A}_{s'}} q_1(s', a') - \max_{a' \in \mathcal{A}_{s'}} q_2(s', a') \right) \right| \\ &\leq \max_{s \in \mathcal{S}} \max_{a \in \mathcal{A}} \gamma \sum_{s' \in \mathcal{S}} p(s'; s, a) \max_{a' \in \mathcal{A}_{s'}} |q_1(s', a') - q_2(s', a')| \\ &\leq \max_{s \in \mathcal{S}} \max_{a \in \mathcal{A}} \gamma \sum_{s' \in \mathcal{S}} p(s'; s, a) \max_{s' \in \mathcal{S}} \max_{a' \in \mathcal{A}_{s'}} |q_1(s', a') - q_2(s', a')| \\ &= \gamma \|q_1 - q_2\|_\infty. \end{aligned}$$

□

We now come to the most important result of this section. The optimal value functions are uniquely characterised as fixed points of Bellman's optimality operators.



Theorem 2.1.27. The optimal value functions are related as

$$V^*(s) = \max_{a \in \mathcal{A}_s} Q^*(s, a)$$

and satisfy the Bellman optimality equations

$$T^*V^* = V^* \quad \text{and} \quad T^*Q^* = Q^*.$$

Proof.

By Theorem 2.1.26 we know that there exist unique fixed points of the Bellman optimality equations, denoted by V^{fix} and Q^{fix} . We will show that these fixed points are given by $V^{\text{fix}} = V^*$ and $Q^{\text{fix}} = Q^*$.

(i) $Q^{\text{fix}} \geq Q^*$: Fix a policy π and note that for all $q : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ and $(s, a) \in \mathcal{S} \times \mathcal{A}$ one has

$$\begin{aligned} (T^*q)(s, a) &= r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'; s, a) \max_{a' \in \mathcal{A}_{s'}} q(s', a') \\ &\geq r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \sum_{a' \in \mathcal{A}_{s'}} p(s'; s, a) \pi(a'; s) q(s', a') = (T^\pi q)(s, a). \end{aligned}$$

Thus,

$$Q^{\text{fix}} = T^*Q^{\text{fix}} \geq T^\pi Q^{\text{fix}},$$

where the inequality holds componentwise. Moreover, for $q_1, q_2 : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ with $q_1 \geq q_2$ componentwise one has $T^*q_1 \geq T^*q_2$. Thus, for $n \in \mathbb{N}$ with $Q^{\text{fix}} \geq (T^\pi)^n Q^{\text{fix}}$ one has

$$Q^{\text{fix}} = (T^*)^{n+1} Q^{\text{fix}} = T^*(T^*)^n Q^{\text{fix}} \geq T^*(T^\pi)^n Q^{\text{fix}} \geq T^\pi (T^\pi)^n Q^{\text{fix}} = (T^\pi)^{n+1} Q^{\text{fix}}.$$

By induction, we get $(T^*)^n Q^{\text{fix}} \geq (T^\pi)^n Q^{\text{fix}}$ for all $n \in \mathbb{N}$. Using Theorem 2.1.23 we get $(T^\pi)^n Q^{\text{fix}} \xrightarrow{n \rightarrow \infty} Q^\pi$, so that $Q^{\text{fix}} \geq Q^\pi$. Taking the supremum over all policies π we get

$$Q^{\text{fix}}(s, a) \geq Q^*(s, a), \quad (s, a) \in \mathcal{S} \times \mathcal{A}.$$

(ii) $Q^{\text{fix}} \leq Q^*$: Define a particular policy π^* as

$$\pi^{\text{fix}}(a; s) = \begin{cases} 1 & : a \in \arg \max_{a'} Q^{\text{fix}}(s, a') \\ 0 & : \text{otherwise} \end{cases}.$$

The policy will later be coined greedy with respect to Q^{fix} . Note that we do not know the matrix Q^{fix} explicitly, so for now this is only a theoretical argument. One has

$$\begin{aligned} Q^{\text{fix}}(s, a) &= (T^* Q^{\text{fix}})(s, a) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'; s, a) \max_{a' \in \mathcal{A}_{s'}} Q^{\text{fix}}(s', a') \\ &= r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \sum_{s'' \in \mathcal{A}_{s'}} p(s'; s, a) \pi^{\text{fix}}(s''; s') Q^{\text{fix}}(s', a'') = (T^\pi Q)(s, a). \end{aligned}$$

Repeating this argument gives $Q^{\text{fix}}(s, a) = (T^{\pi^{\text{fix}}})^n Q^{\text{fix}}(s, a) \xrightarrow{n \rightarrow \infty} Q^{\pi^{\text{fix}}}(s, a)$. Thus,

$$Q^{\text{fix}}(s, a) = Q^{\pi^{\text{fix}}}(s, a) \leq \sup_{\pi} Q^\pi(s, a) = Q^*(s, a).$$

(iii) Bellman equation for Q^* : Combining the two inequalities proved above we obtain the claim for Q^* .

(iv) Relation of V^* and Q^* : The optimal state value function V^* was defined as the supremum over all policies:

$$V^*(s) = \sup_{\pi} V^\pi(s),$$

where for each policy π

$$V^\pi(s) = \sum_a \pi(a; s) Q^\pi(s, a).$$

Since for any fixed s and any policy π we have

$$\sum_a \pi(a; s) Q^\pi(s, a) \leq \max_a Q^\pi(s, a),$$

it follows that

$$V^*(s) = \sup_{\pi} V^\pi(s) \leq \sup_{\pi} \max_a Q^\pi(s, a) = \max_a Q^*(s, a).$$

We proceed with the reverse inequality. To do so we again use the greedy policy $\pi^* := \pi^{\text{fix}}$ introduced above. Let $a^*(s)$ denote an action in $\arg \max_{a \in \mathcal{A}} Q^*(s, a)$. Here is a simple dynamic programming argument. We use a simple fact $T^{\pi^*} Q^* = T^* Q^*$ which follows easily as follows:

$$\begin{aligned} (T^{\pi^*} Q^*)(s, a) &= r(s, a) + \gamma \sum_{s' \in \mathcal{A}_s} \sum_{a' \in \mathcal{A}_{s'}} p(s'; s, a) \pi^*(s'; a') Q^*(s', a') \\ &= r(s, a) + \gamma \sum_{s' \in \mathcal{A}_s} p(s'; s, a) \max_{a' \in \mathcal{A}_{s'}} Q^*(s', a') \\ &= T^* Q^*(s, a), \end{aligned}$$

where for the second equality we used the definition of the greedy policy. Since Q^* solves the Bellman optimality equation we get

$$Q^* = T^* Q^* = T^{\pi^*} Q^*.$$

Uniqueness of Bellman's expectation operator gives $Q^* = Q^{\pi^*}$. Now we can finish the argument:

$$V^*(s) \geq V^{\pi^*}(s) = \sum_a \pi^*(a; s) Q^{\pi^*}(s, a) = \sum_a \pi^*(a; s) Q^*(s, a) = Q^*(s, a^*(s)) = \max_a Q^*(s, a).$$

(v) Bellman equation for V^* : We already know that the optimal state-action value function satisfies the Bellman optimality equation. Taking the maximum over a on both sides, we obtain

$$\max_a Q^*(s, a) = \max_a \left\{ r(s, a) + \gamma \sum_{s'} p(s'; s, a) \max_{a'} Q^*(s', a') \right\}.$$

Using the relation $\max_a Q^*(s, a) = V^*(s)$, we deduce

$$V^*(s) = \max_a \left\{ r(s, a) + \gamma \sum_{s'} p(s'; s, a) V^*(s') \right\}.$$

This is the Bellman equation for the optimal state value function. □

The proof used ideas that are the core of dynamic programming theory. Most importantly, we could see how to find an optimal policy. The optimal policy was just the greedy policy from the optimal state-action function. Unfortunately, a priori Q^* is unknown so we do not know how to find π^* . The main idea behind many optimal control algorithms is to invert the point of view. If the optimal Bellman equation can be solved (somehow) then the greedy policy from this matrix will be optimal. Reusing ideas from the previous proof this will now be turned into theory.



Definition 2.1.28. (Greedy policy)

Given a function $q : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ the deterministic stationary policy defined by

$$\pi_q(a; s) := \begin{cases} 1 & : a = a^*(s) \\ 0 & : \text{otherwise} \end{cases} \quad \text{with } a^*(s) \in \arg \max_{a \in \mathcal{A}_s} q(s, a)$$

is called a greedy policy with respect to q . Sometimes we also write $\text{greedy}(q)$ instead of π_q . If π is greedy with respect to a q_v obtained from the V - Q -transfer operator

$$q_v(s, a) := r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'; s, a)v(s'), \quad s \in \mathcal{S}, a \in \mathcal{A}_s, \quad (2.13)$$

then we also write π_v .

Keep in mind that the greedy policy is very special, only one action is proposed that maximises the given q -function. In case several actions yield the same state-action value a fixed one of them is chosen. Here is a lemma that shows how to deal with greedy policies in the context of Bellman operators.



Lemma 2.1.29. Suppose π_q is a greedy policy obtained from a q , then $T^*q = T^{\pi_q}q$. If q is obtained from a vector v through (2.13) then it also holds that $T^*v = T^{\pi_q}v$.

Proof. Both claims readily follow from the definition of the Bellman operators:

$$\begin{aligned} (T^{\pi_q}q)(s, a) &= r(s, a) + \gamma \sum_{s' \in \mathcal{A}_s} \sum_{a' \in \mathcal{A}_{s'}} p(s'; s, a) \pi_q(s'; a') q(s', a') \\ &= r(s, a) + \gamma \sum_{s' \in \mathcal{A}_s} p(s'; s, a) \max_{a' \in \mathcal{A}_{s'}} q(s', a') = (T^*q)(s, a), \end{aligned}$$

where for the second equality we used the definition of the greedy policy. The same holds if q is obtained from a vector v :

$$\begin{aligned} (T^{\pi_q}v)(s) &= \sum_{a \in \mathcal{A}_s} \pi_q(a; s) \left(r(s, a) + \gamma \sum_{s' \in \mathcal{A}_s} p(s'; s, a)v(s') \right) \\ &= \max_{a \in \mathcal{A}_s} \left(r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'; s, a)v(s') \right) = (T^*v)(s), \end{aligned}$$

where we used the definition of the Bellman operators and the definition of the greedy policy in the second equality. \square

As a consequence we learn how solving the Bellman state-action optimality equation yields an optimal policy.



Theorem 2.1.30. (Dynamic programming algorithm)

An optimal policy π^* always exist and can always be chosen to be stationary and deterministic! Such a policy is given by solving the Bellman state-action optimality



equation and using the greedy policy with respect to the solution.

Alternatively, a solution v to the Bellman optimality equation plugged-into the V - Q -transfer operator

$$q_v(s, a) := r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'; s, a)v(s'), \quad s \in \mathcal{S}, a \in \mathcal{A}_s,$$

yields Q^* and, hence, the greedy optimal policy. Since v is only a vector while q is a matrix it sounds first plausible (up to now) to solve $T^*v = v$ and then transfer to q_v . Solving the non-linear equation is the expensive step, transferring to q not. We will later see that typically we try to avoid transferring from V to Q but in this chapter on stochastic control that sounds like a good idea.

Proof. Suppose we have a solution q of Bellman's state-action optimality equation. By uniqueness q equals Q^* . It remains to show $V^* = V^{\pi_q}$, then π_q is optimal by definition. But this follows from the previous lemma and the fact that $q = Q^*$:

$$Q^* = T^*Q^* = T^{\pi_q}Q^* = T^{\pi_q}Q^*.$$

Uniqueness of Bellman's expectation operator gives $Q^* = Q^{\pi_q}$. Finally, using the definition of the greedy policy and the relations $V^*(s) = \max_{a \in \mathcal{A}_s} Q^*(s, a)$ and $V^\pi(s) = \sum_{a \in \mathcal{A}_s} \pi(a; s)Q^\pi(s, a)$ between V and Q gives

$$V^*(s) = \max_{a \in \mathcal{A}_s} Q^*(s, a) = \max_{a \in \mathcal{A}_s} Q^{\pi_q}(s, a) = \sum_{a \in \mathcal{A}_s} \pi_q(a; s)Q^{\pi_q}(s, a) = V^{\pi_q}(s)$$

We have thus proved that $V^* = V^{\pi_q}$ which shows that π_q is optimal. \square

Let us summarise the findings made so far, the core of everything that will later be called valued-based algorithms.



Solving Bellman's optimality equation (or state-value optimality equation) yields a stationary deterministic policy π^* as the greedy policy obtained from the solution. All algorithms that approximatively solve the Bellman optimality equation give rise to approximation algorithms that find a stationary optimal policy of the stochastic optimal control problem $\sup_{\pi \in \Pi} V^\pi$.

From now on we will only focus on finding stationary optimal policies!

As the remark indicates we will be interested in learning the optimal policies by approximation. Recalling the setting of stochastic bandits (interpreted as one-step MDP) this sounds familiar. In that setting an optimal policy is nothing but a dirac measure on optimal arms, a learning strategy a sequence of policies that learns (approximates) the optimal policy. To make the connection to stochastic bandits let us define the notation of a learning strategy:



Definition 2.1.31. A sequence $(\pi^n)_{n \in \mathbb{N}}$ of policies for a Markov decision model is called a learning strategy, if π^{n+1} only depends on everything seen for the first n rounds of learning.

We keep the definition extremely vague as it will not play any further role in what follows. The aim is to find algorithms that produce learning strategies that converge quickly and efficiently to the optimal strategy π^* . What we mean by convergence will depend on the context, the minimal requirement is convergence of the value function to the optimal value function, i.e. $\|V^{\pi^n} - V^*\|_\infty \rightarrow 0$. In contrast to stochastic bandit theory the regret plays no role in reinforcement learning, the situation is too complex.

There are two typical approaches that we will encounter in different setups:

- value function based learning,
- policy based learning.

For value function the idea is to learn the optimal value function V^* (or optimal state-value function Q^*) and then infer the optimal policy π^* by taking argmax (the greedy policy from Theorem 2.1.30). In contrast, policy learning tries to approximate directly the optimal policy π^* .



The algorithms presented below are called **dynamic programming algorithms**, they belong to a large class of algorithms that break down a problem into (hopefully simpler) subproblems. Here, this means for instance to compute $V^\pi(s)$ from all other $V^\pi(s')$ or $Q^\pi(s, a)$ from all other $Q^\pi(s', a')$. Since for discounted MDP problems the subproblems do not immediately simplify we do not go further into the ideas of dynamic programming but still refer to all algorithms based on Bellman's equations as dynamic programming algorithms.

2.2 Basic tabular value iteration algorithm

We have understood the problem and its ingredients and have learnt that in order to act optimally, that is to know the optimal policy, one only needs to know the optimal value function (or the optimal state-action value function). The approach of this section is to develop methods that approximate the optimal value functions. As seen in the previous section, the Banach fixed point theorem gives not only the existence of a unique fixed point of a contraction mapping T but also the convergence of the sequence $(T^{n+1}v_0)_{n \in \mathbb{N}_0}$ to that fixed point for arbitrary v_0 . We learnt that for T^* the unique fixed point corresponded to V^* . The idea of value iteration is to use this convergence property and turn the Bellman optimality equation into an update procedure. The algorithm is called a tabular algorithm as a table (here the vector V) is updated repeatedly.

Data: Accuracy $\varepsilon > 0$

Result: Approximation $V \approx V^*$, policy $\pi \approx \pi^*$

Initialize $V \equiv 0, V_{\text{new}} \equiv 0$

$\Delta := 1$

while $\Delta > \varepsilon$ **do**

 set $V := V_{\text{new}}$

for $s \in \mathcal{S}$ **do**

$$V_{\text{new}}(s) = \max_{a \in \mathcal{A}_s} \left\{ r(s, a) + \gamma \underbrace{\sum_{s' \in \mathcal{S}} p(s'; s, a) V(s')}_{(T^*V)(s)} \right\}$$

end

$\Delta := \max_{s \in \mathcal{S}} (|V_{\text{new}}(s) - V(s)|)$

end

return $V := V_{\text{new}}$ and the greedy policy with respect to V

Algorithm 7: Value iteration



Theorem 2.2.1. The value iteration Algorithm 7 terminates and the terminal vector satisfies $\|V - V^*\|_\infty \leq \frac{\gamma \varepsilon}{1 - \gamma}$.

Proof. Suppose the sequence of vectors $v_{n+1} = T^*v_n$ is obtained by iteratively applying Bellman's optimality operator and v^* is the limit. The finite-time termination of the algorithm follows immediately from Banach's fixed point theorem as

$$\|v_n - v_{n+1}\|_\infty \stackrel{\Delta}{\leq} \|v_n - v^*\|_\infty + \|v^* - v_{n+1}\|_\infty \rightarrow 0$$

for $n \rightarrow \infty$. Now suppose the algorithm terminated after N steps, i.e. $V = v_N$ so that the termination condition implies $\|v_N - v_{N-1}\|_\infty < \varepsilon$. Using the contraction property of T^* yields, for $m \in \mathbb{N}$,

$$\begin{aligned} \|v_N - v_{N+m}\|_\infty &\stackrel{\Delta}{\leq} \sum_{k=0}^{m-1} \|v_{N+k} - v_{N+k+1}\|_\infty \\ &= \sum_{k=0}^{m-1} \|(T^*)^k v_N - (T^*)^k v_{N+1}\|_\infty \\ &\leq \sum_{k=0}^{m-1} \gamma^k \|v_N - v_{N+1}\|_\infty. \end{aligned}$$

Using continuity of the norm yields

$$\begin{aligned} \|V - V^*\|_\infty &= \lim_{m \rightarrow \infty} \|v_N - v_{N+m}\|_\infty \\ &\leq \lim_{m \rightarrow \infty} \sum_{k=0}^{m-1} \gamma^k \|v_{N+1} - v_N\|_\infty \\ &= \frac{1}{1-\gamma} \|v_{N+1} - v_N\|_\infty \leq \frac{\gamma}{1-\gamma} \|v_n - v_{N-1}\|_\infty. \end{aligned}$$

Now the termination condition $\|v_N - v_{N-1}\|_\infty < \varepsilon$ implies the claim. \square

The approximate value function from Algorithm 7 is clearly not equal to the optimal value function V^* . Hence, the effect on transferring to a policy needs to be analysed.



Definition 2.2.2. For $\varepsilon > 0$ a policy $\pi \in \Pi$ is called ε -optimal if

$$V^*(s) \leq V^\pi(s) + \varepsilon$$

for all $s \in \mathcal{S}$.

Now recall how transfer value-functions into state-value functions. For both the optimal value functions and the value functions for a given policy the vector V is transferred into the Q -matrix as

$$Q(s, a) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'; s, a) V(s'). \quad (2.14)$$

We now do the same and use (2.14) to define from the approximation V of V^* an approximate Q -function. Then the corresponding greedy policy is ε -optimal:



Theorem 2.2.3. Suppose V is the output of Algorithm 7 and Q is obtained from V using the transfer operator (2.14). Then the greedy policy π_Q is $\frac{2\varepsilon\gamma}{1-\gamma}$ -optimal.

Proof. The main point of the argument is a relation of optimality and expectation operator in the case of greedy policies. To emphasise the importance, let us highlight the trick once more:



Bellman's optimality and expectation operators are closely connected for greedy policies!

The crucial computation links the optimality operator with the expectation operator of the greedy policy:

$$\begin{aligned}
T^*V(s) &= \max_{a \in \mathcal{A}_s} \left\{ r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'; s, a) V(s') \right\} \\
&= \max_{a \in \mathcal{A}_s} \{ Q(s, a) \} \\
&= \sum_{a \in \mathcal{A}_s} \pi_Q(a; s) Q(a, s) \\
&= \sum_{a \in \mathcal{A}_s} \pi(a; s) \left(r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'; s, a) V(s') \right) \\
&= T^\pi V(s)
\end{aligned}$$

The rest of the proof is straight forward using fixed points and the contraction property. Suppose the algorithm terminated after N steps, i.e. $V = v_N$ and $\|v_N - v_{N-1}\| < \varepsilon \frac{(1-\gamma)}{2\gamma}$. The identity above yields

$$\begin{aligned}
\|V^\pi - V\|_\infty &= \|T^\pi V^\pi - V\|_\infty \\
&\leq \|T^\pi V^\pi - T^*V\|_\infty + \|T^*V - V\|_\infty \\
&= \|T^\pi V^\pi - T^\pi V\|_\infty + \|T^*V - T^*v_{N-1}\|_\infty \\
&\leq \gamma \|V^\pi - V\|_\infty + \gamma \|V - v_{N-1}\|_\infty.
\end{aligned}$$

Rearranging this equation yields

$$\|V^\pi - V\|_\infty \leq \frac{\gamma}{1-\gamma} \|v_N - v_{N-1}\|_\infty.$$

In the proof of the previous theorem we have already shown that $\|V - V^*\|_\infty \leq \frac{\gamma}{1-\gamma} \|v_N - v_{N-1}\|_\infty$. Finally, using the terminal condition of Algorithm 7 yields

$$\|V^\pi - V^*\|_\infty \leq \|V^\pi - V\|_\infty + \|V - V^*\|_\infty \leq 2 \frac{\gamma}{1-\gamma} \|v_N - v_{N-1}\|_\infty \leq 2\varepsilon \frac{\gamma}{1-\gamma}.$$

This proves the claim. \square

We next analyse the convergenc rates of the value iteration algorithm.



Definition 2.2.4. Suppose $(V, \|\cdot\|)$ is a normed space. For a sequence (y_n) in V with limit y^* we say the convergence is of order $\alpha > 0$ if there exists a constant $K < 1$ for which

$$\|y_{n+1} - y^*\| \leq K \|y_n - y^*\|^\alpha, \quad n \in \mathbb{N}. \quad (2.15)$$

Linear convergence corresponds to $\alpha = 1$.

Since the algorithm is based on Banach's fixed point theorem it is clear that the convergence is at least linear. A simple initialisation shows that the convergence generally cannot be improved.



Theorem 2.2.5. For all initialisations the convergence in Algorithm 7 (without termination) is linear with constant $K = \gamma$. There is an initialisation for which the speed of convergence is exactly linear.

Proof. Let us denote again v_n for the n th update of the iteration $v_n = T^*v_{n-1}$. The linear convergence rate follows directly from the fixed point property of T^* :

$$\|v_{n+1} - V^*\|_\infty = \|T^*v_n - T^*V^*\|_\infty \leq \gamma \|v_n - V^*\|_\infty, \quad n \in \mathbb{N}. \quad (2.16)$$

In fact, the rate cannot be better as the following example shows. If Algorithm 7 is initialised with $v_0 := V^* + k\mathbf{1}$, it holds that

$$\begin{aligned} \|v_1 - V^*\|_\infty &= \|T^*v_0 - V^*\|_\infty = \|T^*(V^* + k\mathbf{1}) - V^*\|_\infty \\ &\stackrel{\text{Def. } T^*}{=} \|V^* + \gamma k\mathbf{1} - V^*\|_\infty \\ &= \gamma \|(V^* + k\mathbf{1}) - V^*\|_\infty = \gamma \|v_0 - V^*\|_\infty. \end{aligned}$$

An induction shows that for this example $\|v_{n+1} - V^*\|_\infty = \gamma \|v_n - V^*\|_\infty$ for all $n \in \mathbb{N}$. \square



Of course we could equally use iterations obtained from Banach's fixed point theorem to directly approximate Q^* instead of approximating V^* and then transferring to Q^* . This will be done later for approximate dynamic programming (Q -learning and SARSA) in which we do not assume explicit knowledge on the transition function, thus, cannot transfer from V to Q . In the explicit case it is more reasonable to work with vector iterations than matrix iterations and only transfer from V to Q once.

2.3 Basic policy iteration (actor-critic) algorithm

In this section we want to explore another method of reaching an optimal value function and hence an optimal policy. The method is not part of the classical methodology of stochastic control but much more common in the reinforcement learning community as it motivates some of the most powerful approximation algorithms. The idea goes as follows. Iterate between the following two steps:

- Policy evaluation: Compute or estimate Q^π (or V^π) for the currently best known policy π .
- Policy improvement: Improve the best known policy by taking $\pi' = \text{greedy}(Q^\pi)$.

In contrast to value iteration the approach is more clever, it uses much more understanding of the optimal control problem. While value iteration is called a value-based method (only the value function is used, the policy is only obtained in the end) policy iteration is a policy-based method, the approach works directly uses the policy.



The approach is called an **actor-critic method** as it alternates between two steps. The critic evaluates the policy (computes the value function) which then the actor uses to improve the quality of the policy.

In the following both steps will be discussed separately and then alternated for the policy iteration algorithm.

2.3.1 Policy evaluation

We are now going to address the question on how to compute $V^\pi(s)$ for a given policy π . There are essentially three direct ideas that one might come up with:

- approximate the expectation by Monte Carlo,
- solve the (linear) Bellman expectation equation using linear algebra (e.g. matrix inversion),
- find the fixed point of the Bellman expectation operator using Banach's fixed point iteration.

The first idea and subtle variants will be topic of the next chapter, the latter two approaches will be address now.

Recall the Bellman expectation operator for a stationary policy $\pi \in \Pi_{\mathcal{S}}$:

$$T^\pi v(s) = \sum_{a \in \mathcal{A}_s} \pi(a; s) \left(r(s, a) + \sum_{s' \in \mathcal{S}} p(s'; s, a) v(s') \right).$$

and that the value function V^π is a fixed point of the Bellman operator T^π .



The one-step reward $r(s, a)$ is a shortcut for $r(s, a) = \sum_{s' \in \mathcal{S}} \sum_{r \in \mathcal{R}} r p(s', r; s, a)$, sometimes we prefer to write the detailed form in order to emphasise more explicitly the dependence on the Markov decision model. The Bellman operator then takes the form

$$(T^\pi v)(s) = \sum_{a \in \mathcal{A}_s} \pi(a; s) \sum_{s' \in \mathcal{S}} \sum_{r \in \mathcal{R}} p(s', r; s, a) [r + \gamma v(s')], \quad s \in \mathcal{S}.$$

The equation looks complicated but (\mathcal{S} is assumed finite) is just a system of linear equations that can be solved by means of linear algebra techniques. This becomes directly evident when we rewrite the fixed point equation in vector notation

$$V^\pi = r^\pi + \gamma P^\pi V^\pi,$$

where

$$P^\pi = \left(\sum_{a \in \mathcal{A}_s} \pi(a; s) p(s'; s, a) \right)_{(s, s') \in \mathcal{S} \times \mathcal{S}}$$

$$r^\pi = \left(\sum_{a \in \mathcal{A}_s} \pi(a; s) r(s, a) \right)_{s \in \mathcal{S}}$$

with $P^\pi \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$ and $r^\pi, V^\pi \in \mathbb{R}^{|\mathcal{S}|}$.



Please check that indeed $T^\pi = r^\pi + \gamma P^\pi$.

Given different use cases, each of these notations may be favorable. The reader may forgive our shifty notation in favor of uncluttered statements.



Proposition 2.3.1. The (affine) linear equation $V = r^\pi + \gamma P^\pi V$ has a unique solution. Hence, V^π can be computed explicitly as

$$V^\pi = (I - \gamma P^\pi)^{-1} r^\pi. \quad (2.17)$$

Proof. This follows immediately as the Bellman expectation operator is a contraction (compare the exercise after Theorem 2.1.26). To compute the solution is straightforward linear algebra operation:

$$\begin{aligned} V^\pi = r^\pi + \gamma P^\pi V^\pi &\Leftrightarrow (I - \gamma P^\pi) V^\pi = r^\pi \\ &\Leftrightarrow V^\pi = (I - \gamma P^\pi)^{-1} r^\pi \end{aligned}$$

The existence of the inverse is guaranteed as the linear equation has a unique solution. \square

As a consequence we see that the evaluation of a policy simply corresponds to the inversion of a matrix. However, although this computation is straightforward, it is a tedious and expensive computation. The complexity of matrix inversion (using Gauss-Jordan elimination) is $\mathcal{O}(n^3)$ and $n = |\mathcal{S}|$ can be huge. Thus we will also explore iterative solution methods.

As seen with value iteration we can use the convergence properties of the iterates of a contraction mapping. We want to do the same for the operator T^π for a stationary policy. Using Banach's fixed point theorem convergence $(T^\pi)^n v_0 \rightarrow V^\pi$ holds for any initialisation v_0 . Implemented as an algorithm we obtain Algorithm 8.

Data: Policy $\pi \in \Pi_S$, $\varepsilon > 0$
Result: Approximation $V \approx V^\pi$
Initialize $V \equiv 0, V_{\text{new}} \equiv 0$
 $\Delta := 1$
while $\Delta > \varepsilon$ **do**
 Set $V = V_{\text{new}}$
 for $s \in \mathcal{S}$ **do**
 $V_{\text{new}}(s) = \sum_{a \in \mathcal{A}_s} \pi(a; s) \sum_{s' \in \mathcal{S}} \sum_{r \in \mathcal{R}} p(s', r; s, a) [r + \gamma V(s')]$
 end
 $\Delta := \max_{s \in \mathcal{S}} |V_{\text{new}}(s) - V(s)|$
end
 $V := V_{\text{new}}$

Algorithm 8: Iterative policy evaluation (Naive)

Theorem 2.3.2. Algorithm 8 terminates and $\|V - V^\pi\|_\infty < \frac{\gamma \varepsilon}{1 - \gamma}$.

Proof. The proof is identical to the proof of Theorem 2.2.1 which is only based on the fact that T^* is a contraction with constant γ . Since T^π is equally a contraction with constant γ the same result holds. \square

Algorithm 8 can be improved in terms of memory. The simple fixed point iteration algorithm requires to occupy memory for $2|\mathcal{S}|$ values since V has to be fully stored in order to compute every value $V_{\text{new}}(s)$. This can be done more efficient by directly using available data, see Algorithm 9. The algorithm does not perform the matrix computation with T^π directly, but row by row, updating coordinate by coordinate instead of all coordinates at once.

Data: Policy $\pi \in \Pi_S$, $\varepsilon > 0$
Result: Approximation $V \approx V^\pi$
Initialize $V \equiv 0$
 $\Delta := \varepsilon$
while $\Delta > \varepsilon$ **do**
 $\Delta := 0$
 for $s \in \mathcal{S}$ **do**
 $v := V(s)$
 $V(s) := \underbrace{\sum_{a \in \mathcal{A}_s} \pi(a; s) \sum_{s' \in \mathcal{S}} \sum_{r \in \mathcal{R}} p(s', r; s, a) [r + \gamma V(s')]}_{T^\pi V(s) = (r^\pi + P^\pi)V(s)}$
 $\Delta := \max(\Delta, |v - V(s)|)$
 end
end

Algorithm 9: Iterative policy evaluation (totally asynchronous updates)

Try to prove convergence of the totally asynchronous policy evaluation algorithm (without termination) to V^π . To do so enumerate \mathcal{S} as s_1, \dots, s_K and define

$$T_s^\pi V(s') = \begin{cases} T^\pi V(s) & : s = s' \\ V(s) & : s \neq s' \end{cases},$$

i.e. T^π is only applied to coordinate s while leaving all other coordinates unchanged. Then the inner loop of the algorithm performs nothing but the composition $\bar{T}^\pi := (T_{s_K}^\pi \circ \dots \circ T_{s_1}^\pi)(V)$, which is not the same as applying T^π ! Show that V^π is a fixed point of the composition and the composition is a contraction on $(U, \|\cdot\|_\infty)$, proceed



step by step using the estimates to show that Bellman operators are contractions. Without the termination the outer loop is nothing but an iteration of \bar{T}^π , hence, without termination the algorithm converges to V^π .

Algorithms in which coordinates s are treated differently are called asynchronous algorithms. In fact, there are other versions of the algorithm in which coordinates are not swept in order but randomly. We will come back to such asynchronous algorithms in the next chapter (e.g. Q -learning is of exactly that kind replacing Bellman's expectation operator by the state-action optimality operator).

2.3.2 Policy improvement

We now aim to improve a given policy, that is, to slightly change it such that its value function takes larger values. Mathematically speaking, for a policy $\pi \in \Pi$ and a value function V^π we aim to find a policy $\pi' \in \Pi$ such that

$$V^\pi(s) \leq V^{\pi'}(s), \quad \forall s \in \mathcal{S}.$$

The improvement is called strict if

$$V^\pi(s) < V^{\pi'}(s) \quad \text{for at least one } s \in \mathcal{S}.$$

If π is not optimal, there always exists a strict improvement, e.g. the optimal policy. We now want to define a procedure to update a non-optimal policy to an improved policy. The key idea is to change the policy at a single state $s \in \mathcal{S}$ to a particular action. For this we look at the action-value function of a stationary policy $\pi \in \Pi_S$. Recalling that

$$V^\pi(s) = \sum_{a \in \mathcal{A}_s} \pi(a; s) Q^\pi(s, a)$$

it becomes apparent that

$$\max_{a \in \mathcal{A}_s} Q^\pi(s, a) \geq V^\pi(s). \quad (2.18)$$

In other words, the inequality above implies that choosing an action $a \in \mathcal{A}$ that maximizes the expected reward in state s and the expected future value of following a policy π is at least as good as following the policy π . Recalling Definition 2.1.28 this suggests to use the greedy policy π_{Q^π} induced by the Q -function of the current policy π which then leads to the simple policy improvement algorithm. The improvement of the greedy policy improvement is a special case of

Data: policy $\pi \in \Pi_S$, Q -function Q^π
Result: improved policy $\pi' = \text{greedy}(Q^\pi)$
 $\pi' := \pi$
for $s \in \mathcal{S}$ **do**
 Choose $a^*(s) \in \arg \max_{a \in \mathcal{A}_s} Q^\pi(s, a)$
 $\pi'(a^*(s); s) := 1$
 for $a \in \mathcal{A}_s \setminus \{a^*(s)\}$ **do**
 $\pi'(a; s) = 0$
 end
end

Algorithm 10: Greedy policy improvement

the policy improvement theorem.


Theorem 2.3.3. (Policy improvement theorem)

Let $\pi, \pi' \in \Pi_S$ be two stationary policies, then the following hold:

(i) If

$$V^\pi(s) \leq \sum_{a \in \mathcal{A}_s} \pi'(a; s) Q^\pi(s, a), \quad \forall s \in \mathcal{S}, \quad (2.19)$$

then π' is an improvement of π .

(ii) If there is a strict inequality in (2.19) for some state s then the improvement is strict.

(iii) For every policy $\pi \in \Pi_S$ the greedy policy obtained from Q^π improves π .

On a heuristic level the theorem is trivial. If a stationary policy π' is better for one step (in all states) then (by the Markov property) the policy will also lead to a larger reward if it is used in all time-steps. But then the value function for π' is bigger than that for π .

Proof. (i) Using definitions, assumptions and the definition of Bellman expectation operators we get

$$\begin{aligned} V^\pi(s) &\leq \sum_{a \in \mathcal{A}_s} \pi'(a; s) Q^\pi(s, a) \\ &= \sum_{a \in \mathcal{A}_s} \pi'(a; s) \left(r(s, a) + \sum_{s' \in \mathcal{S}} p(s'; s, a) V^\pi(s') \right) = T^{\pi'} V^\pi(s). \end{aligned}$$

Monotonicity of Bellman operators allows us to iterate the equation to obtain

$$V^\pi(s) \leq T^{\pi'} V^\pi(s) \leq T^{\pi'} T^{\pi'} V^\pi(s) \leq \dots \leq \lim_{k \rightarrow \infty} (T^{\pi'})^k V^\pi(s).$$

By Banach's fixed point theorem the iterations of $T^{\pi'}$ converge uniformly (thus pointwise) to the unique fixed point of $T^{\pi'}$, which is $V^{\pi'}$. Thus, $V^\pi(s) \leq V^{\pi'}(s)$ for all $s \in \mathcal{S}$.

(ii) For strict inequalities the above steps imply strict inequalities.

(iii) To check that the greedy policy update satisfies the one-step improvement property is easy. With $a^*(s) \in \arg \max_{a \in \mathcal{A}_s} Q^\pi(s, a)$, the greedy policy π' satisfies the improvement condition:

$$V^\pi(s) = \sum_{a \in \mathcal{A}_s} \pi(a; s) Q^\pi(s, a) \leq \max_{a \in \mathcal{A}_s} Q^\pi(s, a) \underbrace{\sum_{a \in \mathcal{A}_s} \pi(a; s)}_{=1} = \sum_{a \in \mathcal{A}_s} \pi'(a; s) Q^\pi(s, a)$$

Hence, (i) implies the claim. \square

For the following section we also want to prove the following lemma that links greedy policy improvement to optimal policies. After all, this is still what we are on the look for.



Proposition 2.3.4. Let $\pi \in \Pi_S$ and π' the greedy policy obtained from Q^π , then

$$V^\pi = V^{\pi'} \text{ (or } Q^\pi = Q^{\pi'}) \implies \pi \text{ and } \pi' \text{ are optimal.}$$

Proof. It follows from Lemma 2.1.20 that equality of the value functions implies equality of the state-action value functions, hence, we work with Q . As in the previous proof we compute, using $a^*(s) \in \arg \max_{a \in \mathcal{A}_s} Q^\pi(s, a)$,

$$\begin{aligned} Q^{\pi'}(s, a) &= r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \sum_{a' \in \mathcal{A}_{s'}} p(s'; s, a) \pi'(a'; s) Q^{\pi'}(s', a') \\ &= r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'; s, a) Q^{\pi'}(s', a^*(s')). \end{aligned}$$

Now suppose that $Q^\pi = Q^{\pi'}$, then the equation becomes

$$Q^\pi(s, a) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'; s, a) \max_{a' \in \mathcal{A}_s} Q^\pi(s', a').$$

Hence, Q^π and $Q^{\pi'}$ both solve Belman's state-action optimality equation. Since this has a unique solution we deduce $Q^{\pi'} = Q^* = Q^\pi$ (and both are optimal). \square



Corollary 2.3.5. For a non-optimal policy $\pi \in \Pi_S$, the greedy policy $\pi' = \text{greedy}(Q^\pi)$ obtained from Q^π is a strict policy improvement.

Proof. The claim follows directly by Lemma 2.3.4. \square

Apart from the greedy policies there is a second interesting application of the policy improvement theorem. For that sake we need a bit of extra notation.



Definition 2.3.6. A policy $\pi \in \Pi_S$ is called

- soft, if it fulfils

$$\pi(a; s) > 0 \quad \forall s \in \mathcal{S}, a \in \mathcal{A}_s,$$

- ε -soft for some $1 \geq \varepsilon > 0$, if it fulfils

$$\pi(a; s) > \frac{\varepsilon}{|\mathcal{A}_s|} \quad \forall s \in \mathcal{S}, a \in \mathcal{A}_s,$$

- ε -greedy with regard to Q , if it selects the greedy action with respect to Q with probability $(1 - \varepsilon)$ and a (uniform) random action with probability ε , i.e.

$$\pi(a; s) = \begin{cases} (1 - \varepsilon) + \frac{\varepsilon}{|\mathcal{A}_s|} & : a = a^*(s) \\ \frac{\varepsilon}{|\mathcal{A}_s|} & : a \text{ otherwise} \end{cases},$$

where $a^*(s) \in \arg \max_a Q(s, a)$.

Let us recall the discussion of the ε -greedy learning strategies for stochastic bandits. Such policies are considered suboptimal as they lead to linear regret, they play suboptimal arms with probability ε . Similarly, ε -soft policies cannot be optimal as suboptimal actions must be played in contrast to greedy policies that only play optimal actions. Hence, ε -greedy policies will typically not improve policies, but they do improve all other ε -soft policies:



Proposition 2.3.7. If $\pi \in \Pi_S$ is ε -soft, then the ε -greedy policy π' obtained from Q^π improves π .

Proof. This follows by checking the condition from the policy improvement theorem:

$$\begin{aligned} V^\pi(s) &= \frac{\varepsilon}{|\mathcal{A}_s|} \sum_{a \in \mathcal{A}_s} Q^\pi(s, a) - \frac{\varepsilon}{|\mathcal{A}_s|} \sum_{a \in \mathcal{A}_s} Q^\pi(s, a) + \sum_a \pi(a; s) Q^\pi(s, a) \\ &= \frac{\varepsilon}{|\mathcal{A}_s|} \sum_{a \in \mathcal{A}_s} Q^\pi(s, a) + (1 - \varepsilon) \sum_a \underbrace{\frac{\pi(a; s) - \frac{\varepsilon}{|\mathcal{A}_s|}}{1 - \varepsilon}}_{\geq 0, \varepsilon\text{-greedy}} Q^\pi(s, a) \\ &\leq \frac{\varepsilon}{|\mathcal{A}_s|} \sum_{a \in \mathcal{A}_s} Q^\pi(s, a) + (1 - \varepsilon) \max_a Q^\pi(s, a) \\ &= \sum_{a \in \mathcal{A}_s} \pi'(a; s) Q^\pi(s, a). \end{aligned}$$

\square

2.3.3 Policy iteration algorithms (tabular actor-critic)

The ingredients developed above can now be combined to obtain the policy iteration algorithm. The idea is simple: alternate policy evaluation and improvement, compute V^π and then improve to π' using the greedy strategy obtained from Q^π . The above results show that every improvement step improves the policy and the limit of the procedure is π^* . Here is an illustration of the procedure:

$$\pi_0 \nearrow V^{\pi_0} \searrow \pi_1 \nearrow V^{\pi_1} \searrow \pi_2 \nearrow \cdots \searrow \pi^*$$

The algorithm obtained this way is called policy iteration algorithm. Since there are many variations how to perform the policy evaluation (exact or approximatedly) we will meet several variants in the chapters below. We will first restrict ourselves to exact evaluations of V^π using

Data: initial policy $\pi_0 \in \Pi_{\mathcal{S}}$
Result: optimal policy $\pi^* \in \Pi_{\mathcal{S}}^D$
 initialise arbitrarily $V_{\text{new}}, \pi_{\text{new}}$
 stop = *False*
while stop = *False* **do**
 Policy evaluation (critic): Obtain V^π by computing (2.17).
 set $Q^\pi(s, a) = \sum_{s' \in \mathcal{S}} \sum_{r \in \mathcal{R}} p(s', r; s, a) [r + \gamma V^\pi(s')]$ for all a, s
 Policy improvement (actor): Obtain the improved greedy policy π_{new} from Q^π
 if $V^{\pi_{\text{new}}} = V^\pi$ **then**
 | stop = *True*
 end
end
 return $\pi^* = \pi_{\text{new}}$

Algorithm 11: Greedy exact policy iteration (actor-critic)

the matrix inversion (2.17).



Theorem 2.3.8. (Greedy exact policy iteration)

Started in any policy the policy iteration Algorithm 12 with exact policy evaluation and greedy policy update for finite MDPs terminates in a finite number of iterations (at most $|\mathcal{A}| \cdot |\mathcal{S}|$) with a solution of the optimality equation and an optimal policy π^* .

Proof. By Corollary 2.3.5 in each iteration there is a strict improvement of the next policy π' (i.e. there exists at least one $s \in \mathcal{S}$ such that $V^\pi(s) < V^{\pi'}(s)$) and the set of deterministic stationary strategies is finite ($|\Pi_{\mathcal{S}}^D| = |\mathcal{S}|^{|\mathcal{A}|} < \infty$) the algorithm has to terminate in finitely many steps. By Lemma 2.3.4 the termination of the algorithm, thus $V^\pi = V^{\pi'}$, implies that π' is optimal. \square

Also for infinite state and action space the exact policy iteration algorithm converges monotonically and in norm to the optimal policy. Since the exact policy evaluation is hardly possible if \mathcal{S} and/or \mathcal{A} are huge we do not go further into the analysis



Typically the value functions V^{π_n} will not be computed explicitly but the policy iteration algorithm will be used in an approximate manner where V^π is estimated. One example is to replace the explicit evaluation of V^π by a few steps of the Banach fixed point iteration corresponding to T^π , compare Algorithms 8 or 9. Other examples will be discussed in Chapter ???. All algorithms alternating between value function estimation a policy improvement (not necessarily greedy)

$$\pi_0 \nearrow \underbrace{\hat{V}^{\pi_0}}_{\approx V^{\pi_0}} \searrow \pi_1 \nearrow \underbrace{\hat{V}^{\pi_1}}_{\approx V^{\pi_1}} \searrow \pi_2 \nearrow \cdots \searrow \pi$$

will be called **generalised policy iteration algorithm**. The aim will be to



generate algorithms that converge as quickly as possible to a policy π which is as close to π^* as possible.

Data: initial policy π

Result: optimal policy π^* (or approximation of π^*)

while *not converged* **do**

 Policy evaluation (critic): Obtain estimates for \hat{V}^π and/or \hat{Q}^π from **some** algorithm.

 Policy improvement (actor): Obtain (hopefully improved) policy π_{new} by **some** algorithm.

 Set $\pi = \pi_{new}$.

end

return π

Algorithm 12: Generalised policy iteration (actor-critic) paradigm

It is not clear at all if a generalised policy iteration algorithm converges to an optimal policy and when to stop the algorithm! The errors made by policy evaluation and improvement might easily build up. It is a very non-trivial task to find and tune approximate algorithms that converge. In Section ?? will come back to generalised policy iteration with ε -greedy policies.



The notion „while not converged“ in algorithm pseudo code will always refer to a non-specified termination condition. In many algorithms we will specify a concrete termination condition.

An interesting version of generalised policy iteration comes in combination with ε -greedy policies that will be useful later to introduce exploitation into some algorithms. Let us extend the notion of optimality to the class of soft policies.



Definition 2.3.9. An ε -soft policy π^* is called ε -soft optimal if

$$V^{\pi^*}(s) = \sup_{\pi \text{ } \varepsilon\text{-soft}} V^\pi(s) =: \tilde{V}^*(s), \quad \forall s \in \mathcal{S}.$$

There is a nice trick how to show convergence of ε -greedy policy iteration. Since ε -greedy policies play some action with probability ε they will never be optimal (except in trivial cases) so they won't converge to an optimal policy. Nonetheless, they do converge to a policy which is optimal among the ε -soft policies.



Theorem 2.3.10. (ε -greedy exact policy iteration)

Started in any ε -soft policy the generalised policy iteration algorithm with exact policy evaluation and ε -greedy policy update converges to an ε -soft optimal policy.

