

# **Autonomous Soaring Policy Initialization through Dynamic Programming**

**Benjamin Rothaupt**

**Masterarbeit**

**2018**



Universität Stuttgart  
**Institut für Flugmechanik und Flugregelung**



# **Erklärung**

Hiermit versichere ich, dass ich diese Masterarbeit selbstständig mit Unterstützung der Betreuer angefertigt und keine anderen als die angegebenen Quellen und Hilfsmittel verwendet habe.

Die Arbeit oder wesentliche Bestandteile davon sind weder an dieser noch an einer anderen Bildungseinrichtung bereits zur Erlangung eines Abschlusses eingereicht worden.

Ich erkläre weiterhin, bei der Erstellung der Arbeit die einschlägigen Bestimmungen zum Urheberschutz fremder Beiträge entsprechend den Regeln guter wissenschaftlicher Praxis<sup>1</sup> eingehalten zu haben. Soweit meine Arbeit fremde Beiträge (z.B. in der Form von Bilder, Zeichnungen, Textpassagen etc.) enthält, habe ich diese Beiträge als solche gekennzeichnet (Zitat, Quellenangabe) und eventuell erforderlich gewordene Zustimmungen der Urheber zur Nutzung dieser Beiträge in meiner Arbeit eingeholt. Mir ist bekannt, dass ich im Fall einer schuldhafoten Verletzung dieser Pflichten die daraus entstehenden Konsequenzen zu tragen habe.

---

Ort, Datum

---

Benjamin Rothaupt

**Betreuer: M. Sc. Stefan Notter**

<sup>1</sup>Nachzulesen in den DFG-Empfehlungen zur „Sicherung guter wissenschaftlicher Praxis“ bzw. in der Satzung der Universität Stuttgart zur „Sicherung der Integrität wissenschaftlicher Praxis und zum Umgang mit Fehlverhalten in der Wissenschaft“



## Abstract

Flying a glider can be regarded as a Markov Decision Process. At each point in time, control variables are set to values that influence how the airplane behaves. Dynamic Programming is a general term for optimization methods that decompose complex problems into subproblems and solve them individually, utilizing the solution of subproblems that have already been solved. A MDP can be decomposed by treating every decision individually and thereby rely on previous decisions and their returns. An advantage of DP over other reinforcement learning algorithms including trust region policy optimization and also optimal control is, that it inherently covers the entire state space and thus can also serve, if the initial state is not known (precisely) *a priori*.

In this thesis, dynamic programming is used to find a time optimal flight path through calm air. The state- and action-spaces are discretized and the discrete problems are solved with DP. Two policy iteration algorithms and one value iteration algorithm are compared. The results are used to train an artificial neural network with supervised learning. The resulting control policy is then optimized by Trust Region Policy Optimization.

The results show, that all three dynamic programming algorithms yield trajectories to the goal from all points in the state space, from where this is physically possible. The trajectories are - however - generally not time-optimal. Optimizing a pre-trained policy with TRPO yields better results than just performing TRPO with a Glorot-initialized neural network as a policy. Most of the policy outputs produced by DP are retained during TRPO. A combination of both can yield a time-optimal flight path from points near the initial state of TRPO, while finding the goal from most other states, and thus outperforms both individual algorithms.

## Kurzfassung

Das Fliegen eines Segelflugzeuges kann als Markov-Entscheidungsprozess (engl.: Markov Decision Process) betrachtet werden. Zu jedem Zeitpunkt werden Steuergrößen auf Werte gesetzt, die das Verhalten des Flugzeugs beeinflussen. Dynamic Programming ist ein Oberbegriff für Optimierungsmethoden, die komplexe Probleme in Teilprobleme zerlegen und individuell lösen, indem sie Erkenntnisse aus bereits gelösten Teilproblemen nutzen. Ein MDP kann zerlegt werden, indem man jede Entscheidung individuell betrachtet und sich dabei auf frühere Entscheidungen und deren Ergebnisse stützt. Ein Vorteil von DP gegenüber anderen Methoden des bestärkenden Lernens einschließlich TRPO und auch einer Optimalsteuerung ist, dass es von Natur aus den gesamten Zustandsraum abdeckt und somit auch verwendbar ist, wenn der Ausgangszustand *a priori* nicht (genau) bekannt ist.