Proof. To prove convergence of exact policy iteration with greedy policy update (Theorem 2.3.8) used the Bellman equation to guarantee that no further improvement means that the current policy is already optimal. This does not work for the ε -greedy policy iteration algorithm. To circumvent this problem we use a trick and “move the ε -softness into the environment”. For that sake let us define a new MDP with transition function

$$\tilde{p}(s', r; s, a) := (1 - \varepsilon)p(s', r; s, a) + \frac{\varepsilon}{|\mathcal{A}_s|} \sum_{b \in \mathcal{A}_s} p(s', r; s, b).$$

This means, that with probability ε , the transition kernel will ignore the selected action and behave as if a uniformly random action was chosen. We can transform stationary ε -soft policies π from the old MDP to stationary policies $\tilde{\pi}$ of the new MDP via

$$\tilde{\pi}(a; s) := \frac{\pi(a; s) - \frac{\varepsilon}{|\mathcal{A}_s|}}{1 - \varepsilon} \geq 0.$$

Let us denote the mapping by $f : \pi \mapsto \tilde{\pi}$. Conversely, for every stationary policy $\tilde{\pi}$ of the transformed MDP we can define the ε -soft policy π by

$$\pi(a; s) := (1 - \varepsilon)\tilde{\pi}(a; s) + \frac{\varepsilon}{|\mathcal{A}_s|},$$

which is the inverse mapping. Therefore f is a bijection between the ε -soft policies in the old MDP and all stationary policies in the new MDP. We now show, that the value functions V^π stay invariant with regard to this mapping. For that sake note that

$$\tilde{p}(s'; s, a) = (1 - \varepsilon)p(s'; s, a) + \frac{\varepsilon}{|\mathcal{A}_s|} \sum_{b \in \mathcal{A}_s} p(s'; s, b)$$

and

$$\tilde{r}(s, a) = (1 - \varepsilon)r(s, a) + \frac{\varepsilon}{|\mathcal{A}_s|} \sum_{b \in \mathcal{A}_s} r(s, b),$$

hence,

$$\begin{aligned} \sum_{a \in \mathcal{A}_s} \tilde{\pi}(a; s) \tilde{r}(s, a) &= \sum_{a \in \mathcal{A}_s} \left(\frac{\tilde{\pi}(a; s) - \frac{\varepsilon}{|\mathcal{A}_s|}}{1 - \varepsilon} \right) \left((1 - \varepsilon)r(s, a) + \frac{\varepsilon}{|\mathcal{A}_s|} \sum_{b \in \mathcal{A}_s} r(s, b) \right) \\ &= \sum_{a \in \mathcal{A}_s} \left(\tilde{\pi}(a; s) - \frac{\varepsilon}{|\mathcal{A}_s|} \right) r(s, a) + \frac{\varepsilon}{|\mathcal{A}_s|} \sum_{b \in \mathcal{A}_s} r(s, b) \underbrace{\sum_{a \in \mathcal{A}_s} \frac{\tilde{\pi}(a; s) - \frac{\varepsilon}{|\mathcal{A}_s|}}{1 - \varepsilon}}_{=1} \\ &= \sum_{a \in \mathcal{A}_s} \tilde{\pi}(a; s) r(s, a). \end{aligned}$$

Similarly:

$$\begin{aligned} \sum_{a \in \mathcal{A}_s} \tilde{\pi}(a; s) \tilde{p}(y; s, a) &= \sum_{a \in \mathcal{A}_s} \left(\frac{\tilde{\pi}(a; s) - \frac{\varepsilon}{|\mathcal{A}_s|}}{1 - \varepsilon} \right) \left((1 - \varepsilon)p(y; s, a) + \frac{\varepsilon}{|\mathcal{A}_s|} \sum_{b \in \mathcal{A}_s} p(y; s, b) \right) \\ &= \sum_{a \in \mathcal{A}_s} \tilde{\pi}(a; s) p(y; s, a) \end{aligned}$$

Combining the above yields

$$\begin{aligned} \tilde{T}^{\tilde{\pi}} V^\pi(s) &= \sum_{a \in \mathcal{A}_s} \tilde{\pi}(a; s) \left[\tilde{r}(s, a) + \gamma \sum_{y \in \mathcal{S}} \tilde{p}(y; s, a) V^\pi(y) \right] \\ &= \sum_{a \in \mathcal{A}_s} \tilde{\pi}(a; s) \left[r(s, a) + \gamma \sum_{y \in \mathcal{S}} p(y; s, a) V^\pi(y) \right] \\ &= T^\pi V^\pi(s) = V^\pi(s). \end{aligned}$$

Since the fixed point is unique it follows that $\tilde{V}^{\tilde{\pi}} = V^\pi$ and, as f is bijective,

$$\sup_{\tilde{\pi} \in \tilde{\Pi}_S} \tilde{V}^{\tilde{\pi}}(s) = \sup_{\pi \text{ } \varepsilon\text{-soft}} V^\pi(s), \quad (2.20)$$

For the Q -functions we obtain

$$\begin{aligned} \tilde{Q}^{\tilde{\pi}}(s, a) &= \tilde{r}(s, a) + \gamma \sum_{y \in \mathcal{S}} \tilde{p}(y; s, a) V^\pi(y) \\ &= (1 - \varepsilon) \left(r(s, a) + \gamma \sum_{y \in \mathcal{S}} p(y; s, a) V^\pi(y) \right) \\ &\quad + \frac{\varepsilon}{|\mathcal{A}_s|} \sum_{b \in \mathcal{A}_s} \left(r(s, b) + \gamma \sum_{y \in \mathcal{S}} p(y; s, b) V^\pi(y) \right) \\ &= (1 - \varepsilon) Q^\pi(s, a) + \frac{\varepsilon}{|\mathcal{A}_s|} \sum_{b \in \mathcal{A}_s} \left(r(s, b) + \gamma \sum_{y \in \mathcal{S}} p(y; s, b) V^\pi(y) \right), \end{aligned}$$

which implies that

$$\arg \max_{a \in \mathcal{A}_s} \tilde{Q}^{\tilde{\pi}}(s, a) = \arg \max_{a \in \mathcal{A}_s} Q^{\pi}(s, a).$$

Therefore greedy with respect to Q^{π} and $\tilde{Q}^{\tilde{\pi}}$ is the same. Let π_n be an ε -soft policy, and let π_{n+1} be ε -greedy with regard to Q^{π_n} . Then $\tilde{\pi}_{n+1} := f(\pi_{n+1})$ is greedy with respect to $\tilde{Q}^{\tilde{\pi}_n}$:

$$\tilde{\pi}_{n+1}(a; s) = \frac{\pi(a; s) - \frac{\varepsilon}{|\mathcal{A}_s|}}{1 - \varepsilon} = \begin{cases} \frac{\left(\frac{(1-\varepsilon) + \frac{\varepsilon}{|\mathcal{A}_s|}}{|\mathcal{A}_s|} - \frac{\varepsilon}{|\mathcal{A}_s|}\right) - \frac{\varepsilon}{|\mathcal{A}_s|}}{1 - \varepsilon} = 1 & : a = a^*(s) \\ \frac{\frac{\varepsilon}{|\mathcal{A}_s|} - \frac{\varepsilon}{|\mathcal{A}_s|}}{1 - \varepsilon} = 0 & : a \text{ otherwise} \end{cases}.$$

Since we proved in Theorem 2.3.8 convergence of exact policy iteration with greedy policy updates the proof can be finished:

$$\sup_{\pi \text{ } \varepsilon\text{-soft}} V^{\pi}(s) \stackrel{(2.20)}{=} \sup_{\tilde{\pi} \in \tilde{\Pi}_S} \tilde{V}^{\tilde{\pi}}(s) = \tilde{V}^{\tilde{\pi}^*}(s) = \lim_{n \rightarrow \infty} \tilde{V}^{\tilde{\pi}_n}(s) = \lim_{n \rightarrow \infty} V^{\pi_n}(s)$$

□

On first view it is completely unclear why there might be any interest in ε -greedy policy iteration if we already have greedy policy iteration that converges to an optimal policy. The reason, as we will discuss in the next chapter, comes from the policy evaluation step. Greedy policies can be very unfavorable to estimate Q -values or the value function if those cannot be computed explicitly. In contrast, ε -greedy policies have pleasant exploration properties, they force the algorithm to look at all actions and not only the ones that are already known to be good. The reader might compare with the exploration-exploitation trade-off for stochastic bandits in Chapter 1.

2.3.4 Relation of policy iteration and value iteration

There is an important simple observation that we should always keep in mind as it reappears later in the discussion of SARSA and Q -learning. Value iteration is a short-cut of policy iteration, it combines the two phases of policy iteration in one. Why is that?



Recall that $T^*v = T^{\pi^*}v$ so that

$$(T^*v)(s) = \max_{a \in \mathcal{A}_s} \left\{ r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'; s, a)v(s') \right\} = (T^{\pi^*}v)(s).$$

It follows that value iteration is exactly the same as policy iteration, where the critic step evaluates V^{π} with just one Bellman expectation update.

One can quite rightly ask whether policy iteration or value iteration is the better algorithm. This comparison shows that value iteration can be more efficient if

2.4 Stochastic control in finite time

So far we discussed stochastic control problems with geometrically distributed random time horizon (equivalently, infinite time horizon with discounted rewards). In this section the time horizon will be a fixed deterministic time T . The geometric time horizon is a simpler problem as the forgetfulness of geometric random variables leads to stationary optimal policies. This is different for deterministic time horizon, optimal policies $(\pi_t^*)_{t \leq T}$ will be non-stationary! ³

³at some point also include a terminal payout (which now is 0)

2.4.1 Setting and dynamic programming

As in the previous section we will consider $(S_t, A_t, R_t)_{t \leq T}$ with the only difference that time is restricted to $D = \{0, \dots, T\}$ for some $T \in \mathbb{N}$ fixed. Action and state spaces as well as transition probabilities p remain unchanged, policies $\pi = (\pi_t : t \in D)$ are defined as before. The finite-time stochastic control problem consists in optimising the state value function $\mathbb{E}_s^\pi[\sum_{t=0}^{T-1} R_t]$ over all policies.

Example 2.4.1. An example to keep in mind is the ice vendor who has a three month summer season (Germans also love ice cream during winter, most others don't). The optimal production strategy of the vendor will depend upon the distance to the terminal day, there is no need in having ice cream left after the last day of the season. An optimal policy of time-dependent problems will henceforth never be stationary!

Compared to infinite time-horizon MDPs crucial differences appear. Most importantly, the idea of dynamic programming (reduce the problem to simpler sub-problems) becomes much clearer. The general idea of dynamic programming is to reduce a problem to a smaller problem and then build a solution to a larger problem from solutions of smaller problems. For infinite horizon control the idea did not become very explicit as the reduced problem (starting one time-step later) is identical to the original problem (multiplied by γ). Thus, the reduction attempt only led to a set of equations. Here, the finite time horizon forces the reduced problem (starting one time-step later) to be simpler. The resulting set of equations will turn out to be a backwards recursion instead of a closed system of equations.

Let us write $\pi \in \Pi_t^T$ to denote that $\pi = (\pi_i)_{i=t}^{T-1}$ is a policy that runs from time t to T , i.e. π consists of $T - t$ Markov kernels. The definition of the state value function for T -step MDPs is as follows.



Definition 2.4.2. For any policy $\pi \in \Pi_t^T$ the functions (vectors) defined by $V_{T,T}^\pi \equiv 0$ and

$$V_{t,T}^\pi : \mathcal{S} \rightarrow \mathbb{R}, \quad s \mapsto \mathbb{E}_s^{\tilde{\pi}} \left[\sum_{i=0}^{T-t-1} R_i \right] =: \mathbb{E}_{S_t=s}^\pi \left[\sum_{i=t}^{T-1} R_i \right],$$

for $t < T$ are called time-state value functions, where $\tilde{\pi}$ is the policy π shifted by t , i.e. $\tilde{\pi}_i = \pi_{t+i}$ for $i = 0, \dots, T - t - 1$. The function $V_{0,T}^\pi$ is called state value function.

Typically the time-state value function is defined as $V_{t,T}^\pi(s) = \mathbb{E}[\sum_{i=t}^{T-1} R_i | S_t = s]$. To avoid discussions of conditioning on zero-sets we shift the starting time to 0 and force the start in s . This is rigorous (not pretty) and we keep in mind that $V_{t,T}^\pi$ is the total reward gained after time t . We also have to redefine the state-action value function



Definition 2.4.3. For any policy $\pi \in \Pi_t^T$ the functions (matrices) defined by $Q_{T,T}^\pi \equiv 0$ and

$$Q_{t,T}^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}, \quad (s, a) \mapsto \mathbb{E}_{s,a}^{\tilde{\pi}} \left[\sum_{i=0}^{T-t-1} R_i \right] =: \mathbb{E}_{S_t=s, A_t=a}^\pi \left[\sum_{i=t}^{T-1} R_i \right],$$

for $t < T$ are called time-state-action value functions. The shifted policy $\tilde{\pi}$ is defined as in the previous definition. The function $Q_{0,T}^\pi$ is called state-action value function.

From now on we will drop the subscript T . We will just write V_t^π and Q_t^π , except in situation in which the emphasise lies on T . As for discounted MDPs the difference between V and Q is

that the first action is fixed to be a for Q . Similarly to Lemma 2.1.20, we also get a lemma that describes the relation between the state value function and the state-action value function for a fixed policy:



Proposition 2.4.4. Given a Markovian policy $\pi = (\pi_t)_{t \leq T}$ and a T -step Markov decision problem. Then the following relation between the state and state-action value function hold

$$V_t^\pi(s) = \sum_{a \in \mathcal{A}_s} \pi_t(a; s) Q_t^\pi(s, a),$$

$$Q_t^\pi(s, a) = r(s, a) + \sum_{s' \in \mathcal{S}} p(s'; s, a) V_{t+1}^\pi(s')$$

for all $t < T$. In particular, the Bellman expectation equations (backwards recursions)

$$V_t^\pi(s) = \sum_{a \in \mathcal{A}} \pi_t(a; s) \left[r(s, a) + \sum_{s' \in \mathcal{S}} p(s'; s, a) V_{t+1}^\pi(s') \right],$$

$$Q_t^\pi(s, a) = r(s, a) + \sum_{s' \in \mathcal{S}} \sum_{a' \in \mathcal{A}_{s'}} p(s'; s, a) \pi_{t+1}(a'; s') Q_{t+1}^\pi(s', a')$$

hold for $t < T$.

Proof. Expressing V_t in terms of Q_t is a matter of definition of \mathbb{E}_s , just as in the infinite horizon case. For $t < T$ we can use the Markov reward property of (S, A, R) and the formula of total probability to derive the recursions. This is exactly the same as in the proof of Proposition 2.1.19. Let's add the proof for completeness. Recall from Proposition 2.1.17⁴ that (S, A, R) is a Markov reward process so that by (2.7)

$$\mathbb{E}_{s,a}^\pi \left[R_1 + R_2 + \dots + R_{T-t-1} \mid S_1 = s', A_1 = a' \right] = \mathbb{E}_{s',a'}^{\tilde{\pi}} \left[R_0 + R_1 + \dots + R_{T-t-2} \right] = Q_{t+1}^\pi(s', a').$$

Thus, using the formula of total probability,

$$\begin{aligned} & Q_t^\pi(s, a) \\ &= \mathbb{E}_{s,a}^{\tilde{\pi}} \left[\sum_{i=0}^{T-t-1} R_i \right] \\ &= \mathbb{E}_{s,a}^{\tilde{\pi}} [R_0] + \mathbb{E}_{s,a}^{\tilde{\pi}} [R_1 + R_2 \dots + R_{T-t-1}] \\ &= r(s, a) + \sum_{s' \in \mathcal{S}, a' \in \mathcal{A}_s} \mathbb{E}_{s,a}^{\tilde{\pi}} \left[R_1 + R_2 + \dots + R_{T-t-1} \mid S_1 = s', A_1 = a' \right] \mathbb{P}_{s,a}^{\tilde{\pi}}(S_1 = s', A_1 = a') \\ &= r(s, a) + \sum_{s' \in \mathcal{S}, a' \in \mathcal{A}_s} \mathbb{P}_{s,a}^{\tilde{\pi}}(S_1 = s', A_1 = a') Q_{t+1}^\pi(s', a') \\ &= r(s, a) + \sum_{s' \in \mathcal{S}} \sum_{a' \in \mathcal{A}_{s'}} p(s'; s, a) \tilde{\pi}_1(a'; s') Q_{t+1}^\pi(s', a') \\ &= r(s, a) + \sum_{s' \in \mathcal{S}} \sum_{a' \in \mathcal{A}_{s'}} p(s'; s, a) \pi_{t+1}(a'; s') Q_{t+1}^\pi(s', a'). \end{aligned}$$

⁴at some point rewrite proposition for Markov policies, works equally

The equation for V follows directly by plugging-in $V_t^\pi(s) = \sum_{a \in \mathcal{A}_s} \pi_t(a; s) Q_t^\pi(s, a)$ twice:

$$\begin{aligned} V_t^\pi(s) &= \sum_{a \in \mathcal{A}_s} \pi_t(a; s) Q_t^\pi(s, a) \\ &= \sum_{a \in \mathcal{A}_s} \pi_t(a; s) \left(r(s, a) + \sum_{s' \in \mathcal{S}} \sum_{a' \in \mathcal{A}_{s'}} p(s'; s, a) \pi_{t+1}(a'; s') Q_{t+1}^\pi(s', a') \right) \\ &= \sum_{a \in \mathcal{A}_s} \pi_t(a; s) \left(r(s, a) + \sum_{s' \in \mathcal{S}} p(s'; s, a) V_{t+1}^\pi(s') \right). \end{aligned}$$

□

It is important to note that the system of equations is different from discounted MDP problems. First, the discounting factor disappears as $\gamma = 1$ and, secondly, the linear systems are actually recursions that gradually simplify the system towards the terminal conditions $Q_T^\pi \equiv 0$, $V_T^\pi \equiv 0$. Or, reversed, starting with the zero-vector (resp. zero-matrix) a backward induction allows to compute V_t and Q_t recursively. There is no system of equations that needs to be solved!

Similarly to the discounted setting the optimal value functions are defined, now depending on the remaining time-horizon:



Definition 2.4.5. For any $t \leq T$ and a given Markov decision problem

- the function $V_t^* : \mathcal{S} \rightarrow \mathbb{R}$ that takes values

$$V_t^*(s) := \sup_{\pi \in \Pi_t^T} V_t^\pi(s), \quad s \in \mathcal{S},$$

is called optimal time-state value function.

- the function $Q_t^* : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ that takes values

$$Q_t^*(s, a) = \sup_{\pi \in \Pi_t^T} Q_t^\pi(s, a), \quad s \in \mathcal{S}, a \in \mathcal{A},$$

is called optimal time-state-action value function.

- a policy π^* that satisfies

$$V_t^\pi(s) = V_t^*(s), \quad s \in \mathcal{S}, t \leq T,$$

is called optimal.

As for discounted MDPs the optimal value functions are the best value functions that are theoretically achievable by a policy. It is far from obvious that there is an optimal policy. As in the discounted setting we get the following relations



Lemma 2.4.6. The following holds for the optimal time-state value function and the optimal time-state-action value function for any $s \in \mathcal{S}$:

- $V_t^*(s) = \max_{a \in \mathcal{A}_s} Q_t^*(s, a)$ for all $t < T$,
- $Q_t^*(s, a) = r(s, a) + \sum_{s' \in \mathcal{S}} p(s'; s, a) V_{t+1}^*(s')$ for all $t < T$

In particular, V^* and Q^* satisfy the following Bellman optimality equations (backwards recursions):

$$V_t^*(s) = \max_{a \in \mathcal{A}_s} \left\{ r(s, a) + \sum_{s' \in \mathcal{S}} p(s'; s, a) V_{t+1}^*(s') \right\}, \quad s \in \mathcal{S},$$



and

$$Q_t^*(s, a) = r(s, a) + \sum_{s' \in \mathcal{S}} p(s'; s, a) \max_{a' \in \mathcal{A}_{s'}} Q_{t+1}^*(s', a'), \quad s \in \mathcal{S}, a \in \mathcal{A}_s,$$

for all $t < T$.

Proof. The proof is the same as in the infinite setting. First write

$$V_t^*(s) = \sup_{\pi} V_t^{\pi} = \sup_{\pi} \sum_{a \in \mathcal{A}_s} \pi_t(a; s) Q_t^{\pi}(s, a)$$

and then try to maximise the righthand side by playing greedy with respect to Q_t^* . \square

For discounted infinite time horizon problems we could now show that it is sufficient to consider stationary greedy policies. **The stationarity is in general not true for finite time horizon MDPs.** Consider for example the ice-vendor MDP and assume that we close our store in the winter. Then the amount of ice cream we want to order depends strongly on the time horizon up to the closing date. It is clear that we would like to order with a different strategy if we can sell the ice cream for 6 weeks rather than just for 1 week. Given this observation it is clear that we can no longer restrict the set of policies to stationary policies and it follows also that we have no fixpoint relation in the value function. We can formulate the following theorem:



Theorem 2.4.7. (Dynamic programming algorithm)

Suppose $v_t : \mathcal{S} \rightarrow \mathbb{R}$, $t = 0, \dots, T$, is a sequence of vectors that fulfill the Bellman optimality equations (backwards recursions)

$$v_t(s) = \begin{cases} 0 & : t = T \\ \max_{a \in \mathcal{A}_s} \left\{ r(s, a) + \sum_{s' \in \mathcal{S}} p(s'; s, a) v_{t+1}(s') \right\} & : t < T \end{cases},$$

then $v_t = V_t^*$ for all $t = 0, \dots, T$. Similarly, if $q_t : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, $t = 0, \dots, T$, is a sequence of vectors that fulfill the Bellman state-action optimality equations (backwards recursions)

$$q_t(s, a) = \begin{cases} 0 & : t = T \\ r(s, a) + \sum_{s' \in \mathcal{S}} p(s'; s, a) \max_{a' \in \mathcal{A}_{s'}} q_{t+1}(s', a') & : t < T \end{cases},$$

then $q_t = Q_t^*$ for all $t = 0, \dots, T$. Most importantly, an optimal (non-stationary) policy is given by the greedy policy

$$\pi_t^q(a; s) = \begin{cases} 1 & : a = a_t^*(s) \\ 0 & : \text{otherwise} \end{cases}, \quad \text{where } a_t^q(s) = \arg \max_{a \in \mathcal{A}_s} q_t(s, a).$$

Proof. Suppose q solves the optimality recursion and π^q is the greedy policy. By uniqueness q equals Q^* . It remains to show $V^* = V^{\pi^q}$, then π^q is optimal by definition. But this follows from the previous lemma and the fact that $q = Q^*$. First for Q :

$$\begin{aligned} Q_t^*(s, a) &= r(s, a) + \sum_{s' \in \mathcal{S}} p(s'; s, a) \max_{a' \in \mathcal{A}_{s'}} Q_{t+1}^*(s', a') \\ &= r(s, a) + \sum_{s' \in \mathcal{S}} p(s'; s, a) \sum_{a' \in \mathcal{A}_{s'}} \pi_{t+1}^q(a'; s') Q_{t+1}^*(s', a'). \end{aligned}$$

This means that also the sequence of matrices Q^* solves the Bellman expectation equation for π^q . Since the recursions have a unique solution it follows that $Q_t^* = Q_t^{\pi^q}$ for all $t \leq T$.

Finally, using the definition of the greedy policy and the relations $V_t^*(s) = \max_{a \in \mathcal{A}_s} Q_t^*(s, a)$ and $V_t^\pi(s) = \sum_{a \in \mathcal{A}_s} \pi_t(a; s) Q_t^\pi(s, a)$ between V and Q gives

$$V^*(s) = \max_{a \in \mathcal{A}_s} Q^*(s, a) = \max_{a \in \mathcal{A}_s} Q^{\pi^q}(s, a) = \sum_{a \in \mathcal{A}_s} \pi_q(a; s) Q^{\pi^q}(s, a) = V^{\pi^q}(s).$$

We have thus proved that $V_t^* = V_t^{\pi^q}$ for all $t \leq T$ which shows that π_q is optimal. \square

The key tool to solve finite time horizon MDPs follows from Theorem 2.4.7 and is called backward induction. If the state-action space is not too large and we have access to the transition probabilities, then the optimal control problem can be solved backwards iteration. Let us recall the ice vendor example on finite time-horizon to get an idea why this is intuitive. In the last timestep, there is no opportunity to sell any more ice cream, the summer season is over. In the optimal situation we obviously have no more stock which corresponds to $V_T^* \equiv 0$. Then we go backwards in time step-by-step and consider what is optimal in every possible state. Suppose one day is left. The recursion gives

$$Q_{T-1}^*(s, a) = \underbrace{r(s, a)}_{\text{last days profit}} + \underbrace{\sum_{s' \in \mathcal{S}} p(s'; s, a) V_T^*(s')}_{=0, \text{ no future to be considered}} = r(s, a)$$

and the optimal policy becomes

$$\pi_{T-1}^*(s) = \arg \max_{a \in \mathcal{A}_s} r(s, a).$$

What does this mean? For the last step we need to chose the action that maximises the expected reward without any thought about the future. No surprise! If the expectations are know, nothing needs to be done. If not, this is nothing but the bandit problem! For the next time-step $T - 2$ the recursion gives

$$Q_{T-2}^*(s, a) = \underbrace{r(s, a)}_{\text{today's profit}} + \underbrace{\sum_{s' \in \mathcal{S}} p(s'; s, a) \max_{a' \in \mathcal{A}_{s'}} Q_{T-1}^*(s', a')}_{\text{next days profit if next day is played optimally}}$$

and optimally the ice vendor orders/produces ice cream according to

$$\pi_{T-2}^*(s) = \arg \max_{a \in \mathcal{A}_s} Q_{T-2}^*(s, a).$$

If one could carry out explicitly the recursion till $t = 0$ the sequence $(\pi_t^*)_{t=1, \dots, T}$ is an optimal time-dependent policy.



Finite-time control problems are similar to the stochastic bandits. If all quantities are known explicitly the problem is essentially trivial, the optimal policy is obtained by choosing largest Q -values and those can be computed without any effort. For discounted infinite-time problems at least the non-linear equation $T^*Q = Q$ must be solved which is not a trivial task. As for stochastic bandits the situation becomes interesting once the Q -values cannot be computed explicitly but must be estimated by interacting with the environment. This will be the subject of the next chapter.

5

2.4.2 Dynamic programming algorithms

In the context of fully available dynamics the dynamic programming algorithms are very simple for finite MDPs. No equations must be solved, only recursion followed backwards must be computed from the terminal condition. Theorem 2.4.7 immediately translates into a simple optimal control algorithm.

⁵Sara: Wuerfelbeispiel

Data: MDP with finite time-horizon T , policy $\pi \in \Pi_0^T$

Result: V_t^π and Q_t^π for all $t = 0, \dots, T-1$

Set backward induction start $V_T^\pi \equiv 0$

for $t = T-1, \dots, 0$ **do**

for $s \in S$ **do**

for $a \in \mathcal{A}_s$ **do**

$$Q_t^\pi(s, a) := r(s, a) + \sum_{s' \in S} p(s'; s, a) V_{t+1}^\pi(s')$$

end

$$V_t^\pi(s) := \sum_{a \in \mathcal{A}_s} \pi_t(a; s) Q_t^\pi(s, a)$$

end

end

Algorithm 13: Policy evaluation for finite-time MDPs

Data: MDP with finite time-horizon T

Result: V_t^* and Q_t^* for all $t = 0, \dots, T$ and optimal policy π^*

Set induction beginning $V_T^* \equiv 0$

for $t = T-1, \dots, 0$ **do**

for $s \in S$ **do**

for $a \in \mathcal{A}_s$ **do**

$$Q_t^*(s, a) := r(s, a) + \sum_{s' \in S} p(s'; s, a) V_{t+1}^*(s')$$

end

$$V_t^*(s) := \sum_{a \in \mathcal{A}} \pi_t^*(a; s) Q_t^*(s, a)$$

$$\pi_t^* := \text{greedy}(Q_t^*)$$

end

end

Algorithm 14: Optimal control for finite-time MDPs

Part II

Non-Tabular Reinforcement Learning

Up to now we studied so-called tabular algorithms under the basic assumption a table for all state-action pairs can be saved. We then developed algorithms that explicitly dealt with all state-action pairs for instance by computing or estimating the matrices Q^π or Q^* . This approach is only feasible if the number of state-action pairs is small enough to at least save such a table (not to think about operating with the tables, i.e. multiplications, inversions, ect.). This assumption is widely unrealistic. It makes sense for gridworld examples but immediately fails for interesting control problems such as playing a computer game. The number of actions is often not the problem but the number of states might be huge (or even infinite). In such situations tabular algorithms are not feasible and further approximations are indispensable. In this part we will discuss a few of the classical approaches:

- **Policy approximation.** The set of all policies is replaced by a parametric family for which the best policy is approximated using (stochastic) gradient descent,
- **Value function approximation.** for which the Q -matrix (or value function) is replaced by an approximate family $(Q_\theta)_{\theta \in \Theta}$ and instead of learning all entries $Q(s, a)$ separately we try to find the best approximating matrix (in some sense) Q_{θ^*} by optimising a suitable errorfunction using (stochastic) gradient descent.
- **Actor-critic.** Mixtures of policy approximation and value function approximation.

All these methods have in common that they somehow involve a (stochastic) gradient algorithm. We thus start with in introductory chapter introducing some of the most classical results.

Chapter 3

A quick dive into gradient ascent/descent methods

We will focus during this section¹ on gradient descent methods for functions $f : \mathbb{R}^d \rightarrow \mathbb{R}$, i.e. we aim to find (local) minima of f . Note that all results hold similar for gradient ascent methods by considering $-f$ instead of f , then all maxima are minima.

Before we get started recall from basic analysis that the negative gradient is the direction of the steepest ascent, i.e.

$$\operatorname{argmax}_{\|d\|=1} \underbrace{f'(x)d}_{\text{slope in direction } d} = \nabla f(x),$$

and the negative gradient the direction of steepest descent, i.e.

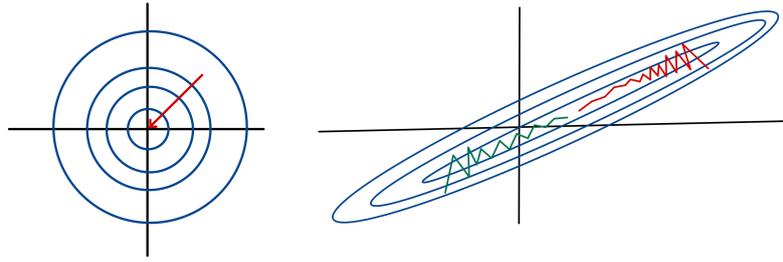
$$\operatorname{argmin}_{\|d\|=1} \underbrace{f'(x)d}_{\text{slope in direction } d} = -\nabla f(x).$$

This motivates the minimisation scheme that goes step-by-step in the direction of the gradient multiplied with the lentheta of each step determined by a step-size:

$$x_{n+1} = x_n - \alpha \nabla f(x_n), \quad x_0 \in \mathbb{R}^d.$$

Since only derivatives of first order are involved such a numerical scheme is called a first order scheme. How to chose the step-size (lentheta of step into gradient direction) is non-trivial. Going too far might also lead to an increase of the function, not going far enough might lead to very slow convergence. To get a first feeling have a look at the illustration of two extreme scenarios. For a quadratic function $f(x) = \|x\|^2$ gradient descent can converge in one step, for functions with narrow valleys the convergence is very slow. The problem in the second example the mixture of flat direction (in the valley) and steep direction (down the ridge) that is reflected in a large ratio of largest and smallest Eigenvalue of the Hessian $\nabla^2 f(x)$. Thus, it is not surprising that good convergence properties hold for functions not too big Eigenvalues (so-called L -smooth functions) and not too small Eigenvalues (so-called strongly convex). We will not talk about Eigenvalues but uses equivalent definitions that are more straightforward to be used in the proofs.

¹For much more please check the lecture notes https://www.wim.uni-mannheim.de/media/Lehrstuehle/wim/doering/OptimL/Sheets/lecture_notes_01.pdf of Simon Weißmann and, of course, the standard text book "Nonlinear programming" of Dimitri P. Bertsekas from which most proofs are taken.



To get a first impression in what follows we will go through the most basic theory, for L -smooth (not too steep) and strongly convex (not too flat) functions. What is more relevant for policy gradient in reinforcement learning is the continuation for functions that satisfy the so-called PL inequality (gradient dominates the function differences).⁹

Let us start with L -smoothness.



Definition 3.0.1. We call a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ L -smooth if f is differentiable and the gradient ∇f is L -Lipschitz continuous, i.e.

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|, \quad \forall x, y \in \mathbb{R}^d. \quad (3.1)$$

L -smoothness is a property in which gradient descent algorithms have a good chance to converge and, luckily, gradients of value functions of MDPs are L -smooth for some choices of policies. To get an idea of what L -smoothness means note that for twice-differentiable f the L -smoothness is equivalent to all eigenvalues of the Hessian $\nabla^2 f(x)$ to be in $[-L, L]$. Thus, L -smooth function do not have very strong curvature.

3.1 Gradient descent for L -smooth functions

In this first section we consider gradient descent with the only assumption of f being L -smooth. In that case gradient descent does not necessarily converge to a (local) minimum but we can show it converges (if it converges) to a stationary point, i.e. a point with vanishing gradient. Stationary points are not necessarily (local) extreme points, saddle points are also stationary.



Lemma 3.1.1. (Descent lemma)

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be L -smooth for some $L > 0$. Then it holds that

$$f(x + y) \leq \underbrace{f(x) + y^T \nabla f(x)}_{\text{tangent at } x} + \frac{L}{2} \|y\|^2 \quad (3.2)$$

for all $x, y \in \mathbb{R}^d$.

Note that small L makes the deviation of ∇f smaller, thus, the function smoother. It will later turn out that convergence of gradient descent is faster for smoother functions.

Proof. We define $\phi(t) = f(x + ty)$ and apply the chain rule in order to derive

$$\phi'(t) = y^T \nabla f(x + ty), \quad t \in [0, 1].$$

By the fundamental theorem of calculus it follows

$$\begin{aligned}
f(x+y) - f(x) &= \phi(1) - \phi(0) = \int_0^1 \phi'(t) dt \\
&= \int_0^1 y^T \nabla f(x+ty) dt \\
&= \int_0^1 y^T \nabla f(x) dt + \int_0^1 y^T (\nabla f(x+ty) - \nabla f(x)) dt \\
&\leq y^T \nabla f(x) + \int_0^1 \|y\| \|\nabla f(x+ty) - \nabla f(x)\| dt \\
&\leq y^T \nabla f(x) + \|y\| \int_0^1 Lt \cdot \|y\| dt \\
&\leq y^T \nabla f(x) + \frac{L}{2} \|y\|^2,
\end{aligned}$$

where we have applied Cauchy-Schwarz followed by L -smoothness. \square

The aim is to show that the gradient descent algorithm converges to a (local) minimum. This is not always the case, the algorithm can for instance also converge to saddle points. What is true without any assumptions, except the minimal differentiability, is the convergence to a stationary point, i.e. a point with vanishing gradient.



Theorem 3.1.2. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be L -smooth and $(x_k)_{k \in \mathbb{N}}$ be defined as

$$x_{k+1} = x_k - \bar{\alpha} \nabla f(x_k), \quad x_0 \in \mathbb{R}^d,$$

with non-negative step-size $\bar{\alpha} \in [\varepsilon, \frac{2-\varepsilon}{L}]$ for some $\varepsilon < \frac{2}{L+1}$. Then every accumulation point \bar{x} of $(x_k)_{k \in \mathbb{N}}$ is a stationary point of f , i.e. $\nabla f(\bar{x}) = 0$.

Proof. Since f is assumed to be L -smooth we can apply the descent Lemma 3.1.1 (with the choice $y \equiv -\bar{\alpha} \nabla f(x_k)$ and $x \equiv x_k$),

$$\begin{aligned}
f(x_k - \bar{\alpha} \nabla f(x_k)) - f(x_k) &\leq (-\bar{\alpha} \nabla f(x_k))^T \nabla f(x_k) + \frac{L}{2} \|\bar{\alpha} \nabla f(x_k)\|^2 \\
&= \bar{\alpha} \|\nabla f(x_k)\|^2 \left(\frac{\bar{\alpha} L}{2} - 1 \right).
\end{aligned} \tag{3.3}$$

In fact, this inequality justifies the name descent lemma. Due to our choice of $\bar{\alpha} \leq \frac{2-\varepsilon}{L}$ we can bound

$$\frac{\bar{\alpha} L}{2} - 1 \leq -\frac{\varepsilon}{2} < 0.$$

We reformulate the inequality (3.3) and obtain

$$f(x_k) - f(x_{k+1}) = f(x_k) - f(x_k - \bar{\alpha} \nabla f(x_k)) \geq \frac{\varepsilon}{2} \bar{\alpha} \|\nabla f(x_k)\|^2 \geq \frac{\varepsilon^2}{2} \|\nabla f(x_k)\|^2. \tag{3.4}$$

Assume that $(x_{n_k})_{k \in \mathbb{N}}$ is a sub-sequence with limit point $\bar{x} \in \mathbb{R}^d$ so that also $\lim_{k \rightarrow \infty} f(x_{n_k}) = f(\bar{x})$ by the continuity of f . By the Cauchy criterion for converging sequence s we also obtain

$$\lim_{k \rightarrow \infty} (f(x_{n_{k+1}}) - f(x_{n_k})) = 0.$$

Then, using the (3.4), it then follows that

$$\|\nabla f(\bar{x})\| = \lim_{k \rightarrow \infty} \|\nabla f(x_k)\| = 0.$$

\square

The proof was simple and the result is very weak. There is no statement about convergence and no way to know if the limit points are minima. For this more assumptions on f are needed.

3.2 Gradient descent for L -smooth, convex functions

We will now study the convergence of gradient descent under the additional assumption of convexity.



Definition 3.2.1. A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is called convex if

$$f(y) \geq f(x) + (y - x)^T \nabla f(x)$$

for all $x, y \in \mathbb{R}^d$.

Recall from basic analysis that the definition has a geometric interpretation. The prime example is a quadratic function $f(x) = x^T A x$ for a positive definite matrix A but also linear functions are convex. At every point the tangent plane lies below the graph of f . For convex and L -smooth f convergence of the gradient descent recursion to a global minimum can be proved if a global minimum exists (linear functions are convex but do not possess minima). In theorems we will always assume a global minimum exists.



All local minima of convex functions must be global minima, in particular, all global minima have same height.

In what follows we will always denote by f_* the height at a (all) global minimum. In order to talk about convergence rates one typically considers an error function $e : \mathbb{R}^d \rightarrow \mathbb{R}$ with the property $e(x) \geq 0$ for all $x \in \mathbb{R}^d$ and $e(x_*) = 0$. A typical choice is $e(x) = f(x) - f_*$. Recall that convergence is called linear (it is actually exponentially fast) if $e(x_{k+1}) \leq c e(x_k)$ for some $c \in (0, 1)$ and sublinear if convergence is slower. For a first order method (only first derivatives are used) convergence will not be linear. For convex and smooth function we can at least prove the following sublinear convergence rate in case a global minimum exists.



Theorem 3.2.2. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex and L -smooth, and let $(x_k)_{k \in \mathbb{N}}$ be generated by

$$x_{k+1} = x_k - \bar{\alpha} \nabla f(x_k), \quad x_0 \in \mathbb{R}^d,$$

with $\bar{\alpha} \leq \frac{1}{L}$. Moreover, we assume that the set of all global minimums of f is non-empty. Then the sequence $(x_k)_{k \in \mathbb{N}}$ converges in the sense that

$$e(x_k) := f(x_k) - f_* \leq \frac{c}{k+1}, \quad k \in \mathbb{N}$$

for some constant $c > 0$ and $f_* = \min_{x \in \mathbb{R}^d} f(x)$.

Note again that since there are no unique minima we do not aim to prove convergence of x_n to x_* but instead convergence of $f(x_n)$ to f_* .

Proof. We again apply the descent lemma 3.1.1 to the recursive scheme ($y = x_{k+1} - x_k$, $x = x_k$) to obtain

$$\begin{aligned} f(x_{k+1}) &\leq f(x_k) - \bar{\alpha} \|\nabla f(x_k)\|^2 + \frac{L\bar{\alpha}^2}{2} \|\nabla f(x_k)\|^2 \\ &= f(x_k) + \bar{\alpha} \left(\frac{L\bar{\alpha}}{2} - 1 \right) \|\nabla f(x_k)\|^2. \end{aligned}$$

Since $\bar{\alpha} \leq \frac{1}{L}$, we have $(\frac{L\bar{\alpha}}{2} - 1) < 0$ and therefore, the sequence $(f(x_k))_{k \in \mathbb{N}}$ is decreasing. Now,

let $x_* \in \mathbb{R}^d$ be a global minimum of f such that due to convexity it holds true that

$$\underbrace{f(x_k) + (x_* - x_k)^T \nabla f(x_k)}_{\text{tangent}} \leq f(x_k).$$

We plug this into the above inequality and obtain

$$\begin{aligned} f(x_{k+1}) &\leq f(x_k) + \bar{\alpha} \left(\frac{L\bar{\alpha}}{2} - 1 \right) \|\nabla f(x_k)\|^2 \\ &\leq f(x_*) - \frac{\bar{\alpha}}{\alpha} (x_* - x_k)^T \nabla f(x_k) + \bar{\alpha} \left(\frac{L\bar{\alpha}}{2} - 1 \right) \|\nabla f(x_k)\|^2 \\ &= f(x_*) + \frac{1}{\bar{\alpha}} \left[\frac{1}{2} \|x_* - x_k\|^2 + \frac{\bar{\alpha}^2}{2} \|\nabla f(x_k)\|^2 - \frac{1}{2} \|(x_* - x_k) + \bar{\alpha} \nabla f(x_k)\|^2 \right] \\ &\quad + \bar{\alpha} \left(\frac{L\bar{\alpha}}{2} - 1 \right) \|\nabla f(x_k)\|^2, \end{aligned}$$

where we have used the polarisation formula $-\langle a, b \rangle = \frac{1}{2} \|a\|^2 + \frac{1}{2} \|b\|^2 - \frac{1}{2} \|a + b\|^2$ for $a, b \in \mathbb{R}^d$. Rearranging the righthand side yields

$$\begin{aligned} f(x_{k+1}) &\leq f(x_*) + \frac{1}{2\bar{\alpha}} (\|x_* - x_k\|^2 - \|x_* - x_{k+1}\|^2) + \bar{\alpha} \left(\frac{L\bar{\alpha}}{2} - \frac{1}{2} \right) \|\nabla f(x_k)\|^2 \\ &\leq f(x_*) + \frac{1}{2\bar{\alpha}} (\|x_* - x_k\|^2 - \|x_* - x_{k+1}\|^2), \end{aligned}$$

where we have used again that $\bar{\alpha} \leq \frac{1}{L}$. Taking the sum over all iterations gives

$$\begin{aligned} \sum_{k=0}^N (f(x_{k+1}) - f(x_*)) &\leq \frac{1}{2\bar{\alpha}} \sum_{k=0}^N (\|x_* - x_k\|^2 - \|x_* - x_{k+1}\|^2) \\ &= \frac{1}{2\bar{\alpha}} (\|x_* - x_0\|^2 - \|x_* - x_{N+1}\|^2) \\ &\leq \frac{1}{2\bar{\alpha}} \|x_* - x_0\|^2, \end{aligned}$$

where we have applied a telescoping sum. With the decrease of $(f(x_k))_{k \in \mathbb{N}}$ it follows that

$$\sum_{k=0}^N f(x_{k+1}) \geq (N+1)f(x_{N+1}),$$

and therefore, the assertion follows with

$$f(x_{N+1}) - f(x_*) \leq \frac{1}{N+1} \sum_{k=0}^N (f(x_{k+1}) - f(x_*)) \leq \frac{1}{N+1} \frac{1}{2\bar{\alpha}} \|x_* - x_0\|^2 =: \frac{c}{N+1}.$$

□



Definition 3.2.3. A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is called μ -strongly convex if

$$f(y) \geq f(x) + \underbrace{(y-x)^T \nabla f(x) + \frac{\mu}{2} \|y-x\|^2}_{\text{tangent at } x \text{ plus quadratic}}$$

for all $x, y \in \mathbb{R}^d$.

Linear functions are not strongly convex, $f(x, y) = x^2$ is not strongly convex as it only bends in the x direction. While f is an example of a convex function with infinitely many global minima (the y -axis) all strongly convex functions have a unique global maxima and no further local minima.



Show that a strongly convex function has a unique global minimum.

If we tighten the assumption on convexity and assume μ -strong convexity instead, we can improve the convergence rate from sublinear to linear convergence as follows.



Theorem 3.2.4. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be μ -strongly convex for some $\mu > 0$ and L -smooth, as let $(x_k)_{k \in \mathbb{N}}$ be generated by

$$x_{k+1} = x_k - \bar{\alpha} \nabla f(x_k), \quad x_0 \in \mathbb{R}^d,$$

with $\bar{\alpha} < \frac{1}{L}$. Then the sequence $(x_k)_{k \in \mathbb{N}}$ converges linearly in the sense that there exists $c \in (0, 1)$ such that

$$e(x_k) := \|x_k - x_*\| \leq c^k \|x_0 - x_*\|, \quad k \in \mathbb{N}$$

where $x_* \in \mathbb{R}^d$ is the unique global minimum of f . The constant can be chosen as $c = (1 + \mu\bar{\alpha})^{-1/2}$.

An example for such a function is $f(x) = x^2 - \cos(x)$. Note that for strongly convex L -smooth functions the convergence also implies linear convergence of $f(x_n)$ to $f(x_*)$ because (Nesterov, page 64)

$$f(x_n) - f(x_*) \leq \frac{1}{2\mu} \|\nabla f(x_n) - \nabla f(x_*)\|^2 \leq \frac{L}{2\mu} \|x_n - x_*\|^2 \leq (c^2)^n \frac{L}{2\mu} \|x_0 - x_*\|^2$$

so that the theorem significantly strengthens the previous theorem on only convex functions.

Proof. Let $x_* \in \mathbb{R}^d$ be the unique global minimum of f with $\nabla f(x_*) = 0$. Since f is assumed to be μ -strongly convex, by Definition 3.2.3 it holds true that

$$\frac{\mu}{2} \|x_{k+1} - x_*\|^2 = (x_{k+1} - x_*)^T \nabla f(x_*) + \frac{\mu}{2} \|x_{k+1} - x_*\|^2 \leq f(x_{k+1}) - f(x_*).$$

Because μ -strongly convex implies convexity we can use the calculations from the previous proof of Theorem 3.2.2, where we have derived that

$$f(x_{k+1}) - f(x_*) \leq \frac{1}{2\bar{\alpha}} (\|x_k - x_*\|^2 - \|x_{k+1} - x_*\|^2)$$

and together we obtain

$$\left(\frac{\mu}{2} + \frac{1}{2\bar{\alpha}}\right) \|x_{k+1} - x_*\|^2 \leq \frac{1}{2\bar{\alpha}} \|x_k - x_*\|^2.$$

The assertion follows via induction using the inequality with $c = \sqrt{\frac{1}{1 + \mu\bar{\alpha}}} \in (0, 1)$. \square

For a choice $\bar{\alpha}$ close to $\frac{1}{L}$ the rate is essentially

$$\|x_k - x_*\| \leq \left(\sqrt{\frac{1}{1 + \frac{\mu}{L}}}\right)^k \|x_0 - x_*\|.$$

This means that the speed of convergence (small c is better) is determined by the ratio of smoothness L and strong convexity μ :

$$c = \sqrt{\frac{1}{1 + \frac{\mu}{L}}}$$

Thus, small L (very smooth function) and large μ (very convex function) yield fast convergence.



Suppose $f(x) = x^T A x$ is a quadratic function. Show that L and μ can be described by the spectrum of A . L is the two times the largest singular value and μ two times the smallest singular value. Thus, convergence is fast if the Eigenvalues are similar while convergence can be slow if the Eigenvalues differ a lot. The latter corresponds to the situation of narrow valleys, the curvatures is big in some direction and small in another direction.

3.3 Gradient descent for L -smooth functions with PL inequality

In the applications for policy gradient methods the assumption of convexity (actually concavity as functions are maximised) is not fulfilled already for simple examples. There is one way to relax this assumption to so called Polyak-Łojasiewicz (PL) conditions.



Definition 3.3.1. A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ satisfies the (strong) PL inequality if for some constant $C > 0$

$$\|\nabla f(x)\|^2 \geq C(f(x) - f_*), \quad x \in \mathbb{R}^d, \quad (3.5)$$

holds for all $x \in \mathbb{R}^d$, where $f_* = \min_{x \in \mathbb{R}^d} f(x) > -\infty$.

In fact, the PL inequality is exactly what gradient descent needs for convergence to a global minimum.



Theorem 3.3.2. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be L -smooth and satisfied the PL condition holds for some $C = 2r$ with $r \in (0, L)$. Then the sequence $(x_k)_{k \in \mathbb{N}}$ generated by

$$x_{k+1} = x_k - \bar{\alpha} \nabla f(x_k), \quad x_0 \in \mathbb{R}^d,$$

with $\bar{\alpha} = \frac{1}{L}$ converges linearly in the sense that

$$e(x_k) := f(x_k) - f_* \leq c^k (f(x_0) - f_*),$$

where $c = 1 - \frac{r}{L} \in (0, 1)$.

Note again that since there are no unique minima we do not aim to prove convergence of x_n to x_* but instead convergence of $f(x_n)$ to $f(x_*)$.

Proof. Using the descent Lemma 3.1.1 as in the proofs above (this only uses L -smooth) yields

$$\begin{aligned} f(x_{k+1}) &\leq f(x_k) - \bar{\alpha} \left(1 - \frac{L\bar{\alpha}}{2}\right) \|\nabla f(x_k)\|^2 \\ &= f(x_k) - \frac{1}{2L} \|\nabla f(x_k)\|^2. \end{aligned}$$

Using the PL-inequality yields

$$f(x_{k+1}) \stackrel{(3.6)}{\leq} f(x_k) - \frac{r}{L} (f(x_k) - f_*).$$

Subtracting f_* from both sides, we get

$$f(x_{k+1}) - f_* \leq \left(1 - \frac{r}{L}\right) (f(x_k) - f_*)$$

and the claim follows by induction. \square

The idea behind inequality (3.6), also known as PL or Polyak-Lojasiewicz inequality (also called gradient domination) is very simple. The inequality ensures that the gradient does not vanish as long as the gradient descent algorithm has not reached x^* . Since the algorithm is based on the gradient the algorithm does not stop moving as long $f(x_k)$ has not reached f_* . Here is another view. For every $x \in \mathbb{R}^d$ the PL inequality lower bounds the norm of the gradient by the difference of $f(x)$ to some global minimum of f . This implies, that a small gradient at x implies a small difference between the function value $f(x)$ and the global optimum. Given that gradient descent for L -smooth functions implies convergence to a stationary point by Theorem 3.4.1 it is not surprising that the PL condition is exactly what is needed for convergence. Estimates of the PL type are called gradient domination inequalities.



Prove that μ -strong convexity and L -smoothness imply the PL condition (3.6). What is C ?

Condition (3.6) can be relaxed a little bit to remain convergence of gradient descent but just with a sublinear convergence rate.



Definition 3.3.3. A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ satisfies the weak PL inequality if for some constant $C > 0$

$$\|\nabla f(x)\| \geq C(f(x) - f_*), \quad x \in \mathbb{R}^d, \quad (3.6)$$

holds for all $x \in \mathbb{R}^d$, where $f_* = \min_{x \in \mathbb{R}^d} f(x) > -\infty$.