In dieser Arbeit wird Dynamic Programming verwendet, um eine zeitoptimale Flugbahn durch ruhige Luft zu finden. Die Zustands- und Aktionsräume werden diskretisiert und die diskreten Probleme mit DP gelöst. Es werden zwei Policy Iteration Algorithmen und ein Value Iteration Algorithmus verglichen. Die Ergebnisse werden verwendet, um ein künstliches neuronales Netz mit überwachtem Lernen zu trainieren. Die daraus resultierende Steuerungs-Policy wird dann durch Trust Region Policy Optimization optimiert.

Die Ergebnisse zeigen, dass alle untersuchten Dynamic Programming Algorithmen Trajektorien ergeben, die von allen Punkten im Zustandsraum, von wo aus dies physikalisch möglich ist, zum Ziel führen. Die Flugbahnen sind jedoch in der Regel nicht zeitoptimal. Die Optimierung einer vortrainierten Policy mit TRPO führt zu besseren Ergebnissen als die Durchführung von TRPO mit einem Glorot-initialisierten neuronalen Netzwerk als Policy. Die meisten der von DP erzeugten Policy-Outputs bleiben während des TRPO erhalten. Eine Kombination aus beidem kann eine zeitoptimale Flugbahn von Punkten nahe dem Ausgangszustand von TRPO ergeben, während das Ziel von den meisten anderen Zuständen ebenfalls gefunden wird, und übertrifft somit beide individuellen Algorithmen.

# Contents

<b>Erklärung</b>	iii
<b>Abstract/Kurzfassung</b>	v
<b>List of Figures</b>	ix
<b>List of Tables</b>	xi
<b>Nomenklatur</b>	xiii
<b>1. Introduction</b>	1
<b>2. Solution Methods for Markov Decision Processes</b>	5
2.1. Markov Decision Process . . . . .	5
2.2. The Principle of Optimality . . . . .	5
2.3. The Agent - Environment System . . . . .	6
2.4. Model Based and Model Free Learning . . . . .	6
2.5. Agent . . . . .	8
2.6. Return . . . . .	8
2.7. Policy . . . . .	9
2.8. Value Functions . . . . .	10
2.9. The Bellman Equation . . . . .	10
2.10. Solving an MDP . . . . .	11
2.11. Dynamic Programming . . . . .	12
2.11.1. The Curses of Dimensionality . . . . .	13
2.11.2. Approximate Dynamic Programming . . . . .	14
2.11.3. Types of Dynamic Programming Algorithms . . . . .	16
2.11.4. The Contraction Mapping Theorem . . . . .	20
2.12. Exploiting the Problem Structure . . . . .	21
<b>3. Function Approximation</b>	25
3.1. Tables . . . . .	25
3.2. Artificial Neural Networks . . . . .	26
3.3. Supervised Learning . . . . .	30
3.4. Optimization Techniques . . . . .	31
3.4.1. Gradient Descent . . . . .	31

## *Contents*

3.4.2. Stochastic Gradient Descent . . . . .	31
3.4.3. The Adam Algorithm . . . . .	32
3.5. Overfitting . . . . .	33
<b>4. Trajectory Optimization with Dynamic Programming</b>	<b>37</b>
4.1. 2D Environment . . . . .	37
4.1.1. Policy Representation . . . . .	37
4.1.2. Equations of Motion . . . . .	37
4.1.3. Reward Function . . . . .	39
4.1.4. Discretization of the State- and Action Space . . . . .	39
<b>5. Results</b>	<b>43</b>
5.1. 2D Dynamic Programming . . . . .	43
5.1.1. Optimal Control Benchmark . . . . .	43
5.1.2. Policy Initialization for TRPO . . . . .	48
<b>6. Discussion</b>	<b>55</b>
<b>A. Appendix</b>	<b>57</b>
A.1. Glider Parameters . . . . .	57
A.2. Computer Configuration and Implementation . . . . .	58
A.2.1. Computer Configuration . . . . .	58
A.2.2. Computational Subtleties of PI and VI . . . . .	59
A.3. Illustration of a One Step Lookahead . . . . .	60
A.4. Optimistic Policy Iteration . . . . .	61
A.5. Table Representation of the Value Function and Policy . . . . .	62
A.6. Scenario 1: Additional Results . . . . .	63
<b>Literaturverzeichnis</b>	<b>65</b>