We will see below that for very special parametrisations of policies the weak PL inequalities holds for the parametrised value function in reinforcement learning. The convergence can still be proved but unfortunately is much slower, weak PL is not similar to strong convexity.



Theorem 3.3.4. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be L -smooth and satisfied the PL condition holds for some $C = 2r$ with $r \in (0, L)$. Then the sequence $(x_k)_{k \in \mathbb{N}}$ generated by

$$x_{k+1} = x_k - \bar{\alpha} \nabla f(x_k), \quad x_0 \in \mathbb{R}^d,$$

with $\bar{\alpha} = \frac{1}{L}$ converges sublinearly in the sense that

$$e(x_k) := f(x_k) - f_* \leq \frac{L}{2r^2(k+1)}.$$

Proof. As in the previous proof we first start with the estimate

$$f(x_{k+1}) \leq f(x_k) - \frac{1}{2L} \|\nabla f(x_k)\|^2.$$

that uses the L -smoothness. Applying the weak PL-inequality gives

$$f(x_{k+1}) \leq f(x_k) - \frac{1}{2L} 4r^2 (f(x_k) - f_*)^2.$$

Subtracting f_* on both sides of the inequality yields the recursion

$$e(x_{k+1}) \leq e(x_k) - \frac{2r^2}{L} e(x_k)^2. \quad (3.7)$$

We will show, that any positive sequence that satisfies the diminishing contraction

$$0 \leq a_{n+1} \leq a_n - qa_n^2 = (1 - qa_n)a_n, \quad n \geq 0,$$

converges to zero with convergence rate

$$a_n \leq \frac{1}{nq + \frac{1}{a_0}} \leq \frac{1}{(n+1)q}.$$

To see this, first note that the positivity implies that $a_n \leq \frac{1}{q}$ as otherwise the inequality gives a contradiction. Now divide the reordered contraction

$$a_n \geq a_{n+1} + qa_n^2$$

by $a_n a_{n+1}$ to obtain

$$\frac{1}{a_{n+1}} \geq \frac{1}{a_n} + q \underbrace{\frac{a_n}{a_{n+1}}}_{\geq 1} \geq \frac{1}{a_n} + q$$

which leads to

$$\frac{1}{a_n} - \frac{1}{a_0} = \sum_{k=0}^{n-1} \left(\frac{1}{a_{k+1}} - \frac{1}{a_k} \right) \geq nq.$$

Reordering we get our claim

$$a_n \leq \frac{1}{nq + \frac{1}{a_0}} \stackrel{a_0 \leq \frac{1}{q}}{\leq} \frac{1}{(n+1)q}.$$

□

One can ask the justified question whether there exists functions which are L -smooth and not convex but satisfy the PL condition. Indeed there are some:



Show that $f(x) = x^2 + 3 \sin^2(x)$ satisfies the PL condition (3.6) and prove that f is not convex. Plot the function to see why gradient descent converges. Hint: The plot can also help to find the parameter r of the PL condition.

3.4 Gradient descent with diminishing step-sizes

Typically, the smoothness parameter L is unknown, especially in a model-free optimization procedure as policy gradient where we only want to assume that we can sample from the MDP. Thus, how could one choose a step-size that satisfies $\bar{\alpha} < \frac{2}{L}$ or $\bar{\alpha} < \frac{1}{L}$? There are plenty of workarounds. Here we show that convergence also works if quickly diminishing step-sizes are used. If step-sizes decay too strongly (being summable) the algorithm might produce a sequence that stops moving away from a stationary point/optimum. We cover the case of convergence to a stationary point by only assuming L -smoothness in the context on diminishing step-sizes. All other results can be transferred to diminishing step-sizes as well.



Theorem 3.4.1. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be L -smooth and $(x_k)_{k \in \mathbb{N}}$ be defined as

$$x_{k+1} = x_k - \alpha_k \nabla f(x_k), \quad x_0 \in \mathbb{R}^d,$$

with non-negative step-size sequence $(\alpha_k)_{k \in \mathbb{N}}$ that fulfills

$$\lim_{k \rightarrow \infty} \alpha_k = 0 \quad \text{and} \quad \sum_{k=0}^{\infty} \alpha_k = \infty.$$

Then for the sequence $(f(x_k))_{k \in \mathbb{N}}$ it holds true that either

$$\lim_{k \rightarrow \infty} f(x_k) = -\infty \quad \text{or} \quad \lim_{k \rightarrow \infty} \nabla f(x_k) = 0.$$



Moreover, every accumulation point \bar{x} of $(x_k)_{k \in \mathbb{N}}$ is a stationary point of f , i.e. $\nabla f(\bar{x}) = 0$.

Note that the step-sizes can also be constant but the convergence also holds for decreasing step-sizes that do not decrease too fast.

Proof. As in all proofs above we use the descent lemma to obtain

$$f(x_{k+1}) \leq f(x_k) - \alpha_k \left(1 - \frac{L\alpha_k}{2}\right) \|\nabla f(x_k)\|^2.$$

Since we have assumed that $\lim_{k \rightarrow \infty} \alpha_k = 0$, there exists a $k_0 \in \mathbb{N}$ such that

$$f(x_{k+1}) \leq f(x_k) - \frac{\alpha_k}{2} \|\nabla f(x_k)\|^2$$

for all $k \geq k_0$. Hence, the sequence $(f(x_k))_{k \geq k_0}$ is decreasing and it either holds that $\lim_{k \rightarrow \infty} f(x_k) = -\infty$ or $\lim_{k \rightarrow \infty} f(x_k) = M$ for some $M \in \mathbb{R}$. Suppose that we are in the case where $\lim_{k \rightarrow \infty} f(x_k) = M$. It follows

$$\sum_{k=k_0}^K \frac{\alpha_k}{2} \|\nabla f(x_k)\|^2 \leq \sum_{k=k_0}^K (f(x_k) - f(x_{k+1})) = f(x_{k_0}) - f(x_K)$$

for ever $K > k_0$. For $K \rightarrow \infty$ this implies

$$\sum_{k=k_0}^{\infty} \frac{\alpha_k}{2} \|\nabla f(x_k)\|^2 \leq f(x_{k_0}) - M < \infty. \quad (3.8)$$

Since $\sum_{k=k_0}^{\infty} \alpha_k = \infty$ it follows that

$$\liminf_{k \rightarrow \infty} \|\nabla f(x_k)\| = 0.$$

In order to prove $\lim_{k \rightarrow \infty} \|\nabla f(x_k)\| = 0$ we will now prove that $\limsup_{k \rightarrow \infty} \|\nabla f(x_k)\| = 0$. Suppose that $\limsup_{k \rightarrow \infty} \|\nabla f(x_k)\| \geq \varepsilon$ for some $\varepsilon > 0$ and consider two sub-sequences $(m_j)_{j \in \mathbb{N}}$ and $(n_j)_{j \in \mathbb{N}}$ with $m_j < n_j < m_{j+1}$ such that

$$\frac{\varepsilon}{3} < \|\nabla f(x_k)\|, \quad \text{for } m_j \leq k < n_j$$

and

$$\|\nabla f(x_k)\| \leq \frac{\varepsilon}{3}, \quad \text{for } n_j \leq k < m_{j+1}.$$

These sequences exists because $\liminf_{k \rightarrow \infty} \|\nabla f(x_k)\| = 0$. Moreover, let $\bar{j} \in \mathbb{N}$ be sufficiently large such that

$$\sum_{k=m_{\bar{j}}}^{\infty} \alpha_k \|\nabla f(x_k)\|^2 \leq \frac{\varepsilon^2}{9L}.$$

Using L -smoothness for $j \geq \bar{j}$ and $m_j \leq m \leq n_j - 1$ it holds true that

$$\begin{aligned} \|\nabla f(x_{n_j}) - \nabla f(x_m)\| &\leq \sum_{k=m}^{n_j-1} \|\nabla f(x_{k+1}) - \nabla f(x_k)\| \\ &\leq L \sum_{k=m}^{n_j-1} \|x_{k+1} - x_k\| \\ &= \frac{3\varepsilon}{3\varepsilon} L \sum_{k=m}^{n_j-1} \alpha_k \|\nabla f(x_k)\| \\ &\leq L \frac{3}{\varepsilon} \sum_{k=m}^{n_j-1} \alpha_k \|\nabla f(x_k)\|^2 \\ &\leq L \frac{3}{\varepsilon} \frac{\varepsilon^2}{9L} = \frac{\varepsilon}{3}, \end{aligned}$$

where we have used that $\|\nabla f(x_k)\| > \frac{\varepsilon}{3}$ for $m_j \leq k \leq n_j - 1$. This implies that

$$\|\nabla f(x_m)\| \leq \|\nabla f(x_{n_j})\| + \|\nabla f(x_{n_j}) - \nabla f(x_m)\| \leq \|\nabla f(x_{n_j})\| + \frac{\varepsilon}{3} \leq \frac{2\varepsilon}{3}$$

and therefore $\|\nabla f(x_m)\| \leq \frac{2\varepsilon}{3}$ for all $m \geq m_j$. This is in contradiction to $\limsup_{k \rightarrow \infty} \|\nabla f(x_k)\| \geq \varepsilon$ and we have proved that

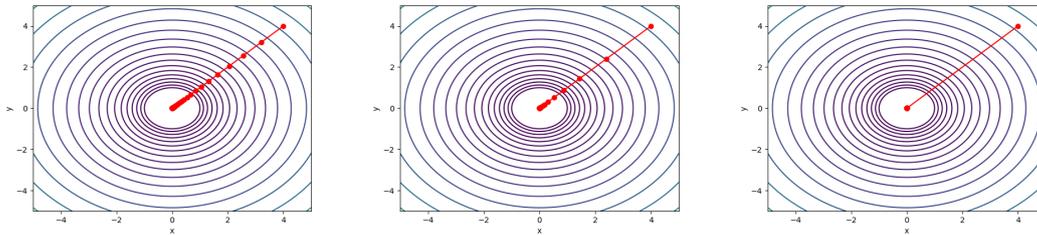
$$\limsup_{k \rightarrow \infty} \|\nabla f(x_k)\| = \liminf_{k \rightarrow \infty} \|\nabla f(x_k)\| = \lim_{k \rightarrow \infty} \|\nabla f(x_k)\| = 0.$$

Finally, let $\bar{x} \in \mathbb{R}^d$ be an accumulating point of $(x_k)_{k \in \mathbb{N}}$. Since $(f(x_k))_{k \geq k_0}$ is decreasing, it follows by continuity that

$$\nabla f(\bar{x}) = \lim_{k \rightarrow \infty} \nabla f(x_k) = 0.$$

□

There is an important point to keep in mind. Choosing step-sizes smaller than possible will slow down convergence a lot. As an example think about $f(x) = \|x\|^2$. The negative gradient $-\nabla f(x, y) = -(2x, 2y)$ points directly to the minimum and gradient descent with step-size $\alpha = \frac{1}{L} = \frac{1}{2}$ converges in just one step to the global minimum. Gradient descent with smaller (or even diminishing step-sizes) might require many steps.



Gradient descent on $f(x) = \|x\|^2$ with constant step-sizes $\alpha = \frac{1}{10}, \frac{1}{5}, \frac{1}{2}$.

3.5 Stochastic gradient descent methods

We consider the following setup. Let $(\Omega, \mathcal{A}, \mathbb{P})$ be the underlying probability space, $Z : \Omega \times \mathbb{R}^d \rightarrow \mathbb{R}^p$ be a random variable on $(\Omega, \mathcal{A}, \mathbb{P})$ for every $x \in \mathbb{R}^d$ with distribution μ_x . We are interested in solving the optimization problem

$$\min_{x \in \mathbb{R}^d} F(x),$$

where the cost function $F : \mathbb{R}^d \rightarrow \mathbb{R}$ is defined as expectation function in form

$$F(x) = \mathbb{E}_{Z \sim \mu_x}[f(x, Z)] = \int_{\mathbb{R}^p} f(x, z) \mu_x(dz), \quad x \in \mathbb{R}^d,$$

for a function $f : \mathbb{R}^d \times \mathbb{R}^p \rightarrow \mathbb{R}$. The mechanism that creates random variables $Z \sim \mu_x$ is often called an oracle. We assume that we can repeatedly ask the oracle for suggestions (samples). Here is a standard set of assumption that typically does not pose any problems:

1. The function $f : \mathbb{R}^d \times \mathbb{R}^p \rightarrow \mathbb{R}$ is $\mathcal{B}(\mathbb{R}^d) \otimes \mathcal{B}(\mathbb{R}^p)/\mathcal{B}(\mathbb{R})$ -measurable.
2. For every $z \in \mathbb{R}^p$ the function $x \mapsto f(x, z)$ is continuously differentiable.
3. For every $x \in \mathbb{R}^d$ we have

$$\mathbb{E}_{Z \sim \mu_x}[|f(x, Z)| + \|\nabla_x f(x, Z)\|] < \infty$$

Data: Initial random variable $X_0 : \Omega \rightarrow \mathbb{R}^d$
Result: Approximation of a stationary point
 Set $k = 0$
while not converged do
 Determine α_k
 Sample Z_{k+1} from μ_{X_k}
 Set $X_{k+1} = X_k - \alpha_k \nabla_x f(X_k, Z_{k+1})$
 $k \mapsto k + 1$
end

Algorithm 15: Plain vanilla stochastic gradient descent method (SGD)

In fact, writing $G(x) = \nabla F(x)$ the algorithm very much looks like stochastic approximation.

$$X_{k+1} = X_k - \alpha_k (G(X_k) + \varepsilon_k)$$

with unbiased error terms $\varepsilon_k = \nabla_x f(X_k, Z_{k+1}) - \nabla F(X_k)$. It is thus not too surprising that the algorithm will converge towards a zero of ∇F if the error terms are sufficiently well behaved (such as bounded variance). Indeed, the following simple variant of SGD is also proved using the Robbins-Siegmund theorem:



Theorem 3.5.1. (SGD almost sure convergence for L -smooth function)

Let $F : \mathbb{R}^d \rightarrow \mathbb{R}$ be L -smooth with $F_* = \inf_{x \in \mathbb{R}^d} F(x) > -\infty$ that satisfies

$$\mathbb{E}_{Z \sim \mu_x} [\|\nabla_x f(x, Z) - \mathbb{E}[\nabla_x f(x, Z)]\|^2] \leq c(1 + (F(x) - F_*)). \quad (3.9)$$

Suppose $(\alpha_k)_{k \in \mathbb{N}}$ are non-negative step-sizes that are deterministic or adapted to the stochastic gradient scheme and satisfy almost surely the Robbins-Monro conditions

$$\sum_{k=0}^{\infty} \alpha_k = \infty \quad \text{and} \quad \sum_{k=0}^{\infty} \alpha_k^2 < \infty.$$

Moreover, let X_0 be a random variable such that $\mathbb{E}[F(X_0)] < \infty$ and $(X_k)_{k \in \mathbb{N}}$ the sequence of random variables generated by the stochastic gradient Algorithm 15. Then $(F(X_k))_{k \in \mathbb{N}}$ converges almost surely to some finite random variable F_∞ and

$$\lim_{k \rightarrow \infty} \|\nabla F(X_k)\|^2 = 0$$

almost surely.

This first theorem on stochastic gradient descent extends Theorem 3.1.2 to the stochastic setting. The assumptions are weak and also the convergence statement is weak. Under stronger assumptions on F (such as convexity, strong convexity, gradient domination) there are also stronger stochastic convergence theorems in the spirit above. There are two points that should be mentioned. Due to the errors step-sizes must tend to zero in all settings and convergence rates in the stochastic setting are much slower.

Proof. For the proof we assume additionally that gradient and expectation can be interchanged (no big problem)

$$\mathbb{E}_{Z \sim \mu_x} [\nabla_x f(x, Z)] = \nabla F(x). \quad (3.10)$$

We define the natural filtration $\mathcal{F} = (\mathcal{F}_k)_{k \in \mathbb{N}}$ through $\mathcal{F}_k = \sigma(X_m, m \leq k) = \sigma(X_0, Z_m, m \leq k)$ and note that $(\alpha_k)_{k \in \mathbb{N}}$ is \mathcal{F} -adapted per construction. Using the descent lemma and the polarisation formula $\langle x, y \rangle = \frac{1}{2}(\|x\|^2 + \|y\|^2 + \|x - y\|^2)$ obtain (pathwise) that

$$\begin{aligned} F(X_{k+1}) &= F(X_k - \alpha_k \nabla_x f(X_k, Z_{k+1})) \\ &\leq F(X_k) - \alpha_k \langle \nabla F(X_k), \nabla_x f(X_k, Z_{k+1}) \rangle + \alpha_k^2 \frac{L}{2} \|\nabla_x f(X_k, Z_{k+1})\|^2 \\ &= F(X_k) - \alpha_k \|\nabla_x F(X_k)\|^2 + \alpha_k \langle \nabla_x F(X_k), \varepsilon_k \rangle \\ &\quad + \alpha_k^2 \frac{L}{2} (\|\nabla_x F(X_k)\|^2 - \langle \nabla_x F(X_k), \varepsilon_k \rangle + \|\varepsilon_k\|^2), \end{aligned}$$

where $\varepsilon_k := \nabla F(X_k) - \nabla_x f(X_k, Z_{k+1})$. Using (3.10) and (3.9) we obtain by construction in Algorithm 15 that

$$\mathbb{E}[\varepsilon_k \mid \mathcal{F}_k] = 0$$

and

$$\mathbb{E}[\|\varepsilon_k\|^2 \mid \mathcal{F}_k] \leq c(1 + (F(X_k) - F_*)).$$

This yields

$$\begin{aligned} \mathbb{E}[\overbrace{F(X_{k+1}) - F_*}^{Z_{k+1} \geq 0} \mid \mathcal{F}_k] &\leq (F(X_k) - F_*) - \alpha_k \left(1 - \frac{L}{2} \alpha_k\right) \|\nabla F(X_k)\|^2 + \frac{L}{2} \alpha_k^2 c(1 + (F(X_k) - F_*)) \\ &= (1 + c \frac{L}{2} \alpha_k^2) \underbrace{(F(X_k) - F_*)}_{Z_k \geq 0} + \underbrace{c \frac{L}{2} \alpha_k^2}_{B_k \geq 0} - \underbrace{\alpha_k \left(1 - \frac{L}{2} \alpha_k\right) \|\nabla F(X_k)\|^2}_{C_k \geq 0}. \end{aligned}$$

Here we used the assumption that F is bounded below to deduce that the sequence (Z_n) is non-negative. We can assume without loss of generality that $\alpha_k \leq (1 - \varepsilon) \frac{2}{L}$ for some $\varepsilon \in (0, 1)$, otherwise let k be sufficiently large, such that $(1 - \frac{L}{2} \alpha_k) \geq \varepsilon > 0$. We can now apply Theorem ?? to deduce that $\lim_{k \rightarrow \infty} F(X_k) - F_*$ exists almost surely and is finite, as well as

$$\varepsilon \sum_{k=0}^{\infty} \alpha_k \|\nabla F(X_k)\|^2 \leq \sum_{k=0}^{\infty} \alpha_k \left(1 - \frac{L}{2} \alpha_k\right) \|\nabla F(X_k)\|^2 < \infty$$

almost surely. Since we have assumed $\sum_{k=0}^{\infty} \alpha_k = \infty$ almost surely, using the same argument as in the proof of Theorem 3.4.1 pathwise we obtain

$$\lim_{k \rightarrow \infty} \|\nabla F(X_k)\|^2 = 0$$

almost surely. \square

There are many other variants of the theorem, with weaker/stronger assumptions on F or the oracle that result in weaker/stronger convergence statements. Here is one other variant that is used in the context of reinforcement learning².

²Yuan, Gower, and Lazaric: "A general sample complexity analysis of vanilla policy gradient", AISTATS 2020



Suppose instead of (3.9) the condition so-called expected smoothness (or ABC) condition

$$\mathbb{E}_{Z \sim \mu_x} [\|\nabla_x f(x, Z)\|^2] \leq 2A(F(x) - F_*) + B\|\nabla F(x)\|^2 + C.$$

Then SGD with constant step-sizes $\eta \in (0, \frac{2}{LB})$ satisfies

$$\min_{k \leq K} \mathbb{E}[\|\nabla F(X_k)\|^2] \leq \frac{2(F(X_0) - F_*)(1 + L\eta^2 A)^K}{\eta K(2 - LB\eta)} + \frac{LC\eta}{2 - LB\eta}$$

and, if $A = 0$, then one can also show

$$\frac{1}{K} \sum_{k=1}^K \mathbb{E}[\|\nabla F(X_k)\|^2] = \mathbb{E}[\|\nabla F(X)\|^2] \leq \frac{2(F(X_0) - F_*)(1)^K}{\eta K(2 - LB\eta)} + \frac{LC\eta}{2 - LB\eta},$$

where X is drawn uniformly from X_1, \dots, X_K . For a prespecified accuracy ε one can then solve for K to achieve the following consequence. For constant step size $\gamma = \min\{\frac{1}{\sqrt{LAK}}, \frac{a}{LB}, \frac{\varepsilon}{2LC}\}$ and a number of iterations $K \geq \frac{12(F(X_0) - F_*)L}{\varepsilon^2} \max\{B, \frac{12A(F(X_0) - F_*)}{\varepsilon^2}, \frac{2C}{\varepsilon^2}\}$ the gradients can be bounded by $\min_{k \leq K} \mathbb{E}[\|\nabla F(X_k)\|^2] \leq \varepsilon$.

3.6 Regularisation

dieses Jahr noch nicht. ³

³discuss regularisation in non-linear optimisation! take previous example as an example

Chapter 4

Policy Gradient methods

The third part of these lecture notes takes a very different view on optimal control problems. While in the previous parts we developed a complete theory of exact and approximate control in the tabular setting the third part of these lecture notes will be less complete. The theory is more complex and much less understood. In the case of state-action spaces that are too large to deal with all Q -values $Q(s, a)$ different further approximations will be discussed. We start with the powerful policy gradient method, where the optimal control problem will be attacked with numerical maximisation of the value function. While (approximate) dynamical programming are methods to find policies that are optimal for all states s (by definition of an optimal policy) we now fix a state s and try to numerically optimise $V(s)$ only. Without much thought an immediate idea is the following:



Fix a subclass $\Pi^\Theta \subseteq \Pi_S$ of parametrised stationary policies $\{\pi^\theta : \theta \in \Theta\}$ that hopefully contains the optimal policy. Maximising the mapping

$$\theta \mapsto J_\mu(\theta) = V^{\pi^\theta}(\mu) = \sum_{s \in \mathcal{S}} \mu(s) V^{\pi^\theta}(s)$$

for a measure μ on \mathcal{S} using a numerical method based on gradient ascent is called a policy gradient method. The best policy is denoted by π^{θ^*} and we hope that $\pi^{\theta^*} \approx \pi^*$. Since $\nabla J_\mu = \sum_{s \in \mathcal{S}} \mu(s) \nabla J_s$ we will only consider the case $J_s(\theta) = V^{\pi^\theta}(s)$, i.e. $\mu = \delta_s$ for a $s \in \mathcal{S}$.

There are many questions to be considered. (1) How to set up an algorithm, and, especially, how to compute the gradients? This is answered with the policy gradient theorems that surprisingly show that gradient estimates can be obtained even in a model free way. (2) How to choose Π^Θ such that $\pi^* \in \Pi^\Theta$? This is a delicate problem, as it results in a trade-off of choosing large and small families of policies. Not much is known in practical applications, typically neural networks are involved and one hopes to get close to π^* . (3) How good is the policy obtained for different starting conditions? A priori the policy gradient approach does not give any information. Typically one will use neural networks and try to optimise them for different states s . In order to use gradient methods it is not surprising that differentiability assumptions are needed.



Definition 4.0.1. Let $\Theta \subset \mathbb{R}^d$, then a set $\{\pi^\theta : \theta \in \Theta\}$ of stationary policies such that $\theta \mapsto \pi^\theta(a; s)$ is differentiable in θ for every $s \in \mathcal{S}$ and $a \in \mathcal{A}$ is called differentiable parameterised family of stationary policies.

We will later discuss policies parametrised by neural networks which must be large enough to include enough (hopefully the optimal) policies but at the same time small enough to be numerically tractable.



If $\Theta = \mathbb{R}^d$ then to reduce the dimension of the (very!) large state-action spaces it is strongly desirable to have $d \ll |\mathcal{A}||\mathcal{S}|$.

There is a simple example to keep in mind which, however, is not achieving the goal of reducing the dimension of the problem. The stationary softmax policies that were already discussed for bandits in Section 1.3.4 readily extend to general Markov decision problems:

Example 4.0.2. We start with the simplest so-called tabular parametrisation. In tabular parametrisations we take as many parameters as there are state-action pairs, i.e. we choose $d = |\mathcal{S}||\mathcal{A}|$. The long vector θ can be written as $\theta = (\theta_{s,a})_{s \in \mathcal{S}, a \in \mathcal{A}}$. There are two parametrisations that come to mind immediately. The direct parametrisation

$$\pi_{\theta}(a; s) = \theta_{s,a}$$

and the softmax parametrisation

$$\pi^{\theta}(a; s) = \frac{e^{\theta_{s,a}}}{\sum_{a' \in \mathcal{A}} e^{\theta_{s,a'}}}.$$

For the direct parametrization the parameter space Θ must be restricted to the simplex where all coordinates are non-negative and sum to 1, the softmax parameterization is defined on the entire \mathbb{R}^d . Instead of the (full) tabular parameterisation one can also choose any differentiable function $\theta \mapsto h_{\theta}(s, a)$ and define the generalised softmax policy

$$\pi^{\theta}(a; s) = \frac{e^{h_{\theta}(s,a)}}{\sum_{a' \in \mathcal{A}} e^{h_{\theta}(s,a')}}.$$

For instance, if $h_{\theta}(s, a) = \theta^T \Phi(s, a)$ then the parametrisation is called linear softmax with features $\Phi(s, a)$ and requires as many dimensions as features are fixed. If there are as many features as state-action pairs and every state-action pair only has its own feature, i.e. the vector $\Phi(s, a)$ is a unit vector with 1 at the position corresponding to (s, a) , then the linear softmax equals the tabular softmax. The feature vector is supposed to reduce complexity, state-action pairs with same feature vector are treated equally. Hence, in practice the number of features should be large enough to separate state-action pairs sufficiently well but small enough to be computationally tractable.

The rest of this chapter will consist of making (practical) sense of gradient ascent algorithms to find optimal policies. Assuming perfect information, i.e. $J_s(\theta) = V^{\pi^{\theta}}(s)$ is known explicitly, then Algorithm 16 is the kind of algorithm we aim for. Unfortunately, there are plenty of problems.

Data: Initial parameter θ_0

Result: Approximate policy $\pi^{\theta} \approx \pi^{\theta^*}$

Set $n = 0$.

while *not converged* **do**

 Choose step-size α .

 Calculate $K = \nabla J(\theta)|_{\theta=\theta_n}$.

 Update $\theta_{n+1} = \theta_n - \alpha K$.

 Set $n = n + 1$.

end

return π^{θ_n} .

Algorithm 16: Plain vanilla policy gradient algorithm (with exact gradients)

Most importantly the following:

- There is no reason $\theta \mapsto J(\theta)$ satisfies typical concavity conditions (it doesn't!) that guarantee convergence of gradient ascent methods to a global optimum. Why should gradient ascent not get stuck in some stationary point, i.e. a suboptimal $\theta \in \Theta$ that satisfies $\nabla J_s(\theta) = 0$?

- The gradient $\nabla J_s(\theta)$ is typically not known. How can we implement a gradient ascent method without direct access to the gradient?
- The first two issues can (partially) be resolved. Still, the simple policy gradient algorithm converges extremely slow. How can the algorithm be speed up?

In the next sections we will shed some light on what is going on. First, classical results for gradient descent/ascent are recalled. Here we focus on what turns out to be the right set-up for gradient ascent in the policy gradient setup. The PL inequality and L -smoothness. To get a certain feeling on how to replace $\nabla J_s(\theta)$ by estimates we quickly recall some results and arguments from stochastic gradient descent. Next, we will discuss the policy gradient theorems and how to combine them with neural network parametrisations. Finally, a couple of modifications are discussed that are used in practice.

4.1 Policy gradient theorems

1

As the name *policy gradient* suggests, the optimisation method we want to use is gradient ascent in order to find a maximum of the objective function J . This requires that J_s is differentiable in θ and a practical formula for the gradient to apply Algorithm 16. To start the discussion we first prove formulas for $\nabla J_s(\theta)$ that go back to Sutton, McAlister, Singh, and Mansour². Their observation was that the gradient $\nabla J_s(\theta)$ can be expressed rather nicely and, most importantly, can be approximated in a model-free way as the appearing expectations do not involve the transition probabilities and can be approximated by sampling only.

For didactic reasons we first deal with finite-time MDPs with undiscounted rewards and additionally consider stationary policies only. The situation is not extremely realistic as most finite-time MDPs do not have optimal strategies that are stationary (think of the ice-vendor where closer towards the final selling day the ice-vendor will store less ice-cream). Nonetheless, there are examples with optimal stationary policies (think of Tic-Tac-Toe) and the computations are simpler to understand the policy gradient theorems. Once the structure of the gradients is understood we proceed with the discounted infinite-time case.

4.1.1 Finite-time undiscounted MDPs, stationary policies

In this didactic section we will derive three different expressions of the gradient. The first theorem will show that differentiability of $\theta \mapsto J_s(\theta)$ follows from the differentiability of the parametrised policy and gives a first formula for the gradient as an expectation. In everything that follows let us write

$$R_t^T = \sum_{k=t}^{T-1} R_{k+1}.$$

for the rewards after time t (this is also called a reward to go) and R_0^T for the total (non-discounted) reward.



Theorem 4.1.1. Assume that (S, A, R) is a T -step MDP with finite state-action spaces and consider a stationary differentiable parameterized family of policies $\{\pi^\theta : \theta \in \Theta\}$. Then the gradient of the value function with respect to θ exists and

¹Notation anpassen. Rewards, Übergangswkeiten haben alte Form. Ueberall auf Anfangsbed μ

²Sutton, McAlister, Singh, Mansour: "Policy Gradient Methods for Reinforcement Learning with Function Approximation", NIPS, 1999



is given by

$$\nabla_{\theta} J_s(\theta) = \mathbb{E}_s^{\pi^{\theta}} \left[R_0^T \sum_{t=0}^{T-1} \nabla_{\theta} (\log \pi^{\theta}(A_t; S_t)) \right].$$

Proof. Consider a trajectory $\tau = (s_0, a_0, r_0, \dots, s_{T-1}, a_{T-1}, r_{T-1}, s_T, r_T)$ of the T -step MDP and recall the MDP path probability

$$\mathbb{P}_s^{\pi^{\theta}}((S, A, R) = \tau) = \delta_s(s_0) \delta_0(r_0) \prod_{t=0}^{T-2} \pi^{\theta}(a_t; s_t) p(s_{t+1}, r_{t+1}; s_t, a_t) \pi^{\theta}(a_T; s_T) p(s_T; s_{T-1}, a_{T-1})$$

from (2.10). Define

$$\mathcal{T} = \{ \tau = (r_0, s_0, a_0, r_1, \dots, s_{T-1}, a_{T-1}, r_T, s_T) : r_t \in \mathcal{R}, s_t \in \mathcal{S}, \forall t \leq T, a_t \in \mathcal{A}, \forall t < T \}$$

as the set of all trajectories. From the definition of the value function for a T -step MDP we deduce that

$$J_s(\theta) = V^{\pi^{\theta}}(s) = \mathbb{E}_s^{\pi^{\theta}} \left[\sum_{t=0}^{T-1} R_{t+1} \right] = \sum_{\tau \in \mathcal{T}} \mathbb{P}_s^{\pi^{\theta}}((S, A, R) = \tau) \sum_{t=0}^{T-1} r_{t+1}.$$

The probabilities are clearly differentiable, thus, for finite state-action states $J(\theta)$ is a finite sum of differentiable functions. This proves the differentiability. Next we use the log-trick we have already seen for policy gradient in bandits. We have

$$\begin{aligned} & \nabla_{\theta} \mathbb{P}_s^{\pi^{\theta}}((S, A, R) = \tau) \\ &= \nabla_{\theta} (\mathbb{P}_s^{\pi^{\theta}}((S, A, R) = \tau)) \frac{\mathbb{P}_s^{\pi^{\theta}}((S, A, R) = \tau)}{\mathbb{P}_s^{\pi^{\theta}}((S, A, R) = \tau)} \\ &= \nabla_{\theta} (\log(\mathbb{P}_s^{\pi^{\theta}}((S, A, R) = \tau))) \mathbb{P}_s^{\pi^{\theta}}((S, A, R) = \tau) \\ &= \nabla_{\theta} \left(\log(\delta_s(s_0)) + \log(\delta_0(r_0)) + \sum_{t=0}^{T-1} (\log(\pi^{\theta}(a_t; s_t)) + \log(p(\{s_{t+1}, r_{t+1}\}; s_t, a_t))) \right) \\ & \quad \times \mathbb{P}_s^{\pi^{\theta}}((S, A, R) = \tau) \\ &= \sum_{t=0}^{T-1} \nabla_{\theta} (\log(\pi^{\theta}(a_t; s_t))) \mathbb{P}_s^{\pi^{\theta}}((S, A, R) = \tau). \end{aligned}$$

Due to the log-trick the product of probabilities factors into a sum, where just the summands π^{θ} depend on θ and so all other summands have derivative 0. Finally we can use the finiteness of the state, action and reward space to see that we can interchange the derivative and the sum, and it follows that $V^{\pi^{\theta}}$ is differentiable. We conclude with

$$\begin{aligned} \nabla_{\theta} J_s(\theta) &= \nabla_{\theta} V^{\pi^{\theta}}(s) \\ &= \sum_{\tau \in \mathcal{T}} \nabla_{\theta} \mathbb{P}_s^{\pi^{\theta}}((S, A, R) = \tau) \sum_{t=0}^{T-1} r_{t+1} \\ &= \sum_{\tau \in \mathcal{T}} \mathbb{P}_s^{\pi^{\theta}}((S, A, R) = \tau) \sum_{t=0}^{T-1} \nabla_{\theta} (\log(\pi^{\theta}(a_t; s_t))) \sum_{t=0}^{T-1} r_{t+1} \\ &= \mathbb{E}_s^{\pi^{\theta}} \left[\sum_{t=0}^{T-1} \nabla_{\theta} (\log(\pi^{\theta}(A_t; S_t))) \sum_{t=0}^{T-1} R_{t+1} \right] \\ &= \mathbb{E}_s^{\pi^{\theta}} \left[\sum_{t=0}^{T-1} \nabla_{\theta} (\log(\pi^{\theta}(A_t; S_t))) R_0^T \right]. \end{aligned}$$

□

The formula for $\nabla J_s(\theta)$ crucially involves $\nabla(\log(\pi^\theta(a; s)))$ that in Section 1.3.4 was already called the score-function of the policy.



Definition 4.1.2. If π is a policy, then the vector $\nabla_\theta(\log(\pi^\theta(a; s)))$ is called the score-function of π^θ .

We have already computed the score-function of the (one-state) softmax policy for stochastic bandits. To get a feeling please redo the computation:



Show for the tabular softmax parametrisation from Example 4.0.2 that

$$\frac{\partial \log(\pi^\theta(a; s))}{\partial \theta_{s', a'}} = \mathbf{1}_{\{s=s'\}}(\mathbf{1}_{\{a=a'\}} - \pi^\theta(a'; s'))$$

and for the linear softmax with features $\Phi(s, a)$

$$\nabla \log(\pi^\theta(a; s)) = \Phi(s, a) - \sum_{a'} \pi^\theta(a'; s) \Phi(s, a').$$

Given the previous theorem, we weight the gradient of the log-probability of the chosen action with the reward of the whole trajectory R_0^T . So the gradient with respect to the t -th summand (and thus the t -th action) is weighted with R_0^T , but due to the Markov property the action in time t is completely independent from the past time points $0, \dots, t-1$. In the next theorem we will see that we can indeed replace R_0^T by the rewards of the future time points $t' \geq t$. We replace the reward of the whole trajectory by the reward to go R_t^T .



Theorem 4.1.3. Assume that (S, A, R) is a T -step MDP with finite state-action spaces and consider a stationary differentiable parameterized family of policies $\{\pi^\theta : \theta \in \Theta\}$. Then the gradient of the value function can also be written as

$$\nabla_\theta J_s(\theta) = \mathbb{E}_s^{\pi^\theta} \left[\sum_{t=0}^{T-1} \nabla_\theta (\log \pi^\theta(A_t; S_t)) R_t^T \right].$$

Proof. From Theorem 4.1.1 we have

$$\begin{aligned} \nabla_\theta J_s(\theta) &= \mathbb{E}_s^{\pi^\theta} \left[\sum_{t=0}^{T-1} \nabla_\theta (\log(\pi^\theta(A_t; S_t))) \sum_{t=0}^{T-1} R_{t+1} \right] \\ &= \sum_{t'=0}^{T-1} \sum_{t^*=0}^{T-1} \mathbb{E}_s^{\pi^\theta} \left[\nabla_\theta (\log(\pi^\theta(A_{t'}; S_{t'}))) R_{t^*+1} \right]. \end{aligned}$$

For every $t^* < t'$ we see

$$\begin{aligned} &\mathbb{E}_s^{\pi^\theta} \left[\nabla_\theta (\log(\pi^\theta(A_{t'}; S_{t'}))) R_{t^*+1} \right] \\ &= \sum_{\tau \in \mathcal{T}} \mathbb{P}_s^{\pi^\theta}((S, A, R) = \tau) \nabla_\theta (\log(\pi^\theta(a_{t'}; s_{t'}))) r_{t^*+1} \\ &= \sum_{a_0} \sum_{s_1} \sum_{r_1} \cdots \sum_{a_{T-1}} \sum_{s_T} \sum_{r_T} \prod_{t=0}^{T-1} \pi^\theta(a_t; s_t) p(\{s_{t+1}, r_{t+1}\}; s_t, a_t) (\nabla_\theta (\log(\pi^\theta(a_{t'}; s_{t'}))) r_{t^*+1}) \\ &= \sum_{a_0} \sum_{s_1} \sum_{r_1} \cdots \sum_{s_{t^*}} \sum_{r_{t^*}} \prod_{t=0}^{t^*} \pi^\theta(a_t; s_t) p(\{s_{t+1}, r_{t+1}\}; s_t, a_t) \sum_{a_{t^*}} \pi^\theta(a_{t^*}; s_{t^*}) \end{aligned}$$

$$\begin{aligned}
& \times \sum_{s_{t^*+1}} \sum_{r_{t^*+1}} p(\{s_{t^*+1}, r_{t^*+1}\}; s_{t^*}, a_{t^*}) r_{t^*+1} \sum_{a_{t^*+1}} \pi^\theta(a_{t^*+1}; s_{t^*+1}) \times \dots \\
& \times \sum_{s_{t'}} \sum_{r_{t'}} p(\{s_{t'}, r_{t'}\}; s_{t'}, a_{t'}) \underbrace{\sum_{a_{t'}} \pi^\theta(a_{t'}; s_{t'}) \nabla_\theta (\log(\pi^\theta(a_{t'}; s_{t'})))}_{= \sum_{a_{t'}} \nabla_\theta \pi^\theta(a_{t'}; s_{t'}) = \nabla_\theta \sum_{a_{t'}} \pi^\theta(a_{t'}; s_{t'}) = \nabla_\theta 1 = 0} \times \overbrace{1}^{\text{sum of probabilities}} \\
& = 0.
\end{aligned}$$

Now the claim follows. \square

In the last step we will replace the reward-to-go with the Q-function.



Theorem 4.1.4. Assume that (S, A, R) is a T -step MDP with finite state-action spaces and consider a stationary differentiable parameterized family of policies $\{\pi^\theta : \theta \in \Theta\}$. Then the gradient of the value function can also be written as

$$\nabla_\theta J_s(\theta) = \mathbb{E}_s^{\pi^\theta} \left[\sum_{t=0}^{T-1} \nabla_\theta (\log \pi^\theta(A_t; S_t)) Q_t^{\pi^\theta}(S_t, A_t) \right].$$

Proof. For the proof we use the property of a Markov property of Markov reward processes:

$$\mathbb{E}_s^{\pi^\theta} [R_t^T | S_t = s, A_t = a] = Q_t^{\pi^\theta}(s, a).$$

Then the claim follows from Theorem 4.1.3:

$$\begin{aligned}
\nabla_\theta J_s(\theta) &= \sum_{t=0}^{T-1} \mathbb{E}_\mu^{\pi^\theta} [\nabla_\theta (\log(\pi^\theta(A_t; S_t))) R_t^T] \\
&= \sum_{t=0}^{T-1} \sum_{s \in \mathcal{S}, a \in \mathcal{A}_s} \mathbb{E}_\mu^{\pi^\theta} [\nabla_\theta (\log(\pi^\theta(A_t; S_t))) R_t^T | S_t = s, A_t = a] \mathbb{P}_\mu^{\pi^\theta}(S_t = s, A_t = a) \\
&= \sum_{t=0}^{T-1} \sum_{s \in \mathcal{S}, a \in \mathcal{A}_s} \nabla_\theta (\log(\pi^\theta(a; s))) \mathbb{E}_\mu^{\pi^\theta} [R_t^T | S_t = s, A_t = a] \mathbb{P}_\mu^{\pi^\theta}(S_t = s, A_t = a) \\
&= \sum_{t=0}^{T-1} \sum_{s \in \mathcal{S}, a \in \mathcal{A}_s} \nabla_\theta (\log(\pi^\theta(a; s))) Q_t^T(s, a) \mathbb{P}_\mu^{\pi^\theta}(S_t = s, A_t = a) \\
&= \mathbb{E}_s^{\pi^\theta} \left[\sum_{t=0}^{T-1} \nabla_\theta (\log(\pi^\theta(A_t; S_t))) Q_t^{\pi^\theta}(S_t, A_t) \right]
\end{aligned}$$

\square

For intuition the Q-function indicates the *expected* reward-to-go and not the reward-to-go for a specific trajectory. More precisely, we have to distinguish between the outer expectation and the expectation in the definition of the Q-function. While the reward-to-go depends on the trajectory of the outer expectation, the Q-function calculates its own expectation. For the moment one might wonder whether there is any use in this representation, Q^{π^θ} is not known. The miracle why this can be useful will be solved later when we discuss actor-critic methods for which a second parametrisation is used for Q .



Definition 4.1.5. In general optimization methods that aim to maximise f and the gradient of f takes the form of an expectation that can be sampled are typically



called stochastic gradient methods. The usual form of the gradients is

$$\nabla f(\theta) = \mathbb{E}[g(X, \theta)],$$

for policy gradient the form is more delicate

$$\nabla f(\theta) = \mathbb{E}^\theta[g(X, \theta)].$$

If (unbiased) samples $\tilde{\nabla} f(\theta)$ for the gradients are available (samples of the expectation), the update scheme

$$\theta \leftarrow \theta + \alpha \tilde{\nabla} f(\theta)$$

is called stochastic gradient algorithm. If more than one sample is used to estimate the gradient as a Monte Carlo average the algorithm is called a batch stochastic gradient ascent algorithm. The number of samples averaged is called the batch size.

Let us come back to the goal of performing a gradient ascent algorithm to maximise V^{π^θ} with initial distribution μ . Using the policy gradient representations the gradient ascent update can be written as

$$\theta_{t+1} = \theta_t + \alpha \mathbb{E}_\mu^{\pi^\theta} \left[\sum_{t=0}^{T-1} \nabla_\theta (\log \pi^\theta(A_t; S_t)) Q_t^{\pi^\theta}(S_t, A_t) \right]$$

or

$$\theta_{t+1} = \theta_t + \alpha \mathbb{E}_\mu^{\pi^\theta} \left[\sum_{t=0}^{T-1} \nabla_\theta (\log \pi^\theta(A_t; S_t)) \sum_{t'=t}^{T-1} R_{t'+1} \right].$$

Unfortunately, the exact expectations are typically unknown and must be estimated. Similar to the approximate dynamic programming chapter we rewrote the gradient as an expectation and thus can estimate the expectation with samples. As earlier the first obvious idea to estimate the gradients is the model-free Monte Carlo method. As the Q -function appearing inside the expectation is an unknown value (and an expectation itself) it is most common in applications to use the second policy gradient theorem instead of the third. This can be interpreted as estimating the Q -function by the unbiased estimator R_t^T , the reward-to-go. An estimator for the gradient is then given by

$$\tilde{\nabla} J_\mu(\theta) = \frac{1}{K} \sum_{i=1}^K \left[\sum_{t=0}^{T-1} \nabla_\theta (\log \pi^\theta(a_t^i; s_t^i)) \sum_{t'=t}^{T-1} r_{t'+1} \right], \quad (4.1)$$

with K trajectories $(s_0^i, a_0^i, s_1^i, r_1^i, \dots, a_{T-1}^i, s_T^i, r_T^i)$ sampled according to the current policy π^θ with initial distribution μ .

This famous algorithm, even not in exactly that form, is due to Ronald Williams³. To use the algorithm in practice a stopping condition needs to be fixed. Typically, either the number of iterations is fixed or the algorithm is terminated when the norm of the gradient reaches a small threshold.



Sampling a trajectory sounds harmless for a Mathematician. But it is not at all! Imagine we learn how to steer a car or a certain robotic task and π^θ is a policy. Then sampling means running the car (or the robot) K times with the current (possibly very bad) policy while evaluating the rewards! It is thus easy to imagine pitfalls

³Williams: "Simple Statistical Gradient-Following Algorithms for Connectionist Reinforcement Learning", Machine Learning, 8, 229-256, 1992

Data: Initial parameter θ_0 , $K \geq 1$, initial distribution μ

Result: Approximate policy $\pi^{\theta_L} \approx \pi^{\theta^*}$

$l = 1$.

while *not converged* **do**

for $i = 1, \dots, K$ **do**

 Sample initial condition s_0^i from μ .

 Sample trajectory $(s_0^i, a_0^i, s_1^i, r_1^i, \dots, a_{T-1}^i, s_T^i, r_T^i)$ using policy $\pi^{\theta_{l-1}}$.

end

 Determine step-size α .

$$\tilde{\nabla} J(\theta_{l-1}) = \frac{1}{K} \sum_{i=1}^K \left[\sum_{t=0}^{T-1} \nabla_{\theta} (\log \pi^{\theta_{l-1}}(a_t^i; s_t^i)) \sum_{t'=t}^{T-1} r_{t'+1}^i \right]$$

$$\theta_l = \theta_{l-1} + \alpha \tilde{\nabla} J(\theta_{l-1})$$

end

Set $l = l + 1$.

Algorithm 17: REINFORCE: (Batch-)Stochastic policy gradient algorithm



for practical applications. Sample efficiency (using as little samples as possible) in learning is essential and also sampling from very bad policies might be a costly endeavour.

4.1.2 Infinite-time MDPs with discounted rewards

In this section we will consider MDPs with infinite time horizon and discounted rewards. Recall that in this setting the use of stationary policies is justified. As before we aim to rewrite the gradient of the value function as an expectation, but now the infinite time horizon makes it much more complex to calculate the gradient of the value function. The objective function using a parametrised policy π^{θ} is given by

$$J_s(\theta) = V^{\pi^{\theta}}(s) = \mathbb{E}_s^{\pi^{\theta}} \left[\sum_{t=0}^{\infty} \gamma^t R_t \right] = \sum_{a \in \mathcal{A}_s} \pi^{\theta}(a; s) Q^{\pi^{\theta}}(s, a). \quad (4.2)$$

The infinite sum makes the computation of the gradient more challenging.



Lemma 4.1.6. Under the assumption that J_s is differentiable for every start state $s \in \mathcal{S}$ we have a closed form of the gradient

$$\nabla J_s(\theta) = \sum_{s' \in \mathcal{S}} \sum_{a \in \mathcal{A}_{s'}} \rho_s^{\pi^{\theta}}(s') \nabla \pi^{\theta}(a; s') Q^{\pi^{\theta}}(s', a),$$

where $\rho_s^{\pi}(s') = \sum_{t=0}^{\infty} \gamma^t \mathbb{P}_{\mu}^{\pi}(S_t = s') = \mathbb{E}_{\mu}^{\pi}[\sum_{t=0}^{\infty} \gamma^t \mathbf{1}_{S_t=s'}]$.

Proof. The trick is essentially simple but tedious to write down. One applies the gradient to the right-hand side of (4.2) and applies the product rule. The claim follows from an induction, repeating this process.