# List of Figures

2.1. The Agent-Environment-System [12] . . . . .	7
2.2. MDP solution methods, adopted from [2]) . . . . .	12
2.3. Dynamic programming as a reinforcement learning method (taken from [17])	15
2.4. The policy iteration algorithm (adopted from [17, section 4.6]) . . . . .	19
3.1. A neuron in an artificial neural network . . . . .	26
3.2. Activation functions (adopted from [20]) . . . . .	27
3.3. The policy MLP [20] . . . . .	28
3.4. Overfitting in case of a high order polynomial and noisy data . . . . .	34
3.5. Stop criterion to avoid overfitting . . . . .	35
4.1. Reference frame transformation [12] . . . . .	38
4.2. The discretized state space in all 2D scenarios. . . . .	40
5.1. Results of synchronous value iteration and optimal control in scenario 1 .	46
5.2. Results of value iteration and optimal control in scenario 1 with a fine grid	47
5.3. Results of value iteration and optimal control in scenario 2 . . . . .	48
5.4. Average return per iteration with and without policy initialization . . .	49
5.5. Maximum return per iteration with and without policy initialization . .	49
5.6. VI - trajectories from various initial states with zero velocity . . . . .	51
5.7. VI - trajectories from various initial states with ${}_g u_K = -20 \frac{m}{s}$ . . . . .	51
5.8. TRPO - trajectories from various initial states with zero velocity . . . .	52
5.9. TRPO - trajectories from various initial states with ${}_g u_K = -20 \frac{m}{s}$ . . . .	52
5.10. VI and TRPO - trajectories from various initial states with zero velocity	53
5.11. VI and TRPO - trajectories from various initial states with ${}_g u_K = -20 \frac{m}{s}$	53
A.1. A one step lookahead in the 2D environment. . . . .	60
A.2. A four-dimensional table. . . . .	62
A.3. Results of synchronous generalized policy iteration and optimal control in scenario 1 with a coarse grid . . . . .	63
A.4. Results of synchronous optimistic policy iteration and optimal control in scenario 1 with a coarse grid . . . . .	64



# List of Tables

4.1. Grid parameters for trajectory optimization . . . . .	42
5.1. Some information about the investigated scenarios . . . . .	43
5.2. Scenario 1, comparison of synchronous algorithms on a coarse grid . . . . .	45
5.3. Scenario 1, comparison of asynchronous algorithms on a coarse grid . . . . .	45
5.4. Scenario 2, performance of asynchronous VI on a fine grid . . . . .	45
5.5. Flight- and computation-times for OC and asynchronous VI in scenario 2	46
5.6. Flight times of agents with three different policy initialization techniques and different start states . . . . .	50
A.1. glider data . . . . .	57
A.2. computer specifications . . . . .	58



# Nomenclature

## Abbreviations

ANN	Artificial Neural Network
CPU	central processing unit
DP	Dynamic Programming
GPI	Generalized Policy Iteration
GPU	graphics processing unit
iFR	Institute for Flight Mechanics and Control
MDP	Markov Decision Process
MLP	multi layer perceptron
MSE	mean squared error
OC	optimal control
OPI	Optimistic Policy Iteration
PE	Policy Evaluation
PI	Policy Iteration
RAM	random access memory
RL	Reinforcement Learning
SL	Supervised Learning
VI	Value Iteration
VRAM	video random access memory