With the notation $p(s \rightarrow s'; t, \pi) = \mathbb{P}_s^{\pi}(S_t = s')$ we first show the following identity by induction:



For all $n \in \mathbb{N}_0$ it holds that

$$\begin{aligned} \nabla J_s(\theta) &= \sum_{t=0}^n \sum_{s' \in \mathcal{S}} \gamma^t p(s \rightarrow s'; t, \pi^{\theta}) \sum_{a \in \mathcal{A}_{s'}} \nabla \pi^{\theta}(a; s') Q^{\pi^{\theta}}(s', a) \\ &\quad + \sum_{s' \in \mathcal{S}} \gamma^{n+1} p(s \rightarrow s'; n+1, \pi^{\theta}) \nabla J_{s'}(\theta). \end{aligned} \quad (4.3)$$

By the chain rule we have that

$$\begin{aligned}
\nabla J_s(\theta) &= \sum_{a \in \mathcal{A}_s} \left(\nabla \pi^\theta(a; s) Q^{\pi^\theta}(s, a) + \pi^\theta(a; s) \nabla Q^{\pi^\theta}(s, a) \right) \\
&= \sum_{a \in \mathcal{A}_s} \left(\nabla \pi^\theta(a; s) Q^{\pi^\theta}(s, a) + \pi^\theta(a; s) \nabla (r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'; s, a) V^{\pi^\theta}(s')) \right) \\
&= \sum_{a \in \mathcal{A}_s} \left(\nabla \pi^\theta(a; s) Q^{\pi^\theta}(s, a) + \pi^\theta(a; s) \gamma \sum_{s' \in \mathcal{S}} p(s'; s, a) \nabla J_{s'}(\theta) \right).
\end{aligned}$$

This is the induction for $n = 0$. For the inductive step suppose the statement holds for some $n \in \mathbb{N}_0$. Then, a similar computation gives

$$\begin{aligned}
\nabla J_s(\theta) &= \sum_{t=0}^n \sum_{s' \in \mathcal{S}} \gamma^t p(s \rightarrow s'; t, \pi^\theta) \sum_{a \in \mathcal{A}_{s'}} \nabla \pi^\theta(a; s') Q^{\pi^\theta}(s', a) \\
&\quad + \sum_{s' \in \mathcal{S}} \gamma^{n+1} p(s \rightarrow s'; n+1, \pi^\theta) \nabla J_{s'}(\theta) \\
&= \sum_{t=0}^n \sum_{s' \in \mathcal{S}} \gamma^t p(s \rightarrow s'; t, \pi^\theta) \sum_{a \in \mathcal{A}_{s'}} \nabla \pi^\theta(a; s') Q^{\pi^\theta}(s', a) \\
&\quad + \sum_{s' \in \mathcal{S}} \gamma^{n+1} p(s \rightarrow s'; n+1, \pi^\theta) \\
&\quad \times \left(\sum_{a' \in \mathcal{A}_{s'}} \left(\nabla \pi^\theta(a'; s') Q^{\pi^\theta}(s', a') + \pi^\theta(a'; s') \gamma \sum_{s'' \in \mathcal{S}} p(s''; s', a') \nabla J_{s''}(\theta) \right) \right) \\
&= \sum_{t=0}^{n+1} \sum_{s' \in \mathcal{S}} \gamma^t p(s \rightarrow s'; t, \pi^\theta) \sum_{a \in \mathcal{A}_{s'}} \nabla \pi^\theta(a; s') Q^{\pi^\theta}(s', a) \\
&\quad + \sum_{s' \in \mathcal{S}} \gamma^{n+1} p(s \rightarrow s'; n+1, \pi^\theta) \cdot \left(\sum_{a' \in \mathcal{A}_{s'}} \pi^\theta(a'; s') \gamma \sum_{s'' \in \mathcal{S}} p(s''; s', a') \nabla J_{s''}(\theta) \right) \\
&= \sum_{t=0}^{n+1} \sum_{s' \in \mathcal{S}} \gamma^t p(s \rightarrow s'; t, \pi^\theta) \sum_{a \in \mathcal{A}_{s'}} \nabla \pi^\theta(a; s') Q^{\pi^\theta}(s', a) \\
&\quad + \sum_{s'' \in \mathcal{S}} \gamma^{n+2} p(s \rightarrow s''; n+2, \pi^\theta) \nabla J_{s''}(\theta).
\end{aligned}$$



The remainder term in (4.3) vanishes for $n \rightarrow \infty$.

$$\left| \sum_{s'' \in \mathcal{S}} \gamma^{n+1} p(s \rightarrow s''; n+1, \pi^\theta) \nabla J_{s''}(\theta) \right| \leq \gamma^{n+1} \max_{s \in \mathcal{S}} |\nabla J_s(\theta)| \rightarrow 0, \quad n \rightarrow \infty.$$

Hence, we proved that

$$\begin{aligned}
\nabla J_s(\theta) &= \sum_{t=0}^{\infty} \sum_{s' \in \mathcal{S}} \gamma^t p(s \rightarrow s'; t, \pi^\theta) \sum_{a \in \mathcal{A}_{s'}} \nabla \pi^\theta(a; s') Q^{\pi^\theta}(s', a) \\
&= \sum_{s' \in \mathcal{S}} \rho_s^{\pi^\theta}(s') \sum_{a \in \mathcal{A}_{s'}} \nabla \pi^\theta(a; s') Q^{\pi^\theta}(s', a).
\end{aligned}$$

The exchange of limits and sums is justified as we assume \mathcal{S} to be finite. □



For episodic MDPs (the MDP terminates almost surely under all policies π_θ), we can get rid of the assumption of the existence of $\nabla J_s(\theta)$. Go through the proof of Theorem 4.1.6 and argue why it is enough to assume the existence of $\nabla \pi_\theta(\cdot; s)$ for all $s \in \mathcal{S}$.

Taking a closer look at the expression from the gradient, it is obvious that the infinite sum $\rho^\pi(s') = \mathbb{E}_s^\pi[\sum_{t=0}^{\infty} \gamma^t \mathbf{1}_{S_t=s'}]$ is unpleasant for implementations. There is an alternative way to rewrite the policy gradient theorem using the discounted state visitation measure:



Definition 4.1.7. The discounted state-visitation measure under policy π for starting state $s \in \mathcal{S}$ is given by

$$d_\mu^\pi(s) := \frac{\rho_\mu^\pi(s)}{\sum_{s'} \rho_\mu^\pi(s')}.$$

First note that, using Fubini's theorem,

$$\sum_{s' \in \mathcal{S}} \rho_\mu^\pi(s') = \sum_{s' \in \mathcal{S}} \mathbb{E}_\mu^\pi \left[\sum_{t=0}^{\infty} \gamma^t \mathbf{1}_{S_t=s'} \right] = \frac{1}{1-\gamma}.$$

Hence, we can write $d^\pi(s) = (1-\gamma)\rho^\pi(s)$ for the discounted state-visitation measure. With this definition the gradient can also be written as follows:



Theorem 4.1.8. Assume that $J_s(\theta)$ is differentiable for every start state $s \in \mathcal{S}$ and $\pi^\theta(a; s) \neq 0$ for all $\theta \in \Theta, s \in \mathcal{S}$ and $a \in \mathcal{A}_s$. Then, the gradient can be expressed as

$$\nabla J_s(\theta) = \frac{1}{1-\gamma} \mathbb{E}_{S \sim d_s^{\pi^\theta}, A \sim \pi^\theta(\cdot; S)} [\nabla \log(\pi^\theta(A; S)) Q^{\pi^\theta}(S, A)].$$

In essence the theorem says the following. Draw a state-action pair (s, a) according to their relevance measured in terms of expected discounted visitation frequency, compute the direction using the score function, and follow that direction according to the expected reward.

Proof. By Theorem 4.1.6 we have

$$\begin{aligned} \nabla J_s(\theta) &= \sum_{s' \in \mathcal{S}} \rho^{\pi^\theta}(s') \sum_{a \in \mathcal{A}_{s'}} \nabla \pi^\theta(a; s') Q^{\pi^\theta}(s', a) \\ &= \frac{1}{1-\gamma} \sum_{s' \in \mathcal{S}} d_s^{\pi^\theta}(s') \sum_{a \in \mathcal{A}_{s'}} \pi^\theta(a; s') \nabla \log(\pi^\theta(a; s')) Q^{\pi^\theta}(s', a) \\ &= \frac{1}{1-\gamma} \mathbb{E}_{S \sim d_s^{\pi^\theta}, A \sim \pi^\theta(\cdot; S)} [\nabla \log(\pi^\theta(A; S)) Q^{\pi^\theta}(S, A)], \end{aligned}$$

where we used the definition of the state-visitation measure and the log-trick. The second equality gives the first claim, the final equality the second claim. \square

For practical use it will be important to estimate the gradient (an expectation) using Monte Carlo. Thus, it will be necessary to sample from the discounted occupation measures and then compute score-function and Q -value. On first sight sampling from the d^π looks infeasible as an infinite sum is involved. In fact, a sample can be obtained by following a rollout up to an independent geometric random variable, counting the number of visitations and then drawing from this empirical distribution. This is justified by a quick computation:



Lemma 4.1.9. If $T \sim \text{Geo}(1 - \gamma)$ is independent of the MDP, then

$$\rho_\mu^\pi(s) = \mathbb{E}_\mu^\pi \left[\sum_{t=0}^T \mathbf{1}_{S_t=s} \right] \quad \text{and} \quad d_\mu^\pi(s) = \mathbb{E}_\mu^\pi [\mathbf{1}_{S_T=s}].$$

In particular, sampling from d_μ^π can be carried out as follows.

1. Draw $T \sim \text{Geo}(1 - \gamma)$, independently of the MDP.
2. Run the MDP under policy π from initial distribution μ until time T .
3. Output the terminal state S_T .

Alternatively, $(1 - \gamma)\hat{\rho}(s)$ provides an unbiased estimator of $d_\mu^\pi(s)$, where

$$\hat{\rho}(s) := \sum_{t=0}^T \mathbf{1}_{S_t=s}.$$

Proof. This is easily seen as

$$\begin{aligned} \mathbb{E}_\mu^\pi \left[\sum_{t=0}^T \mathbf{1}_{S_t=s} \right] &= \mathbb{E}_\mu^\pi \left[\sum_{t=0}^{\infty} \mathbf{1}_{t \leq T} \mathbf{1}_{S_t=s} \right] \\ &\stackrel{\text{Fubini}}{=} \sum_{t=0}^{\infty} \mathbb{E}_\mu^\pi [\mathbf{1}_{t \leq T} \mathbf{1}_{S_t=s}] \\ &\stackrel{\text{ind.}}{=} \sum_{t=0}^{\infty} \mathbb{P}(t \leq T) \mathbb{E}_\mu^\pi [\mathbf{1}_{S_t=s}] \\ &= \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_\mu^\pi [\mathbf{1}_{S_t=s}] \\ &= \mathbb{E}_\mu^\pi \left[\sum_{t=0}^{\infty} \gamma^t \mathbf{1}_{S_t=s} \right] = \rho_\mu^\pi(s) \end{aligned}$$

and

$$\mathbb{E}_s^\pi [\mathbf{1}_{S_T=s}] = \sum_{t=0}^{\infty} \mathbb{P}(T = t) \mathbb{P}^\pi(S_t = s) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}^\pi(S_t = s) = d^\pi(s).$$

□



Lemma 4.1.10. Let $F : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ be a bounded function. Then,

$$\mathbb{E}_\mu^\pi \left[\sum_{t=0}^{\infty} \gamma^t F(S_t, A_t) \right] = \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \rho_\mu^\pi(s) \pi(a; s) F(s, a).$$

Proof.

$$\begin{aligned} \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \rho_\mu^\pi(s) \pi(a; s) F(s, a) &= \sum_{t=0}^{\infty} \gamma^t \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \mathbb{P}_\mu^\pi(S_t = s) \pi(a; s) F(s, a) \\ &= \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_\mu^\pi [F(S_t, A_t)] \\ &= \mathbb{E}_\mu^\pi \left[\sum_{t=0}^{\infty} \gamma^t F(S_t, A_t) \right]. \end{aligned}$$

The final exchange of sum and expectation is justified by dominated convergence and the boundedness of F . \square

Here is another representation in the spirit of the finite-time policy gradient theorem. The representation avoids the sampling from the discounted occupation measure but instead uses full trajectories. It follows directly from Lemma 4.1.6 and Lemma 4.1.10.



Theorem 4.1.11. Suppose that $(s, a) \mapsto \nabla_{\theta}(\log \pi^{\theta}(a; s))Q^{\pi^{\theta}}(s, a)$ is bounded. Then

$$\nabla_{\theta} J_s(\theta) = \mathbb{E}_s^{\pi^{\theta}} \left[\sum_{t=0}^{\infty} \gamma^t \nabla_{\theta}(\log \pi^{\theta}(A_t; S_t)) Q^{\pi^{\theta}}(S_t, A_t) \right].$$

⁴ The assumption does not hold for arbitrary parametrisation and might also be hard to check. For softmax policies the score-function can be computed and is bounded for instance for bounded feature vector. If the rewards are additionally bounded (which implies the Q -function is bounded) the assumption of the theorem holds. A dirty version of the REINFORCE algorithm is then obtained by running trajectories up to some large T and to use the truncated series as an estimator for the gradient. Since the truncated estimator is not unbiased smarter approaches have been proposed. As for the discounted occupation measure one might ask if geometric times can be used to truncate the infinite time-horizon. Under certain assumptions on the policy parametrisation this can be done.



Assumption 4.1.12. The policy π^{θ} is differentiable with respect to θ and $\nabla \log(\pi^{\theta}(a; s))$ exists and is L_{Θ} -smooth and has bounded norm for any $(s, a) \in \mathcal{S} \times \mathcal{A}$, i.e.

$$\|\nabla \log(\pi^{\theta_1}(a; s)) - \nabla \log(\pi^{\theta_2}(a; s))\| \leq L_{\Theta} \|\theta_1 - \theta_2\|, \quad \text{for any } \theta_1, \theta_2 \in \Theta,$$

$$\|\nabla \log(\pi^{\theta}(a; s))\| \leq B_{\Theta}, \quad \text{for any } \theta \in \Theta,$$

for some $L_{\Theta}, B_{\Theta} > 0$.

This assumption is for example fulfilled by tabular and linear softmax parametrisations or by Gaussian policies.



Show that the tabular and linear softmax parametrisation fulfills Assumption 4.1.12.

Using a computation similar to the one for the occupation measure gives the following theorem⁵.



Proposition 4.1.13. Suppose that the policy parametrisation fulfills Assumption 4.1.12 and $T \sim \text{Geo}(1 - \gamma)$, $T' \sim \text{Geo}(1 - \gamma^{1/2})$ are independent of each other and the MDP, then

$$\nabla J(\theta) = \frac{1}{1 - \gamma} \mathbb{E}_{\mu}^{\pi^{\theta}} \left[\nabla \log(\pi^{\theta}(S_T; A_T)) \sum_{t=T}^{T+T'} \gamma^{(t-T)/2} R_t \right].$$

The advantage of this policy gradient theorem is clear, unbiased estimates can be obtained by running the MDP for a finite (random) number of steps.

⁴brauche auch die version mit reward to go, beweis wie oben mit reinbedingen

⁵Zhang, Koppel, Zhu, Basar: "Global Convergence of Policy Gradient Methods to (Almost) Locally Optimal Policies", SIAM Journal on Control and Optimization, 2020

Proof. As before we prove the claim for $\mu = \delta_s$ for some $s \in \mathcal{S}$ and in the end sum-out μ to get the general claim.



The term $\hat{Q}^{\pi^\theta}(s, a) := \sum_{t'=0}^{T'} \gamma^{t'/2} R_{t'}$ is an unbiased estimator of the Q-function:

$$\mathbb{E}_{s,a} \left[\sum_{t'=0}^{T'} \gamma^{t'/2} R_{t'} \right] = Q^{\pi^\theta}(s, a)$$

First note that by the bounded reward assumption

$$\mathbb{E}_{s,a} \left[\sum_{t=0}^N \mathbf{1}_{0 \leq t \leq T'} \gamma^{t/2} R_{t+1} \right] \leq R_* \mathbb{E}_{T'} \left[\sum_{t=0}^N \mathbf{1}_{0 \leq t \leq T'} \gamma^{t/2} \right].$$

If we write $E_s^{\pi^\theta}$ we mean the expectation with respect to all appearing random variables, i.e. also of T' . When we bound R by R_* , only the expectation with respect to T' is left, which is why we denote it with $\mathbb{E}_{T'}$. In the following we will write $\mathbb{E}_{T'}[\mathbb{E}_s^{\pi^\theta}[\cdot]]$ to highlight the independence of T' to the MDP. As the RHS of the inequality is monotonically increasing and the limit exists, we can use the monotone convergence theorem to show that

$$\begin{aligned} \mathbb{E}_{s,a}[\hat{Q}^{\pi^\theta}(s, a)] &= \mathbb{E}_{T'} \left[\mathbb{E}_{s,a} \left[\sum_{t'=0}^{T'} \gamma^{t'/2} R_{t'} \right] \right] \\ &= \mathbb{E}_{T'} \left[\mathbb{E}_{s,a} \left[\lim_{N \rightarrow \infty} \sum_{t'=0}^N \mathbf{1}_{0 \leq t' \leq T'} \gamma^{t'/2} R_{t'} \right] \right] \\ &= \lim_{N \rightarrow \infty} \sum_{t'=0}^N \mathbb{E}_{T'} \left[\mathbb{E}_{s,a} \left[\mathbf{1}_{0 \leq t' \leq T'} \gamma^{t'/2} R_{t'} \right] \right] \\ &= \lim_{N \rightarrow \infty} \sum_{t'=0}^N \mathbb{E}_{s,a} \left[\mathbb{E}_{T'}[\mathbf{1}_{0 \leq t' \leq T'}] \gamma^{t'/2} R_{t'} \right] \\ &= \lim_{N \rightarrow \infty} \sum_{t'=0}^N \mathbb{E}_{s,a} \left[\gamma^{t'} R_{t'} \right] \\ &= \lim_{N \rightarrow \infty} \mathbb{E}_{s,a} \left[\sum_{t'=0}^N \gamma^{t'} R_{t'} \right] \\ &= Q^{\pi^\theta}(s, a), \end{aligned}$$

where we have used that T' is independent of the MDP and that $\mathbb{E}_{T'}[\mathbf{1}_{0 \leq t' \leq T'}] = \mathbb{P}(T' \geq t) = \gamma^{t/2}$ by the $\text{Geo}(1 - \gamma^{1/2})$ distribution. In the last equation we used the dominated convergence theorem to bring the limit back inside the expectation.

For the second step we define the estimator $\hat{\nabla} J_s(\theta) = \frac{1}{1-\gamma} \nabla \log(\pi^\theta(A_T; S_T)) \hat{Q}^{\pi^\theta}(S_T, A_T)$ and show



$$\frac{1}{1-\gamma} \mathbb{E}_s^{\pi^\theta} \left[\nabla \log(\pi^\theta(S_T; A_T)) \sum_{t=T}^{T+T'} \gamma^{(t-T)/2} R_t \right] = \mathbb{E}_s^{\pi^\theta} \left[\hat{\nabla} J_s(\theta) \right]$$

This follows from the (strong) Markov property applied at T :

$$\begin{aligned} &\mathbb{E}_s^{\pi^\theta} \left[\nabla \log(\pi^\theta(S_T; A_T)) \sum_{t=T}^{T+T'} \gamma^{(t-T)/2} R_t \right] \\ &= \mathbb{E}_s^{\pi^\theta} \left[\nabla \log(\pi^\theta(S_T; A_T)) \mathbb{E}_s^{\pi^\theta} \left[\sum_{t=T}^{T+T'} \gamma^{(t-T)/2} R_t \mid S_T, A_T \right] \right] \end{aligned}$$

$$\begin{aligned}
&= \mathbb{E}_s^{\pi^\theta} \left[\nabla \log(\pi^\theta(S_T; A_T)) \mathbb{E}_{S_T, A_T}^{\pi^\theta} \left[\sum_{t=0}^{T'} \gamma^{t/2} R_t \right] \right] \\
&= \mathbb{E}_s^{\pi^\theta} \left[\nabla \log(\pi^\theta(S_T; A_T)) \hat{Q}^{\pi^\theta}(S_T, A_T) \right]
\end{aligned}$$

For the final equality we used the first step.

It remains to show the following claim:



$$\mathbb{E}_s^{\pi^\theta} \left[\widehat{\nabla} J_s(\theta) \right] = \nabla J_s(\theta)$$

As T and T' are independent and independent of the MDP, we have directly that

$$\begin{aligned}
\mathbb{E}_s^{\pi^\theta} \left[\widehat{\nabla} J_s(\theta) \right] &= \frac{1}{1-\gamma} \mathbb{E}_{s, T, T'}^{\pi^\theta} \left[\nabla \log(\pi^\theta(S_T; A_T)) \hat{Q}^{\pi^\theta}(S_T, A_T) \right] \\
&= \frac{1}{1-\gamma} \mathbb{E}_{s, T}^{\pi^\theta} \left[\nabla \log(\pi^\theta(S_T; A_T)) \mathbb{E}_{T', S_T, A_T}^{\pi^\theta} [\hat{Q}(S_T, A_T)] \right] \\
&= \frac{1}{1-\gamma} \mathbb{E}_{s, T}^{\pi^\theta} \left[\nabla \log(\pi^\theta(S_T; A_T)) Q(S_T, A_T) \right],
\end{aligned}$$

where the last equality is due to step 1. As Q is bounded by $\frac{R_*}{1-\gamma}$ and $\|\nabla \log(\pi^\theta(a; s))\| \leq B_\Theta$ by Assumption 4.1.12, we can again use the monotone convergence theorem as in step 1 to show that

$$\begin{aligned}
\mathbb{E}_s^{\pi^\theta} \left[\widehat{\nabla} J_s(\theta) \right] &= \frac{1}{1-\gamma} \mathbb{E}_{s, T}^{\pi^\theta} \left[\nabla \log(\pi^\theta(S_T; A_T)) Q(S_T, A_T) \right] \\
&= \frac{1}{1-\gamma} \mathbb{E}_{s, T}^{\pi^\theta} \left[\sum_{t=0}^{\infty} \mathbf{1}_{t=T} \nabla \log(\pi^\theta(S_t; A_t)) Q(S_t, A_t) \right] \\
&= \lim_{N \rightarrow \infty} \sum_{t=0}^N \mathbb{E}_T[\mathbf{1}_{t=T}] \frac{1}{1-\gamma} \mathbb{E}_s^{\pi^\theta} \left[\nabla \log(\pi^\theta(S_t; A_t)) Q(S_t, A_t) \right] \\
&= \lim_{N \rightarrow \infty} \sum_{t=0}^N \gamma^t \mathbb{E}_s^{\pi^\theta} \left[\nabla \log(\pi^\theta(S_t; A_t)) Q(S_t, A_t) \right] \\
&= \mathbb{E}_s^{\pi^\theta} \left[\sum_{t=0}^{\infty} \gamma^t \nabla \log(\pi^\theta(S_t; A_t)) Q(S_t, A_t) \right] \\
&= \nabla J_s(\theta),
\end{aligned}$$

where we used the independence of T to the MDP and that $\mathbb{E}_T[\mathbf{1}_{t=T}] \frac{1}{1-\gamma} = \frac{\mathbb{P}(T=t)}{1-\gamma} = \gamma^t$. The last equation is due to the policy gradient theorem 4.1.11 \square

Using this proposition results in the following REINFORCE algorithm with geometric rollouts. The representation is not only useful from a practical point of view. The random variables can be shown to have bounded variance, an ingredient that was crucial for the proof of convergence to a stationary point for stochastic gradient schemes. We will see this in the subsequent section. Combined with an additional PI inequality (only known for tabular softmax parametrisation) then also yields convergence to an optimal policy.

4.1.3 Convergence of REINFORCE

Under certain assumptions we can prove convergence of REINFORCE to stationary points using Theorem 3.5.1. First, we will need L -smoothness, which can be shown under Assumption 4.1.12. Secondly, we need an unbiased estimator with bounded variance of the gradient $\nabla J(\theta)$ for any

Data: Initial parameter θ_0 , $K \geq 1$, initial state s

Result: Approximate policy $\pi^{\theta_L} \approx \pi^{\theta^*}$

$l = 1$

while *not converged* **do**

for $i = 1, \dots, K$ **do**

 Sample $T_i \sim \text{Geo}(1 - \gamma)$

 Sample trajectory $(s_0^i = s, a_0^i, s_1^i, r_1^i, \dots, a_{T_i-1}^i, s_{T_i}^i, r_{T_i}^i, a_{T_i}^i)$ using policy $\pi^{\theta_{l-1}}$.

 Sample $T'_i \sim \text{Geo}(1 - \gamma^{\frac{1}{2}})$

 Set $\tilde{s}_0^i = s_{T_i}^i$ and $\tilde{a}_0^i = a_{T_i}^i$

 Sample trajectory $(\tilde{s}_0^i, \tilde{a}_0^i, \tilde{s}_1^i, \tilde{r}_1^i, \dots, \tilde{a}_{T'_i-1}^i, \tilde{s}_{T'_i}^i, \tilde{r}_{T'_i}^i, \tilde{a}_{T'_i}^i)$ using policy $\pi^{\theta_{l-1}}$.

end

 Determine step-size α .

$$\widehat{\nabla} J_s(\theta_{l-1}) = \frac{1}{1-\gamma} \frac{1}{K} \sum_{i=1}^K \left[\nabla_{\theta} (\log \pi^{\theta_{l-1}}(a_{T_i}^i; s_{T_i}^i)) \sum_{t'=0}^{T'_i-1} \gamma^{t'/2} \tilde{r}_{t'+1}^i \right]$$

$$\theta_l = \theta_{l-1} - \alpha \widehat{\nabla} J_s(\theta_{l-1})$$

end

Set $l = l + 1$.

Algorithm 18: Mini-batch REINFORCE for infinite time horizon

$\theta \in \mathbb{R}^d$. This can be achieved using the geometric rollouts discussed in the previous section. The unbiasedness has already been shown in Proposition 4.1.13 and the bounded variance is discussed below. The algorithm we consider is the REINFORCE Algorithm 18.

First we prove L -smoothness of the objective.



Lemma 4.1.14. Under Assumption 4.1.12, the objective $J_{s_0}(\theta)$ is L -smooth with $L = \frac{R_* L_{\Theta}}{(1-\gamma)^2} + \frac{(1+\gamma)R_* B_{\Theta}^2}{(1-\gamma)^3}$, where R_* is the maximal reward of the bounded reward assumption.

Proof. From the policy gradient theorem 4.1.8 we obtain for an initial state s_0 that

$$\nabla J_{s_0}(\theta) = \frac{1}{1-\gamma} \sum_{s \in \mathcal{S}, a \in \mathcal{A}} d_{s_0}^{\pi^{\theta}}(s) \pi^{\theta}(a; s) \nabla \log(\pi^{\theta}(a; s)) Q^{\pi^{\theta}}(s, a).$$

Now note that we can rewrite Q to be

$$Q^{\pi^{\theta}}(s, a) = \sum_{t'=0}^{\infty} \gamma^{t'} \sum_{s' \in \mathcal{S}, a' \in \mathcal{A}} p((s, a) \rightarrow s'; t', \pi^{\theta}) \pi^{\theta}(a'; s') r(s', a'),$$

where $p((s, a) \rightarrow s'; t', \pi^{\theta}) = \mathbb{P}_s^{\pi_a^{\theta}}(S_{t'} = s')$. Using the definition of the state-visitation measure and this form of $Q^{\pi^{\theta}}$ we obtain for the gradient

$$\begin{aligned} & \nabla J_{s_0}(\theta) \\ &= \sum_{s \in \mathcal{S}, a \in \mathcal{A}} \sum_{t=0}^{\infty} \gamma^t p(s_0 \rightarrow s; t, \pi^{\theta}) \pi^{\theta}(a; s) \nabla \log(\pi^{\theta}(a; s)) Q^{\pi^{\theta}}(s, a) \\ &= \sum_{t=0}^{\infty} \gamma^t \sum_{s \in \mathcal{S}, a \in \mathcal{A}} p(s_0 \rightarrow s; t, \pi^{\theta}) \pi^{\theta}(a; s) \nabla \log(\pi^{\theta}(a; s)) \sum_{t'=0}^{\infty} \gamma^{t'} \sum_{s' \in \mathcal{S}, a' \in \mathcal{A}} p((s, a) \rightarrow s'; t', \pi^{\theta}) \pi^{\theta}(a'; s') r(s', a') \\ &= \sum_{t=0}^{\infty} \sum_{t'=0}^{\infty} \gamma^{t+t'} \sum_{s \in \mathcal{S}, a \in \mathcal{A}, s' \in \mathcal{S}, a' \in \mathcal{A}} p(s_0 \rightarrow s; t, \pi^{\theta}) \pi^{\theta}(a; s) p((s, a) \rightarrow s'; t', \pi^{\theta}) \pi^{\theta}(a'; s') \nabla \log(\pi^{\theta}(a; s)) r(s', a'). \end{aligned}$$

For notation simplicity we define

$$f_{t,t'}^{s_0,\theta}(s, a, s', a') := p(s_0 \rightarrow s; t, \pi^\theta) \pi^\theta(a; s) p((s, a) \rightarrow s'; t', \pi^\theta) \pi^\theta(a'; s')$$

and get

$$\nabla J_{s_0}(\theta) = \sum_{t=0}^{\infty} \sum_{t'=0}^{\infty} \gamma^{t+t'} \sum_{s \in \mathcal{S}, a \in \mathcal{A}, s' \in \mathcal{S}, a' \in \mathcal{A}} f_{t,t'}^{s_0,\theta}(s, a, s', a') \nabla \log(\pi^\theta(a; s)) r(s', a').$$

For any θ_1 and θ_2 we can analyse the difference of gradients using this representation as follows.

$$\begin{aligned} & \|\nabla J_{s_0}(\theta_1) - \nabla J_{s_0}(\theta_2)\| \\ &= \left\| \sum_{t=0}^{\infty} \sum_{t'=0}^{\infty} \gamma^{t+t'} \cdot \left\{ \sum_{s \in \mathcal{S}, a \in \mathcal{A}, s' \in \mathcal{S}, a' \in \mathcal{A}} f_{t,t'}^{s_0,\theta}(s, a, s', a') (\nabla \log(\pi^{\theta_1}(a; s)) - \nabla \log(\pi^{\theta_2}(a; s))) r(s', a') \right. \right. \\ & \quad \left. \left. + \sum_{s \in \mathcal{S}, a \in \mathcal{A}, s' \in \mathcal{S}, a' \in \mathcal{A}} (f_{t,t'}^{s_0,\theta_1}(s, a, s', a') - f_{t,t'}^{s_0,\theta_2}(s, a, s', a')) \nabla \log(\pi^{\theta_2}(a; s)) r(s', a') \right\} \right\| \\ &\leq \sum_{t=0}^{\infty} \sum_{t'=0}^{\infty} \gamma^{t+t'} \cdot \left\{ \sum_{s \in \mathcal{S}, a \in \mathcal{A}, s' \in \mathcal{S}, a' \in \mathcal{A}} f_{t,t'}^{s_0,\theta}(s, a, s', a') \|\nabla \log(\pi^{\theta_1}(a; s)) - \nabla \log(\pi^{\theta_2}(a; s))\| |r(s', a')| \right. \\ & \quad \left. + \sum_{s \in \mathcal{S}, a \in \mathcal{A}, s' \in \mathcal{S}, a' \in \mathcal{A}} |f_{t,t'}^{s_0,\theta_1}(s, a, s', a') - f_{t,t'}^{s_0,\theta_2}(s, a, s', a')| \|\nabla \log(\pi^{\theta_2}(a; s))\| |r(s', a')| \right\}. \end{aligned}$$

Using L_Θ -smoothness of the log-policy and the bounded reward assumption we have for the first term, that

$$\sum_{s \in \mathcal{S}, a \in \mathcal{A}, s' \in \mathcal{S}, a' \in \mathcal{A}} f_{t,t'}^{s_0,\theta}(s, a, s', a') \|\nabla \log(\pi^{\theta_1}(a; s)) - \nabla \log(\pi^{\theta_2}(a; s))\| |r(s', a')| \leq R_* L_\Theta \|\theta_1 - \theta_2\|. \quad (4.4)$$

For the second term we have to work a little harder to bound the difference of the probability densities $|f_{t,t'}^{s_0,\theta_1}(s, a, s', a') - f_{t,t'}^{s_0,\theta_2}(s, a, s', a')|$. Therefore, we first denote by \mathcal{T}_t all trajectories from 0 to t , i.e.

$$\mathcal{T}_t = \{\tau = \{s_0, a_0, \dots, s_t, a_t\} | s_0, a_0 \in \mathcal{A}_0, \dots, s_t \in \mathcal{S}, a_t \in \mathcal{A}_{s_t}\}$$

and rewrite $f_{t,t'}^{s_0,\theta}(s, a, s', a')$ as follows

$$\begin{aligned} f_{t,t'}^{s_0,\theta}(s, a, s', a') &= p(s_0 \rightarrow s; t, \pi^\theta) \pi^\theta(a; s) p((s, a) \rightarrow s'; t', \pi^\theta) \pi^\theta(a'; s') \\ &= \sum_{\tau \in \mathcal{T}_{t+t'}} \mathbf{1}_{s_t=s, a_t=a, s_{t'}=s', a_{t'}=a'} \prod_{n=0}^{t+t'} \pi^\theta(a_n; s_n) \prod_{n=0}^{t+t'-1} p(s_{n+1}; s_n, a_n). \end{aligned}$$

So,

$$\begin{aligned} & f_{t,t'}^{s_0,\theta_1}(s, a, s', a') - f_{t,t'}^{s_0,\theta_2}(s, a, s', a') \\ &= \sum_{\tau \in \mathcal{T}_{t+t'}} \mathbf{1}_{s_t=s, a_t=a, s_{t'}=s', a_{t'}=a'} \left(\prod_{n=0}^{t+t'} \pi^{\theta_1}(a_n; s_n) - \prod_{n=0}^{t+t'} \pi^{\theta_2}(a_n; s_n) \right) \prod_{n=0}^{t+t'-1} p(s_{n+1}; s_n, a_n). \end{aligned}$$

Then, using Taylor expansion (or also the more dimensional mean value or Lagrange theorem) of

$\theta \mapsto \prod_{n=0}^{t+t'} \pi^\theta(a_n; s_n)$ we obtain the existence of a $\tilde{\theta} \in [\theta_1, \theta_2]$ such that

$$\begin{aligned} \left| \prod_{n=0}^{t+t'} \pi^{\theta_1}(a_n; s_n) - \prod_{n=0}^{t+t'} \pi^{\theta_2}(a_n; s_n) \right| &\leq \left| (\theta_1 - \theta_2)^T \nabla_\theta \left(\prod_{n=0}^{t+t'} \pi^\theta(a_n; s_n) \right) \Big|_{\theta=\tilde{\theta}} \right| \\ &\leq \|\theta_1 - \theta_2\| \left\| \sum_{n=0}^{t+t'} \nabla \pi^{\tilde{\theta}}(a_n; s_n) \prod_{m=0, m \neq n}^{t+t'} \pi^{\tilde{\theta}}(a_m; s_m) \right\| \\ &\leq \|\theta_1 - \theta_2\| \sum_{n=0}^{t+t'} \left\| \nabla \log(\pi^{\tilde{\theta}}(a_n; s_n)) \right\| \prod_{m=0}^{t+t'} \pi^{\tilde{\theta}}(a_m; s_m) \\ &\leq \|\theta_1 - \theta_2\| (t + t' + 1) B_\Theta \prod_{m=0}^{t+t'} \pi^{\tilde{\theta}}(a_m; s_m). \end{aligned}$$

Finally, we follow for the second term, that

$$\begin{aligned} &\sum_{s \in \mathcal{S}, a \in \mathcal{A}, s' \in \mathcal{S}, a' \in \mathcal{A}} |f_{t,t'}^{s_0, \theta_1}(s, a, s', a') - f_{t,t'}^{s_0, \theta_2}(s, a, s', a')| \|\nabla \log(\pi^{\theta_2}(a; s))\| |r(s', a')| \\ &\leq \sum_{\tau \in \mathcal{T}_{t+t'}} \left| \prod_{n=0}^{t+t'} \pi^{\theta_1}(a_n; s_n) - \prod_{n=0}^{t+t'} \pi^{\theta_2}(a_n; s_n) \right| \prod_{n=0}^{t+t'-1} p(s_{n+1}; s_n, a_n) B_\Theta R_* \\ &\leq \|\theta_1 - \theta_2\| (t + t' + 1) B_\Theta^2 R_* \underbrace{\sum_{\tau \in \mathcal{T}_{t+t'}} \prod_{m=0}^{t+t'} \pi^\theta(a_m; s_m) \prod_{n=0}^{t+t'-1} p(s_{n+1}; s_n, a_n)}_{=1} \\ &= \|\theta_1 - \theta_2\| (t + t' + 1) B_\Theta^2 R_*. \end{aligned}$$

All in all, we obtain

$$\begin{aligned} \|\nabla J_{s_0}(\theta_1) - \nabla J_{s_0}(\theta_2)\| &\leq \sum_{t=0}^{\infty} \sum_{t'=0}^{\infty} \gamma^{t+t'} (R_* L_\Theta \|\theta_1 - \theta_2\| + \|\theta_1 - \theta_2\| (t + t' + 1) B_\Theta^2 R_*) \\ &\leq L \|\theta_1 - \theta_2\|, \end{aligned}$$

for $L = \frac{R_* L_\Theta}{(1-\gamma)^2} + \frac{(1+\gamma)R_* B_\Theta^2}{(1-\gamma)^3}$, where we used that

$$\sum_{t=0}^{\infty} \sum_{t'=0}^{\infty} \gamma^{t+t'} (t + t' + 1) = \frac{1 + \gamma}{(1 - \gamma)^3}.$$

□

Secondly, we show that the estimator used in Algorithm 18 has bounded variance.



Lemma 4.1.15. The estimator $\widehat{\nabla} J_s(\theta)$ proposed in Algorithm 18 for $K = 1$ has bounded variance. More precisely,

$$\|\widehat{\nabla} J_s(\theta)\| \leq \frac{B_\Theta R_*}{(1-\gamma)(1-\gamma^{1/2})}$$

and therefore,

$$\mathbb{E}[\|\widehat{\nabla} J_s(\theta) - \nabla J(\theta)\|^2] \leq C$$

$$\text{for } C = R_*^2 B_\Theta^2 \left(\frac{1}{(1-\gamma)^2 (1-\gamma^{1/2})^2} + \frac{2}{(1-\gamma)^3 (1-\gamma^{1/2})} + \frac{1}{(1-\gamma)^4} \right).$$

Proof. By the definition of $\widehat{\nabla} J_s(\theta)$ we have that

$$\begin{aligned} \|\widehat{\nabla} J_s(\theta)\| &= \left\| \frac{1}{1-\gamma} \nabla_{\theta} \log(\pi^{\theta}(A_T; S_T)) \sum_{t'=0}^{T'-1} \gamma^{t'/2} R(S_{T+t'}, A_{T+t'}) \right\| \\ &= \frac{1}{1-\gamma} \|\nabla_{\theta} \log(\pi^{\theta}(A_T; S_T))\| \sum_{t'=0}^{T'-1} \gamma^{t'/2} |R(S_{T+t'}, A_{T+t'})| \\ &\leq \frac{B_{\Theta} R_*}{(1-\gamma)} \sum_{t'=0}^{T'-1} \gamma^{t'/2} \\ &= \frac{B_{\Theta} R_*}{(1-\gamma)(1-\gamma^{1/2})}. \end{aligned}$$

For the second claim, first note that

$$\begin{aligned} \|\nabla J_s(\theta)\| &= \left\| \frac{1}{1-\gamma} \mathbb{E}_{S \sim d_s^{\pi^{\theta}}, A \sim \pi^{\theta}(\cdot; S)} [\nabla \log(\pi^{\theta}(A; S)) Q^{\pi^{\theta}}(S, A)] \right\| \\ &\leq \frac{1}{1-\gamma} \mathbb{E}_{S \sim d_s^{\pi^{\theta}}, A \sim \pi^{\theta}(\cdot; S)} [\|\log(\pi^{\theta}(A; S))\| Q^{\pi^{\theta}}(S, A)] \\ &\leq \frac{1}{1-\gamma} \mathbb{E}_{S \sim d_s^{\pi^{\theta}}, A \sim \pi^{\theta}(\cdot; S)} [B_{\Theta} \frac{R_*}{1-\gamma}] \\ &= \frac{B_{\Theta} R_*}{(1-\gamma)^2}. \end{aligned}$$

Hence, we have that

$$\begin{aligned} &\mathbb{E}[\|\widehat{\nabla} J_s(\theta) - \nabla J_s(\theta)\|^2] \\ &\leq \mathbb{E}[\|\widehat{\nabla} J_s(\theta)\|^2] + 2\mathbb{E}[\|\widehat{\nabla} J_s(\theta)\| \|\nabla J_s(\theta)\|] + \|\nabla J_s(\theta)\|^2 \\ &\leq \frac{R_*^2 B_{\Theta}^2}{(1-\gamma)^2 (1-\gamma^{1/2})^2} + 2 \frac{R_* B_{\Theta}}{(1-\gamma)(1-\gamma^{1/2})} \frac{R_* B_{\Theta}}{(1-\gamma)^2} + \frac{R_*^2 B_{\Theta}^2}{(1-\gamma)^4} \\ &= R_*^2 B_{\Theta}^2 \left(\frac{1}{(1-\gamma)^2 (1-\gamma^{1/2})^2} + \frac{2}{(1-\gamma)^3 (1-\gamma^{1/2})} + \frac{1}{(1-\gamma)^4} \right). \end{aligned}$$

Define $C = R_*^2 B_{\Theta}^2 \left(\frac{1}{(1-\gamma)^2 (1-\gamma^{1/2})^2} + \frac{2}{(1-\gamma)^3 (1-\gamma^{1/2})} + \frac{1}{(1-\gamma)^4} \right)$ proves the claim. \square



Theorem 4.1.16. Assume the objective function $\theta \mapsto J(\theta)$ is L -smooth and the policy parametrisation fulfills Assumption 4.1.12. Consider the stochastic process $(\theta_n)_{n \geq 0}$ generated in Algorithm 18 for $K = 1$, where the step-sizes $(\alpha_l)_{l \in \mathbb{N}}$ (deterministic or \mathcal{F} -adapted) satisfy

$$\alpha_l > 0, \quad \sum_{l=0}^{\infty} \alpha_l = \infty \quad \text{and} \quad \sum_{l=0}^{\infty} \alpha_l^2 < \infty$$

(almost surely). Then,

$$\lim_{l \rightarrow \infty} \|\nabla J_s(\theta_l)\|^2 = 0$$

almost surely.

Proof. First note that we consider a maximisation problem instead of a minimisation problem and therefore have to consider $-J$ in the SGD Theorem 3.5.1. By the bounded reward assumption we know that $J_{s,*} = \sup_{\theta \in \mathbb{R}^d} J_s(\theta) \leq \frac{R_*}{1-\gamma} < \infty$, i.e. $\inf_{\theta} -J_s(\theta) > -\infty$. Furthermore, we know that

J_s (and thus also $-J_s$) is L -smooth. We have by $\widehat{\nabla} J_s(\theta)$ from Algorithm 18 an estimator which fulfills the equations (3.10) and (3.9) shown in Proposition 4.1.13 and Lemma 4.1.15. Hence, defining

- $F = -J_s$,
- $x = \theta$ and Z takes the roll of the MDP rollouts up to the geometric length ℓ , i.e. $Z(x) = (T, T', S_T, A_T, S_{T+1}, \dots, S_{T+T'}, A_{T+T'})$ under policy π^θ ,
- $\nabla_x f(x, Z) = \frac{1}{1-\gamma} \nabla_\theta \log(\pi^\theta(S_T; A_T)) \sum_{t=T}^{T+T'} \gamma^{(t-T)/2} R(S_t, A_t)$,

leads to almost sure convergence:

$$\lim_{l \rightarrow \infty} \|\nabla J_s(\theta_l)\|^2 = 0.$$

□

4.1.4 Baseline

Traditionally it is believed that reducing the variance in estimating the gradient $\nabla J(\theta)$ is crucial. Recalling the discussion from Section 1.3.4 there are also other effects when using the variance reduction by baselines for the softmax policy gradient approach to stochastic bandits. Nonetheless, reducing variances is certainly a good idea.

Baseline

Let us generalise the baseline idea from stochastic bandits to general MDP policy gradient, formulated in the finite MDP setting.



Theorem 4.1.17. Assume that (S, A, R) is a T -step MDP with finite state-action spaces and consider a stationary differentiable parametrized family of policies $\{\pi^\theta : \theta \in \Theta\}$. For any $b \in \mathbb{R}$ the gradient of $J(\theta)$ can also be written as

$$\nabla_\theta J(\theta) = \mathbb{E}_\mu^{\pi^\theta} \left[\sum_{t=0}^{T-1} \nabla_\theta (\log \pi^\theta(A_t; S_t)) \underbrace{(Q_t^{\pi^\theta}(S_t, A_t) - b)}_{\text{or } (R_t^T - b) \text{ or } (R_0^T - b)} \right].$$

Proof. The computation is very similar to the bandit case, we redo the log-trick and see that for every $t \leq T-1$

$$\begin{aligned} \mathbb{E}_\mu^{\pi^\theta} [\nabla_\theta (\log \pi^\theta(A_t; S_t)) b] &= b \sum_{s_t \in \mathcal{S}} \sum_{a_t \in \mathcal{A}_s} \mathbb{P}_\mu^{\pi^\theta}(S_t = s_t) \pi^\theta(a_t; s_t) \nabla_\theta (\log \pi^\theta(a_t; s_t)) \\ &= b \sum_{s_t \in \mathcal{S}} \mathbb{P}_\mu^{\pi^\theta}(S_t = s_t) \sum_{a_t \in \mathcal{A}_s} \nabla_\theta \pi^\theta(a_t; s_t) \\ &= b \sum_{s_t \in \mathcal{S}} \mathbb{P}_\mu^{\pi^\theta}(S_t = s_t) \nabla_\theta \underbrace{\sum_{a_t \in \mathcal{A}} \pi^\theta(a_t; s_t)}_{=1} = 0. \end{aligned}$$

Thus, the additional term vanishes in all three representations of the gradient so that $-b$ can be added. □



Show that the constant baseline b can be replaced by any deterministic state-dependent baseline $b : \mathcal{S} \rightarrow \mathbb{R}$, i.e.

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\mu}^{\pi^{\theta}} \left[\sum_{t=0}^{T-1} \nabla_{\theta} (\log \pi^{\theta}(A_t; S_t)) (Q_t^{\pi^{\theta}}(S_t, A_t) - b(S_t)) \right].$$

We will come back to the most important state-dependent baseline $b(s) = V^{\pi^{\theta}}(s)$ when we discuss actor-critic methods.



Write down and prove the baseline gradient representation for infinite discounted MDPs.

In the bandit case we have already seen that the baseline trick is a variance reduction trick. The same is true in the MDP setting.

Importance sampling trick/likelihood ratio method

There is a nice trick from stochastic numerics that is useful in many situations for reinforcement learning, the so-called importance sampling trick (or likelihood ratio method). Suppose $m = \mathbb{E}[g(X)]$ is to be approximated using a Monte Carlo method. The likelihood ratio trick is used in two situations:

- If it is hard to sample X but easy to simulate a random variable Y with similar density.
- The variance $g(X)$ should be reduced.

The first advantage of the method is what is closer to the interest of reinforcement learning. Before showing why let us recall the trick. Suppose X has density f_X and Y is another random variable with (positive) density f_Y . Then

$$\mathbb{E}[g(X)] = \int g(x) f_X(x) dx = \int g(y) \frac{f_X(y)}{f_Y(y)} f_Y(y) dy = \mathbb{E}[\tilde{g}(Y)]$$

with $\tilde{g}(y) = g(y) \frac{f_X(y)}{f_Y(y)}$. Suppose that f_X is close to a known density f_Y but somehow wiggles around. It might then be much easier to simulate Y instead and correct the Monte Carlo estimator with the likelihood ratios. Alternatively, if $\tilde{g}(Y)$ has smaller variance than $g(X)$, then the Monte Carlo convergence using Y would converge faster. This is typically the case if Y has more mass on the important values (justifying the name importance sampling) giving the most contribution to the expectation.

We now turn towards reinforcement learning and, quite surprisingly, find the trick useful to construct off-policy estimators for the policy gradients! These are estimators that do not use samples produced from π^{θ} . Here is the main trick:



Suppose $\tau = (s_0, a_0, \dots, s_{T-1}, a_{T-1}, s_T)$ is a policy and π, π^b are two policies. Using the path probabilities gives

$$\frac{\mathbb{P}_{\mu}^{\pi}(\tau)}{\mathbb{P}_{\mu}^{\pi^b}(\tau)} = \frac{\mu(s_0) \delta_0(r_0) \prod_{i=0}^{T-1} \pi(a_i; s_i) p(s_{i+1}; s_i, a_i)}{\mu(s_0) \delta_0(r_0) \prod_{i=0}^{T-1} \pi^b(a_i; s_i) p(s_{i+1}; s_i, a_i)} = \prod_{i=0}^{T-1} \frac{\pi(a_i; s_i)}{\pi^b(a_i; s_i)}.$$

What is crucial here is the cancellation of the transition p , the change from π to π^b is model-free!

The likelihood ratio trick now does the following. Write down the policy gradient estimator, replace π^{θ} by some behavior policy π^b and compensate by inserting the likelihood ratio. Here is the version of the finite-time non-discounted case:


Proposition 4.1.18. (Off-policy actor-critic policy gradient theorem)

Suppose π^b is a behaviour policy (with strictly positive weights), then

$$\nabla_{\theta} J_{\mu}(\theta) = \mathbb{E}_{\mu}^{\pi^b} \left[\sum_{t=0}^{T-1} \gamma^t \prod_{i=0}^t \frac{\pi^{\theta}(A_i; S_i)}{\pi^b(A_i; S_i)} \nabla_{\theta} (\log \pi^{\theta}(A_t; S_t)) A_t^{\pi^{\theta}}(S_t, A_t) \right],$$

where $A_t^{\pi^{\theta}}(s, a) = Q_t^{\pi^{\theta}}(s, a) - V_t^{\pi^{\theta}}(s)$ is the advantage function.

Proof. The claim follows from the definition of discrete expectations $\mathbb{E}[g(X)] = \sum_k g(k)\mathbb{P}(X = k)$ and the MDP likelihood trick:

$$\begin{aligned} \nabla J(\theta) &= \sum_{t=0}^{T-1} \gamma^t \mathbb{E}_{\mu}^{\pi^{\theta}} \left[\nabla_{\theta} (\log \pi^{\theta}(A_t; S_t)) Q_t^{\pi^{\theta}}(S_t, A_t) \right] \\ &= \sum_{t=0}^{T-1} \gamma^t \sum_{\tau: \text{trajectory}} \left(\nabla_{\theta} (\log \pi^{\theta}(a_t; s_t)) Q_t^{\pi^{\theta}}(s_t, a_t) \right) \mathbb{P}_{\mu}^{\pi^{\theta}}(\tau) \\ &= \sum_{t=0}^{T-1} \gamma^t \sum_{\tau: \text{trajectory}} \left(\frac{\mathbb{P}_{\mu}^{\pi^{\theta}}(\tau)}{\mathbb{P}_{\mu}^{\pi^b}(\tau)} \nabla_{\theta} (\log \pi^{\theta}(a_t; s_t)) Q_t^{\pi^{\theta}}(s_t, a_t) \right) \mathbb{P}_{\mu}^{\pi^b}(\tau) \\ &= \sum_{t=0}^{T-1} \gamma^t \sum_{\tau: \text{trajectory}} \left(\prod_{i=0}^t \frac{\pi^{\theta}(a_i; s_i)}{\pi^b(a_i; s_i)} \nabla_{\theta} (\log \pi^{\theta}(a_t; s_t)) Q_t^{\pi^{\theta}}(s_t, a_t) \right) \mathbb{P}_{\mu}^{\pi^b}(\tau) \\ &= \sum_{t=0}^{T-1} \gamma^t \mathbb{E}_{\mu}^{\pi^b} \left[\prod_{i=0}^t \frac{\pi^{\theta}(A_i; S_i)}{\pi^b(A_i; S_i)} \nabla_{\theta} (\log \pi^{\theta}(A_t; S_t)) Q_t^{\pi^{\theta}}(S_t, A_t) \right] \end{aligned}$$

Note that the product only runs up to t , all other factors sum out when computing the probabilities. Finally, we use the value function as a state-dependent baseline and we obtain the off-policy actor-critic representation. \square

There are several reasons why the trick can be useful. It can be used to reduce the variance of policy gradients (the original idea of the trick). Somehow this seems not to be very successful. Alternatively, the trick can be used to simulate from only one policy to approximate all policy gradients $\nabla_{\theta} J_{\mu}(\theta)$, i.e. run an off-policy policy gradient ascent, see for instance⁶. This approach turns out to be delicate as small values of π^b can strongly increase the variance. From the theoretical point of view the likelihood trick is very useful as the REINFORCE algorithm turns into a true stochastic gradient descent algorithm of a function $f(\theta) = \mathbb{E}[h(X, \theta)]$, the dependence of θ is removed from the sampled law.

4.2 PPO

We start with REINFORCE with cyclic importance sampling updates in order to reuse sample rollouts for a view gradient steps. Since the variance is big we use a clever gradient estimator, that has a bias but a smaller variance. This is what PPO does. There is only one difference: We do not clip importance sampling weights in our analysis, we do early stopping. The idea of the early stopping is to stop the importance sampling updates if the variance of the gradient estimator gets too big.

4.2.1 Bias of the PPO surrogate gradient

We have a fixed behaviour policy π_{old} from which we can sample and want to use this to estimate the gradient $\nabla_{\theta} J(\pi_{\theta})$.

⁶Metelli, Papini, Faccio, Restelli: "Policy Optimization via Importance Sampling", NIPS 2018


Theorem 4.2.1. (PPO surrogate bias)

Suppose π_{old} is a behaviour policy (with strictly positive weights) and define the surrogate gradient by

$$\widehat{\nabla_{\theta} J}(\theta) = \mathbb{E}^{\pi_{\text{old}}} \left[\sum_{t=0}^{T-1} \gamma^t \frac{\nabla_{\theta} \pi^{\theta}(A_t; S_t)}{\pi_{\text{old}}(A_t; S_t)} A_t^{\pi^{\theta}}(S_t, A_t) \right]$$

and additionally assume that $\|A_t^{\pi_{\text{old}}}(s, \cdot)\|_{\infty} \leq A_{\text{max}}$. It then holds that

$$\begin{aligned} & \|\nabla_{\theta} J(\pi_{\theta}) - \widehat{\nabla_{\theta} J}(\pi_{\theta})\|_{\infty} \\ & \leq 2\sqrt{2}A_{\text{max}} \sqrt{\max_{s \in \mathcal{S}} \text{KL}(\pi_{\text{old}}(\cdot; s) \parallel \pi_{\theta}(\cdot; s))} \|\nabla_{\theta} \log \pi_{\theta}\|_{\infty} \sum_{t=0}^{T-1} \gamma^t t. \end{aligned}$$

For the proof we need a couple of preparations.

Let us start with a well-known identity, also known as the performance difference lemma (see for instance [?, Lemma 6.1] and [?, Eqn. (1)] for the infinite time-horizon setting).


Lemma 4.2.2. (Performance difference lemma)

For two policies $\pi, \tilde{\pi}$,

$$J(\tilde{\pi}) - J(\pi) = \mathbb{E}^{\tilde{\pi}} \left[\sum_{t=0}^{T-1} \gamma^t A_t^{\tilde{\pi}}(S_t, A_t) \right] = \sum_{t=0}^{T-1} \gamma^t \mathbb{E}^{\tilde{\pi}} [A_t^{\tilde{\pi}}(S_t, A_t)] =: \sum_{t=0}^{T-1} \gamma^t \Delta_t. \quad (4.5)$$

Proof. We have $Q_t^{\pi}(s, a) = \mathbb{E}_{S_t=s, A_t=a}^{\pi} [R_t + \gamma V_{t+1}^{\pi}(S_{t+1})]$ and $\mathbb{E}_{S \sim d_t^{\pi}, A \sim \pi(\cdot; s)} [A_t^{\pi}(S, A)] = 0$. Thus,

We bound the KL-divergence between the state marginals for different policies in terms of the KL-divergence between the policies, which we can actually control. The following lemma is essentially an application of the formula $\text{KL}(\otimes_i \mathbb{P}_i \parallel \otimes_i \mathbb{Q}_i) = \sum_i \text{KL}(\mathbb{P}_i \parallel \mathbb{Q}_i)$ for product measures [?, p. 85]:



Lemma 4.2.3. Let \mathbb{P}_t^{π} and $\mathbb{P}_t^{\tilde{\pi}}$ be the distributions of $(S_0, A_0, R_0, \dots, S_t)$ under π and $\tilde{\pi}$, respectively. Then

$$\text{KL}(\mathbb{P}_t^{\tilde{\pi}} \parallel \mathbb{P}_t^{\pi}) = \sum_{i=0}^{t-1} \mathbb{E}^{\tilde{\pi}} [\text{KL}(\tilde{\pi}(\cdot; S_i) \parallel \pi(\cdot; S_i))].$$

Proof. By factorising \mathbb{P}_t^{π} and $\mathbb{P}_t^{\tilde{\pi}}$ using the MDP dynamics, we find

$$\log \frac{\mathbb{P}_t^{\tilde{\pi}}(s_0, a_0, r_0, \dots, s_t)}{\mathbb{P}_t^{\pi}(s_0, a_0, r_0, \dots, s_t)} = \sum_{i=0}^{t-1} \log \frac{\tilde{\pi}(a_i; s_i)}{\pi(a_i; s_i)}$$

for any $(s_0, a_0, r_0, \dots, s_t)$. Using the definition of KL and Fubini gives

$$\begin{aligned} \text{KL}(\mathbb{P}_t^{\tilde{\pi}} \parallel \mathbb{P}_t^{\pi}) &= \mathbb{E}^{\tilde{\pi}} \left[\sum_{i=0}^{t-1} \log \frac{\tilde{\pi}(A_i; S_i)}{\pi(A_i; S_i)} \right] = \sum_{i=0}^{t-1} \mathbb{E}_{S \sim d_i^{\tilde{\pi}}} \left[\mathbb{E}_{A \sim \tilde{\pi}(\cdot; S)} \left[\log \frac{\tilde{\pi}(A; S)}{\pi(A; S)} \right] \right] \\ &= \sum_{i=0}^{t-1} \mathbb{E}_{S \sim d_i^{\tilde{\pi}}} [\text{KL}(\tilde{\pi}(\cdot; S) \parallel \pi(\cdot; S))]. \end{aligned}$$

□



Lemma 4.2.4. For two policies $\pi, \tilde{\pi}$ with $\max_s \text{KL}(\tilde{\pi}(\cdot; s) \parallel \pi(\cdot; s)) \leq \varepsilon$, we have

$$\text{KL}(d_t^{\tilde{\pi}} \parallel d_t^{\pi}) \leq t\varepsilon,$$

where $d_t^{\pi} = \mathbb{P}^{\pi}(S_t \in \cdot)$ and $d_t^{\tilde{\pi}} = \mathbb{P}^{\tilde{\pi}}(S_t \in \cdot)$.

Proof. Because d_t^{π} and $d_t^{\tilde{\pi}}$ are marginal distributions of \mathbb{P}_t^{π} and $\mathbb{P}_t^{\tilde{\pi}}$, respectively, the chain rule for the KL-divergence⁷ implies

$$\text{KL}(d_t^{\tilde{\pi}} \parallel d_t^{\pi}) \leq \text{KL}(\mathbb{P}_t^{\tilde{\pi}} \parallel \mathbb{P}_t^{\pi}) = \sum_{i=0}^{t-1} \mathbb{E}^{\tilde{\pi}}[\text{KL}(\tilde{\pi}(\cdot; S_i) \parallel \pi(\cdot; S_i))] \leq t\varepsilon.$$

□

We next give an adaptation and improvement of Lemma 2 from the TRPO paper for the finite-time horizon setting.



Lemma 4.2.5. If $\max_s \text{KL}(\pi_{\text{old}}(\cdot; s) \parallel \pi_{\theta}(\cdot; s)) \leq \varepsilon$ and $\|A_t^{\pi_{\text{old}}}(s, \cdot)\|_{\infty} \leq A_{\text{max}}$, then

$$\|g_t^{\theta}\|_{\infty} \leq \sqrt{2\varepsilon}A_{\text{max}},$$

where $g_t^{\theta}(s) := \mathbb{E}_{A \sim \pi_{\theta}(\cdot; s)}[A_t^{\pi_{\text{old}}}(s, A)]$, $s \in \mathcal{S}$.

Proof. Because $\mathbb{E}_{A \sim \pi_{\text{old}}(\cdot; s)}[A_t^{\pi_{\text{old}}}(s, A)] = 0$, a characterisation of the total variation distance⁸ and Pinsker's inequality⁹ imply

$$\begin{aligned} \|g_t^{\theta}(s)\|_{\infty} &= \|\mathbb{E}_{A \sim \pi_{\theta}(\cdot; s)}[A_t^{\pi_{\text{old}}}(s, A)] - \mathbb{E}_{A \sim \pi_{\text{old}}(\cdot; s)}[A_t^{\pi_{\text{old}}}(s, A)]\|_{\infty} \\ &\leq 2\text{TV}(\pi_{\theta}(\cdot; s), \pi_{\text{old}}(\cdot; s))A_{\text{max}} \\ &\leq \sqrt{2\varepsilon}A_{\text{max}}. \end{aligned}$$

□

We can now give the proof of Theorem 5.3.1.

Proof of Theorem 5.3.1: First note that Fubini allows us to rewrite

$$\widehat{\nabla_{\theta} J} = \sum_{t=0}^{T-1} \gamma^t \nabla_{\theta} \hat{\Delta}_t$$

with $\hat{\Delta}_t := \mathbb{E}^{\pi_{\text{old}}}[g_t^{\theta}(S_t)]$, $\Delta_t := \mathbb{E}^{\pi_{\theta}}[g_t^{\theta}(S_t)]$, and g_t^{θ} as defined above. For a single coordinate θ_j

⁷For instance Theorem 2.5.3 of T. Cover "Elements of information theory"

⁸see Prop. 4.5 of Levin, Peres "Mixing times of Markov chains"

⁹see for instance A. Tsybakov, "Introduction to Nonparametric Estimation"

of θ we obtain

$$\begin{aligned}
\partial_{\theta_j}(\Delta_t - \hat{\Delta}_t) &= \partial_{\theta_j} \sum_{s \in \mathcal{S}} (d_t^{\pi_\theta}(s) - d_t^{\pi_{\text{old}}}(s)) g_t^\theta(s) \\
&= \sum_{s \in \mathcal{S}} \sum_{\tau: s_t=s} \mathbb{P}_t^{\pi_\theta}(\tau) \sum_{i=0}^{t-1} \partial_{\theta_j} \log \pi_\theta(a_i; s_i) g_t^\theta(s) + \sum_{s \in \mathcal{S}} (d_t^{\pi_\theta}(s) - d_t^{\pi_{\text{old}}}(s)) \partial_{\theta_j} g_t^\theta(s) \\
&\leq \|g_t^\theta\|_\infty \mathbb{E}_{\tau \sim \mathbb{P}_t^{\pi_\theta}} \left[\sum_{i=0}^{t-1} \partial_{\theta_j} \log \pi_\theta(A_i; S_i) \right] + \sum_{s \in \mathcal{S}} (d_t^{\pi_\theta}(s) - d_t^{\pi_{\text{old}}}(s)) \underbrace{\sum_{a \in \mathcal{A}} \partial_{\theta_j} \pi_\theta(a; s) A_t^{\pi_{\text{old}}}(s, a)}_{= \mathbb{E}_{A \sim \pi_\theta(\cdot; s)} [\partial_{\theta_j} \log \pi_\theta(A; s) A_t^{\pi_{\text{old}}}(s, A)]} \\
&=: \|g_t^\theta\|_\infty \mathbb{E}_{\tau \sim \mathbb{P}_t^{\pi_\theta}} \left[\sum_{i=0}^{t-1} \partial_{\theta_j} \log \pi_\theta(A_i; S_i) \right] + (\mathbb{E}_{S \sim d_t^{\pi_\theta}} [f(S)] - \mathbb{E}_{S \sim d_t^{\pi_{\text{old}}}} [f(S)]).
\end{aligned}$$

For the second equality we used

$$\partial_{\theta_j} d_t^{\pi_\theta}(s) = \sum_{\tau: s_t=s} \partial_{\theta_j} \mathbb{P}_t^{\pi_\theta}(\tau) = \sum_{\tau: s_t=s} \mathbb{P}_t^{\pi_\theta}(\tau) \partial_{\theta_j} \log \mathbb{P}_t^{\pi_\theta}(\tau) = \sum_{\tau: s_t=s} \mathbb{P}_t^{\pi_\theta}(\tau) \sum_{i=0}^{t-1} \partial_{\theta_j} \log \pi_\theta(a_i; s_i).$$

The second summand is now estimated as

$$\begin{aligned}
|\mathbb{E}_{S \sim d_t^{\pi_\theta}} [f(S)] - \mathbb{E}_{S \sim d_t^{\pi_{\text{old}}}} [f(S)]| &\stackrel{\text{TV property}}{\leq} 2\|f\|_\infty TV(d_t^{\pi_\theta}, d_t^{\pi_{\text{old}}}) \\
&\stackrel{\text{Pinsker}}{\leq} \sqrt{2}\|f\|_\infty \sqrt{\text{KL}(d_t^{\pi_\theta} \| d_t^{\pi_{\text{old}}})} \\
&\stackrel{\text{Ass.}}{\leq} A_{\max} \|\nabla_\theta \log \pi_\theta\|_\infty \sqrt{2t\epsilon},
\end{aligned}$$

while the first summand more directly from 5.3.5. Combining both yields the upper bound

$$\leq \sqrt{2\epsilon} A_{\max} t \|\nabla_\theta \log \pi_\theta\|_\infty + A_{\max} \|\nabla_\theta \log \pi_\theta\|_\infty \sqrt{t} \sqrt{2\epsilon}.$$

The result follows by summing over t . □

Chapter 5

Proximal policy optimization (PPO)

In this chapter we describe one of most important RL algorithms, proximal policy optimization (PPO). We take an unusual point of view and derive the algorithms as a more sample efficient variant of the REINFORCE algorithm. We derive (and analyse) PPO in four steps.

- We look at REINFORCE with cyclic off-policy sampling. This is supposed to improve sample efficiency, using sampled rollouts for multiple gradient steps.
- We replace the reward to go with the advantage, reducing the variance a lot. In particular, we discuss generalized advantage estimation (GAE) for finite-time MDPs.
- To reduce the variance of the gradient estimator we introduce a bias to the gradient that allows to sample from transitions instead of rollouts. We analyse the variance and the bias of the surrogate gradient. Essentially, this is early stopping PPO.
- **Simon/Sebastian:** Finally, we develop a convergence theory for early stopping PPO. For that sake we use the smoothness of the value function (under certain policy assumptions) and throw bias/variance estimates into standard machinery. Suitable early stopping and step-size assumptions yield convergence to stationary points of the true value function.
- **Simon/Sebastian** We improve the analysis to a local convergence, and under stronger assumptions on model and policy also to convergence to an optimal policy.

Before we start, let us recall standard REINFORCE, which is based on standard policy gradient. Note: All we do is for finite-time MDPs, this is what is implemented anyway. If J denotes the parametrised value function (start distribution is fixed), then

$$\nabla_{\theta} J_s(\theta) = \mathbb{E}^{\pi^{\theta}} \left[\sum_{t=0}^{T-1} \nabla_{\theta} (\log \pi^{\theta}(A_t; S_t)) R_t^T \right].$$

Thus, an SGD scheme uses sampled entire rollouts to estimate the gradient. Keeping in mind

$$\text{Cov} \left(\sum_{t=0}^{T-1} X_t \right) = \sum_{t=0}^{T-1} \text{Cov}(X_t) + \sum_{t \neq t'} \text{Cov}(X_t, X_{t'})$$

it becomes clear that the rollout estimator has large variance. If a reward to go R_t^T is large, then also $R_{t'}^T$ is large for $t \approx t'$. In what follows, we will see how to reduce the variances a lot.

5.1 REINFORCE with cyclic off-policy updates

Idea: Sample trajectories and use them for a few gradient steps (this is a cycle), using the importance sampling correction. Then sample again. For this approach we use the importance sampling variant of the gradient:

$$\nabla_{\theta} J(\theta) = \mathbb{E}^{\pi^{\theta_{\text{old}}}} \left[\sum_{t=0}^{T-1} \gamma^t \prod_{i=0}^t \frac{\pi^{\theta}(A_i; S_i)}{\pi^{\theta_{\text{old}}}(A_i; S_i)} \nabla_{\theta} (\log \pi^{\theta}(A_t; S_t)) R_t^T \right],$$

Here is pseudocode for the resulting algorithm. Note: This is how the cycling is typically done:

- At the beginning of a cycle sample a number of trajectories \mathcal{D} . Then split \mathcal{D} into disjoint K minibatches.
- Do K gradient updates, using the K minibatches. This is called an epoch.
- Now perform several epochs (aka pass through the data) on the same batch of rollouts. Shuffle the full batch, form new minibatches and start a new epoch.

It might well be that in this simple situation it would be optimal to run full-batch updates only.

Require: number of cycles C ; episodes per cycle N ; number of passes per cycle E ; number of minibatches per pass K ; learning rate α

- 1: Initialise π_{old}
- 2: **for** cycle = 1, 2, ..., C **do**
- 3: Set sampling (behavior) policy $\mu = \pi_{\theta_{\text{old}}}$
- 4: Initialize empty dataset $\mathcal{D} = \emptyset$
 {— Data Collection Phase —}
- 5: **for** $n = 1, \dots, N$ **do**
- 6: Initialize episode: observe initial state $s_0 \sim p(s_0)$
- 7: Rollout trajectory $\tau_n = (s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_T)$ under μ :

$$a_t \sim \mu(\cdot | s_t), \quad s_{t+1} \sim p(\cdot | s_t, a_t)$$

- 8: Compute return to go for each timestep:

$$R_t^T = \sum_{t'=t}^T \gamma^{t'-t} r_{t'}$$

- 9: Store full trajectory τ_n in dataset \mathcal{D}
- 10: **end for**
 {— Optimization Phase —}
- 11: **for** epoch = 1, ..., E **do**
- 12: Randomly shuffle \mathcal{D}
- 13: Split \mathcal{D} into K disjoint minibatches: $\{\mathcal{B}_1, \dots, \mathcal{B}_K\}$
- 14: **for** $k = 1, \dots, K$ **do**
- 15: Estimate gradient:

$$\hat{g}_k = \frac{1}{|\mathcal{B}_k|} \sum_{\tau \in \mathcal{B}_k} \left[\sum_{t=0}^T w_t(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) G_t \right]$$

where importance ratio:

$$w_t(\tau) = \prod_{i=0}^t \frac{\pi_{\theta}(a_i | s_i)}{\mu(a_i | s_i)}$$

- 16: Update policy parameters:

$$\theta \leftarrow \theta + \alpha \hat{g}_k$$

- 17: **end for**
- 18: **end for**
- 19: Set $\theta_{\text{old}} \leftarrow \theta$ for next cycle
- 20: **end for**

Algorithm 19: Off-Policy REINFORCE with multi-batching for MDPs with finite time-horizon (Rollout-Based)

No theory provided here. The algorithm is useless anyway.

5.2 Generalized advantage estimation, actor-critic

The first step to reduce the variance of the gradient estimator is to take conditional expectations insight (reward to go only depends on the future transitions) and subtract a baseline. If we take

the current value function as a base-line, this is called advantage:

$$\begin{aligned}
 \nabla_{\theta} J_s(\theta) &= \mathbb{E}^{\pi^{\theta}} \left[\sum_{t=0}^{T-1} \nabla_{\theta} (\log \pi^{\theta}(A_t; S_t)) R_t^T \right] \\
 &= \mathbb{E}^{\pi^{\theta}} \left[\sum_{t=0}^{T-1} \nabla_{\theta} (\log \pi^{\theta}(A_t; S_t)) Q_t^T(S_t, A_t) \right] \\
 &= \mathbb{E}^{\pi^{\theta}} \left[\sum_{t=0}^{T-1} \nabla_{\theta} (\log \pi^{\theta}(A_t; S_t)) \underbrace{(Q_t^{\pi^{\theta}}(S_t, A_t) - V_t^{\pi^{\theta}}(S_t))}_{=: A_t^{\pi^{\theta}}(S_t, A_t)} \right]
 \end{aligned}$$

We are now entering the world of actor-critic policy gradient schemes. In order to improve the policy (the actor) we need to evaluate the current policy (done by the critic). Why? Because the gradient now involves the current value function.

Problem: It is generally not a good idea to estimate Q AND V separately. Thus, we typically use that

$$A^{\pi}(s, a) = \mathbb{E}_{s,a}^{\pi} \left[\sum_{t=0}^{n-1} \gamma^t R_t + \gamma^n V^{\pi}(S_n) \right] - V^{\pi}(s) \approx \sum_{t=0}^{n-1} \gamma^t R_t + \gamma^n V^{\pi}(S_n) - V^{\pi}(s).$$

The equality is due to the n -step Bellman-operator (Markov property after n -steps, instead of 1 step as usually). The A2C algorithm uses this simple advantage estimator using the current path as

$$A^{\pi_{\theta}}(S_t, A_t) = \sum_{t'=0}^{n-1} \gamma^{t'} R_{t+t'} + \gamma^n V^{\pi_{\theta}}(S_{t+n}) - V^{\pi_{\theta}}(S_t),$$

where $V^{\pi_{\theta}}$ is estimated in the critic step from the current data. Here is pseudo-code. Note that A2C is fully on-policy, there are not off-policy cycles. The pseudocode below is not exactly A2C, but it uses the simple A2C advantage estimator.

Require: Policy (actor) π_θ ; value function (critic) V_ϕ ; number of cycles C ; episodes per cycle N ; number of passes per cycle E ; number of minibatches per pass K ; learning rates α_π, α_V

- 1: Initialize actor parameters θ_0 and critic parameters ϕ_0
- 2: **for** cycle = 1, 2, ..., C **do**
- 3: Set sampling (behavior) policy $\mu = \pi_{\theta_{\text{old}}}$
- 4: Initialize empty dataset $\mathcal{D} = \emptyset$
- 5: — Collection of rollouts from current policy —
- 6: **for** $n = 1, \dots, N$ **do**
- 7: Initialize episode: observe initial state $s_0 \sim p(s_0)$
- 8: Rollout trajectory $\tau_n = (s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_{T_\tau})$ under μ :

$$a_t \sim \mu(\cdot|s_t), \quad s_{t+1} \sim p(\cdot|s_t, a_t)$$

- 9: Store (s_t, a_t, r_t, s_{t+1}) for all $t = 0, \dots, T_\tau$ in \mathcal{D}
- 10: **end for**
- 11: — Estimation of the value function, update of previous network —
- 12: **for** epoch = 1, ..., E_V **do**
- 13: Randomly shuffle \mathcal{D} , split \mathcal{D} into K critic minibatches: $\{\mathcal{B}_1^V, \dots, \mathcal{B}_K^V\}$
- 14: **for** $k = 1, \dots, K$ **do**
- 15: Compute bootstrapped targets (e.g., TD(0)):

$$y_t = r_t + \gamma V_\phi(s_{t+1})$$

- 16: Update critic (value function estimator) by minimizing squared error by mixing over all transitions:

$$\phi \leftarrow \phi - \alpha_V \nabla_\phi \frac{1}{|\mathcal{B}_k^V|} \sum_{(s_t, a_t, r_t, s_{t+1}) \in \mathcal{B}_k^V} (V_\phi(s_t) - y_t)^2$$

- 17: **end for**
- 18: **end for**
- 19: — Computation of all advantages estimators $A_t(s_t, a_t)$ —
- 20: **for** each trajectory $\tau = (s_0, a_0, r_0, \dots, s_{T_\tau})$ in \mathcal{D} **do**
- 21: **for** $t = 0, \dots, T_\tau$ **do**
- 22: Compute some advantage estimator, for instance

$$\hat{A}_t(s_t, a_t) = r_t + \gamma V_\phi(s_{t+1}) \mathbf{1}_{s_{t+1} \neq \text{terminal}} - V_\phi(s_t)$$

- 23: **end for**
- 24: **end for**
- 25: — Computation of the gradients and update of policy network —
- 26: **for** epoch = 1, ..., E **do**
- 27: Randomly shuffle \mathcal{D} at the trajectory level
- 28: Split \mathcal{D} into K actor minibatches: $\{\mathcal{B}_1^A, \dots, \mathcal{B}_K^A\}$
- 29: **for** $k = 1, \dots, K$ **do**
- 30: Estimate policy gradient per minibatch:

$$\hat{g}_k = \frac{1}{|\mathcal{B}_k^A|} \sum_{\tau \in \mathcal{B}_k^A} \sum_{t=0}^{T_\tau} w_t(\tau) \nabla_\theta \log \pi_\theta(a_t|s_t) \hat{A}_t(s_t, a_t)$$

$$\text{where } w_t(\tau) = \prod_{i=0}^t \frac{\pi_\theta(a_i|s_i)}{\mu(a_i|s_i)}$$

- 31: Update actor parameters:

$$\theta \leftarrow \theta + \alpha_\pi \hat{g}_k$$

- 32: **end for**
- 33: **end for**
- 34: Set $\theta_{\text{old}} \leftarrow \theta$ for next cycle
- 35: **end for**

Algorithm 20: REINFORCE with cyclic updates, using A2C critic estimation with $n = 1$

5.2.1 Infinite horizon GAE

As noted before, $Q^\pi(s, a) = \mathbb{E}_{s,a}^\pi[\sum_{t=0}^{n-1} \gamma^t R_t + \gamma^n V^\pi(S_n)]$ allows us to work with different advantage estimators. Here is an important point. Small n uses stronger bootstrapping of already learned values, large n mostly uses new samples.

- small n has less variance, but larger V -estimation bias
- large n has more variance, but less V -estimation bias

What is GAE? GAE is just a mixture of the n -step advantage estimators. Let's write this up:

To construct estimators of the advantage function in the actor-critic framework, we rely on the notion of temporal-difference (TD) errors. Given a value function estimate \hat{V} , the one-step TD-error at time t is defined as

$$\hat{\delta}_t := R_t + \gamma \hat{V}(S_{t+1}) - \hat{V}(S_t).$$

If $\hat{V} = V^\pi$, the expected TD-error satisfies

$$\begin{aligned} \mathbb{E}[\delta_t \mid S_t, A_t] &= \mathbb{E}[R_t + \gamma V^\pi(S_{t+1}) - V^\pi(S_t) \mid S_t, A_t] \\ &= Q^\pi(S_t, A_t) - V^\pi(S_t) = A^\pi(S_t, A_t), \end{aligned}$$

due to the Bellman-equation, making δ_t an unbiased estimator of the true advantage. Using the TD-errors, we define *k-step advantage estimators* that accumulate information from k future steps before bootstrapping with V :

$$\hat{A}_t^{(k)} \stackrel{\text{telescop}}{:=} \sum_{\ell=0}^k \gamma^\ell R_{t+\ell} + \gamma^{k+1} \hat{V}(S_{t+k+1}) - \hat{V}(S_t) = \sum_{\ell=0}^k \gamma^\ell \delta_{t+\ell}.$$

Again, if $\hat{V} = V^\pi$, the k -step advantage estimator provides an unbiased estimator of the true advantage.

The (infinite horizon) Generalized Advantage Estimator introduced in [?] combines all k -step estimators using exponentially decaying weights controlled by a parameter $\lambda \in [0, 1]$:

$$\hat{A}_t^{\text{GAE}(\lambda)} = \hat{A}_t^\infty := (1 - \lambda) \sum_{k=0}^{\infty} \lambda^k \hat{A}_t^{(k)}.$$

Since each estimator is unbiased, also the convex combination is unbiased. Averaging geometrically over all k -step advantage estimators, it implicitly assumes access to trajectories of infinite length.

5.2.2 Finite time-horizon GAE (new!)

To account for the finite rollout length T used in practice, we define the fixed-time GAE, which replaces the infinite geometric weighting in standard GAE by a normalized, finite geometric sum up to the last available timestep T . For a given timestep $t < T$, the estimator is defined as

$$\hat{A}_t^T := \frac{1 - \lambda}{1 - \lambda^{T-t}} \sum_{k=0}^{T-t-1} \lambda^k \hat{A}_t^{(k)},$$

where again $\hat{A}_t^{(k)}$ denotes the k -step advantage estimator $\hat{A}_t^{(k)} = \sum_{\ell=0}^k \gamma^\ell \delta_{t+\ell}$. The normalization factor $\frac{1-\lambda}{1-\lambda^{T-t}}$ ensures that the geometric weights sum to one. This formulation yields a consistent fixed-horizon analogue of GAE that aligns with the data available from truncated rollouts.

5.2.3 Termination-time GAE (new)

In practical reinforcement learning settings, trajectories often terminate before reaching the fixed rollout horizon T , because an absorbing state is reached. In such cases, it is no longer meaningful to average over k -step estimators beyond the point of termination, as those future returns cannot be observed. To address this, we restrict the geometric averaging to the range of steps actually available before termination. This leads to an estimator that depends on a random horizon, given by the episode's termination time. We denote by $\tilde{\tau}$ the termination time, i.e.

$$\tilde{\tau} := \inf\{t \geq 0 \mid S_t \text{ is terminal}\},$$

which is a stopping time with respect to the natural filtration. We assume that $\tilde{\tau}$ is a.s. finite. Moreover, let us define $\tau := \tilde{\tau} \wedge T$. Conditioning on τ allows us to construct an advantage estimator that mirrors the fixed-time formulation but adapts to the random episode length. Specifically, we define

$$\hat{A}_t^\tau := \frac{1 - \lambda}{1 - \lambda^{\tau-t}} \sum_{k=0}^{\tau-t-1} \lambda^k \hat{A}_t^{(k)},$$

where the sum now runs only over k -step estimators that only observe information before termination. When $\tau = T$, the estimator coincides with the finite-horizon version \hat{A}_t^T , but for early termination it automatically adjusts to the available rollout length, avoiding the use of unavailable future returns.

5.2.4 The backwards induction

5.2.5 Relationship of the finite and the stopping time GAE

Throughout this section, we assume that we terminate in finite time and that the value function estimator satisfies: $\hat{V}(S_t) = \hat{V}(S_\tau) = V^\pi(S_\tau) = 0$ for all $t \geq \tau$. By separating the geometric sum in the definition of the finite-horizon estimator at the (random) termination time τ , we can express \hat{A}_t^T as a combination of the stopping-time estimator \hat{A}_t^τ and the last available k -step estimate $\hat{A}_t^{(\tau-t-1)}$:



Proposition 5.2.1. Suppose that we terminate in finite time. For $t = 0, \dots, \tau - 1$, we have

$$\hat{A}_t^T = \frac{1 - \lambda^{\tau-t}}{1 - \lambda^{T-t}} \hat{A}_t^\tau + \frac{\lambda^{\tau-t} - \lambda^{T-t}}{1 - \lambda^{T-t}} \hat{A}_t^{(\tau-t-1)}. \quad (5.1)$$

Proof. Let $t \in \{0, \dots, \tau - 1\}$. Then,

$$\begin{aligned} \hat{A}_t^T &= \frac{1 - \lambda}{1 - \lambda^{T-t}} \left(\sum_{k=0}^{T-t-1} \lambda^k \hat{A}_t^{(k)} \right) \\ &= \frac{1 - \lambda}{1 - \lambda^{T-t}} \left(\sum_{k=0}^{\tau-t-1} \lambda^k \hat{A}_t^{(k)} + \sum_{k=\tau-t}^{T-t-1} \lambda^k \hat{A}_t^{(k)} \right). \end{aligned}$$

Note, that $V^\pi(S_\tau) = 0$. Furthermore, once we have reached a terminal state S_τ , we stay in the terminal state for the rest of the rollout, i.e. $S_t = S_\tau$, for all $t = \tau, \dots, T$ and do not receive any

new reward. Thus, $\hat{A}_t^{(k)} = \hat{A}_t^{(\tau-t-1)}$ for all $k \in \{\tau-t, T-t-1\}$, which allows us to rewrite

$$\begin{aligned}
\hat{A}_t^T &= \frac{1-\lambda}{1-\lambda^{T-t}} \left(\sum_{k=0}^{\tau-t-1} \lambda^k \hat{A}_t^{(k)} + \hat{A}_t^{(\tau-t-1)} \sum_{k=\tau-t}^{T-t-1} \lambda^k \right) \\
&= \frac{1-\lambda}{1-\lambda^{T-t}} \left(\sum_{k=0}^{\tau-t-1} \lambda^k \hat{A}_t^{(k)} + \hat{A}_t^{(\tau-t-1)} \lambda^{\tau-t} \sum_{k=0}^{T-1-\tau} \lambda^k \right) \\
&= \frac{1-\lambda}{1-\lambda^{T-t}} \left(\sum_{k=0}^{\tau-t-1} \lambda^k \hat{A}_t^{(k)} + \hat{A}_t^{(\tau-t-1)} \lambda^{\tau-t} \frac{1-\lambda^{T-\tau}}{1-\lambda} \right) \\
&= \frac{1-\lambda^{\tau-t}}{1-\lambda^{T-t}} \hat{A}_t^\tau + \frac{1-\lambda}{1-\lambda^{T-t}} \left(\hat{A}_t^{(\tau-t-1)} \lambda^{\tau-t} \frac{1-\lambda^{T-\tau}}{1-\lambda} \right) \\
&= \frac{1-\lambda^{\tau-t}}{1-\lambda^{T-t}} \hat{A}_t^\tau + \frac{\lambda^{\tau-t} - \lambda^{T-t}}{1-\lambda^{T-t}} \hat{A}_t^{(\tau-t-1)}. \quad \square
\end{aligned}$$

The relation shows that the finite-horizon estimator \hat{A}_t^T can be interpreted as a reweighted version of the stopping-time estimator \hat{A}_t^τ . Once termination occurs, all k -step advantage estimates beyond the termination time coincide, i.e. $\hat{A}_t^{(k)} = \hat{A}_t^{(\tau-t-1)}$ for all $k \geq \tau-t$, because no further rewards or transitions are observed after entering an absorbing state. Consequently, the geometric averaging in \hat{A}_t^T effectively reuses the last available k -step estimate multiple times for the remaining weights $\lambda^{\tau-t}, \lambda^{\tau-t+1}, \dots, \lambda^{T-t-1}$. This causes the second term above, which redistributes more weighting mass onto the final nontrivial estimate $\hat{A}_t^{(\tau-t-1)}$. In contrast, the stopping-time estimator \hat{A}_t^τ averages only over the actually observed k -step estimates up to termination, assigning zero weight to unobserved future terms.

Relation between the stopping-time and the standard GAE estimator. Analogously, we can relate the stopping-time estimator \hat{A}_t^τ to the standard infinite-horizon GAE estimator \hat{A}_t^∞ . By separating the geometric sum at the termination time τ , we obtain: Suppose that we terminate in finite time. For $t = 0, \dots, \tau-1$, we have

$$\hat{A}_t^\infty = (1-\lambda^{\tau-t}) \hat{A}_t^\tau + \lambda^{\tau-t} \hat{A}_t^{(\tau-t-1)}.$$

Proof. From the definition of the infinite-horizon estimator,

$$\hat{A}_t^\infty = (1-\lambda) \sum_{k=0}^{\infty} \lambda^k \hat{A}_t^{(k)}.$$

We decompose this sum at the termination time τ :

$$\hat{A}_t^\infty = (1-\lambda) \left(\sum_{k=0}^{\tau-t-1} \lambda^k \hat{A}_t^{(k)} + \sum_{k=\tau-t}^{\infty} \lambda^k \hat{A}_t^{(k)} \right).$$

After termination, again all future k -step estimators coincide with the last nontrivial one. Hence,

$$\begin{aligned}
\hat{A}_t^\infty &= (1-\lambda) \left(\sum_{k=0}^{\tau-t-1} \lambda^k \hat{A}_t^{(k)} + \hat{A}_t^{(\tau-t-1)} \sum_{k=\tau-t}^{\infty} \lambda^k \right) \\
&= (1-\lambda) \left(\sum_{k=0}^{\tau-t-1} \lambda^k \hat{A}_t^{(k)} + \hat{A}_t^{(\tau-t-1)} \lambda^{\tau-t} \frac{1}{1-\lambda} \right) \\
&= (1-\lambda^{\tau-t}) \hat{A}_t^\tau + \lambda^{\tau-t} \hat{A}_t^{(\tau-t-1)}. \quad \square
\end{aligned}$$

This expression shows that the infinite-horizon GAE estimator \hat{A}_t^∞ can be interpreted as a reweighted extension of the stopping-time estimator \hat{A}_t^τ . After termination, all k -step advantage

estimates collapse to the final value $\hat{A}_t^{(\tau-t-1)}$, which is repeatedly weighted in the infinite geometric sum. The coefficient $\lambda^{\tau-t}$ thus represents the corresponding cumulative geometric weight. When the episode length is long or λ is small, the difference between \hat{A}_t^∞ and \hat{A}_t^τ becomes negligible.

Example: Relation between the stopping-time and the standard GAE estimator.

$$\begin{aligned}\hat{A}_{\tau-1}^\tau &= \hat{A}_{\tau-1}^{(0)} = \delta_{\tau-1} \\ \hat{A}_{\tau-2}^\tau &= \frac{1-\lambda}{1-\lambda^2}(\hat{A}_{\tau-2}^{(0)} + \lambda\hat{A}_{\tau-2}^{(1)}) \\ &= \frac{1-\lambda}{1-\lambda^2}(\delta_{\tau-2} + \lambda(\delta_{\tau-2} + \gamma\delta_{\tau-1})) = \delta_{\tau-2} + \lambda\gamma\frac{1-\lambda}{1-\lambda^2}\hat{A}_{\tau-1}^\tau \\ \hat{A}_{\tau-3}^\tau &= \frac{1-\lambda}{1-\lambda^3}(\hat{A}_{\tau-3}^{(0)} + \lambda\hat{A}_{\tau-3}^{(1)} + \lambda^2\hat{A}_{\tau-3}^{(2)}) = \dots \\ &= \delta_{\tau-3} + \gamma\lambda\frac{1-\lambda^2}{1-\lambda^3}\hat{A}_{\tau-2}^\tau.\end{aligned}$$

$$\begin{aligned}\hat{A}_{\tau-1}^\infty &= \hat{A}_{\tau-1}^{(0)} = \delta_{\tau-1} \\ \hat{A}_{\tau-2}^\infty &= (1-\lambda)(\hat{A}_{\tau-2}^{(0)} + \hat{A}_{\tau-2}^{(1)} \sum_{l=1}^{\infty} \lambda^l) \\ &= (1-\lambda)(\delta_{\tau-2} + (\delta_{\tau-2} + \gamma\delta_{\tau-1})\frac{\lambda}{1-\lambda}) = \delta_{\tau-2} + \lambda\gamma\hat{A}_{\tau-1}^\infty \\ \hat{A}_{\tau-3}^\infty &= (1-\lambda)(\hat{A}_{\tau-3}^{(0)} + \lambda\hat{A}_{\tau-3}^{(1)} + \sum_{l=2}^{\infty} \lambda^l \hat{A}_{\tau-3}^{(2)}) \\ &= (1-\lambda)((\delta_{\tau-3} + \lambda(\delta_{\tau-3} + \gamma\delta_{\tau-2})) \\ &\quad + \sum_{l=2}^{\infty} \lambda^l (\delta_{\tau-3} + \gamma\delta_{\tau-2} + \gamma^2\delta_{\tau-1})) \\ &= \delta_{\tau-3} + \gamma\lambda\hat{A}_{\tau-2}^\infty.\end{aligned}$$

5.2.6 Unbiasedness of the estimators

In this section, we establish that if the true value function V^π of the policy π is used to compute the temporal-difference errors, then all proposed advantage estimators are unbiased estimators of the true advantage function $A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$.

Theorem 5.2.2. Suppose that the true value function V^π of π is used in the TD-errors:

$$\delta_t := R_t + \gamma V^\pi(S_{t+1}) - V^\pi(S_t), \quad \gamma \in (0, 1).$$

Then, for any finite horizon T :

$$\mathbb{E}[\hat{A}_t^T \mid S_t, A_t] = A^\pi(S_t, A_t) = Q^\pi(S_t, A_t) - V^\pi(S_t).$$

Proof. For a fixed starting time t , recall that the fixed-time estimator is defined by

$$\hat{A}_t^T = \frac{1-\lambda}{1-\lambda^{T-t}} \sum_{k=0}^{T-t-1} \lambda^k \hat{A}_t^{(k)}.$$

For any deterministic $k \leq T - t - 1$, we have due to the Bellman equation

$$\begin{aligned} \mathbb{E}[\hat{A}_t^{(k)} | S_t, A_t] &= \mathbb{E}\left[\sum_{\ell=0}^k \gamma^\ell R_{t+\ell} + \gamma^{k+1} V^\pi(S_{t+k+1})\right. \\ &\quad \left. - V^\pi(S_t) | S_t, A_t\right] \\ &= A^\pi(S_t, A_t). \end{aligned}$$

As the fixed-time estimator is a normalized, geometrically weighted average of the k -step estimators using aboves identity, yields

$$\begin{aligned} \mathbb{E}[\hat{A}_t^T | S_t, A_t] &= \frac{1 - \lambda}{1 - \lambda^{T-t}} \sum_{k=0}^{T-t-1} \lambda^k \mathbb{E}[\hat{A}_t^{(k)} | S_t, A_t] \\ &= \frac{1 - \lambda}{1 - \lambda^{T-t}} \sum_{k=0}^{T-t-1} \lambda^k A^\pi(S_t, A_t) \\ &= A^\pi(S_t, A_t), \end{aligned}$$

since the geometric weights sum to one. \square

[Unbiasedness of the last term in Stopping-Time Estimator] Asume that we use the true Value function V^π of π for our TD Errors:

$$\delta_t := R_t + \gamma V^\pi(S_{t+1}) - V^\pi(S_t), \quad \gamma \in (0, 1).$$

Then, for $t \leq \tau - 1$:

$$\mathbb{E}[\hat{A}_t^{(\tau-t-1)} | S_t, A_t] = A^\pi(S_t, A_t).$$

Proof. Define

$$\mathcal{G}_k := \sigma(S_t, A_t, \dots, S_{t+k}, A_{t+k}, S_{t+k+1}) \quad (k \geq 0),$$

and the process

$$\begin{aligned} M_k &:= \hat{A}_t^{(k)} = \sum_{\ell=0}^k \gamma^\ell \delta_{t+\ell} \\ &= \sum_{l=0}^k \gamma^l R_{t+l} + \gamma V^\pi(S_{t+k+1}) - V^\pi(S_t). \end{aligned}$$

Suppose (M_k) is uniformly integrable (e.g., bounded rewards imply bounded V^π , hence $|\delta_s|$ is bounded, hence (M_k) is UI).

(i) **Martingale property.** By definition (M_k) is \mathcal{G}_k -adapted and integrable. Moreover,

$$\begin{aligned} &\mathbb{E}[M_{k+1} - M_k | \mathcal{G}_k] \\ &= \gamma^{k+1} \mathbb{E}[\delta_{t+k+1} | \mathcal{G}_k] \\ &= \gamma^{k+1} \mathbb{E}[R_{t+k+1} + \gamma V^\pi(S_{t+k+2}) - V^\pi(S_{t+k+1}) | \mathcal{G}_k] \\ &= \gamma^{k+1} \mathbb{E}[R_{t+k+1} + \gamma V^\pi(S_{t+k+2}) - V^\pi(S_{t+k+1}) | \mathcal{S}_{t+k+1}] \\ &= 0, \end{aligned}$$

where we have used the Bellman Equation for V^π .

(ii) **Define a stopping time.** Let $N := \max\{\tau - t - 1, 0\}$. For any k , $\{N \leq k\} = \{\tau \leq t + k + 1\} \in \mathcal{G}_k$ because \mathcal{G}_k contains the entire state information up to S_{t+k+1} . Thus, N is a stopping time for (\mathcal{G}_k) .