## Latin Letters

$\bar{e}$	mean error
$\mathbf{b}$	bias
$\mathbf{W}$	weights
$\Delta t$	time step
$\mathcal{A}$	action space
$\mathcal{N}$	Gaussian Distribution
$\mathcal{P}$	State Transition Probability

## List of Tables

$\mathcal{S}$	state space
$\rho$	$[\frac{\text{kg}}{\text{m}^3}]$ air density
$a$	action
$e$	error
$m$	[kg] glider mass
$o$	observation
$Q$	Action Value Function
$q$	$[\frac{\text{N}}{\text{m}^2}]$ dynamic pressure
$r$	reward
$s$	state
$V$	State Value Function
D	drag
L	lift

## Greek Letters

$\alpha$	[rad]	angle of attack
$\theta$		parameter vector
$\gamma$		discount factor
$\Lambda$		aspect ratio
$\mu$		mean
$\phi$	[rad]	flight path angle
$\pi$		policy
$\sigma$		standard deviation
$e$		oswald efficiency factor

## Indices

*	optimal
$\pi$	w.r.t. policy $\pi$
$T$	terminal
$t$	time step
D	drag
L	lift

# 1. Introduction

The goal of this work is to find an optimal trajectory of a glider moving through calm air and reaching a given distance in minimal time. The resulting optimal policy can be used as a starting point for training the glider on scenarios with different wind conditions utilizing other algorithms like Reinforcement Learning (RL). The basic idea is to use the knowledge about the glider dynamics to pre-train a policy. After that, an RL algorithm is used to train the policy on scenarios with the agent facing a stochastic environment.

In soaring competitions, there are two main goals that pilots must be able to achieve in order to be successful: exploiting updrafts to gain potential energy and covering a given distance in minimal time. These are two conflicting tasks. The longer a pilot decides to stay within an updraft, the more energy can be harvested. The time spent there however increases the total flight time which interferes with the other goal. Generally, the locations where updrafts occur are not known a-priori. So a pilot has to react spontaneously to changes in his vicinity in order to harvest as much energy as possible, which takes a lot of experience. Automation of the optimal-flying task is therefore challenging.

Another topic, where using energy from updrafts can be helpful, is electric flying. Electric aircraft have been an important research topic in recent years. They can operate environmentally friendly and are easy to maintain. Their range and endurance is - however - limited by the energy density of current batteries. This means they either have to land frequently to recharge or be equipped with solar panels on their wings which weigh them down and are difficult to maintain. Electric UAVs suffer from the same deficiency. They can theoretically be used for a wide range of applications from gathering scientific data to surveillance with battery capacity being their only limiting factor. The exploitation of thermal updrafts is one way to increase the endurance of an aircraft without changing anything on the airframe itself, which makes it relatively cheap to implement.

Reacting to changes in the environment is a topic that is inherently covered by various Machine Learning algorithms. In fact, the most challenging and therefore interesting and powerful aspect of *learning* is gaining the ability to generalize and behave reasonably in new situations, such as unknown updraft distributions. This sets Machine Learning apart from traditional optimal control algorithms that rely on precise knowledge of the underlying problem.

Finding an optimal flight path with respect to varying mission goals in unknown atmosphere is a challenging task. While the dynamics of the aircraft itself are known and understood, the occurrence of updrafts cannot be predicted with sufficient accuracy. RL algorithms are very good at finding an optimal path through an unknown environment. They require no prior knowledge and learn from experience they gather through inter-

## 1. Introduction

acting with the environment. The stochastics subject to the likelihood of occurrence of thermal updrafts and their specific characteristics are hard to model. Sampling from experience constitutes a way to learn how to act, while facing (hidden) stochastics within the environment. This way, RL algorithms can deal with (apparently) random updrafts efficiently.

In [14], the authors combined simulations of atmospheric flow with reinforcement learning algorithms to find strategies for gliders to cope with and exploit turbulent thermal updrafts. They utilized model-free algorithms to train the agents. The resulting policies were able to find thermal updrafts in a 3D environment and exploit them more or less conservatively, depending on the intensity of the turbulence.

Markus Zuern applied Trust Region Policy Optimization (TRPO) and Asynchronous Advantage Actor Critic (A3C) to a trajectory optimization problem [20] (some of the topics of [20] are also covered in [12]). In his work, both algorithms deal with the complete environment, i.e. the glider dynamics and the stochastic wind distribution, at the same time. They require no knowledge about the environment, nor do they try to learn a model. Instead, they directly optimize a policy. In TRPO, there does not even exist a value function that represents known information about the MDP. Although the results in [20] are remarkable, they do not exploit the fact that one part of the trajectory optimization MDP is known *a priori*: the glider dynamics when flying through calm air.

The goal of this work is to utilize this knowledge by applying Dynamic Programming to calculate a time-optimal trajectory through calm air. This is achieved by first finding an optimal policy and following it. The resulting policy is then trained with unknown wind distributions. This is similar to the behavior one would expect from a human pilot. If he/she is confronted with unknown updrafts, he/she would at first also rely on his experience from prior flights, be it through calm air or not. Dynamic Programming (DP) is one of the most popular algorithms for MDPs with a state space that is not too large. Typically, MDPs with up to a few millions of states can be solved with DP algorithms.

Dynamic Programming is a general term that summarizes a certain kind of algorithms. These algorithms can be used to solve problems that can be split into subproblems. The optimal solution of each subproblem can then be used to recompose the optimal solution to the whole problem. The underlying principle is called the *principle of optimality*. The term was first used by Richard E. Bellman in 1954 [3]. He proved that various optimization problems could be solved efficiently by making use of the principle of optimality instead of traditional optimal control theory. In his book *Applied Dynamic Programming*, Bellman covered numerous examples of optimal control problems and shows how they can be solved by Dynamic Programming [4]. If the subproblems are similar, DP is especially effective. DP however suffers from the "curse of dimensionality", that is in large or continuous state or action spaces it quickly becomes inefficient [13]. The obvious approach is to discretize the state space and deal with a smaller number of states. If done properly, this yields an approximation of the exact optimal solution. The achieved approximate solution is close to the exact one if the discretization is sufficiently fine. Like in many technical topics, there is a tradeoff between precision and cost.

Dynamic Programming algorithms have already been used to calculate optimal trajectories for airborne agents. In [19], dynamic programming is used to calculate an optimal trajectory for a solar powered UAV across the atlantic ocean in consideration of meteorological conditions. (mehr Bsp.)

In this work, a time optimal flight path through calm air is calculated with dynamic programming. First, the state and action space are discretized in section 4.1.4. A policy iteration (PI) algorithm is then used to calculate optimal paths in the vertical plane and in a 3D environment. The usefulness of the results from the 3D scenario is obvious. If airspace restrictions or topography restrict the horizontal trajectory, any 3D scenario degenerates to a 2D optimization of a vertical trajectory along the predefined ground trace. Therefore, the results of optimization in the vertical plane shown in section ?? are also relevant.

The first part of this thesis presents the general concepts and terms of Reinforcement Learning (cf. chapter 2) and the theory of Dynamic Programming and the algorithms used in this work (cf. chapter 2.11). Chapter 3 gives a brief overview of artificial neural networks (ANNs) and how training an ANN works. The trajectory optimization problem is described in chapter 4 and the results are presented and compared to an optimal control result in chapter 5. The last chapter contains a brief discussion of the results and how using an ANN, that has been pre-trained by DP, for TRPO affects the results.



## 2. Solution Methods for Markov Decision Processes

### 2.1. Markov Decision Process

Path planning through calm air can be done analyzing how a glider behaves in calm air. Flying a glider can be regarded as a decision process where the pilot decides what control surface position is suitable for the current situation. By varying, for example, the elevator position, the pilot can control the speed of the glider and thereby the flight path angle. In reality, this process is time-continuous. Computers - however - can only deal with discrete processes. Therefore it is common practice to introduce a time step  $\Delta t$ . At each moment in time  $t_k = t_0 + k\Delta t$ , an elevator position is chosen and is kept constant until  $t_{k+1}$ . The smaller the time step, the better the approximation of a continuous process. Each state  $s_t$  in this decision process has the Markov property, i.e. the information needed to calculate  $s_{t+1}$  is fully contained in  $s_t$ . In other words, the probability of reaching a state  $s_{t+1}$  by taking an action  $a_t$  only depends on the state  $s_t$ . The predecessor states of  $s_t$  are irrelevant for  $P(s_{t+1}|s_t, a_t)$ . A decision process where the Markov property applies to all states is called a Markov Decision Process (MDP).