(iii) **Optional stopping.** Applying Optional Stopping Theorem, yields:

$$\mathbb{E}[M_N | \mathcal{G}_0] = \mathbb{E}[\hat{A}_t^{(\tau-t-1)} | \mathcal{G}_0] = M_0 = \hat{A}_t^{(0)} = \delta_t. \quad \text{a.s.}$$

Conditioning on (S_t, A_t) and using the tower property, yields

$$\begin{aligned} & \mathbb{E}[\hat{A}_t^{(\tau-t-1)} | S_t, A_t] \\ &= \mathbb{E}[\mathbb{E}[\hat{A}_t^{(\tau-t-1)} | S_t, A_t, S_{t+1}] | S_t, A_t] \\ &= \mathbb{E}[\mathbb{E}[\hat{A}_t^{(\tau-t-1)} | \mathcal{G}_0] | S_t, A_t] \\ &= \mathbb{E}[\delta_t | S_t, A_t] \\ &= \mathbb{E}[R_t + \gamma V^\pi(S_{t+1}) - V^\pi(S_t) | S_t, A_t] \\ &= Q^\pi(S_t, A_t) - V^\pi(S_t) \\ &= A^\pi(S_t, A_t). \end{aligned} \quad \square$$

Theorem 5.2.3. Suppose that the true value function V^π of π is used in the TD-errors:

$$\delta_t := R_t + \gamma V^\pi(S_{t+1}) - V^\pi(S_t), \quad \gamma \in (0, 1).$$

Then,

$$\mathbb{E}[\hat{A}_t^\tau | S_t, A_t] = A^\pi(S_t, A_t) = Q^\pi(S_t, A_t) - V^\pi(S_t).$$

Proof. Todo. □

5.3 PPO as importance sampling variance reduction method

Here is the point of PPO. So far we computed gradients with advantage estimation. We also used importance sampling to improve sample efficiency. The problem is, importance sampling makes the variances explode. We next introduce two tricks. First, early stopping to avoid explosion of the variance if sampling and current policy are too far apart. Second, we flatten rollouts to a transition level, this reduces variances of the gradient estimator as it reduces correlations through rollouts. To make the second trick work, we need the PPO trick. Let us first recall the gradient representation with importance sampling and advantage:

$$\nabla_\theta J(\theta) = \mathbb{E}^{\pi^{\text{old}}} \left[\sum_{t=0}^{T-1} \gamma^t \prod_{i=0}^t \frac{\nabla_\theta \pi^\theta(A_i; S_i)}{\pi^{\text{old}}(A_i; S_i)} A_t^{\pi^\theta}(S_t, A_t) \right],$$

Instead of using entire sampled rollouts we would like to use transitions (S_i, A_i) that are sampled from independent rollouts. Why? Because it reduces the variance of the sum, inter-time correlations disappear. To do this, we modify the gradient estimator to

$$\widehat{\nabla_\theta J(\theta)} = \mathbb{E}^{\pi^{\text{old}}} \left[\sum_{t=0}^{T-1} \gamma^t \frac{\nabla_\theta \pi^\theta(A_t; S_t)}{\pi^{\text{old}}(A_t; S_t)} A_t^{\pi^\theta}(S_t, A_t) \right].$$

Doing this introduces a bias to the gradient but allows us to either sample with weight proportional to γ^t transitions that contribute with $\frac{\nabla_\theta \pi^\theta(A_t; S_t)}{\pi^{\text{old}}(A_t; S_t)} A_t^{\pi^\theta}(S_t, A_t)$ to the gradient estimate, or, alternatively, sample uniformly in time with contribution $\frac{\nabla_\theta \pi^\theta(A_t; S_t)}{\pi^{\text{old}}(A_t; S_t)} \gamma^t A_t^{\pi^\theta}(S_t, A_t)$.

5.3.1 Bias of the PPO surrogate gradient

We have a fixed behaviour policy π_{old} from which we can sample and want to use this to estimate the gradient $\nabla_\theta J(\pi_\theta)$.


Theorem 5.3.1. (PPO surrogate bias (no clipping))

Suppose π_{old} is a behaviour policy (with strictly positive weights) and define the surrogate gradient by

$$\widehat{\nabla_{\theta} J}(\theta) = \mathbb{E}^{\pi_{\text{old}}} \left[\sum_{t=0}^{T-1} \gamma^t \frac{\nabla_{\theta} \pi^{\theta}(A_t; S_t)}{\pi_{\text{old}}(A_t; S_t)} A_t^{\pi_{\text{old}}}(S_t, A_t) \right]$$

and additionally assume that $\|A_t^{\pi_{\text{old}}}(s, \cdot)\|_{\infty} \leq A_{\text{max}}$. Furthermore, we abbreviate

$$\overline{\text{KL}}(\pi_{\text{old}}, \pi_{\theta}) := \mathbb{E}_{S \sim d^{\pi_{\text{old}}}} [\text{KL}(\pi_{\text{old}}(\cdot; S) \parallel \pi_{\theta}(\cdot; S))],$$

where $d^{\pi_{\text{old}}} := \frac{1}{T} \sum_{t=0}^{T-1} d_t^{\pi_{\text{old}}}$. Then

$$\|\nabla_{\theta} J(\pi_{\theta}) - \widehat{\nabla_{\theta} J}(\pi_{\theta})\|_{\infty} \leq \sqrt{16} \sqrt{\overline{\text{KL}}(\pi_{\text{old}}, \pi_{\theta})} T^{5/2} \|\nabla_{\theta} \log \pi_{\theta}\|_{\infty} A_{\text{max}}.$$

We actually prove

$$\left\| \nabla_{\theta} J(\pi_{\theta}) - \widehat{\nabla_{\theta} J}(\pi_{\theta}) \right\|_{\infty} \leq C \left(T^2 \sqrt{\varepsilon} + \sqrt{T} \varepsilon (\sqrt{2T(T-1)} + T) \right),$$

with $C = \sqrt{2} \|\nabla_{\theta} \log \pi_{\theta}\|_{\infty} A_{\text{max}}$ but write the estimate in a more compact form.

For the proof we need a couple of preparations.

Let us start with a well-known identity, also known as the performance difference lemma (see for instance [?, Lemma 6.1] and [?, Eqn. (1)] for the infinite time-horizon setting).


Lemma 5.3.2. (Performance difference lemma)

For two policies $\pi, \tilde{\pi}$,

$$J(\tilde{\pi}) - J(\pi) = \mathbb{E}^{\tilde{\pi}} \left[\sum_{t=0}^{T-1} \gamma^t A_t^{\tilde{\pi}}(S_t, A_t) \right] = \sum_{t=0}^{T-1} \gamma^t \mathbb{E}^{\tilde{\pi}} [A_t^{\tilde{\pi}}(S_t, A_t)] =: \sum_{t=0}^{T-1} \gamma^t \Delta_t. \quad (5.2)$$

Proof. We have $Q_t^{\pi}(s, a) = \mathbb{E}_{S_t=s, A_t=a}^{\pi} [R_t + \gamma V_{t+1}^{\pi}(S_{t+1})]$ and $\mathbb{E}_{S \sim d_t^{\pi}, A \sim \pi(\cdot; S)} [A_t^{\pi}(S, A)] = 0$. Thus,

$$\begin{aligned} J(\tilde{\pi}) - J(\pi) &= \sum_{t=0}^{T-1} \gamma^t (\mathbb{E}^{\tilde{\pi}} [R_t] - \mathbb{E}^{\pi} [R_t]) \\ &= \sum_{t=0}^{T-1} \gamma^t \left(\mathbb{E}_{S \sim d_t^{\tilde{\pi}}, A \sim \tilde{\pi}(\cdot; S)} \left[\underbrace{\mathbb{E}_{S_t=S, A_t=A}^{\pi} [R_t]}_{=Q_t^{\pi}(S, A) - \gamma \mathbb{E}_{S' \sim p(\cdot; S, A)} [V_{t+1}^{\pi}(S')]} \right] - \mathbb{E}_{S \sim d_t^{\pi}, A \sim \pi(\cdot; S)} \left[\mathbb{E}_{S_t=S, A_t=A}^{\pi} [R_t] \right] \right) \\ &= \sum_{t=0}^{T-1} \gamma^t \left(\mathbb{E}_{S \sim d_t^{\tilde{\pi}}, A \sim \tilde{\pi}(\cdot; S)} [A_t^{\tilde{\pi}}(S, A) - \gamma \mathbb{E}_{S' \sim p(\cdot; S, A)} [V_{t+1}^{\pi}(S')] + V_t^{\pi}(S)] \right. \\ &\quad \left. - \mathbb{E}_{S \sim d_t^{\pi}, A \sim \pi(\cdot; S)} [A_t^{\pi}(S, A) - \gamma \mathbb{E}_{S' \sim p(\cdot; S, A)} [V_{t+1}^{\pi}(S')] + V_t^{\pi}(S)] \right) \\ &= \sum_{t=0}^{T-1} \gamma^t \mathbb{E}_{S \sim d_t^{\tilde{\pi}}, A \sim \tilde{\pi}(\cdot; S)} [A_t^{\tilde{\pi}}(S, A)] \\ &\quad + \sum_{t=0}^{T-1} \gamma^t \left(\gamma \mathbb{E}_{S \sim d_t^{\tilde{\pi}}, A \sim \tilde{\pi}(\cdot; S)} \left[\mathbb{E}_{S' \sim p(\cdot; S, A)} [V_{t+1}^{\pi}(S')] \right] - \mathbb{E}_{S \sim d_t^{\pi}, A \sim \pi(\cdot; S)} [V_t^{\pi}(S)] \right) \\ &\quad - \sum_{t=0}^{T-1} \gamma^t \left(\gamma \mathbb{E}_{S \sim d_t^{\pi}, A \sim \pi(\cdot; S)} \left[\mathbb{E}_{S' \sim p(\cdot; S, A)} [V_{t+1}^{\pi}(S')] \right] - \mathbb{E}_{S \sim d_t^{\tilde{\pi}}, A \sim \tilde{\pi}(\cdot; S)} [V_t^{\pi}(S)] \right). \end{aligned}$$

For the second sum, we calculate, using $V_T^\pi \equiv 0$,

$$\begin{aligned}
& \sum_{t=0}^{T-1} \gamma^t (\gamma \mathbb{E}_{S \sim d_t^\pi, A \sim \pi(\cdot; S)} [\mathbb{E}_{S' \sim p(\cdot; S, A)} [V_{t+1}^\pi(S')]] - \mathbb{E}_{S \sim d_t^\pi, A \sim \pi(\cdot; S)} [V_t^\pi(S)]) \\
&= \sum_{t=0}^{T-1} \gamma^t (\gamma \mathbb{E}^\pi [V_{t+1}^\pi(S_{t+1})] - \mathbb{E}^\pi [V_t^\pi(S_t)]) \\
&= \gamma^T \mathbb{E}^\pi [V_T^\pi(S_T)] - \mathbb{E}^\pi [V_0^\pi(S_0)] + \sum_{t=0}^{T-2} \gamma^{t+1} \mathbb{E}^\pi [V_{t+1}^\pi(S_{t+1})] - \sum_{t=1}^{T-1} \gamma^t \mathbb{E}^\pi [V_t^\pi(S_t)] \\
&= -\mathbb{E}^\pi [V_0^\pi(S_0)],
\end{aligned}$$

and analogously for the third sum. Because $d_0^\pi = \mu = d_0^{\tilde{\pi}}$, we have $\mathbb{E}^\pi [V_0^\pi(S_0)] = \mathbb{E}^{\tilde{\pi}} [V_0^\pi(S_0)]$, meaning these sums cancel, which finishes the proof. \square

We bound the KL-divergence between the state marginals for different policies in terms of the KL-divergence between the policies, which we can actually control. The following lemma is essentially an application of the formula $\text{KL}(\otimes_i \mathbb{P}_i \parallel \otimes_i \mathbb{Q}_i) = \sum_i \text{KL}(\mathbb{P}_i \parallel \mathbb{Q}_i)$ for product measures [?, p. 85]:



Lemma 5.3.3. Let \mathbb{P}_t^π and $\mathbb{P}_t^{\tilde{\pi}}$ be the distributions of $(S_0, A_0, R_0, \dots, S_t)$ under π and $\tilde{\pi}$, respectively. Then

$$\text{KL}(\mathbb{P}_t^{\tilde{\pi}} \parallel \mathbb{P}_t^\pi) = \sum_{i=0}^{t-1} \mathbb{E}^{\tilde{\pi}} [\text{KL}(\tilde{\pi}(\cdot; S_i) \parallel \pi(\cdot; S_i))].$$

Proof. By factorising \mathbb{P}_t^π and $\mathbb{P}_t^{\tilde{\pi}}$ using the MDP dynamics, we find

$$\log \frac{\mathbb{P}_t^{\tilde{\pi}}(s_0, a_0, r_0, \dots, s_t)}{\mathbb{P}_t^\pi(s_0, a_0, r_0, \dots, s_t)} = \sum_{i=0}^{t-1} \log \frac{\tilde{\pi}(a_i; s_i)}{\pi(a_i; s_i)}$$

for any $(s_0, a_0, r_0, \dots, s_t)$. Using the definition of KL and Fubini gives

$$\begin{aligned}
\text{KL}(\mathbb{P}_t^{\tilde{\pi}} \parallel \mathbb{P}_t^\pi) &= \mathbb{E}^{\tilde{\pi}} \left[\sum_{i=0}^{t-1} \log \frac{\tilde{\pi}(A_i; S_i)}{\pi(A_i; S_i)} \right] = \sum_{i=0}^{t-1} \mathbb{E}_{S \sim d_i^{\tilde{\pi}}} \left[\mathbb{E}_{A \sim \tilde{\pi}(\cdot; S)} \left[\log \frac{\tilde{\pi}(A; S)}{\pi(A; S)} \right] \right] \\
&= \sum_{i=0}^{t-1} \mathbb{E}_{S \sim d_i^{\tilde{\pi}}} [\text{KL}(\tilde{\pi}(\cdot; S) \parallel \pi(\cdot; S))].
\end{aligned}$$

\square



Lemma 5.3.4. For two policies $\pi, \tilde{\pi}$ with $\max_s \text{KL}(\tilde{\pi}(\cdot; s) \parallel \pi(\cdot; s)) \leq \varepsilon$, we have

$$\text{KL}(d_t^{\tilde{\pi}} \parallel d_t^\pi) \leq t\varepsilon,$$

where $d_t^\pi = \mathbb{P}^\pi(S_t \in \cdot)$ and $d_t^{\tilde{\pi}} = \mathbb{P}^{\tilde{\pi}}(S_t \in \cdot)$.

Proof. Because d_t^π and $d_t^{\tilde{\pi}}$ are marginal distributions of \mathbb{P}_t^π and $\mathbb{P}_t^{\tilde{\pi}}$, respectively, the chain rule for the KL-divergence¹ implies

$$\text{KL}(d_t^{\tilde{\pi}} \parallel d_t^\pi) \leq \text{KL}(\mathbb{P}_t^{\tilde{\pi}} \parallel \mathbb{P}_t^\pi) = \sum_{i=0}^{t-1} \mathbb{E}^{\tilde{\pi}} [\text{KL}(\tilde{\pi}(\cdot; S_i) \parallel \pi(\cdot; S_i))] \leq t\varepsilon.$$

\square

¹For instance Theorem 2.5.3 of T. Cover "Elements of information theory"

We next give an adaptation and improvement of Lemma 2 from the TRPO paper for the finite-time horizon setting.



Lemma 5.3.5. If $\max_s \text{KL}(\pi_{\text{old}}(\cdot; s) \parallel \pi_\theta(\cdot; s)) \leq \varepsilon$ and $\|A_t^{\pi_{\text{old}}}(s, \cdot)\|_\infty \leq A_{\text{max}}$, then

$$\|g_t^\theta\|_\infty \leq \sqrt{2\varepsilon}A_{\text{max}},$$

where $g_t^\theta(s) := \mathbb{E}_{A \sim \pi_\theta(\cdot; s)}[A_t^{\pi_{\text{old}}}(s, A)]$, $s \in \mathcal{S}$.

Proof. Because $\mathbb{E}_{A \sim \pi_{\text{old}}(\cdot; s)}[A_t^{\pi_{\text{old}}}(s, A)] = 0$, a characterisation of the total variation distance² and Pinsker's inequality³ imply

$$\begin{aligned} \|g_t^\theta(s)\|_\infty &= \|\mathbb{E}_{A \sim \pi_\theta(\cdot; s)}[A_t^{\pi_{\text{old}}}(s, A)] - \mathbb{E}_{A \sim \pi_{\text{old}}(\cdot; s)}[A_t^{\pi_{\text{old}}}(s, A)]\|_\infty \\ &\leq 2\text{TV}(\pi_\theta(\cdot; s), \pi_{\text{old}}(\cdot; s))A_{\text{max}} \\ &\leq \sqrt{2\varepsilon}A_{\text{max}}. \end{aligned}$$

□

We can now give the proof of Theorem 5.3.1.

Proof of Theorem 5.3.1: First note that Fubini allows us to rewrite

$$\widehat{\nabla}_\theta J = \sum_{t=0}^{T-1} \gamma^t \nabla_\theta \hat{\Delta}_t$$

with $\hat{\Delta}_t := \mathbb{E}^{\pi_{\text{old}}}[g_t^\theta(S_t)]$, $\Delta_t := \mathbb{E}^{\pi_\theta}[g_t^\theta(S_t)]$, and g_t^θ as defined above. For a single component θ_j of θ , we calculate

$$\partial_{\theta_j} d_t^{\pi_\theta}(s) = \sum_{\tau: s_t=s} \partial_{\theta_j} \mathbb{P}_t^{\pi_\theta}(\tau) = \sum_{\tau: s_t=s} \mathbb{P}_t^{\pi_\theta}(\tau) \partial_{\theta_j} \log \mathbb{P}_t^{\pi_\theta}(\tau) = \sum_{\tau: s_t=s} \mathbb{P}_t^{\pi_\theta}(\tau) \sum_{i=0}^{t-1} \partial_{\theta_j} \log \pi_\theta(a_i; s_i)$$

and

$$|\partial_{\theta_j} g_t^\theta(s)| = \left| \sum_{a \in \mathcal{A}} \partial_{\theta_j} \pi_\theta(a; s) A_t^{\pi_{\text{old}}}(s, a) \right| = |\mathbb{E}_{A \sim \pi_\theta(\cdot; s)}[\partial_{\theta_j} \log \pi_\theta(A; s) A_t^{\pi_{\text{old}}}(s, A)]| \leq \|\nabla_\theta \log \pi_\theta\|_\infty A_{\text{max}}.$$

Combining this with Lemma 5.3.5, we have

$$\begin{aligned} |\partial_{\theta_j}(\Delta_t - \hat{\Delta}_t)| &= \left| \sum_{s \in \mathcal{S}} \sum_{\tau: s_t=s} \mathbb{P}_t^{\pi_\theta}(\tau) \sum_{i=0}^{t-1} \partial_{\theta_j} \log \pi_\theta(a_i; s_i) g_t^\theta(s) + \sum_{s \in \mathcal{S}} (d_t^{\pi_\theta}(s) - d_t^{\pi_{\text{old}}}(s)) \partial_{\theta_j} g_t^\theta(s) \right| \\ &\leq t \|\nabla_\theta \log \pi_\theta\|_\infty \mathbb{E}_{\tau \sim \mathbb{P}_t^{\pi_\theta}}[|g_t^\theta(s_t)|] + 2\text{TV}(d_t^{\pi_{\text{old}}}, d_t^{\pi_\theta}) \|\nabla_\theta g_t^\theta\|_\infty \\ &= t \|\nabla_\theta \log \pi_\theta\|_\infty \mathbb{E}_{S \sim d_t^{\pi_\theta}}[|g_t^\theta(S)|] + 2\text{TV}(d_t^{\pi_{\text{old}}}, d_t^{\pi_\theta}) \|\nabla_\theta g_t^\theta\|_\infty \\ &\leq \underbrace{\sqrt{2} \|\nabla_\theta \log \pi_\theta\|_\infty A_{\text{max}}}_{=: C} (t \mathbb{E}_{S \sim d_t^{\pi_\theta}}[\sqrt{\text{KL}(\pi_{\text{old}}(\cdot; S) \parallel \pi_\theta(\cdot; S))}] + \sqrt{\text{KL}(d_t^{\pi_{\text{old}}} \parallel d_t^{\pi_\theta})}). \end{aligned}$$

Using

$$\mathbb{E}_{S \sim d_t^{\pi_\theta}}[\sqrt{\text{KL}(\pi_{\text{old}}(\cdot; S) \parallel \pi_\theta(\cdot; S))}] \leq \mathbb{E}_{S \sim d_t^{\pi_{\text{old}}}}[\sqrt{\text{KL}(\pi_{\text{old}}(\cdot; S) \parallel \pi_\theta(\cdot; S))}] + \sqrt{2\text{KL}(d_t^{\pi_{\text{old}}} \parallel d_t^{\pi_\theta})}$$

gives

$$|\partial_{\theta_j}(\Delta_t - \hat{\Delta}_t)| \leq C(t \mathbb{E}_{S \sim d_t^{\pi_{\text{old}}}}[\sqrt{\text{KL}(\pi_{\text{old}}(\cdot; S) \parallel \pi_\theta(\cdot; S))}] + (t\sqrt{2} + 1)\sqrt{\text{KL}(d_t^{\pi_{\text{old}}} \parallel d_t^{\pi_\theta})}). \quad (5.3)$$

²see Prop. 4.5 of Levin, Peres "Mixing times of Markov chains"

³see for instance A. Tsybakov, "Introduction to Nonparametric Estimation"

Proceeding as in the proof of Lemma 5.3.4, we have

$$\text{KL}(d_t^{\pi_{\text{old}}} \| d_t^{\pi_\theta}) \leq \sum_{i=0}^{T-1} \mathbb{E}_{S \sim d_i^{\pi_{\text{old}}}} [\text{KL}(\pi_{\text{old}}(\cdot; S) \| \pi_\theta(\cdot; S))] = T \overline{\text{KL}}.$$

Plugging this into (5.3) and summing over t gives

$$\begin{aligned} |\partial_{\theta_j} J(\pi_\theta) - \widehat{\nabla_{\theta_j} J(\pi_\theta)}| &\leq C \sum_{t=0}^{T-1} (t \mathbb{E}_{S \sim d_t^{\pi_{\text{old}}}} [\sqrt{\text{KL}(\pi_{\text{old}}(\cdot; S) \| \pi_\theta(\cdot; S))}] + (t\sqrt{2} + 1)\sqrt{T \overline{\text{KL}}}) \\ &\leq C(T^2 \mathbb{E}_{S \sim d^{\pi_{\text{old}}}} [\sqrt{\text{KL}(\pi_{\text{old}}(\cdot; S) \| \pi_\theta(\cdot; S))}] + \sqrt{T \overline{\text{KL}}}(\sqrt{2T(T-1)} + T)) \\ &\leq C(T^2 \sqrt{\overline{\text{KL}}} + \sqrt{T \overline{\text{KL}}}(\sqrt{2T(T-1)} + T)). \end{aligned}$$

□

We now continue with the PPO gradient estimator, including the clipping. We denote by

$$\widetilde{\nabla_{\theta} J}(\theta) = \mathbb{E}^{\pi_{\text{old}}} \left[\sum_{t=0}^{T-1} \gamma^t \frac{\nabla_{\theta} \pi^{\theta}(A_t; S_t)}{\pi_{\text{old}}(A_t; S_t)} \mathbf{1}_{\left| \frac{\pi^{\theta}(A_t; S_t)}{\pi_{\text{old}}(A_t; S_t)} - 1 \right| < \epsilon} A_t^{\pi_{\text{old}}}(S_t, A_t) \right]$$



WAS IST UEBERHAUPT DAS DING MIT TILDE

Theorem 5.3.6. (Clipped PPO Gradient Approximation)

Under the above assumptions and for sufficiently small $\text{KL} \leq \epsilon$, the clipped PPO gradient approximates the true REINFORCE gradient as

$$\begin{aligned} &\|\nabla J(\theta) - \widetilde{\nabla J}(\theta)\|_{\infty} \\ &\leq C_1 T^2 \sqrt{\epsilon} + C_2 \|\nabla_{\theta} \log \pi_{\theta}\|_{\infty} A_{\max} \frac{\sqrt{\epsilon}}{\epsilon} + O\left(\sqrt{T\epsilon}(\sqrt{2T(T-1)} + T)\right) = \sqrt{\overline{\text{KL}}}(\dots). \end{aligned}$$

The proof is a triangle argument. We already estimated the bias of the PPO-type gradients without clipping, now we compare the gradients with and without clipping. The idea is simple: We distinguish two events, the one in which clipping occurs and the one without clipping. Clipping only occurs when the likelihood ratios are large, but this means the KL-divergence would be large.



Lemma 5.3.7.

$$\mathbb{E}^{\pi_{\text{old}}} \left[\left| \frac{\pi_{\theta}(A_t; S_t)}{\pi_{\text{old}}(A_t; S_t)} - 1 \right| \right] \leq \sqrt{2 \mathbb{E}^{\pi_{\text{old}}} [\text{KL}(\pi_{\text{old}}(\cdot | s) \| \pi_{\theta}(\cdot | s))]}.$$

and in particular

$$\mathbb{P}^{\pi_{\text{old}}} \left(\left| \frac{\pi_{\theta}(A_t; S_t)}{\pi_{\text{old}}(A_t; S_t)} - 1 \right| > \epsilon \right) \leq \frac{\sqrt{\mathbb{E}^{\pi_{\text{old}}} [2 \text{KL}(\pi_{\text{old}}(\cdot | S_t) \| \pi_{\theta}(\cdot | S_t))]}{\epsilon}$$

Proof. Fix a state $s \in \mathcal{S}$, and denote $P := \pi_{\text{old}}(\cdot | s)$ and $Q := \pi_{\theta}(\cdot | s)$. Then

$$\mathbb{E}_{A \sim \pi_{\text{old}}(\cdot | s)} \left[\left| \frac{\pi_{\theta}(A; s)}{\pi_{\text{old}}(A; s)} - 1 \right| \right] = \sum_a P(a) \left| \frac{Q(a)}{P(a)} - 1 \right| = \sum_a |Q(a) - P(a)| = 2 \text{TV}(P, Q), \quad (5.4)$$

where $\text{TV}(P, Q)$ is the total variation distance

$$\text{TV}(P, Q) = \frac{1}{2} \sum_a |P(a) - Q(a)|.$$

Pinsker's inequality states that for any two probability measures P, Q on \mathcal{A} ,

$$\text{TV}(P, Q) \leq \sqrt{\frac{1}{2} \text{KL}(P\|Q)}. \quad (5.5)$$

We thus obtain

$$\mathbb{E}_{A \sim \pi_{\text{old}}(\cdot|s)} \left[\left| \frac{\pi_{\theta}(A; s)}{\pi_{\text{old}}(A; s)} - 1 \right| \right] \leq \sqrt{2 \text{KL}(\pi_{\text{old}}(\cdot|s) \| \pi_{\theta}(\cdot|s))}. \quad (5.6)$$

Integrating out the distribution of S_t with the use of Fubini and using Jensen's inequality yields

$$\begin{aligned} \mathbb{E}^{\pi_{\text{old}}} \left[\left| \frac{\pi_{\theta}(A_t; S_t)}{\pi_{\text{old}}(A_t; S_t)} - 1 \right| \right] &\leq \mathbb{E}^{\pi_{\text{old}}} \left[\sqrt{2 \text{KL}(\pi_{\text{old}}(\cdot|S_t) \| \pi_{\theta}(\cdot|S_t))} \right] \\ &\leq \sqrt{\mathbb{E}^{\pi_{\text{old}}} \left[2 \text{KL}(\pi_{\text{old}}(\cdot|S_t) \| \pi_{\theta}(\cdot|S_t)) \right]}. \end{aligned}$$

□

We can now proceed with the main argument.

Proof of Theorem 5.3.1: Using that the true and the simplified PPO surrogate gradients are equal when there is no clipping, the difference is

$$\begin{aligned} |\widehat{\nabla_{\theta} J}(\theta) - \widetilde{\nabla_{\theta} J}| &= \left| \mathbb{E}^{\pi_{\text{old}}} \left[\sum_{t=0}^{T-1} \gamma^t \frac{\nabla_{\theta} \pi^{\theta}(A_t; S_t)}{\pi_{\text{old}}(A_t; S_t)} \mathbf{1}_{\left| \frac{\pi^{\theta}(A_t; S_t)}{\pi_{\text{old}}(A_t; S_t)} - 1 \right| > \epsilon} A_t^{\pi_{\text{old}}}(S_t, A_t) \right] \right| \\ &= \left| \mathbb{E}^{\pi_{\text{old}}} \left[\sum_{t=0}^{T-1} \gamma^t \frac{\pi^{\theta}(A_t; S_t)}{\pi_{\text{old}}(A_t; S_t)} \nabla_{\theta} \log(\pi_{\theta}(A_t, S_t) \mathbf{1}_{\left| \frac{\pi^{\theta}(A_t; S_t)}{\pi_{\text{old}}(A_t; S_t)} - 1 \right| > \epsilon} A_t^{\pi_{\text{old}}}(S_t, A_t)) \right] \right| \\ &\leq C \mathbb{E}^{\pi_{\text{old}}} \left[\sum_{t=0}^{T-1} \gamma^t r_t \mathbf{1}_{|r_t - 1| > \epsilon} \right] \\ &= \dots??? \\ &= \mathbb{E}^{\pi_{\text{old}}} \left[\sum_{t=0}^{T-1} \gamma^t (1 + \epsilon) \mathbf{1}_{|r_t - 1| > \epsilon} A_t^{\pi_{\text{old}}}(S_t, A_t) \right] \\ &\leq A_{\max} (1 + \epsilon) \sum_{t=0}^{T-1} \gamma^t \mathbb{P}^{\pi_{\text{old}}} (|r_t - 1| > \epsilon) \\ &\leq A_{\max} (1 + \epsilon) \sum_{t=0}^{T-1} \gamma^t \frac{\sqrt{\mathbb{E}^{\pi_{\text{old}}} \left[2 \text{KL}(\pi_{\text{old}}(\cdot|S_t) \| \pi_{\theta}(\cdot|S_t)) \right]}}{\epsilon} \\ &\leq A_{\max} (1 + \epsilon) \frac{1 - \gamma^T}{1 - \gamma} \frac{\sqrt{\frac{1 - \gamma}{1 - \gamma^T} \sum_{t=0}^{T-1} \gamma^t \mathbb{E}^{\pi_{\text{old}}} \left[2 \text{KL}(\pi_{\text{old}}(\cdot|S_t) \| \pi_{\theta}(\cdot|S_t)) \right]}}{\epsilon} \\ &= A_{\max} (1 + \epsilon) \frac{1 - \gamma^T}{1 - \gamma} \frac{\sqrt{KL}}{\epsilon} \end{aligned}$$

In the last step we used Jensen's inequality. Finally, the triangle inequality yields

.....

□

The preceding bound treated the clipping condition as symmetric in r_t , i.e. activated whenever $|r_t - 1| > \epsilon$. In the actual PPO objective, clipping is *asymmetric*: only the upper side $r_t > 1 + \epsilon$ is active for positive advantages $\hat{A}_t > 0$, and the lower side $r_t < 1 - \epsilon$ is active for negative advantages $\hat{A}_t < 0$. Assuming the advantages take both signs with similar frequency, the probability of clipping on the relevant side of the ratio is at most half of the symmetric bound, yielding a factor 2 in the estimates. The true PPO estimator potentially has a slightly smaller bias than estimated.

We consider discrete distributions P and Q . We define the "indistinguishable set" \mathcal{A} where the likelihood ratio is small:

$$\mathcal{A} = \left\{ x : \ln \frac{p(x)}{q(x)} \leq \ln(1 + \epsilon) \right\}.$$

Let $\gamma = \ln(1 + \epsilon)$ be our threshold. Our goal is to bound $P(\mathcal{A})$ using the Kullback–Leibler divergence.

1. The Log-Likelihood Random Variable

We define the random variable L as the log-likelihood ratio evaluated on data generated by P :

$$L = \ln \frac{p(X)}{q(X)}, \quad \text{where } X \sim P.$$

The expectation of L is exactly the KL divergence:

$$\mathbb{E}_P[L] = \sum_x p(x) \ln \frac{p(x)}{q(x)} = D_{\text{KL}}(P\|Q).$$

The event $x \in \mathcal{A}$ corresponds to the event $L \leq \gamma$. If $D_{\text{KL}}(P\|Q) > \gamma$, then the event \mathcal{A} represents a "tail event" (a deviation below the mean).

2. Rigorous Bound: Cantelli's Inequality

Since KL divergence is only the first moment (mean), it cannot bound tail probabilities alone without assuming something about the spread of the distribution. We introduce the variance of the log-likelihood (sometimes called the v-entropy):

$$\sigma^2 = \text{Var}_P(L) = \mathbb{E}_P \left[\left(\ln \frac{p(X)}{q(X)} - D_{\text{KL}}(P\|Q) \right)^2 \right].$$

We use **Cantelli's Inequality** (a one-sided Chebyshev bound), which states that for a random variable L with mean μ and variance σ^2 , and for any $k > 0$:

$$P(L \leq \mu - k) \leq \frac{\sigma^2}{\sigma^2 + k^2}.$$

Setting $\mu = D_{\text{KL}}(P\|Q)$ and $\mu - k = \gamma$, we define the "gap" $k = D_{\text{KL}}(P\|Q) - \ln(1 + \epsilon)$. Assuming $D_{\text{KL}}(P\|Q) > \ln(1 + \epsilon)$, we obtain the bound:

$$P(\mathcal{A}) \leq \frac{\sigma^2}{\sigma^2 + (D_{\text{KL}}(P\|Q) - \ln(1 + \epsilon))^2}.$$

This bound is rigorous and uses KL explicitly, though the decay is polynomial ($1/k^2$) rather than exponential.

3. Exponential Approximation (Local Analysis)

If we are willing to assume that P is "close" to Q (the local regime), the random variable L is approximately Gaussian.

In the local regime, it is a known property of information geometry that the variance is approximately twice the KL divergence:

$$\sigma^2 \approx 2D_{\text{KL}}(P\|Q).$$

If we assume L follows a Gaussian distribution $\mathcal{N}(D_{\text{KL}}, 2D_{\text{KL}})$, we can apply the standard Chernoff bound for a Gaussian variable. The probability of deviating from the mean by k is bounded by:

$$P(L \leq \mu - k) \leq \exp\left(-\frac{k^2}{2\sigma^2}\right).$$

Substituting $k \approx D_{\text{KL}}$ (assuming the threshold $\gamma \approx 0$) and $\sigma^2 \approx 2D_{\text{KL}}$:

$$P(\mathcal{A}) \lesssim \exp\left(-\frac{(D_{\text{KL}})^2}{2(2D_{\text{KL}})}\right) = \exp\left(-\frac{1}{4}D_{\text{KL}}(P\|Q)\right).$$

This explains why KL-based exponential bounds often appear with a factor of 1/4 or 1/2 in the exponent in asymptotic literature.

4. Summary

To use KL divergence for a single-shot bound, we must account for the variance:

1. **Rigorous:** Use Cantelli's inequality. The bound depends on the ratio of the variance σ^2 to the squared KL distance.
2. **Heuristic:** If $P \approx Q$, the probability decays as roughly $\exp(-\frac{1}{4}D_{\text{KL}}(P\|Q))$.

1. The PPO Clipping Mechanism

In Proximal Policy Optimization (PPO), the surrogate objective is defined as:

$$L^{CLIP}(\theta) = \mathbb{E}_t[\min(r_t(\theta)A_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)A_t)]$$

where $r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}$.

The gradient $\nabla_\theta L^{CLIP}$ is non-zero **only** when the ratio is unclipped (or when the advantage sign dictates movement towards the clip, but we focus on the vanishing region). The "Active Region" \mathcal{A} where the policy update is fully active is:

$$\mathcal{A} = \left\{ (s, a) : 1 - \epsilon < \frac{\pi_\theta(a|s)}{\pi_{\theta_{old}}(a|s)} < 1 + \epsilon \right\}.$$

If the probability measure of this set $P(\mathcal{A}) \rightarrow 0$, the gradient contribution vanishes.

2. Log-Ratio Formulation

We analyze the log-ratio random variable $L = \ln r_t(\theta)$. The active region corresponds to:

$$\ln(1 - \epsilon) < L < \ln(1 + \epsilon).$$

Let $\gamma_+ = \ln(1 + \epsilon)$ and $\gamma_- = \ln(1 - \epsilon)$. Since ϵ is small (usually 0.2), this interval is centered near 0.

However, the expected value of L under the current policy π_θ is the KL divergence:

$$\mathbb{E}_{\pi_\theta}[L] = \mathbb{E}_{\pi_\theta} \left[\ln \frac{\pi_\theta}{\pi_{\theta_{old}}} \right] = D_{\text{KL}}(\pi_\theta \| \pi_{\theta_{old}}).$$

Note: We analyze the expectation under π_θ because we are asking "Where does the current policy put its mass relative to the old policy?"

3. The "Gradient Silencing" Bound

Let $\mu = D_{\text{KL}}(\pi_\theta \| \pi_{\theta_{old}})$. As the policy updates, if π_θ becomes very different from $\pi_{\theta_{old}}$, then μ becomes large and positive.

The active region requires $L < \ln(1 + \epsilon)$. If $\mu \gg \ln(1 + \epsilon)$, the active region is far in the "left tail" of the distribution of L .

Using **Cantelli's Inequality** (one-sided Chebyshev):

$$P(L \leq \ln(1 + \epsilon)) = P(L \leq \mu - k) \leq \frac{\sigma^2}{\sigma^2 + k^2}$$

where the gap is:

$$k = D_{\text{KL}}(\pi_\theta \| \pi_{\theta_{old}}) - \ln(1 + \epsilon).$$

And $\sigma^2 = \text{Var}_{\pi_\theta}(L)$.

4. Conclusion: Convergence to Zero

The probability that a sample contributes to the gradient is bounded by:

$$P(\text{Active}) \leq \frac{\text{Var}(L)}{\text{Var}(L) + (D_{\text{KL}}(\pi_\theta \| \pi_{\theta_{old}}) - \ln(1 + \epsilon))^2}$$

As the policies diverge ($D_{\text{KL}} \rightarrow \infty$):

1. The "gap" k grows linearly with the KL divergence.
2. The denominator grows quadratically with the KL divergence.
3. Consequently, $P(\text{Active}) \rightarrow 0$.

This proves that PPO naturally "silences" gradients from the current policy if it drifts too far from the sampling distribution, preventing catastrophic updates based on off-distribution samples.



ACHTUNG, WIE KOMMT REIN, DASS DIE RECHTE SEITE NICHT NEGATIVE IST

Theorem 5.3.8. For bounded advantage and score function we have

$$\|\widetilde{\nabla J(\theta)}\|_\infty \leq A_{max} C(1 + \epsilon) T \left(1 - \frac{\overline{KL}}{C_\epsilon}\right),$$

where $C_\epsilon := \log(1 + \epsilon) + |\log(1 - \epsilon)|$ and as above $\overline{KL} = \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}^{\pi_{old}} \left[\text{KL}(\pi_{old}(\cdot | S_t) \| \pi_\theta(\cdot | S_t)) \right]$.

Proof. Now the trick is that the gradient is zero for the summands in which the ratio is clipped to 1, nothing else depends on θ . Puh, ist da ein Problem mit Differenzierbarkeit??? MACHT ES

UEBERHAUPT SINN, DEN GRADIENTEN RAUSZUEIHEN? GRADIENTEN VERSTEHEN, WENN $\pi_\theta = p^{old}$ ist. WANN KANN MAN REIN UND RAUSZIEHEN?

$$\begin{aligned}
\|\widetilde{\nabla J(\theta)}\|_\infty &\leq \|\nabla_\theta \mathbb{E}^{\pi^{\text{old}}} \left[\sum_{t=0}^{T-1} \gamma^t \frac{\pi^\theta(A_t; S_t)}{\pi^{\text{old}}(A_t; S_t)} \mathbf{1}_{|r_t-1| < \epsilon} A_t^{\pi^{\text{old}}}(S_t, A_t) \right]\|_\infty \\
&= \|\mathbb{E}^{\pi^{\text{old}}} \left[\sum_{t=0}^{T-1} \gamma^t \frac{\pi^\theta(A_t; S_t)}{\pi^{\text{old}}(A_t; S_t)} \mathbf{1}_{|r_t-1| < \epsilon} \nabla_\theta \log(\pi_\theta(A_t, S_t)) A_t^{\pi^{\text{old}}}(S_t, A_t) \right]\|_\infty \\
&\leq A_{\max} C(1 + \epsilon) \sum_{t=1}^{T-1} \gamma^t \mathbb{P}^{\pi^{\text{old}}}(|r_t - 1| < \epsilon) \\
&\leq A_{\max} C(1 + \epsilon) \sum_{t=1}^{T-1} \mathbb{P}^{\pi^{\text{old}}}(|r_t - 1| < \epsilon) \\
&\leq A_{\max} C(1 + \epsilon) \sum_{t=0}^{T-1} \left(1 - \frac{\mathbb{E}^{\pi^{\text{old}}} \left[\text{KL}(\pi_{\text{old}}(\cdot | S_t) \| \pi_\theta(\cdot | S_t)) \right]}{C_\epsilon} \right) \\
&= A_{\max} C(1 + \epsilon) T \left(1 - \frac{\overline{\text{KL}}}{C_\epsilon} \right)
\end{aligned}$$

The last inequality is linearity, using the properly normalised sum is 1. \square

Was wuerde das ganze Bedeuten? Wir wissen, wie nah wir am Gradienten sind und gleichzeitig wissen wir, dass wir irgendwann einfach stehen bleiben, wenn $KL \approx C_\epsilon$.

Impact on gradient norm.

5.4 Effect of Large Mean KL on Gradient Magnitude

We now show that when the *mean* policy KL divergence across rollout states becomes large, the magnitude of the PPO gradient is necessarily small with high probability. Together with the converse (Section ??), this yields a self-consistent stability mechanism: a large policy shift suppresses further gradient updates.

5.4.1 Setup and Definitions

Let $\pi_{\text{old}} = \pi_{\theta_{\text{old}}}$ denote the behavior policy and π_{θ} the updated policy. For each sample (S, A) collected under π_{old} , define the likelihood ratio

$$r(S, A) := \frac{\pi_{\theta}(A|S)}{\pi_{\text{old}}(A|S)}. \quad (5.7)$$

The clipped PPO objective for one sample is

$$L^{\text{clip}}(S, A) = \min\left(r(S, A) \hat{A}(S, A), \text{clip}(r(S, A), 1 - \epsilon, 1 + \epsilon) \hat{A}(S, A)\right), \quad (5.8)$$

and its gradient with respect to θ satisfies

$$g(S, A) := \nabla_{\theta} L^{\text{clip}}(S, A) = \nabla_{\theta} \log \pi_{\theta}(A|S) \tilde{A}(S, A) \quad (5.9)$$

with

$$\tilde{A}(S, A) = \begin{cases} \hat{A}(S, A), & |r(S, A) - 1| \leq \epsilon, \\ 0, & \text{otherwise (appropriate clipping sign)}. \end{cases} \quad (5.10)$$

Hence the gradient is nonzero only for ratio values inside the unclipped interval $[1 - \epsilon, 1 + \epsilon]$.

5.4.2 Mean KL Divergence and Ratio Deviations

For each fixed state s , denote the old and new action distributions $P_s(a) = \pi_{\text{old}}(a|s)$, $Q_s(a) = \pi_{\theta}(a|s)$, and $r_s(a) = Q_s(a)/P_s(a)$. The per-state KL divergence is

$$\text{KL}(s) = \mathbb{E}_{A \sim P_s} \left[\log \frac{P_s(A)}{Q_s(A)} \right] = -\mathbb{E}_{A \sim P_s} [\log r_s(A)]. \quad (5.11)$$

Let

$$\bar{\text{KL}} := \mathbb{E}_{S \sim d^{\pi_{\text{old}}}} [\text{KL}(S)] = -\mathbb{E}_{S, A \sim \pi_{\text{old}}} [\log r(S, A)]$$

be the *mean* KL divergence across rollout states.

Bound on the measure of large ratios. Fix $\epsilon > 0$ and define per-state sets $\Omega_{s,+} = \{a : r_s(a) > 1 + \epsilon\}$ and $\Omega_{s,-} = \{a : r_s(a) < 1 - \epsilon\}$. By Jensen's and monotonicity of the logarithm,

$$\begin{aligned} \text{KL}(s) &= \mathbb{E}_{A \sim P_s} [-\log r_s(A)] \\ &\geq (-\log(1 + \epsilon)) P_s(\Omega_{s,+}) + |\log(1 - \epsilon)| P_s(\Omega_{s,-}) \\ &= C_{\epsilon} P_s(|r_s(A) - 1| > \epsilon), \end{aligned} \quad (5.12)$$

where

$$C_{\epsilon} := \log(1 + \epsilon) + |\log(1 - \epsilon)|. \quad (5.13)$$

Averaging (5.12) over rollout states $S \sim d^{\pi_{\text{old}}}$ gives

$$\mathbb{E}_{S, A \sim \pi_{\text{old}}} [\mathbb{I}\{|r(S, A) - 1| > \epsilon\}] \leq \frac{\bar{\text{KL}}}{C_{\epsilon}}. \quad (5.14)$$

Let

$$p_{\text{clip(avg)}} := \mathbb{E}_{S,A \sim \pi_{\text{old}}} [\mathbb{I}\{|r(S,A) - 1| > \epsilon\}] \quad (5.15)$$

be the fraction of ratio values outside the unclipped region. Equation (5.14) guarantees that

$$p_{\text{clip(avg)}} \geq \frac{\bar{\text{KL}}}{C_\epsilon}. \quad (5.16)$$

5.4.3 Expected Gradient Magnitude

Assume uniform bounds on the log-policy gradients and advantages:

$$\|\nabla_\theta \log \pi_\theta(A|S)\| \leq G_{\max}, \quad |\hat{A}(S,A)| \leq A_{\max}.$$

Then the sample gradient magnitude is bounded as

$$\|g(S,A)\| \leq G_{\max} A_{\max} \mathbb{I}\{|r(S,A) - 1| \leq \epsilon\}. \quad (5.17)$$

Taking expectation with respect to $(S,A) \sim \pi_{\text{old}}$,

$$\begin{aligned} \mathbb{E}\|g(S,A)\| &\leq G_{\max} A_{\max} \Pr(|r(S,A) - 1| \leq \epsilon) \\ &= G_{\max} A_{\max} (1 - p_{\text{clip(avg)}}) \\ &\leq G_{\max} A_{\max} \left(1 - \frac{\bar{\text{KL}}}{C_\epsilon}\right). \end{aligned} \quad (5.18)$$

Thus, as the mean KL divergence $\bar{\text{KL}}$ increases, the expected gradient magnitude decreases linearly, vanishing when $\bar{\text{KL}} \approx C_\epsilon$.

5.4.4 High-Probability Small-Gradient Bound

The random variable $\|g(S,A)\|$ lies in $[0, G_{\max} A_{\max}]$. For any $t > 0$, Markov's inequality yields

$$\Pr(\|g(S,A)\| \geq t) \leq \frac{\mathbb{E}\|g(S,A)\|}{t}. \quad (5.19)$$

Choosing $t = G_{\max} A_{\max} \left(1 - \frac{\bar{\text{KL}}}{2C_\epsilon}\right)$ and substituting (5.18) gives

$$\begin{aligned} \Pr\left(\|g(S,A)\| > G_{\max} A_{\max} \left(1 - \frac{\bar{\text{KL}}}{2C_\epsilon}\right)\right) &\leq \frac{1 - \bar{\text{KL}}/C_\epsilon}{1 - \bar{\text{KL}}/(2C_\epsilon)} \\ &\leq \frac{\bar{\text{KL}}}{C_\epsilon}. \end{aligned} \quad (5.20)$$

Therefore, with probability at least $1 - \bar{\text{KL}}/C_\epsilon$, the sample gradient norm satisfies

$$\|g(S,A)\| \leq G_{\max} A_{\max} \left(1 - \frac{\bar{\text{KL}}}{2C_\epsilon}\right). \quad (5.21)$$

5.4.5 Conclusion

Equations (5.18)–(5.21) formally show that when the mean policy KL $\bar{\text{KL}}$ increases, an increasing fraction of state–action pairs fall into the clipped region (probability at least $\bar{\text{KL}}/C_\epsilon$). For those samples, the surrogate objective is flat, and the gradient vanishes. Hence the expected and typical (high-probability) gradient magnitudes decrease linearly with $\bar{\text{KL}}$. This yields the qualitative mechanism:

small KL \Rightarrow active update, large KL \Rightarrow gradient suppression.

The PPO algorithm therefore maintains policy stability even without an explicit trust-region constraint.

5.4.6 Variance of the PPO surrogate gradient with GAE estimator

We now bound the variance of the surrogate gradient under following assumptions:

1. **Bounded rewards and values:** $|R_t| \leq R_{\max}$, $|V(S_t)| \leq V_{\max}$. Hence TD-errors $\delta_t = R_t + \gamma V(S_{t+1}) - V(S_t)$ satisfy $|\delta_t| \leq \delta_{\max} = R_{\max} + 2\gamma V_{\max}$.
2. **Bounded score function:** $\|\nabla_{\theta} \log \pi_{\theta}(a|s)\|_{\infty} \leq G_{\max}$.
3. **KL constraint:** $\mathbb{E}_s[D_{\text{KL}}(\pi_{\text{old}}(\cdot|s) \|\pi_{\theta}(\cdot|s))] \leq \varepsilon$, giving $\mathbb{E}[r_t^m] \leq e^{m\varepsilon}$ for all m .

ACHTUNG: PROBLEM MIT KL KOMMT NOCH

Define the λ -return (GAE) and per-step gradient term

$$\hat{A}_t^{(\lambda)} = \sum_{l=0}^{T-t-1} (\gamma\lambda)^l \delta_{t+l}, \quad g_t = \gamma^t r_t \nabla_{\theta} \log \pi_{\theta}(A_t|S_t) \hat{A}_t^{(\lambda)}.$$

Let $G_T = \sum_{t=0}^{T-1} g_t$. Our aim is to bound $\text{Var}[G_T]$.

Let us first show how to deal with the time-correlations that are induced by the advantage estimators:



Lemma 5.4.1. For any $0 \leq t < T$ and $k \geq 0$ with $t+k < T$,

$$|\hat{A}_t^{(\lambda)}, \hat{A}_{t+k}^{(\lambda)}| \leq \frac{4\delta_{\max}^2}{1 - (\gamma\lambda)^2} (\gamma\lambda)^{2k}.$$

Proof. First recall the definition of GAE

$$\hat{A}_t = \sum_{l=0}^{T-t-1} (\gamma\lambda)^l \delta_{t+l}, \quad \hat{A}_{t+k} = \sum_{l'=0}^{T-t-k-1} (\gamma\lambda)^{l'} \delta_{t+k+l'},$$

so that

$$\hat{A}_t \hat{A}_{t+k} = \sum_{l, l'} (\gamma\lambda)^{l+l'} \delta_{t+l} \delta_{t+k+l'}.$$

Because $|\delta_t| \leq \delta_{\max}$ by assumption, we get

$$|\hat{A}_t \hat{A}_{t+k}| \leq 2\delta_{\max}^2 \sum_{l=k}^{T-t-1} \sum_{l'=0}^{T-t-k-1} (\gamma\lambda)^{l+l'} \leq \frac{4\delta_{\max}^2}{1 - (\gamma\lambda)^2} (\gamma\lambda)^{2k}.$$

□



Theorem 5.4.2. (Finite-horizon variance bound for one rollout)

Define

$$C_{\star} = \frac{8e^{2\varepsilon} G_{\max}^2 \delta_{\max}^2}{(1 - \gamma^2)[1 - \gamma(\gamma\lambda)^2][1 - (\gamma\lambda)^2]} (1 - \gamma^{2T}).$$

Then for a single trajectory, $\text{Tr}[\text{Var}[G_T]] \leq C_{\star}$.

Proof. Compute

$$\text{Var}[G_T] = \sum_{t, t'} \text{Cov}(g_t, g_{t'}).$$

Using the bounds on score functions (G_{\max}), ratios ($\mathbb{E}[r_t^4 r_{t'}^4] \leq e^{4\varepsilon}$), and advantages ($|\text{Cov}(\hat{A}_t, \hat{A}_{t'})| \leq \frac{4\delta_{\max}^2}{1 - (\gamma\lambda)^2} (\gamma\lambda)^{2|t-t'|}$), one obtains

$$|\text{Cov}(g_t, g_{t+k})| \leq 2e^{2\varepsilon} G_{\max}^2 \frac{\delta_{\max}^2}{1 - (\gamma\lambda)^2} \gamma^{t+t+k} (\gamma\lambda)^{2k}.$$

Summing all (t, k) pairs yields C_{\star} .

□

!!! Ratios ist falsch!!!

We now turn to the mini-batch situation in which there is a finite set of N rollouts of length T and $M \leq NT$ transitions are sampled randomly from the transition buffer. The idea is the following. Transitions from independent rollouts are independent, thus do not contribute correlations to the variance. For two correlations from the same rollout we can easily estimate the correlations. Thus, we need to count how many transitions we have from the same rollouts and integrate out their time.



Theorem 5.4.3. Suppose we have N rollouts of length T and sample the gradient with M randomly chosen transitions from the buffer. Then

$$\text{Tr} [\text{Var}[\tilde{G}_M]] \leq \frac{C_T^*}{MNT^2} \left(1 - \frac{M}{NT}\right)$$

where

$$C_T^* := \frac{8e^{2\varepsilon} G_{\max}^2 \delta_{\max}^2}{(1 - \gamma^2) [1 - \gamma(\gamma\lambda)^2] [1 - (\gamma\lambda)^{23}] [1 - \gamma^{2T}]}$$

Warning: The constant depends on T .

Require: Policy (actor) π_θ ; value function (critic) V_ϕ ; number of cycles C ; transitions per cycle N ; number of passes per cycle E ; number of minibatches per pass K ; learning rates α_π, α_V ; clipping range ϵ

```

1: Initialize policy parameters  $\theta_0$  and value function parameters  $\phi_0$ 
2: for cycle = 1, 2, ...,  $C$  do
3:   Set behavior policy  $\mu = \pi_{\theta_{\text{old}}}$ 
4:   Initialize empty transition buffer  $\mathcal{D} = \emptyset$ 
5:   — Collection of transitions under current behavior policy —
6:   while  $|\mathcal{D}| < N$  do
7:     Observe current state  $s_t$ 
8:     Sample action  $a_t \sim \mu(\cdot|s_t)$  and execute in environment
9:     Observe reward  $r_t$  and next state  $s_{t+1}$ 
10:    Store  $(s_t, a_t, r_t, s_{t+1})$  in  $\mathcal{D}$ 
11:   end while
12:   — Critic (Value Function) Update —
13:   for epoch = 1, ...,  $E_V$  do
14:     Randomly shuffle  $\mathcal{D}$ 
15:     Split  $\mathcal{D}$  into  $K$  critic minibatches  $\{\mathcal{B}_1^V, \dots, \mathcal{B}_K^V\}$ 
16:     for  $k = 1, \dots, K$  do
17:       Compute TD targets:

```

$$y_t = r_t + \gamma V_\phi(s_{t+1})$$

```

18:       Update critic by minimizing value loss:

```

$$\phi \leftarrow \phi - \alpha_V \nabla_\phi \frac{1}{|\mathcal{B}_k^V|} \sum_{(s_t, a_t, r_t, s_{t+1}) \in \mathcal{B}_k^V} (V_\phi(s_t) - y_t)^2$$

```

19:     end for
20:   end for
21:   — Advantage Estimation —
22:   for each transition  $(s_t, a_t, r_t, s_{t+1}) \in \mathcal{D}$  do
23:     Compute:

```

$$\hat{A}_t = r_t + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$

```

24:   end for
25:   — Policy (Actor) Update with PPO Clipping —
26:   for epoch = 1, ...,  $E$  do
27:     Randomly shuffle  $\mathcal{D}$ 
28:     Split  $\mathcal{D}$  into  $K$  actor minibatches  $\{\mathcal{B}_1^A, \dots, \mathcal{B}_K^A\}$ 
29:     for  $k = 1, \dots, K$  do
30:       Compute per-transition importance ratio (using only current step):

```

$$r_t = \frac{\pi_\theta(a_t|s_t)}{\mu(a_t|s_t)}$$

```

31:       Compute clipped policy gradient estimator:

```

$$\hat{g}_k = \frac{1}{|\mathcal{B}_k^A|} \sum_{(s_t, a_t, r_t, s_{t+1}) \in \mathcal{B}_k^A} \nabla_\theta \log \pi_\theta(a_t|s_t) \min(r_t \hat{A}_t, \text{clip}(r_t, 1-\epsilon, 1+\epsilon) \hat{A}_t)$$

```

32:       Update actor parameters:

```

$$\theta \leftarrow \theta + \alpha_\pi \hat{g}_k$$

```

33:     end for
34:   end for
35:   Set  $\theta_{\text{old}} \leftarrow \theta$ 
36: end for

```

Algorithm 21: Proximal Policy Optimization (PPO) with Transition Buffer and Multi-Batching

And here is the early stopping version without clipping. This is what we analyse so far.

Require: Policy (actor) π_θ ; value function (critic) V_ϕ ; number of cycles C ; transitions per cycle N ; number of passes per cycle (actor) E ; number of minibatches per pass K ; learning rates α_π, α_V ; KL threshold δ_{KL}

```

1: Initialize policy parameters  $\theta_0$  and value function parameters  $\phi_0$ 
2: for cycle = 1, 2, ...,  $C$  do
3:   Set behavior policy  $\mu = \pi_{\theta_{\text{old}}}$ 
4:   Initialize empty transition buffer  $\mathcal{D} = \emptyset$ 
5:   — Collection of transitions under current behavior policy —
6:   while  $|\mathcal{D}| < N$  do
7:     Observe current state  $s_t$ 
8:     Sample action  $a_t \sim \mu(\cdot|s_t)$  and execute in environment
9:     Observe reward  $r_t$  and next state  $s_{t+1}$ 
10:    Store  $(s_t, a_t, r_t, s_{t+1})$  in  $\mathcal{D}$ 
11:  end while
12:  — Critic (Value Function) Update —
13:  for epoch = 1, ...,  $E_V$  do
14:    Randomly shuffle  $\mathcal{D}$ 
15:    Split  $\mathcal{D}$  into  $K$  critic minibatches  $\{\mathcal{B}_1^V, \dots, \mathcal{B}_K^V\}$ 
16:    for  $k = 1, \dots, K$  do
17:      Compute TD targets:  $y_t = r_t + \gamma V_\phi(s_{t+1})$ 
18:      Update critic parameters:

```

$$\phi \leftarrow \phi - \alpha_V \nabla_\phi \frac{1}{|\mathcal{B}_k^V|} \sum_{(s_t, a_t, r_t, s_{t+1}) \in \mathcal{B}_k^V} (V_\phi(s_t) - y_t)^2$$

```

19:    end for
20:  end for
21:  — Advantage Estimation —
22:  for each transition  $(s_t, a_t, r_t, s_{t+1}) \in \mathcal{D}$  do
23:    Compute advantage  $\hat{A}_t(s_t, a_t) = r_t + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$ 
24:  end for
25:  — Policy (Actor) Update with KL Early Stopping —
26:  Initialize flag stop_updates  $\leftarrow$  False
27:  for epoch = 1, ...,  $E$  do
28:    if stop_updates = True then
29:      break
30:    end if
31:    Randomly shuffle  $\mathcal{D}$ , Split  $\mathcal{D}$  into  $K$  actor minibatches  $\{\mathcal{B}_1^A, \dots, \mathcal{B}_K^A\}$ 
32:    Initialize variable avg_KL  $\leftarrow$  0
33:    for  $k = 1, \dots, K$  do
34:      Compute policy gradient estimator (no clipping):

```

$$\hat{g}_k = \frac{1}{|\mathcal{B}_k^A|} \sum_{(s_t, a_t, r_t, s_{t+1}) \in \mathcal{B}_k^A} \frac{\nabla \pi_\theta(a_t|s_t)}{\mu(a_t|s_t)} \hat{A}_t(a_t, s_t)$$

```

35:      Update actor:  $\theta \leftarrow \theta + \alpha_\pi \hat{g}_k$ 
36:      Compute approximate KL between new and old policies on minibatch:

```

$$\text{KL}_k = \frac{1}{|\mathcal{B}_k^A|} \sum_{(s_t, a_t) \in \mathcal{B}_k^A} \left[\log \mu(a_t|s_t) - \log \pi_\theta(a_t|s_t) \right]$$

```

37:      Accumulate avg_KL  $\leftarrow$  avg_KL +  $\text{KL}_k$ 
38:    end for
39:    Compute mean KL across all minibatches: avg_KL  $\leftarrow$   $\frac{\text{avg\_KL}}{K}$ 
40:    if avg_KL >  $\delta_{\text{KL}}$  then
41:      Set stop_updates  $\leftarrow$  True
42:    break
43:  end if
44: end for
45:  Set  $\theta_{\text{old}} \leftarrow \theta$ 
46: end for

```

Algorithm 22: Proximal Policy Optimization (PPO) with KL-Based Early Stopping and

This is the algorithm that we analyse, but with a better advantage estimator.

5.6 Theory

5.6.1 Improved GAE

Explain forwards and backwards TD. Then explain that in finite time we get a bias because we should condition the geometric to be smaller than T. Write down unbiased version, and also the stopping time version. Try to prove something on the variance.

5.6.2 Convergence

Reminder: What we know on the true value function J

We know smoothness, PL under strong assumptions.

Bounding the gradients

In what follows we consider the true value function

$$J(\theta) = \mathbb{E}^{\pi_\theta} \left[\sum_{t=0}^{T-1} \gamma^t R_t \right]$$

with the gradient

$$\begin{aligned} \nabla J(\theta) &= \mathbb{E}^{\pi_\theta} \left[\sum_{t=0}^T \gamma^t \nabla \log \pi_\theta(a_t, s_t) A_t(s_t, a_t) \right] \\ &= \mathbb{E}^\mu \left[\sum_{t=0}^T \prod_{i=0}^t \frac{\pi_\theta(a_i, s_i)}{\mu(a_i, s_i)} \gamma^t \nabla \log \pi_\theta(a_t, s_t) A_t(s_t, a_t) \right] \\ &= \mathbb{E}^\mu \left[\sum_{t=0}^T \prod_{i=0}^t \frac{\nabla \pi_\theta(a_i, s_i)}{\mu(a_i, s_i)} \gamma^t A_t(s_t, a_t) \right] \end{aligned}$$

SK: The last equality should not be correct. Doesn't matter since we throwing it away anyways.

and the surrogate gradient (throwing away factors in the product)

$$\widetilde{J}(\theta) = \mathbb{E}^\mu \left[\sum_{t=0}^T \frac{\nabla \pi_\theta(a_t, s_t)}{\mu(a_t, s_t)} \gamma^t A_t(s_t, a_t) \right]$$

Recall: The reason we use the surrogate is that we can produce samples based on sampled transitions instead of entire rollouts. The sum becomes part of the sampling.

Chapter 6

Reinforcement learning with function approximation

Tabular methods generally assumed small state- and action-spaces. Small in the sense that in principle all Q -values $Q(s, a)$ could be considered separately in a table. In reality the number of states/actions is typically way too big. Additionally, it seems unpractical to estimate (learn) all Q -values separately as small changes in states might only have little impact on $Q(s, a)$. For instance, if s is a picture such as the state of an Atari game. In that case improved learning of $Q(s, a)$ should also effect $Q(s', a)$ for s' similar (in some sense) to s . As an illustration one might think of a function $f : \mathbb{N}^d \rightarrow \mathbb{R}$ that should be learned in some learning context. If $f(n) = a \cdot n + b$ is (affine) linear it would be completely unreasonable to learn all values $f(n)$ separately, only the vectors a and b would have to be estimated. If linear functions are a good approximation to the true function f then it might still be much more reasonable to estimate two parameters for an approximation instead of approximating all true values $f(n)$. Exactly this happens in RL with function approximation, the state-value function V or the state-action value function Q is approximated using a parametric class of functions V_w (resp. Q_w) and instead of finding estimates $\hat{V}(s)$ for all $s \in S$ (or $\hat{Q}(s, a)$ for all s, a) a parameter-vector w is estimated that best approximates $V(s)$ for all s (or $Q(s, a)$ for all s, a). In this chapter we will mostly deal with linear function approximation, in practice V_w or Q_w are neural network functions and w is the vector with all weights. As in dynamic programming there are different possible tasks:

- policy evaluation with function approximation, i.e. find V_w from the approximating parametric class that best approximates V^π (or Q_w that best approximates Q^π),
- policy iteration with function approximation, i.e. alternate between policy evaluation with function approximation and policy improvement, which can either be model-based or model-free using samples,
- value iteration with function approximation or a sample based Q -learning with function approximation,
- policy gradient with function approximation.

The good news is that even though it's complicated to make it work, in many settings dynamic programming or policy gradient with function approximation works. Bad news are lack of mathematical understanding. In very restricted settings proofs can be given but for the most relevant use of function approximation there are still large gaps in the understanding.

6.1 Policy evaluation with function approximation

In this section we deal with the evaluation of a fixed policy π , i.e. the computation of the state-value function V^π and the state-action-value function Q^π . Of course, we already know

exact and approximative tabular solution methods that use the theory of dynamic programming and stochastic approximation to obtain estimates of V^π and Q^π .

The aim of this section is to adapt the algorithms from the tabular case to obtain more general methods that do not require us to store estimates of every entry of V^π and Q^π . Estimates of V^π and Q^π will not be stored as lookup-tables (or vectors), where estimates for each of the entries of $V^\pi(s)$ and/or $Q^\pi(a, s)$ individually. Instead, we will use a parametrized family $\{f_w : w \in \mathbb{R}^d\}$. Typically the number d of parameters will be much smaller than the number of states (or state-action pairs) which means that we will have to tune far fewer parameters than in our previous tabular algorithms. While this is definitely an advantage, we cannot expect convergence to the value functions V^π and Q^π once approximations are introduced. One obvious reason for that is that they may not be within the set of functions that we can represent exactly using our chosen parametrization. In this new setting, the following questions appear naturally:

- (i) Does a given algorithm converge to something?
- (ii) If the algorithm does convergence, is the limit close (in some sense) to the true value function V^π (resp. Q^π)?
- (iii) If the algorithm does not converge, does it at least oscillate within some small region around V^π (resp. Q^π)?

6.1.1 Function approximation

We will broadly classify the parametrized approximation architectures as either linear or non-linear. Linear architectures approximate the state-value function as follows:

$$f_w(s) = w \cdot \phi(s) = \sum_{i=1}^d w_i \phi_i(s),$$

where $\phi : \mathcal{S} \rightarrow \mathbb{R}^d$ and $\phi_i(s)$ denotes the i -th entry of $\phi(s)$. Here ϕ is a fixed mapping from \mathcal{S} to \mathbb{R}^d called the feature extractor, ϕ should extract the most important features of a state. Feature vectors can be anything, a good choice for concrete problems is clearly a tricky task. Essentially, using feature extraction the state-space is reduced to a subset of \mathbb{R}^d and we approximate a function $\mathbb{R}^d \rightarrow \mathbb{R}$ by linear functions only. The linear architecture for Q^π is entirely analogous:

$$f_w(s, a) = w \cdot \phi(s, a) = \sum_{i=1}^d w_i \phi_i(s, a),$$

with $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$. In this section we will only deal with policy evaluation for the state-value function, the action-values can be handled similarly as the Bellman operator is essentially the same. As we will see soon, convergence of many algorithms can only be guaranteed in the case of linear approximation architectures. Still, non-linear parametrizations like neural network functions are very popular in practical applications of RL (in this particular example, w would denote the weights and biases of the network). Fortunately, linear methods still allow for a good deal of flexibility as we will discuss below.

Example 6.1.1. Let us partition \mathcal{S} into m disjoint subsets \mathcal{S}_i , i.e. $\mathcal{S} = \dot{\bigcup}_{i=1}^m \mathcal{S}_i$ with $\mathcal{S}_i \cap \mathcal{S}_j = \emptyset$ for all $i \neq j$. Define

$$\phi : \mathcal{S} \rightarrow \mathbb{R}^m \text{ such that } \phi_i : s \mapsto \mathbf{1}_{\mathcal{S}_i}(s).$$

Then states from a fixed partition set are treated equally and the linear function approximation simplifies to $f_w(s) = w_i$ for all $s \in \mathcal{S}_i$. This particularly simple linear function approximation is referred to as state aggregation. State aggregation can for instance be useful in a computer game where many of the possible images on the screen may have near identical state values because many features (for instance the sun in the background of super mario) of the image are irrelevant for the game.

It is not surprising that many rigorous results exist for state aggregation. In a way it first simplifies the problem and then treats the remaining problem as a tabular problem.



Tabular reinforcement learning can be seen as reinforcement learning with function approximation. Tabular is nothing but state aggregation with the finest partition $\mathcal{S}_i = \{s_i\}$.

Example 6.1.2. From the more familiar tasks of regression and interpolation we already that polynomials are good at approximating real-valued functions. It is not so surprising, then, that polynomials can be used in tasks where the states can be expressed as numbers, e.g. positions or angles in a robotics task. Suppose, for example, $\mathcal{S} \subseteq \mathbb{R}^2$ and let

$$\phi(s) = (1, s_1, s_2, s_1 s_2) \quad \text{so that} \quad f_w(s) = w_1 + w_2 s_1 + w_3 s_2 + w_4 s_1 s_2,$$

where the last term takes into account interactions between the dimensions. Of course, we could also use higher-dimensional polynomials to model more complex interactions, say $\phi(s) = (1, s_1, s_2, s_1 s_2, s_1^2, s_1^2 s_2^3, s_1^3, s_2^3)$ if we expect V^π to behave like a cubic function. This already highlights how we should incorporate prior knowledge about the nature of the task into the construction of useful feature vectors.

6.1.2 Sample based policy evaluation with function approximation

All function approximation algorithms must use some notion of error between functions to evaluate the quality of the function approximation. To give simple algorithms the notion should be useful and simple. The most classical one is the weighted mean-squared error:



Definition 6.1.3. Let μ a (discrete) probability measure on \mathcal{S} and $\{f_w : w \in \mathbb{R}^d\}$ a set of approximation architectures. We define the mean square value error as

$$E(w) = \frac{1}{2} \sum_{s \in \mathcal{S}} \mu(s) (V^\pi(s) - f_w(s))^2 = \frac{1}{2} \mathbb{E}_{S \sim \mu} [(V^\pi(S) - f_w(S))^2].$$

The weights μ are supposed to indicate the importance of the respective states, thus, a natural choice is the (discounted) state visitation measure. Ideally, we want to find w such that $E(w)$ attains its global minimum. If our parametrized function class contains the true value function the value of the minimum will, of course, be 0. But there is no reason that should be the case, and typically it won't. It should be noted that it is not clear that E is the right performance objective for reinforcement learning. Our reason to even learn the value function is (at least most of the time) to obtain a better policy. However, the best value function approximation need not be the one that minimizes E . Even so, it is not obvious what a sensible alternative goal would be.

The following discussion will be kept somewhat informal to illustrate the important points, a formal treatment follows below. Least-squares problem as minimising E appear in many contexts, solution methods are quite standard. If the approximation architecture is linear the minimisation problem can be solved explicitly but involves non-trivial matrix inversions. In general, as long as the approximation architecture is differentiable the minimum (at least a local minimum) can be obtained using gradient descent on the parameters as the gradient can be obtained easily from the chain rule. This yields

$$\begin{aligned} \nabla_w E(w) &= \frac{1}{2} \sum_{s \in \mathcal{S}} \mu(s) \nabla_w (V^\pi(s) - f_w(s))^2 \\ &= - \sum_{s \in \mathcal{S}} \mu(s) (V^\pi(s) - f_w(s)) \nabla_w f_w(s) \\ &= \mathbb{E}_{S \sim \mu} [(f_w(S) - V^\pi(S)) \nabla_w f_w(S)]. \end{aligned} \tag{6.1}$$

What we see is that the gradient takes precisely the form needed for stochastic gradient algorithms (SGD). In order for SGD algorithms to work, a minimal necessary requirement is that we should use unbiased estimates of the gradient of the potential function. For the sake of argument, assume just for now that we somehow have access to the value function. As suggested by the calculations above, a sample based SGD algorithm for calculating the best function approximation would be given by

$$w_{n+1} = w_n + \alpha_n (V^\pi(s_{n+1}) - f_{w_n}(s_{n+1})) \nabla_w f_{w_n}(s_{n+1}), \quad w_0 \in \mathbb{R}^d,$$

where the states s_1, \dots are sampled independently according to μ and step-sizes that satisfy decay rates depending on the problem (most of the time the Robbins-Monro conditions work). Obviously, we do not have access to the value function in practice. Therefore, we will have to replace $V^\pi(s_n)$ by estimates U_t which leads to the following general algorithm:

$$w_{n+1} = w_n + \alpha_n (U_{n+1} - f_{w_n}(s_{n+1})) \nabla_w f_{w_n}(s_{n+1}), \quad w_0 \in \mathbb{R}^d. \quad (6.2)$$

If f is a linear architecture, then the algorithm simplifies to

$$w_{n+1} = w_n + \alpha_n (U_{n+1} - w_n \cdot \phi(s_{n+1})) \phi(s_{n+1}), \quad w_0 \in \mathbb{R}^d.$$

We can get creative here and chose U_t as one of the estimates of the value function from sample-based dynamic programming. In theory the simplest choice of U is a Monte Carlo estimator where independent samples are generated:

Example 6.1.4. Setting $U_n = \sum_{t=0}^{\infty} \gamma^t R(S_{t+1}^n, A_{t+1}^n)$ with rollouts $(S_t^n, A_t^n)_{t \geq 0}$ sampled according to the policy π with starting state $S_0^n = s_n$, the Monte Carlo version of policy evaluation with function approximation is

$$w_{n+1} = w_n + \alpha_n \left(\underbrace{\sum_{t=0}^{\infty} \gamma^t R(S_{t+1}^{n+1}, A_{t+1}^{n+1})}_{\text{Theoretical Monte Carlo estimator for } V^\pi(s_{n+1})} - f_{w_n}(s_{n+1}) \right) \nabla_w f_{w_n}(s_{n+1}), \quad w_0 \in \mathbb{R}^d.$$

This is somewhat a double stochastic gradient scheme, sampling independently from states and total reward distributions. We will see below that for linear function architectures convergence follows directly from standard stochastic gradient theorems.

Similarly to sample-based dynamic programming the Monte Carlo approach is inefficient, there is no bootstrapping of rollouts. Rollouts are only used once, with the advantage that there is a lot of independence that facilitates convergence guarantees. Reusing rollouts by including the already estimated value function yields a one-step temporal difference approach (TD(0)). Recall that in TD(0) a value function estimation is improved by only resampling one times-step and then reusing the old estimate of V^π :

$$V^\pi(s) \leftarrow R(s, a) + \gamma V^\pi(s'),$$

where $a \sim \pi(\cdot; s)$ and $s' \sim p(\cdot; s, a)$ is a one-step sample. While this sounds plausible, a formal justification is through the Bellman operator written in terms of expectations. Since V^π is unknown it is plausible to replace V^π by the best-known current approximation of V^π . The corresponding update formula is as follows:

Example 6.1.5. The TD(0) update rule for policy evaluation with function approximation is

$$w_{n+1} = w_n + \alpha_n \left(\underbrace{R(s_{n+1}, a) + \gamma f_{w_n}(s')}_{\text{TD(0) estimator for } V^\pi(s_{n+1})} - f_{w_n}(s_{n+1}) \right) \nabla_w f_{w_n}(s_{n+1}), \quad w_0 \in \mathbb{R}^d, \quad (6.3)$$

where $a \sim \pi(\cdot; s_{n+1})$ and $s' \sim p(\cdot; s_{n+1}, a)$.

Analogously to sample based dynamic programming one-step temporal differences are also replaced by n -steps or a TD(λ) version:

Example 6.1.6. The N -step temporal difference update rule for policy evaluation with function approximation is

$$w_{n+1} = w_n + \alpha_n \underbrace{\left(\sum_{t=0}^{N-1} \gamma^t R(S_t^{n+1}, A_t^{n+1}) + \gamma^N f_{w_n}(S_N^n) - f_{w_n}(s_{n+1}) \right)}_{N\text{-step TD estimator for } V^\pi(s_{n+1})} \nabla_w f_{w_n}(s_{n+1}), \quad (6.4)$$

where $A_0^n, S_0^n, \dots, S_N^n$ are MDP rollouts started in s_n .



Using approximators f_w to estimate V^π destroys the gradient descent property of the algorithm. The update rule is not anymore obtained from the chain rule of the least-square error! In reinforcement learning such algorithms are often called semi-gradient, they look like gradient descent updates but aren't. As a matter of fact, most algorithms do not converge, and if so, the analysis does not follow from gradient descent theorems but from the very specific form of the iteration.

Up to now the discussion was informal and, in fact, problematic. There is a big problem with the chain rule. If, as we did for the TD(0) update, the value function is replaced by a function approximation f_w , then the dependence of w yields a different derivative and the gradient descent update is not (6.4). Nonetheless, the algorithm could still be reasonable and (at least for linear function approximation) actually is.



Theorem 6.1.7. (Stochastic Gradient Descent)

Suppose $(\Omega, \mathcal{A}, \mathcal{F}_n)$ is a filtered probability with \mathcal{F}_{n+1} -measurable errors ε_n and define the recursion

$$r_{n+1} = r_n - \alpha_n (\nabla g(r_n) + \varepsilon_n)$$

where $g : \mathbb{R}^n \rightarrow \mathbb{R}_{\geq 0}$ is a cost function. We assume that the stepsizes α_n are nonnegative, \mathcal{F}_n -measurable and satisfy

$$\sum_{n=0}^{\infty} \alpha_n = \infty \quad \text{and} \quad \sum_{n=0}^{\infty} \alpha_n^2 < \infty.$$

g is assumed to be L -smooth, i.e. g is differentiable and the gradient ∇g is L -Lipschitz continuous. In addition, we assume that the noise terms are \mathcal{F}_{n+1} -measurable and conditionally unbiased in the sense that

$$\mathbb{E}[\varepsilon_n | \mathcal{F}_n] = 0$$

and have sufficiently small variances

$$\mathbb{E}[\|\varepsilon_n\|^2 | \mathcal{F}_n] \leq A + B \|\nabla g(r_n)\|^2$$

for some constants $A, B \in \mathbb{R}$. Under these assumptions the sequence $(g(r_n))_{n \in \mathbb{N}}$ converges almost surely and $\lim_{n \rightarrow \infty} \nabla g(r_n) = 0$. Moreover, every limit point of $(r_n)_{n \in \mathbb{N}}$ is a stationary point of g .

Note that the theorem only requires L -smoothness on g but no convexity assumption. Hence, the result is only convergence to a critical point but not convergence to a global minimum¹. To prove at least something we can show without much effort the convergence of the gradient descent based value prediction algorithm with linear function approximation in the case the value function is estimated without bias (e.g. with Monte Carlo estimation):

¹The proof of this theorem can be found in Bertsekas' and Tsitsiklis' book *Neuro-Dynamic Programming*



Theorem 6.1.8. Suppose $(w_n)_{n \in \mathbb{N}}$ is the sequence from Example 6.1.4 initialised in some $w_0 \in \mathbb{R}^d$, i.e. approximate policy evaluation with Monte Carlo estimates. If features and rewards are bounded and the approximation architecture is linear, then $(w_n)_{n \in \mathbb{N}}$ converges almost surely to the global minimum of E .

The linear least-squared problem can also be solved explicitly by matrix inversion. In contrast to gradient descent methods (with provable convergence or not) the direct approach does not extend to non-linear approximation architectures. To give a first impression we discuss a simple convergence proof for the simple linear situation with Monte Carlo estimates.

Proof. Since the f_w are assumed linear the mean-squared error is a quadratic function. Thus, it is L -smooth (the derivative is linear) and has a unique global minimum (which is the only critical point). Furthermore, the gradient is

$$\nabla_w E(w) \stackrel{(6.1)}{=} \mathbb{E}_\mu [(f_w(S) - V^\pi(S)) \nabla_w f_w(S)].$$

To link the update scheme formally to the stochastic gradient scheme we write

$$w_{n+1} = w_n + \alpha_n [(U_{n+1} - f_w(S_{n+1})) \nabla_w f_w(S_{n+1})] = w_n + \alpha_n [\nabla_w E(w_n) + \varepsilon_n],$$

with

$$\varepsilon_n := (f_{w_n}(S_{n+1}) - U_{n+1}) \nabla_w f_{w_n}(S_{n+1}) - \nabla_w E(w_n),$$

where U_n is a Monte Carlo sample of $V^\pi(S_n)$ generated independently of $S_n \sim \mu$. Next, denote by \mathcal{F}_n the filtration generated by all random variables of the algorithm to produce w_n . These are all α_k with $k < n$ and all rollouts (S^k, A^k) that result in U_k with $k < n$. Then U_n is independent of \mathcal{F}_n and \mathcal{F}_{n+1} -measurable (thus, ε_n is \mathcal{F}_{n+1} -measurable). Furthermore, S_n is independent of \mathcal{F}_n . Hence, the tower property and measurability/independence properties of conditional expectation yield

$$\begin{aligned} & \mathbb{E}[\varepsilon_n \mid \mathcal{F}_n] \\ &= \mathbb{E}[f_{w_n}(S_{n+1}) \nabla_w f_{w_n}(S_{n+1}) \mid \mathcal{F}_n] - \mathbb{E}[U_{n+1} \nabla_w f_{w_n}(S_{n+1}) \mid \mathcal{F}_n] - \nabla_w E(w_n) \\ &= \mathbb{E}[f_{w_n}(S_{n+1}) \nabla_w f_{w_n}(S_{n+1}) \mid \mathcal{F}_n] - \mathbb{E}[\mathbb{E}[U_{n+1} \nabla_w f_{w_n}(S_{n+1}) \mid \sigma(\mathcal{F}_n, S_{n+1})] \mid \mathcal{F}_n] - \nabla_w E(w_n) \\ &= \mathbb{E}[f_{w_n}(S_{n+1}) \nabla_w f_{w_n}(S_{n+1}) \mid \mathcal{F}_n] - \mathbb{E}[\mathbb{E}[U_{n+1} \mid \sigma(\mathcal{F}_n, S_{n+1})] \nabla_w f_{w_n}(S_{n+1}) \mid \mathcal{F}_n] - \nabla_w E(w_n) \\ &= \mathbb{E}[f_{w_n}(S_{n+1}) \nabla_w f_{w_n}(S_{n+1}) \mid \mathcal{F}_n] - \mathbb{E}[\mathbb{E}[U_{n+1} \mid S_{n+1}] \nabla_w f_{w_n}(S_{n+1}) \mid \mathcal{F}_n] - \nabla_w E(w_n) \\ &= \mathbb{E}[f_{w_n}(S_{n+1}) \nabla_w f_{w_n}(S_{n+1}) \mid \mathcal{F}_n] - \mathbb{E}[V^\pi(S_{n+1}) \nabla_w f_{w_n}(S_{n+1}) \mid \mathcal{F}_n] - \nabla_w E(w_n) \\ &\stackrel{\text{ind.}}{=} \mathbb{E}_{S_{n+1} \sim \mu} [(f_{w_n}(S_{n+1}) - V^\pi(S_{n+1})) \nabla_w f_{w_n}(S_{n+1})] - \nabla_w E(w_n) \\ &= 0. \end{aligned}$$

² Finally, the second moments satisfy

$$\begin{aligned} & \mathbb{E}[\|\varepsilon_n\|^2 \mid \mathcal{F}_n] \\ &= \mathbb{E}[\|2(f_{w_n}(S_{n+1}) - U_{n+1}) \nabla_w f_{w_n}(S_{n+1}) - \nabla_w E(w_n)\|^2 \mid \mathcal{F}_n] \\ &\leq \mathbb{E}[\|(2\|U_n \nabla_w f_{w_n}(S_{n+1})\|_2 + \|2f_{w_n}(S_{n+1}) \nabla_w f_{w_n}(S_{n+1}) - \nabla_w E(w_n)\|)^2 \mid \mathcal{F}_n] \\ &= \mathbb{E}[4U_n^2 \|\nabla_w f_{w_n}(S_{n+1})\|^2 \mid \mathcal{F}_n] \\ &\quad + \mathbb{E}[2|U_n| \cdot \|\nabla_w f_{w_n}(S_{n+1})\|_2 \cdot \|2f_{w_n}(S_{n+1}) \nabla_w f_{w_n}(S_{n+1}) - \nabla_w E(w_n)\| \mid \mathcal{F}_n] \\ &\quad + \mathbb{E}[\|2f_{w_n}(S_{n+1}) \nabla_w f_{w_n}(S_{n+1}) - \nabla_w E(w_n)\|^2 \mid \mathcal{F}_n] \\ &\leq C_1 \cdot \mathbb{E}[U_{n+1}^2 \mid \mathcal{F}_n] + C_2 \cdot \mathbb{E}[|U_{n+1}| \mid \mathcal{F}_n] + C_3 \\ &= C_1 \cdot \mathbb{E}[\mathbb{E}[U_{n+1}^2 \mid S_n] \mid \mathcal{F}_n] + C_2 \cdot \mathbb{E}[\mathbb{E}[|U_{n+1}| \mid S_n] \mid \mathcal{F}_n] + C_3 \\ &\leq C_1 K_1 + C_2 K_2 + C_3 \end{aligned}$$

²warum sind die gradienten beschränkt? stimmt nicht

where $C_1, C_2, C_3 \in \mathbb{R}$ exist due to the finiteness of the state space and $K_1, K_2 \in \mathbb{R}$ exist due to our assumption on the variance of U_t . By the SGD theorem, it immediately follows that w_t converges to a stationary point of E . Since E is a finite sum of convex functions if we use a linear architecture, it is a convex function itself. Therefore, the only stationary point is a global minimum which proves the theorem. \square

The proof did not require much. The linear architecture was used to ensure the least-square error is an L -smooth (here quadratic) function, and that the error has a global minimum. The algorithm might still converge to a stationary point for other estimators and other function approximations. Unfortunately, the sampling will usually have a bias as the true value function is replaced by the current estimate. For linear architectures with TD(λ) estimators we refer to the famous paper of Tsitsiklis and van Roy³.

What we did in this section works equally well for approximate evaluation of the Q -function for a given policy. The minimisation goal becomes

$$E(w) = \sum_{s \in \mathcal{S}, a \in \mathcal{A}} \mu(s) \pi(a; s) (Q^\pi(s, a) - f_w(s, a))^2.$$

The generic gradient descent motivated minimisation algorithm is

$$w_{n+1} = w_n + \alpha_n (U_{n+1} - f_w(s_{n+1}, a_{n+1})) \nabla_w f_w(s_{n+1}, a_{n+1}), \quad w_0 \in \mathbb{R}^d,$$

where s_n is drawn from the reference measure μ on \mathcal{S} , a_n according to $\pi(\cdot; s_n)$, and U_n are estimates for $Q^\pi(s_n, a_n)$. In the simplest situation of independent Monte Carlo estimates and linear function approximations $f_w(s, a) = w \cdot \phi(s, a)$ the same stochastic gradient argument shows convergence of $(w_n)_{n \in \mathbb{N}}$ to the best linear approximation of Q^π .

6.2 Approximate policy improvement

Explain how policy evaluation is done approximately and what is known.

- Replace for large action space $\operatorname{argmax}_a Q(s, a)$ by an approximation such as softmax or Gumbel softmax. Explain why this is useful and prove what is known.
- Make link to policy gradient, seeing a gradient step as approximation to taking best action. Also make connection to natural gradient.

6.3 Generic bounds for policy iteration with approximations

So far we discussed both ingredients for policy iteration (policy evaluation, policy improvement) with function approximation, i.e. replacing the true targets by targets that are simpler to obtain. Since both steps involve errors it is a priori not clear how strongly the errors accumulate while iterating. In this section we give a generic result on how errors propagate.



Proposition 6.3.1. Consider a policy improvement algorithm that produces a sequence of stationary policies π_k and a corresponding sequence of approximate state-value functions \widehat{V}^{π_k} with given accuracies

- Policy evaluation error: $\|\widehat{V}^{\pi_k} - V^{\pi_k}\|_\infty \leq \varepsilon$ for all $k \in \mathbb{N}$
- Policy improvement error: $\|T^{\pi_{k+1}} \widehat{V}^{\pi_k} - T^* \widehat{V}^{\pi_k}\|_\infty \leq \delta$ for all $k \in \mathbb{N}$

³J. Tsitsiklis, B. van Roy: "An Analysis of Temporal-Difference Learning with Function Approximation", 1997, IEEE Transactions on Automatic Control, Vol 42(5)



It then holds that

$$\limsup_{k \rightarrow \infty} \|V^* - V^{\pi_k}\|_\infty \leq \frac{\delta + 2\gamma\varepsilon}{(1 - \gamma)^2}.$$

Before we turn to the proof, let us discuss what this statement really means. In approximative policy iteration there are two error sources:

- errors obtained in policy evaluation (estimation errors and approximation errors),
- errors obtained in finding the greedy policy corresponding to the estimated value functions \widehat{V}^{π_k} .

While the first error is clear from the statement, the second error is more hidden. It becomes clearer recalling an important fact from dynamic programming.



A policy π is greedy with respect to some function $V \in \mathbb{R}^n$ if and only if $T^*V = T^\pi V$.

Thus, δ measures, in some sense, the deviation of policies π_{k+1} from the greedy policy obtained from \widehat{V}^{π_k} . But why do we not just estimate the policies in some sense as measures, for instance in total deviation? Using the definitions the error measure used in the proposition turns out to be the deviation from the best policy (the greedy one) weighted by the future impact of actions when playing the action:

$$\begin{aligned} & (T^* \widehat{V}^{\pi_k})(s) - (T^{\pi_{k+1}} \widehat{V}^{\pi_k})(s) \\ &= \max_{a \in \mathcal{A}_s} \left\{ r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'; s, a) \widehat{V}^{\pi_k}(s') \right\} - \sum_{a \in \mathcal{A}_s} \pi_{k+1}(a; s) \left(r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'; s, a) \widehat{V}^{\pi_k}(s') \right) \\ &= \sum_{a \in \mathcal{A}_s} \underbrace{\left(\text{greedy}(\widehat{V}^{\pi_k})(a; s) - \pi_{k+1}(a; s) \right)}_{\text{deviation from best policy}} \underbrace{\left(r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'; s, a) \widehat{V}^{\pi_k}(s') \right)}_{\text{currently estimated Q-value}} \end{aligned}$$

The proposition is very general and sheds some light to different situations.

(i) If all errors are zero, an explicit policy iteration can be performed, then the proposition yields another proof of convergence.

(ii) Suppose the model is known, everything is computable but policy evaluation is performed with function approximation. Then the Q -functions can be computed and the greedy policies obtained without error. In that case $\delta = 0$ and ε measures the approximation error of the value function using the approximation class. The proposition shows how badly the errors can at most accumulate.

What might be possible reasons why $\delta \neq 0$? One way this can happen is if there is no model of the system available. Without knowledge of the transition probabilities, we cannot compute Q^{π_k} and consequently cannot compute the greedy policy of V^{π_k} either. Another reason could be that we do indeed have a model available, but we do not want to use a strictly greedy policy in the next step to ensure sufficient exploration (e.g. an ε -greedy method).



As a side note, by letting $\varepsilon, \delta = 0$ we recover the convergence result for the exact policy iteration procedure.

The above proposition does not only apply in the case of function approximation but in the tabular setting as well, where it might still happen that $\varepsilon \neq 0$ if we use an approximate method for the policy evaluation like Monte-Carlo. Is the bound provided by the proposition any good? Clearly, the term $(1 - \gamma)^2$ in the denominator is annoying if γ is close to 1, but the good news is that the bound only depends linearly on ε and δ .

Two properties of Bellman's expectation operator T^π (and T^*) will be used:

- (i) *Monotonicity*: $V(s) \geq U(s)$ for all $s \in \mathcal{S}$ implies $T^\pi V(s) \geq T^\pi U(s)$ for all $s \in \mathcal{S}$.
(ii) Let $c \in \mathbb{R}$ and let $\mathbf{1} \in \mathbb{R}^{|\mathcal{S}|}$ the vector with 1 in every component. Then,

$$T^\pi(V + c\mathbf{1})(s) = T^\pi V(s) + \gamma c, \quad \forall s \in \mathcal{S}.$$

Both properties follow immediately by the definition or by writing $T^\pi V(s) = \mathbb{E}_s^\pi[R(s, A_0) + \gamma V(S_1)]$. The proposition is proved in two steps. A first lemma deals with a one-step version of the proposition, which is then extended to n steps. Taking limits will then prove the proposition.



Lemma 6.3.2. (generalised policy improvement lemma)

Let π and $\bar{\pi}$ be some policies and let $V \in \mathbb{R}^{|\mathcal{S}|}$ be a vector. If

- $\|V - V^\pi\|_\infty \leq \varepsilon$ for some $\varepsilon > 0$,
- $\|T^*V - T^{\bar{\pi}}V\|_\infty \leq \delta$ for some $\delta > 0$,

then

$$\max_{s \in \mathcal{S}} (V^\pi(s) - V^{\bar{\pi}}(s)) \leq \frac{\delta + 2\gamma\varepsilon}{1 - \gamma}.$$

Why do we call the lemma generalised policy improvement? If $V = V^\pi$ and $\pi = \bar{\pi}$ then the choices $\varepsilon = \delta = 0$ yield the policy improvement lemma for the greedy policy obtained from V^π because T^*V^π is nothing but the value function of the greedy policy obtained from V^π . The generalisation tells us that value weakening can be controlled in the approximate policy iteration step. We will later apply the lemma to $\pi = \pi_k$, $\bar{\pi} = \pi_{k+1}$ and $V = \widehat{V}^{\pi_k}$ for fixed $k \in \mathbb{N}$. The assumptions of the proposition imply the assumption of the lemma for this particular choice.

Proof. Define $\xi = \max_{s \in \mathcal{S}} (V^\pi(s) - V^{\bar{\pi}}(s))$ so that

$$V^{\bar{\pi}}(s) + \xi \geq V^\pi(s), \quad \forall s \in \mathcal{S}.$$

Using the assumption we obtain

$$0 \leq -(T^*V)(s) + (T^{\bar{\pi}}V)(s) + \delta \leq -(T^\pi V)(s) + (T^{\bar{\pi}}V)(s) + \delta$$

for all s . We have used our assumption on $\bar{\pi}$ in the first inequality and the general inequality $T^*V(s) \geq T^\pi V(s)$ in the second one. We can now derive the main inequality that we need,

$$\begin{aligned} V^\pi(s) - V^{\bar{\pi}}(s) &= V^\pi(s) - T^{\bar{\pi}}V^{\bar{\pi}}(s) \\ &\leq (V^\pi(s) - T^{\bar{\pi}}V^\pi(s) + \gamma\xi) + (-(T^\pi V)(s) + (T^{\bar{\pi}}V)(s) + \delta) \\ &\leq \|V^\pi - T^{\bar{\pi}}V^\pi\|_\infty + \|T^{\bar{\pi}}V - T^{\bar{\pi}}V^\pi\|_\infty + \gamma\xi + \delta \\ &= \|T^\pi V^\pi - T^{\bar{\pi}}V^\pi\|_\infty + \|T^{\bar{\pi}}V - T^{\bar{\pi}}V^\pi\|_\infty + \gamma\xi + \delta \\ &\leq 2\gamma\|V^\pi - V\|_\infty + \gamma\xi + \delta \\ &\leq 2\gamma\varepsilon + \gamma\xi + \delta, \end{aligned}$$

where we have used additionally the fixedpoint and contraction properties of the Bellman expectation operators. The definition of ξ implies that $\xi \leq 2\gamma\varepsilon + \gamma\xi + \delta$ which is equivalent to

$$\xi \leq \frac{\delta + 2\gamma\varepsilon}{1 - \gamma}$$

which proves the statement. □

Next, define

$$\xi_k = \|V^{\pi_k} - V^{\pi_{k+1}}\|_\infty \quad \text{and} \quad \zeta_k = \max_{s \in \mathcal{S}} (V^*(s) - V^{\pi_k}(s)).$$

We will now derive a bound for ζ_k which involves in the upper bound the term estimated in the previous lemma for the right choices of $\pi, \bar{\pi}$.



Lemma 6.3.3. Suppose the assumptions from Proposition 6.3.1 hold, then

$$\zeta_{k+1} \leq \gamma\zeta_k + \gamma\xi_k + \delta + 2\gamma\varepsilon$$

for all $k \in \mathbb{N}$.

Proof. First note that

$$V^*(s) - \zeta_k \leq V^{\pi_k}(s), \quad \forall k \in \mathbb{N}.$$

Monotonicity and linearity of T^* as well as the fixed point property yield

$$T^*V^{\pi_k}(s) \geq T^*(V^* - \zeta_k\mathbf{1})(s) = T^*V^*(s) - \gamma\zeta_k = V^*(s) - \gamma\zeta_k, \quad \forall s \in \mathcal{S}.$$

Using the assumption and monotonicity of T^{π_k} this yields

$$\begin{aligned} T^{\pi_{k+1}}V^{\pi_k}(s) &\geq T^{\pi_{k+1}}(\widehat{V}^{\pi_k} - \varepsilon\mathbf{1})(s) \\ &= T^{\pi_{k+1}}\widehat{V}^{\pi_k}(s) - \gamma\varepsilon \\ &\geq T^*\widehat{V}^{\pi_k}(s) - \delta - \gamma\varepsilon \\ &\geq T^*(V^{\pi_k} - \varepsilon\mathbf{1})(s) - \delta - \gamma\varepsilon \\ &= T^*V^{\pi_k}(s) - \delta - 2\gamma\varepsilon \\ &\geq T^*(V^* - \zeta_k\mathbf{1})(s) - \delta - 2\gamma\varepsilon \\ &= V^*(s) - \gamma\zeta_k - \delta - 2\gamma\varepsilon \end{aligned}$$

for all $s \in \mathcal{S}$, where we have used our assumptions on the approximation error of \widehat{V}^{π_k} and the computational error of the corresponding greedy policy. It follows that

$$\begin{aligned} V^{\pi_{k+1}}(s) &= T^{\pi_{k+1}}V^{\pi_{k+1}}(s) \\ &\geq T^{\pi_{k+1}}(V^{\pi_k} - \xi_k\mathbf{1})(s) \\ &= T^{\pi_{k+1}}V^{\pi_k}(s) - \gamma\xi_k \\ &\geq V^*(s) - \gamma\zeta_k - \delta - 2\gamma\varepsilon - \gamma\xi_k \end{aligned}$$

for all $s \in \mathcal{S}$. Rearranging this yields

$$0 \leq V^*(s) - V^{\pi_{k+1}}(s) \leq \gamma\zeta_k + \gamma\xi_k + \delta + 2\gamma\varepsilon$$

for all $s \in \mathcal{S}$ from which we conclude the result. \square

We can now easily prove the original proposition by using a familiar trick from introductory analysis courses.

Proof of Proposition 6.3.1: The second lemma yields

$$\zeta_{k+1} - \gamma\zeta_k \leq \gamma\xi_k + \delta + 2\gamma\varepsilon$$

which, combined with the first lemma (using $\pi = \pi_k$, $\bar{\pi} = \pi_{k+1}$, $V = \widehat{V}^{\pi_k}$ and the assumption of the proposition to check the assumptions of the lemma), gives

$$\zeta_{k+1} - \gamma\zeta_k \leq \gamma \frac{\delta + 2\gamma\varepsilon}{1 - \gamma} + \delta + 2\gamma\varepsilon.$$

Taking the limit superior on both sides yields

$$(1 - \gamma) \limsup_{k \rightarrow \infty} \zeta_k \leq \gamma \frac{\delta + 2\gamma\varepsilon}{1 - \gamma} + \delta + 2\gamma\varepsilon.$$

Simplifying the righthand side and dividing by $(1 - \gamma)$ proves the claim. \square

6.4 Q-learning with linear function approximation

Let us quickly recall what we learnt in dynamic programming and the sample based analogue. Dynamic programming lead to two kind of algorithms:

- Policy iteration algorithms based on Banach's fixed point theorem. Iterating T^π or T^* leads to iterative algorithms that approximate the value function V^π (resp. Q^π) or the optimal value-function V^* (resp. Q^*). If V^* or Q^* is known, then an optimal policy is obtained by taking the greedy policy.
- Using that T^π and T^* for Q (not for V !) can be written as expectations the stochastic approximation theorem for $\|\cdot\|_\infty$ -contractions gives sample based algorithms to approximate V^π , Q^π , and Q^* . Since the mathematical structure is the same (stochastic approximation of fixedpoint) the algorithms looked very similar, using the respective operator.

So far we discussed policy evaluation with function approximation and gave a rigorous justification of the simplest update rule

$$w_{n+1} = w_n + \alpha_n \underbrace{(V^\pi(s_{n+1}) - f_{w_n}(s_{n+1}))}_{\approx U_{n+1}} \nabla_w f_{w_n}(s_{n+1}), \quad w_0 \in \mathbb{R}^d,$$

in terms of stochastic gradient descent. The update minimises (at least for linear approximation architectures) the least-squared error. Similarly, for Q -learning with function approximation the goal is to minimise

$$E(w) = \frac{1}{2} \sum_{s \in \mathcal{S}, a \in \mathcal{A}_s} \mu(s, a) (Q^*(s, a) - f_w(s, a))^2 = \frac{1}{2} \mathbb{E}_{(S, A) \sim \mu} [(Q(S, A) - f_w(S, A))^2]$$

for some reference measure μ . Just as before a stochastic gradient algorithm to minimise E is

$$w_{n+1} = w_n + \alpha_n (Q^*(s_{n+1}, a_{n+1}) - f_{w_n}(s_{n+1}, a_{n+1})) \nabla_w f_{w_n}(s_{n+1}, a_{n+1}), \quad w_0 \in \mathbb{R}^d.$$

As for the temporal difference approach to policy evaluation with function approximation we use the Bellman operator. Recalling that Q^* is the unique fixedpoint of $T^*Q(s, a) = \mathbb{E}_s^{\pi^a} [R(s, a) + \gamma \max_{a' \in \mathcal{A}_{S_1}} Q(S_1, a')]$ we replace the unknown Q^* by a sample of the expectation and the unknown Q by the best known approximation f_w .



Definition 6.4.1. The update scheme

$$w_{n+1} = w_n + \alpha_n (R(s_{n+1}, a_{n+1}) + \max_a f_{w_n}(s_{n+1}, a) - f_{w_n}(s_{n+1}, a_{n+1})) \nabla_w f_{w_n}(s_{n+1}, a_{n+1})$$

is called Q -learning with function approximation. If $\{f_w = w^T \Phi : w \in \mathbb{R}^d\}$ one speaks of Q -learning with linear function approximation, if $\{f_w : w \in \mathbb{R}^d\}$ are neural network functions one speaks of deep Q -learning (DQN). The choice of state-action pairs (s_n, a_n) is handled differently in different versions of the algorithm.

6.5 Policy gradient with linear function approximation

Similarly to stochastic approximation

$$w_{n+1} = w_n + \alpha_n \left(R(s_n, a_n) + \gamma \max_{a'} f_w(s', a') - f_w(s_n, a_n) \right) \nabla_w f_w(s_n, a_n),$$

6.6 A quick dive into neural networks

6.6.1 What are neural network functions?

To get the discussion started let us look at the simplest example, that of a single neuron. A neuron in our brain is an electrically excitable cell that fires electric signals called action potentials. In mathematical term, a neuron is a function that takes a vector (incoming information) and returns a number (signal, e.g. 0 or 1). A network of billions of communicating neurons passes incoming signals to outgoing action signals. Inspired from neuroscience a century of research in computer science has led to incredible results in artificial intelligence. Let us first describe the simplest model of a neuron, and then continue with artificial networks of neurons (so-called neural networks or perceptrons). The simplest model⁴ is

$$f(x) = \mathbf{1}_{[0, \infty)} \left(\sum_{i=1}^d w_i x_i + b \right).$$

If the combined strength of the signals (weighted according to neurons judgement of importance for the signals) exceeds a certain level the neuron fires, otherwise not. More generally, for the definition of an artificial neuron the indicator is typically replaced by a generic function σ , the so-called activation function.



Definition 6.6.1. For $\sigma : \mathbb{R} \rightarrow \mathbb{R}$, $w \in \mathbb{R}^d$, and $b \in \mathbb{R}$, the function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ defined by

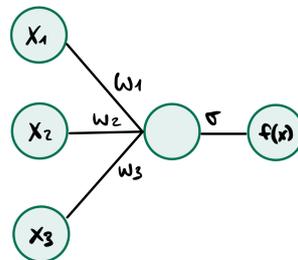
$$f(x) = \sigma \left(\sum_{i=1}^d w_i x_i + b \right)$$

is called an artificial neuron with activation function σ , weights w_1, \dots, w_d for the signals x_1, \dots, x_d and bias b .

The weights can be interpreted as importance of the signal for the neuron, the bias as susceptibility of the neuron.

Example 6.6.2. Here are some typical activation functions:

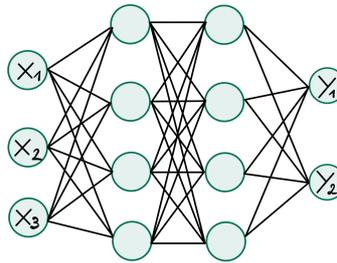
- $\sigma(s) = \mathbf{1}_{[0, \infty)}(s)$, the Heavyside function, yields a neuron that either fires or not,
- $\sigma(s) = s$ yields a linear neuron which is not really used,
- $\sigma(s) = \max\{0, s\}$ is called rectified linear unit (ReLU), a very typical practical choice,
- $\sigma(s) = \frac{1}{1+e^{-s}} = \frac{e^s}{1+e^s}$ is called logistic activation function. The logistic activation function is nothing but the softmax probability for one of two possible actions.



Graphical representation of an artificial neuron

⁴McCulloch, Pitts: "A logical calculus of ideas immanent in nervous activity", Bull. Math. Biophys., 5:115-133, 1943

As mentioned the brain consists of billions of neurons that form a network in some very complicated manner. We will do the same, transfer the outputs of neurons into new neurons.



Graphical representation of a neural network with 2 hidden layers, 3 input dimensions, 2 output dimensions

Here is a mathematical formulation of a simple network of neurons using the same activation function, such networks of artificial neurons are also called multilayer perceptrons⁵:



Definition 6.6.3. (Feedforward neural networks)

Let $d, L \in \mathbb{N}$, $L \geq 2$, and $\sigma : \mathbb{R} \rightarrow \mathbb{R}$. Then a multilayer perceptron (or fully connected feedforward neural network) with d -dimensional input, L layers, and activation function σ is a function $F : \mathbb{R}^d \rightarrow \mathbb{R}^k$ that can be written as

$$F(x) = T_L(\sigma(T_{L-1}(\dots \sigma(T_1(x))))),$$

where $L - 1$ is the number of hidden layers with N_l units (number of neurons per layer), $N_0 = d$ and $N_L = k$, and $T_l(x) = W_l x + b_l$ are affine functions with

- weight matrices $W_l \in \mathbb{R}^{N_l \times N_{l-1}}$,
- biases $b_l \in \mathbb{R}^{N_l}$.

Note that the function σ is applied coordinate-wise to the vectors. A neural network is called shallow if $L = 1$, i.e. there is only one hidden layer and deep if the number of layers is large. Finally, if the output is supposed to be a probability vector one applies a final normalisation with the softmax function $\Phi(y)_i = \frac{e^{y_i}}{\sum_{j=1}^k e^{y_j}}$ to the output vector:

$$F(x) = \Phi(T_L(\sigma(T_{L-1}(\dots \sigma(T_1(x)))))).$$

⁶ The multilayer perceptron is quite special in its structure. All neurons of a fixed layer receive the same input vector (the outputs of the neurons before) which is then multiplied with the weight-vector of the neuron, shifted, and plugged-into the activation function to yield the signal. To visualise the mapping one can for instance write the real-valued weights on the arrow pointing into a neuron. If a weight is forced to be zero the arrow is usually removed from the graphical representation. This neural network was made up by thinking, not by brute-force training on mass data. The network produces a vector of features of an image, such as "number of circles", "number of edges", ...

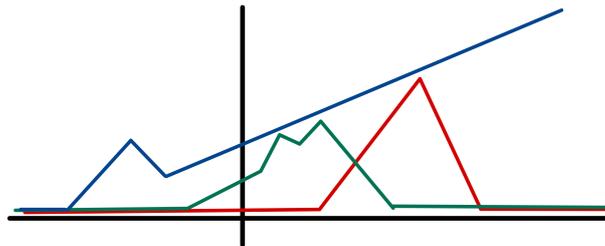
⁵F. Rosenblatt: "The perceptron: a probabilistic model for information storage and organization in the brain" Psychological review, 65(6):386, 1958

⁶include drawings



Definition 6.6.4. The architecture of a feedforwards neural network is the chosen number of layers, units, and fixed zero-weights.

The class of neural network functions has been investigated for decades. While the first approach was to understand what particular architecture leads to desired behavior the modern approach is to use mass data to learn the weights and biases by optimising a carefully chosen loss-function. Keeping in mind that the graphical representation of the architecture is only a visualisation of the weights it can be instructive to look at plots of very small-dimensional neural network functions with ReLU activation function:



Three functions obtained from ReLU neural networks

To get an idea why neural networks with the simple ReLU activation and linear output function can create such functions two observations are useful:

- One can construct addition, shifting, and linear transformations of functions obtained from a neural network through a neural network with same activation function. To do so write the neural networks on top of each other without connections (i.e. set the linking weights to 0), and use the output weights to obtain the desired transformation.
- The sum $\text{ReLU}(-1) - 2\text{ReLU}(0) + \text{ReLU}(1)$ gives a spike at 0, with $\text{ReLU}(b) = \max\{0, b\}$.

There are other classes of artificial neural networks that are much more relevant in applications. Neural networks that do not only pass signals forwards but also have arrows pointing to previously used neurons are called recurrent neural networks. Those are particularly useful to model time-dependent behavior. More recently, the so-called transformer network architecture has received enormous attention and allowed massive progress in large language models. Since this is an introductory course to reinforcement learning we will only discuss two important features of neural networks, function approximation and differentiation by backwards propagation through the network.

6.6.2 Approximation properties

We will show that already the simplest neural network functions are useful enough to approximate important classes of functions, so-called Boolean functions and continuous functions. The proofs are not complicated but only lead to rather uninteresting existence results⁷.



Theorem 6.6.5. (Approximation of Boolean functions)

Every Boolean function $f : \{0, 1\}^d \rightarrow \{0, 1\}$ can be represented by a neural network with activation function $\sigma(x) = \mathbf{1}_{[0, \infty)}$.

⁷Check for instance these lecture notes of Philipp Christian Petersen for more.

Proof. First note that for all Boolean variables $a, b \in \{0, 1\}$ it holds that $2ab - a - b \leq 0$ with equality if and only if $a = b$ (check the cases). Thus, for $u \in \{0, 1\}^d$,

$$\mathbf{1}_{x=u} = \mathbf{1}_{[0,\infty)}\left(\sum_{i=1}^d (2x_i u_i - x_i - u_i)\right) = \sigma\left(\sum_{i=1}^d (2x_i u_i - x_i - u_i)\right).$$

Now denoting by $A := \{u \in \mathbb{R}^d : f(u) = 1\}$ it follows that

$$f(x) = \sigma\left(-1 + \sum_{u \in A} \mathbf{1}_{x=u}\right) = \sigma\left(-1 + \sum_{u \in A} \sigma\left(\sum_{i=1}^d (2x_i u_i - x_i - u_i)\right)\right).$$

This, in fact, this is nothing but a neural network with two hidden layers. The first with $|A|$ units, the second with one unit. □



Definition 6.6.6. A subset \mathcal{F} of neural networks has the universal approximation property if for any $K \subset \mathbb{R}^d$, any $\varepsilon > 0$, and any continuous function $g : K \rightarrow \mathbb{R}$ there exists a neural network $f \in \mathcal{F}$ with $\|f - g\|_\infty < \varepsilon$.

In fact, for many activation functions the universal approximation property holds. An activation function is called sigmoid (*S*-formed) if it is bounded with limits $\lim_{x \rightarrow +\infty} \sigma(x) = 1$ and $\lim_{x \rightarrow -\infty} \sigma(x) = 0$. The terminology comes from the logistic activation function that really looks like an *S*. An important result on neural networks is that shallow (one hidden layer) feedforward neural networks with sigmoid activation function have the universal approximation property.



Theorem 6.6.7. (Universal approximation theorem)

If the activation function σ is sigmoid then the set of all shallow feedforward neural networks has the universal approximation property.

Proof. Let us denote by \mathcal{F}_p the subset of shallow neural networks with p neurons. We first consider $d = 1$ and show that Lipschitz functions on $[0, 1]$ can be approximated. To do so we first approximate h with piecewise constant functions and then approximate the piecewise constant function with a shallow neural network function. For that sake let $x_i = \frac{i}{p}$ and set $h_p(x) = \sum_{i=1}^p h(x_i) \mathbf{1}_{[x_{i-1}, x_i)}(x)$.



There is $f \in \mathcal{F}_p$ with $\|f - h\|_\infty \leq C \omega(h, 1/p)$ with C independent of h .

In the estimate the so-called modulus of continuity appears:

$$\omega(h, \delta) = \sup_{x, y \in \mathbb{R}: |x-y| \leq \delta} |h(x) - h(y)|$$

For a function $h : \mathbb{R} \rightarrow \mathbb{R}$ the modulus of continuity is a very rough measure for the flatness of h . The smaller ω the flatter the function. For Lipschitz continuous functions with Lipschitz constant L it follows immediately that $\omega(h, \delta) \leq L\delta$, the fluctuation of f over intervals of length δ are at most δL . The estimate shows that Lipschitz functions with smaller constant L can be well-approximated with $1/L$ many units, thus, rougher functions need more units to be approximated. We will only use the statement for the exponential function which is clearly Lipschitz continuous on compact sets (bounded derivative on compact sets). To prove the first claim note that

$$\sup_{x \in [0,1]} |h(x) - h_p(x)| = \sup_{x \in [0,1]} \left| h(x) - \sum_{i=1}^p h(x_i) \mathbf{1}_{[x_{i-1}, x_i)}(x) \right| \leq \omega(h, 1/p)$$

by the definition of h_p and ω . With a telescopic sum we rewrite

$$h_p(x) = h(x_1) + \sum_{i=1}^{\lfloor px \rfloor} (h(x_{i+1}) - h(x_i)).$$

Next, we define a the shallow neural network function with p units as

$$f_p(x) = h(x_1)\sigma(\alpha) + \sum_{i=1}^{p-1} (h(x_{i+1}) - h(x_i))\sigma(\alpha(px - i)) \in \mathcal{F}_p$$

for some $\alpha \in \mathbb{R}$ that is specified next. Fix $\varepsilon > 0$ and choose α large enough such that $|\sigma(z) - \mathbf{1}_{[0,\infty)}(z)| \leq \frac{\varepsilon}{p}$ whenever $|z| \geq \alpha$. Since σ is a sigmoid function this is possible. The choice of α shows that

$$|\sigma(\alpha(px - i)) - \mathbf{1}_{i \leq \lfloor px \rfloor}| \leq \frac{\varepsilon}{p}$$

holds for all $i \notin \{\lfloor px \rfloor, \lfloor px \rfloor + 1\}$ because then $|\alpha(px - i)| > \alpha$. It then follows that

$$\begin{aligned} & |f_p(x) - h_p(x)| \\ &= \left| h(x_1)(\sigma(\alpha) - 1) + \sum_{i=1}^{p-1} (h(x_{i+1}) - h(x_i))(\sigma(\alpha(px - i)) - \mathbf{1}_{i \leq \lfloor px \rfloor}) \right| \\ &\leq \frac{\varepsilon}{p} (|h(x_1)| + (p-2)\omega(h, 1/p)) + |h(x_{\lfloor px \rfloor+1}) - h(x_{\lfloor px \rfloor})| |\sigma(\alpha(px - \lfloor px \rfloor)) - 1| \\ &\quad + |h(x_{\lfloor px \rfloor+2}) - h(x_{\lfloor px \rfloor+1})| |\sigma(\alpha(px - \lfloor px \rfloor - 1)) - 1|. \end{aligned}$$

The first summand can be made arbitrarily small, the other two can both be estimated by $\omega(h, 1/p)(1 + \|\sigma\|_\infty)$. Combining the approximation of h by h_p and the approximation of h_p by f_p it follows that for all $\delta > 0$ there is some p such that

$$\begin{aligned} \sup_{x \in [0,1]} |h(x) - f_p(x)| &\stackrel{\Delta}{\leq} \sup_{x \in [0,1]} |h(x) - h_p(x)| + \sup_{x \in [0,1]} |h_p(x) - f_p(x)| \\ &\leq \omega(h, 1/p) + \delta + 2(\|\sigma\|_\infty + 1)\omega(h, 1/p). \end{aligned}$$

Now the claim follows. Covering a compact set $K \subseteq \mathbb{R}$ by finitely many intervals it then follows that for all $h : K \rightarrow \mathbb{R}$ and all $\varepsilon > 0$ there is a neural network function f such that $\sup_{x \in K} |h(x) - f(x)| < \varepsilon$.

It remains to extend the arguments to the d -dimensional case. To reduce to the one-dimensional case a simple dense family is used:

The set

$$\mathcal{E} := \left\{ g : K \rightarrow \mathbb{R} \mid g(x) = \sum_{i=1}^N s_i e^{(v^i, x)}, x \in \mathbb{R}^d, N \in \mathbb{N}, v^i \in \mathbb{R}^d, s_i \in \{-1, 1\} \right\}$$

is dense in $(C(K), \|\cdot\|_\infty)$.

Using Stone-Weierstrass from functional analysis one only needs to check that \mathcal{E} is an algebra that separates points. Easy.

Next, we approximate \mathcal{E} by neural network functions with one hidden layer:

For every $g(x) = \sum_{i=1}^N s_i e^{(v^i, x)} \in \mathcal{E}$ and every $\varepsilon > 0$ there is a neural network function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ with one hidden layer such that $\sup_{x \in K} |g(x) - f(x)| < \varepsilon$.

Define $g_{v^i}(x) = e^{\langle v^i, x \rangle}$. The trick is simple. Let \bar{K}_i the linear transformation of K_i under $x \mapsto \langle v^i, x \rangle$. Then $\bar{K}_i \subseteq \mathbb{R}$ is compact again. Next, apply the first step to chose from \mathcal{F}_p an

$$f_i(x) = \sum_{l=1}^p w_l \sigma(x - b)$$

with $\sup_{x \in \bar{K}_i} |f_i(x) - e^x| < \frac{\varepsilon}{N}$. Then define the neural network function $\bar{f}_i : \mathbb{R}^d \rightarrow \mathbb{R}$ as

$$\bar{f}_i(x) = s_i f_i \left(\sum_{k=1}^d v_k^i x_k \right)$$

so that

$$\begin{aligned} \sup_{x \in K} |\bar{f}_i(x) - s_i e^{\langle v^i, x \rangle}| &= \sup_{x \in K} |s_i f_i(\langle v^i, x \rangle) - s_i e^{\langle v^i, x \rangle}| \\ &= \sup_{x \in \bar{K}_i} |s_i f_i(x) - s_i e^x| < \frac{\varepsilon}{N}. \end{aligned}$$

To see that \bar{f}_i is a neural network function with one hidden layer note that

$$\bar{f}_i(x) = \sum_{l=1}^p s_i w_l \sigma \left(\sum_{k=1}^d v_k^i x_k - b_l \right).$$

Recalling that sums of neural network functions are neural network functions with the same number of hidden layers (stacking neural networks without connections) a neural network approximation of g is given by $f = \sum_{i=1}^N \bar{f}_i$:

$$\sup_{x \in K} |f(x) - g(x)| \leq \sum_{i=1}^N \sup_{x \in K} |\bar{f}_i(x) - s_i e^{\langle v^i, x \rangle}| < \varepsilon.$$



The universal approximation property is satisfied for \mathcal{F}_p .

Suppose $h \in (C(K), \|\cdot\|_\infty)$ and $\varepsilon > 0$. Using the two steps before we can choose $g \in \mathcal{E}$ with $\sup_{x \in K} |g(x) - h(x)| < \frac{\varepsilon}{2}$ and a neural network function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ with $\sup_{x \in K} |f(x) - h(x)| < \frac{\varepsilon}{2}$. Then $\sup_{x \in K} |h(x) - f(x)| \leq \sup_{x \in K} |h(x) - g(x)| + \sup_{x \in K} |g(x) - f(x)| < \varepsilon$. This proves the theorem. \square

6.6.3 Differentiation properties

There are a couple of reasons why the use of deep neural networks is extremely successful in reinforcement learning. We will highlight here a differentiation property that led to the breakthroughs in image classification. Suppose a set of labeled images is given: $\{(x_i, y_i) : i \leq N\}$. We assume the images are already transformed into a pixel vector, say \mathbb{R}^d , and they are labeled as a k -dimensional feature vector. A feature vector is an abstraction of an image used to characterize and numerically quantify the contents of an image. Simply put, a feature vector is a list of numbers used to represent an image. The pixel vector itself is a possible feature vector that fully describes an image but the aim is extract as little abstract information as possible that describes a images as good as possible. This could be counting geometric forms or more specifically the number of persons on the images.



A main goal of image analysis is to define feature vectors that encode the images well and then find mappings that are simple to evaluate and map an image to its feature vector.

The typical approach consists learning both tasks on a large training set $\{(x_i, y_i)\}$ of images x_i for which the labels y_i were coded manually by humans. There are plenty of such examples, for instance imageNet⁸ with more than fourteen million labeled images. There are two amazing features why neural networks are used to map images to feature vectors:

- Differentiation: neural network structures allow to differentiate well for the weights, thus, allowing gradient descent methods to minimise the error between outputs of the neural network and the labels. Since neural networks used have dozens of millions of parameters this is crucial. Imagine for every gradient step to compute a gradient of fifty million partial derivatives. This is only possible if many of the computations can be reused, a property that holds for neural network functions due to the so-called back propagation algorithm.
- Generalisation: for some magic reasons that are not fully understood, yet, neural networks generalise extremely well. This means that the mapping from image to labels obtained by gradient descent on a finite labeled set generalises well to further images. Overfitting is a smaller problem than one might imagine (at least there are efficient tricks such as dropout) if for instance sixty million parameters are used to classify one million images. For imageNet one typically uses a training set of 1.2 million images and then compares the generalisation property for the other images.

Since this is not a course on deep learning we only discuss a very nice differentiation trick. Algorithm 23 summarises the most famous one, the delta-method that can be used to train a neural network on given input/output pairs or to just compute the gradient $\nabla_w F(x)$.



Theorem 6.6.8. The δ -method from Algorithm 23 is nothing but stochastic gradient descent for $F_w := \sum_{i=1}^N \frac{1}{2} \|F_w(x)_i - y_i\|^2$ with (x_i, y_i) chosen uniformly from the training set.

⁹ In the algorithm and the proof we allowed a generic function Φ to normalise the outputs so the chain of concatenations is

$$F(x) = T_L(\sigma(T_{L-1}(\dots \sigma(T_1(x))))),$$

This is needed if the network should output a probability vector.

Proof. In what follows we use the definitions of h, V appearing in the algorithm. Please check Figure ??¹⁰ for a graphical representation to ease the understanding. ¹¹ First of all, note that the optimisation is equivalent to optimising $\frac{1}{N} F_w$, thus, F_w can be written as expectation $F_w = \mathbb{E}[\|F_w(x) - y\|^2]$, where the randomness comes by uniformly choosing the pair (x, y) from the training set. Hence, a stochastic gradient algorithm first chooses a pair (x, y) uniformly and then computes the gradient $\nabla_w \|F_w(x) - y\|^2$. In what follows we show all partial derivatives $\frac{\partial}{\partial w_{i,j}^l}$ of the gradient are indeed computed by the δ -rule.

Output layer: To compute the partial derivative with respect to the final weights $w_{i,j}^L$ we

⁸imageNet website

⁹ Φ muss weg weil es nicht koordinatenweise wirkt, dazu passen die indices nicht

¹⁰bildchen noch malen

¹¹put drawing

proceed as follows. Writing the norm and the definition of the neural network yields

$$\begin{aligned}
\frac{\partial}{\partial w_{i,j}^L} \frac{1}{2} \|F_w(x) - y\|^2 &= \frac{\partial}{\partial w_{i,j}^L} \sum_{s=1}^{N_L} \frac{1}{2} \left(\Phi \left(\sum_{t=1}^{N_{L-1}} w_{t,s}^L V_t^{L-1} \right) - y_s \right)^2 \\
&= \frac{\partial}{\partial w_{i,j}^L} \frac{1}{2} \left(\Phi \left(\sum_{t=1}^{N_{L-1}} w_{t,j}^L V_t^{L-1} \right) - y_j \right)^2 \\
&= \left(\Phi \left(\sum_{t=1}^{N_{L-1}} w_{t,j}^L V_t^{L-1} \right) - y_j \right) \Phi'(h_j^L) \frac{\partial}{\partial w_{i,j}^L} \sum_{t=1}^{N_{L-1}} w_{t,j}^L V_t^{L-1} \\
&= (V_j^L - y_j) \Phi'(h_j^L) V_i^{L-1} \\
&=: \delta_j^L V_i^{L-1}.
\end{aligned}$$

There are two crucial points involved in the second and fourth equality. For the second equality we use that only the j th term depends on $w_{i,j}^L$ (compare the graphical representation of the network) so that all other derivatives equal zero. Similarly, for the fourth equality only the term $t = i$ depends on $w_{i,j}^L$ so that all other derivatives vanish and the t th term simply has the derivative V_i^{L-1} .

Last hidden layer: To compute the partial derivative with respect to the final weights $w_{i,j}^{L-1}$ we proceed similarly. Writing the norm and the definition of the neural network yields

$$\begin{aligned}
&\frac{\partial}{\partial w_{i,j}^{L-1}} \frac{1}{2} \|F_w(x) - y\|^2 \\
&= \frac{\partial}{\partial w_{i,j}^{L-1}} \sum_{s=1}^{N_L} \frac{1}{2} \left(\Phi \left(\sum_{t=1}^{N_{L-1}} w_{t,s}^L \sigma \left(\sum_{r=1}^{N_{L-2}} w_{r,t}^{L-1} V_r^{L-2} \right) \right) - y_s \right)^2 \\
&= \sum_{s=1}^{N_L} (V_s^L - y_s) \Phi'(h_s^L) \frac{\partial}{\partial w_{i,j}^{L-1}} \sum_{t=1}^{N_{L-1}} w_{t,s}^L \sigma \left(\sum_{r=1}^{N_{L-2}} w_{r,t}^{L-1} V_r^{L-2} \right) \\
&= \sum_{s=1}^{N_L} (V_s^L - y_s) \Phi'(h_s^L) \sum_{t=1}^{N_{L-1}} w_{t,s}^L \sigma' \left(\sum_{r=1}^{N_{L-2}} w_{r,t}^{L-1} V_r^{L-2} \right) \frac{\partial}{\partial w_{i,j}^{L-1}} \sum_{r=1}^{N_{L-2}} w_{r,t}^{L-1} V_r^{L-2} \\
&= \sum_{s=1}^{N_L} (V_s^L - y_s) \Phi'(h_s^L) \sum_{t=1}^{N_{L-1}} w_{t,s}^L \sigma'(h_j^{L-2}) \frac{\partial}{\partial w_{i,j}^{L-1}} \sum_{r=1}^{N_{L-2}} w_{r,t}^{L-1} V_r^{L-2}.
\end{aligned}$$

All derivatives of the last term disappear, except for $r = i$ and $t = j$. The remaining derivative equals 1, thus, the expression simplifies to

$$\begin{aligned}
\frac{\partial}{\partial w_{i,j}^{L-1}} \frac{1}{2} \|F_w(x) - y\|^2 &= V_i^{L-2} \sigma'(h_j^{L-2}) \sum_{s=1}^{N_L} (V_s^L - y_s) \Phi'(h_s^L) w_{j,s}^L \\
&= V_i^{L-2} \sigma'(h_j^{L-2}) \sum_{s=1}^{N_L} \delta_j^L w_{j,s}^L \\
&=: V_i^{L-2} \delta_j^{L-1}.
\end{aligned}$$

Well, anyone who likes playing with indices is warmly invited to turn this into a clean induction :).¹² \square

Since the derivatives of σ repeatedly appear it is clearly desirable to chose activation functions that are easy to differentiate. The logistic function $\sigma(x) = \frac{1}{1+e^{-x}}$ is a typical example with a reasonable derivative $\sigma'(x) = \frac{e^x}{(1+e^x)^2}$.

¹²write induction



Go through the previous arguments to check that the gradient $\nabla_w F_w(x)$ is computed analogously with $\delta_j^L = \Phi'(h_j^L)$ and $\delta_i^l = \sigma'(h_i^{l-1}) \sum_{j=1}^{N_i} w_{i,j}^l \delta_j^l$.

Data: labeled training set (x_i, y_i)

Result: weights w to minimise training error $\sum_i \|F_w(x_i) - y_i\|^2$

Set $n = 0$.

Initialise all weights (for instance iid Gaussian).

while *not converged* **do**

 Chose (x, y) uniformly from the labeled training set.

Forwards propagation: Starting from the input x compute forwards through the network $V_i^0 := x_i$ and

- $h_i^1 := \sum_{k=1}^{N_{l-1}} w_{k,l}^1 V_k^{l-1}, i = 1, \dots, N_l,$

- $V_i^l := \sigma(h_i^l), i = 1, \dots, N_l,$

for $l < L$ and

- $h_i^L := \sum_{k=1}^{N_L} w_{k,L}^L V_k^{L-1}$

- $V_i^L := \Phi(h_i^L)$

Back propagation:

- compute $\delta_j^L = (V_j^L - y_j) \Phi'(h_j^L).$

- compute $\delta_i^l = \sigma'(h_i^{l-1}) \sum_{j=1}^{N_{l+1}} w_{i,j}^{l+1} \delta_j^{l+1}$

 Chose a step-size α .

 Update all weights as $w_{i,j}^l = w_{i,j}^l + \alpha \delta_j^l V_i^l$.

end

Algorithm 23: δ -method to train a neural network

6.6.4 Using neural networks to approximate value functions

6.6.5 Using neural networks to parametrise policies

The success of neural networks in supervised learning was not overseen in the reinforcement learning community where neural networks are used in all branches with increasing success. We give a very brief overview of the use of neural networks for policy gradient schemes but also Q -learning. In later sections we will see how to combine both (actor-critic) and how to proceed much further in deep Q -learning.

Linear softmax with neural networks representations:

First recall the simplest parametrisation in the tabular setting. The tabular softmax policy was defined as $\pi^\theta(a; s) = \frac{e^{\theta_{s,a}}}{\sum_{a'} e^{\theta_{s,a'}}$ with a single parameter $\theta_{s,a}$ tuning the probability of action a in state s . Since this is only reasonable for small state-action spaces the policy has no practical relevance. More relevant are linear softmax policies of the form

$$\pi^\theta(a; s) = \frac{e^{\theta^T \cdot \Phi(s,a)}}{\sum_{a'} e^{\theta^T \cdot \Phi(s,a')}, \quad \theta \in \mathbb{R}^d,$$

or

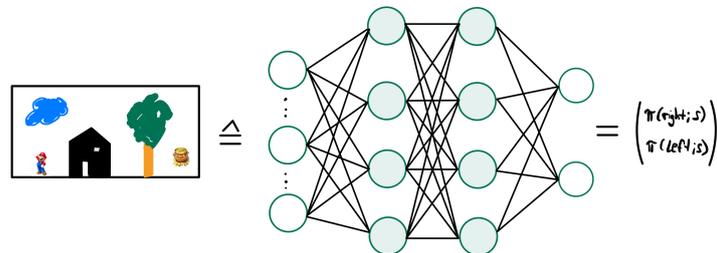
$$\pi^\theta(a; s) = \frac{e^{\theta_a^T \cdot \Phi(s)}}{\sum_{a'} e^{\theta_{a'}^T \cdot \Phi(s)}}, \quad \theta = (\theta_{a_1}, \dots, \theta_{a_k}) \in \mathbb{R}^{ad},$$

where Φ is a representation (feature vector) of the state (resp. state-action pair). In order to use such policies a reasonable representation Φ must be obtained. One way is to work with representations Φ obtained in supervised learning. For a concrete example let us think of learning to play Atari games. The state-space is large, it consists of all possible Atari frames of 210x160 pixel. Without any prior knowledge the state-space is 210×160 dimensional. Using neural networks trained on large libraries of images (such as imageNet) there are large neural networks that give a reasonably large representation of the states, say 1000-dimensional. The function Φ takes an input Atari frame and returns the last hidden layer of the neural network. Essentially, the precise structure of the state s is replaced by the essential information $\Phi(s)$. Using the linear softmax policy we could now run a policy gradient algorithm with the linear softmax policies, which is a 1000-dimensional stochastic gradient algorithm.

Direct use of neural network

The linear softmax has the advantage to provide a the simple and explicit score function $\Phi(s, a) = \sum_{a'} \pi^\theta(a'; s) \Phi(s, a')$ which is useful to compute the policy gradient. A neural network can also be used directly to write down a parametrised policy $\pi^\theta(s, a)$, where θ is the vector of all weights. There are two approaches with the same disadvantage: the parameter vector θ is the entire weight vector of the neural network, which, for many applications, will be huge.

State as input, action probabilities as output: Thinking of Atari games (many states, few actions) once more here is a visulisation of the idea. An image (which is nothing but a long pixel vector) is fed into a neural network, the output vector are the probabilities $\pi(a_1; s), \dots, \pi(a_k; s)$. In that case a softmax function transforms the final hidden layer into probabilities.



Direct policy parametrisation using neural networks

State-action as input, action probability as output: Alternatively, one might feed a state-action pair into a neural network and obtain a one-dimensional output $\pi(a; s)$.

6.7 Deep Q-learning (DQN)

6.8 Deep policy gradient (actor critic methods)

Rethinking the plain vanilla REINFORCE algorithm it comes as no surprise that more advanced algorithms have been developed. REINFORCE is nothing but stochastic gradient descent with non-iid data using the policy gradient theorem to evaluate the expectations. Since stochastic gradient descent is known to be slow and for reinforcement learning we are interested in very large problems there is generally not much hope that REINFORCE works in practice. Indeed, it doesn't but combined with further idea it does. In this section we will discuss some of the most important ideas for practical use of REINFORCE.

6.8.1 Simple actor-critic (AC) and advantage actor-critic (A2C)

So far we have discussed two classes of approaches to solve optimal control problems:

- Policy-based algorithms that learn the optimal policy using gradient descent. A possible drawback of such methods is that the gradient estimators may have a large variance, no bootstrapping of old values is used in the gradient updates.
- Value-based algorithms that evaluate the value-function (or Q -function). Such methods are indirect, they do not try to optimise directly over a set of policies.

Algorithms that directly deal with policies are called actor-only methods while algorithms that improve by criticising the outcome (value-function) of a policy that someone else obtained are called critic-only algorithms. Mixing both approaches is called actor-critic. The typical approach is as follows. Using the policy gradient theorems to perform the REINFORCE algorithm incorporates the Q -functions $Q^{\pi^{\theta_n}}$ or, since the Q -function is typically unknown, estimates of the Q -functions (for instance rewards to go) that introduce variance in the gradient update. An alternative approach is to fix a (simpler) parametrised family ($Q_w^{\pi^{\theta_n}}$) to approximate $Q^{\pi^{\theta_n}}$. Usually the critic's family ($Q_w^{\pi^{\theta_n}}$) is parametrised by a neural network and w is the vector of weights. Actor-critic then works by alternating actor and critic ^{13,14}:

- critic step: approximate $Q^{\pi^{\theta_n}}$ by some $Q_w^{\pi^{\theta_n}}$,
- actor step: perform gradient update using the policy gradient theorem with $Q_w^{\pi^{\theta_n}}$ instead of $Q^{\pi^{\theta_n}}$.

Compared to the direct REINFORCE algorithm the actor-critic approach has advantages and disadvantages. On the plus side actor-critic algorithms can reduce the variance of the policy gradient by using the critic's value function as a baseline. Less variance means less steps for convergence. Next, actor-critic algorithms can incorporate ideas from temporal-difference learning to force bootstrapping of samples. Actor-critic algorithms also have some disadvantages compared to plain vanilla policy gradient algorithms. For instance, they introduce a trade-off between bias and variance, as the approximation in the critic step introduces a bias. Additionally, the approximation increases complexity and computational cost for the algorithm as they require the actor steps and typically the training of a neural network for the parametrisation of the critic's value function. We will next show a result that shows that actor-critic with linear function approximations theoretically converges but using non-linear function approximations (such as neural network functions) there is no theoretical justification for convergence and in practice it is a delicate matter to make such an algorithm converge. Let us start with the plain vanilla actor-critic in a very theoretical setup of linear function approximation. ¹⁵



Proposition 6.8.1. If \mathbf{Q}^{π^θ} satisfies the compatibility with the policy (score function)

$$\nabla_w \mathbf{Q}_w^{\pi^\theta}(s, a) = \nabla_\theta \log(\pi^\theta(s, a))$$

and the approximation property

$$\sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}_s} d^{\pi^\theta}(s) \pi^\theta(a; s) (Q^{\pi^\theta}(s, a) - \mathbf{Q}_w^{\pi^\theta}(s, a)) \nabla_w \mathbf{Q}_w^{\pi^\theta}(s, a) = 0,$$

then

$$\nabla J_s(\theta) = \frac{1}{1-\gamma} \sum_{s' \in \mathcal{S}} \sum_{a \in \mathcal{A}_{s'}} d^{\pi^\theta}(s') \pi^\theta(a; s') \nabla \log(\pi^\theta(a; s')) \mathbf{Q}_w^{\pi^\theta}(s', a).$$

¹³R. Sutton, D. McAllester, S. Singh, Y. Mansour: "Policy Gradient Methods for Reinforcement Learning with Function Approximation", NIPS 1999

¹⁴V. Konda, J. Tsitsiklis: "Actor-Critic Algorithms", NIPS 1999

¹⁵ueberall μ in gradient einarbeiten, auch fuer d^π

¹⁶ The exact same formulas also hold with the other policy gradient representations from Section 4.1.2.

Proof. By the chain rule the first condition is equivalent to $\nabla_w \mathbf{Q}_w^{\pi^\theta}(s, a) \pi^\theta(s, a) = \nabla_\theta \pi^\theta(s, a)$. Plugging-into the second yields

$$\sum_{s' \in \mathcal{S}} d^{\pi^\theta}(s') \sum_{a \in \mathcal{A}_{s'}} \nabla_\theta \pi^\theta(a; s) (Q^{\pi^\theta}(s', a) - \mathbf{Q}_w^{\pi^\theta}(s', a)) = 0$$

Using the policy gradient theorem (version of Theorem 4.1.6) and the above yields the claim:

$$\begin{aligned} \nabla J_s(\theta) &= \frac{1}{1-\gamma} \sum_{s' \in \mathcal{S}} \sum_{a \in \mathcal{A}_{s'}} d^{\pi^\theta}(s') \pi^\theta(a; s) \nabla \log(\pi^\theta(a; s')) Q^{\pi^\theta}(s', a) \\ &\quad - \frac{1}{1-\gamma} \sum_{s' \in \mathcal{S}} \sum_{a \in \mathcal{A}_{s'}} d^{\pi^\theta}(s') \pi^\theta(a; s) \nabla_\theta \log(\pi^\theta(s, a)) (Q^{\pi^\theta}(s', a) - \mathbf{Q}_w^{\pi^\theta}(s', a)) \\ &= \frac{1}{1-\gamma} \sum_{s' \in \mathcal{S}} \sum_{a \in \mathcal{A}_{s'}} d^{\pi^\theta}(s') \pi^\theta(a; s) \nabla \log(\pi^\theta(a; s')) \mathbf{Q}_w^{\pi^\theta}(s', a) \end{aligned}$$

□

On first sight the theorem does not look useful at all as the conditions on \mathbf{Q} are very special. Here are two examples that show that the assumptions are not completely artificial. They are satisfied for linear function approximations.

Example 6.8.2. Recall the linear softmax policy

$$\pi^\theta(a; s) = \frac{e^{\theta^T \cdot \Phi(s, a)}}{\sum_{a' \in \mathcal{A}} e^{\theta^T \cdot \Phi(s, a')}}$$

with feature vectors $\Phi(s, a)$ for the state-action pairs. The linear softmax has score function

$$\Phi(s, a) - \sum_{a'} \pi^\theta(a'; s) \Phi(s, a').$$

thus, requiring linear function approximation

$$\mathbf{Q}_w^{\pi^\theta}(s, a) = w^T \cdot \left(\Phi(s, a) - \sum_{a'} \pi^\theta(a'; s) \Phi(s, a') \right).$$

What does it mean? \mathbf{Q} must be a linear function with the same features as the policy, except normalized to be mean zero for each state.

Example 6.8.3. A similar example can be obtained for a continuous action space. Suppose π^θ is a linear Gaussian policy, i.e.

$$\pi^\theta(a; s) \sim \mathcal{N}(\theta^T \cdot \Phi(s), \sigma^2)$$

for feature vectors $\Phi(s)$ and a constant σ . For the compatibility we first need to compute the score function:

$$\nabla \log \pi^\theta(s, a) = \nabla_\theta \frac{-(\theta^T \cdot \Phi(s))^2}{2\sigma^2} = \frac{\theta^T \cdot \Phi(s)}{\sigma^2} \Phi(s)$$

For f to satisfy the compatibility, as in the previous example \mathbf{Q} must be linear:

$$\mathbf{Q}_w^{\pi^\theta}(s, a) = w^T \cdot \frac{\theta^T \cdot \Phi(s)}{\sigma^2} \Phi(s).$$

Data: Initial parameter θ_0 and approximation class $\mathbf{Q}_w^{\pi^\theta}$.
Result: Approximate policy $\pi^\theta \approx \pi^{\theta^*}$
Set $n = 0$.
while *not converged* **do**
 Policy evaluation of critic: find a critical point w of the weighted approximation error,
 i.e.

$$\nabla_w \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}_s} d^{\pi^{\theta_n}}(s) \pi^{\theta_n}(a; s) (Q^{\pi^{\theta_n}}(s, a) - \mathbf{Q}_w^{\pi^\theta}(s, a))^2 = 0.$$

 Policy improvement of actor: chose step-size α and set

$$\theta_{n+1} = \theta_n + \alpha \frac{1}{1-\gamma} \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}_s} d^{\pi^{\theta_n}}(s) \pi^{\theta_n}(a; s) \nabla \log(\pi^{\theta_n}(a; s')) \mathbf{Q}_w^{\pi^\theta}(s, a).$$

end

Algorithm 24: Theoretical simple actor-critic algorithm

Taking account to the new representation of the policy gradient here is a first version of an actor-critic algorithm with provable convergence. We call the algorithm theoretical actor critic as the connection of critic computation and actor update is nothing else then a regular update of the policy gradient algorithm (Proposition 6.8.1). Thus, whenever the exact gradient ascent algorithm converges to a maximum (or a stationary point) the simple actor-critic will do the same. If implemented there are a couple of approximations entering the scene:

- functions Q_w will be used that do not satisfy the needed conditions,
- the approximation of Q^{π^θ} will produce an approximation error in every round,
- the gradient will be approximated as in the REINFORCE algorithm replacing the Q -function by the estimated Q -function $\mathbf{Q}_w^{\pi^\theta}$.

Each approximation will generate an error and a priori it seems difficult to make the algorithm converge.

In our previous discussions of baselines for policy gradients a good guess for the baseline seemed to be the value function as it reinforces actions that are better than the average and negatively reinforces actions that are worse then the mean. So far we did not follow upon the idea as the value function V^π is unknown, so how could it serve as a baseline of an implementable algorithm? Let's ignore this fundamental issue for a moment.



Definition 6.8.4. The advantage function of a policy is defined as $A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$.

The advantage function measures the advantage (or disadvantage) of first playing action a . Recalling that baselines can be state-dependent the policy gradient can be written as

$$\begin{aligned} \nabla J_s(\theta) &= \frac{1}{1-\gamma} \sum_{s' \in \mathcal{S}} \sum_{a \in \mathcal{A}_{s'}} d^{\pi^\theta}(s') \pi^\theta(a; s) \nabla \log(\pi^\theta(a; s')) (Q^{\pi^\theta}(s', a) - V^{\pi^\theta}(s')) \\ &= \frac{1}{1-\gamma} \sum_{s' \in \mathcal{S}} \sum_{a \in \mathcal{A}_{s'}} d^{\pi^\theta}(s') \pi^\theta(a; s) \nabla \log(\pi^\theta(a; s')) A^{\pi^\theta}(s', a). \end{aligned}$$

In what follows we will discuss what is known as A2C (advantage-actor-critic) algorithms in which similarly to the approximation of Q by linear functions the advantage function is approximated. All A2C algorithms consist of two alternating steps:

- Use the current policy π^θ to produce samples from which an approximation \hat{A}^{π^θ} of the unknown A^{π^θ} is derived.

¹⁶Bedeutung beider Eigenschaften ausführen

- Plug \hat{A}^{π^θ} into the policy gradient formula to obtain an estimate of the gradient which is used to update θ .

There are plenty of approaches to turn this into an algorithm. All of them use ideas from temporal difference learning. The simplest approach is as follows. As for 1-step temporal difference we write the advantage function as

$$A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s) = r(s, a) + V^\pi(s') - V^\pi(s), \quad s' \sim p(\cdot; s, a).$$

This shows that only the value function needs to be estimated. Now suppose there is a parametric family $V_w : \mathcal{S} \rightarrow \mathbb{R}$ that we intend to use to approximate all value functions V^{π^θ} . Typically, V_w will be given by a neural network with weight vector w . So how do we fit V_w to V^{π^θ} for some policy π^θ ? We approximate $V^{\pi^\theta}(s)$ using samples (Monte Carlo) and then minimise the error. More precisely, suppose y_1, \dots, y_n are the discounted rewards of n rollouts under policy π^θ , then we solve the regression problem

$$\min_w \sum_{i=1}^n \|V_w(s_0^i) - y_i\|^2.$$

If V_w is given by neural networks the regression problem is solved (approximatively) by back propagation. The solution V_w is then used to give an estimator \hat{A}^{π^θ} of A^{π^θ} from the gradient step is computed.

Data: Initial parameter θ and approximation functions V_w

Result: Approximate policy $\pi^\theta \approx \pi^{\theta^*}$

while not converged do

	Sample N starting points s_0^i from d^{π^θ} and using policy π^θ rollouts (s_0^i, a_0^i, \dots) .
	minimise $\sum_i \ V_w(s_0^i) - \sum_{t=0}^{\infty} \gamma^t R(s_{t+1}^i, a_{t+1}^i)\ ^2$ over w .
	set $\hat{A}(s_0^i, a_0^i) = R(s_0^i, a_0^i) + \hat{V}_w^\pi(s_1^i) - \hat{V}_w^\theta(s_0^i)$.
	set $\hat{\nabla}_\theta J(\theta) = \frac{1}{N} \sum_{i=1}^N \frac{1}{1-\gamma} \nabla \log(\pi^\theta(a_0^i; s_0^i)) \hat{A}(s_0^i, a_0^i)$.
	update $\theta = \theta + \alpha \nabla \pi^\theta(a; s')$.

end

Algorithm 25: A2C with Monte Carlo approximation

More advanced versions use n -step or TD(λ) approximations for V^{π^θ} to bootstrap samples. It is far from obvious if any of these versions converges to an optimal policy. ¹⁷

6.8.2 Soft actor-critic (SAC)

6.8.3 Proximal policy optimisation (PPO)

¹⁷bespreche sampling aus d^π durch stoppen nach geometrischer zeit. und bedeutung von $d^\pi \pi$

Chapter 7

Monte Carlo Tree Search

Chapter 8

Policy-Gradient Methods