$$P(s', r|s, a) = \mathbb{P}[s_{t+1} = s', r_{t+1} = r | s_t = s, a_t = a] \quad (2.1)$$

$$= \mathbb{P}[s_{t+1} | s_t, a_t, s_{t-1}, a_{t-1}, s_{t-2}, a_{t-2}, \dots] \quad (2.2)$$

### 2.2. The Principle of Optimality

According to Bellman's principle of optimality, the solution of some optimization problems can be put together from the solution of subproblems. This is the basis for all dynamic programming algorithms. If an agent chooses the optimal action in a state  $s_t$  and all states it visits afterwards ( $s_k, k = t + 1, t + 2, \dots, T$ ), the resulting trajectory is the optimal one from  $s_t$  to the terminal state  $s_T$ .

The principle of optimality can be expressed by the Bellman Equation, which is shown in section 2.9. Prior to introducing the Bellman Equation, a few necessary concepts and terms are explained in the following sections.

(Bsp. Fibonacci Folge und/oder shortest path)

## 2.3. The Agent - Environment System

In Reinforcement Learning scenarios, an agent interacts with its environment and thereby learns how to maximize some sort of scalar reward signal. The agent aims to maximize its total reward by choosing - at each time step - the best possible action  $a$  from a set of actions  $\mathcal{A}$ . The agent has no prior knowledge about the environment and about how to behave in an optimal way. Instead, it learns by interacting with its environment and the information it receives during training. Before each step, the agent is in a state  $s_t$ , chooses an action  $a_t$  and receives a scalar reward  $r_t$  and an observation  $o_t$  representing the next state  $s_{t+1}$ . The process is repeated, then. Figure 2.1 illustrates what happens at each time step. Note that, unlike in this thesis, a timestep is denoted by an index  $k$  instead of  $t$ . The observation  $o_t$  can just be  $s_{t+1}$ . However, in partially observable environments, it might only be a part of  $s_{t+1}$  or any information regarding the true state  $s_{t+1}$ . In this work,  $o_t$  is equal to  $s_{t+1}$  (i.e. the agent has complete knowledge about his state at any time).

The reward  $r_t$  the agent receives at each step is determined by a reward function  $\mathcal{R}(s, a)$ . The reward function maps from states (and sometimes actions) to a scalar reward and characterizes the mission target. The choice of a reward function is crucial for the success of any RL-based trajectory optimization algorithm. The only way to "tell" the agent, what his goal is, is by defining the reward function.

Generally, a reward function should guide the agent to the target. There are, however, restrictions on what a reward function should contain. If too much information about the problem is put into the reward function, this might lead to a suboptimal result. If, for example, a user trying to solve an MDP thinks a specific terminal state might be better than the others, he might be tempted to increase the final reward for reaching this - seemingly optimal - terminal state. But any (unnecessary) assumption, that is put in the reward function, can restrict the agent's ability to find an optimal trajectory on his own. A false assumption can sometimes even jeopardize convergence of the whole algorithm. Therefore, a good reward function should only give a reward for reaching a target state and penalize every step the agent does, if time is crucial. It can also penalize any step that leads to an undesired terminal state (in this work, touching the ground is such a case). This way, the agent can find a good trajectory without restrictions. More on the reward function used here can be found in section 4.1.3.

## 2.4. Model Based and Model Free Learning

MDP solution methods can be divided into model-free and model-based methods. In model-based methods, the agent uses a model to predict the reaction of the environment to any of his actions. An environment model typically takes a state-action pair and returns the next state and next reward. If the environment is discrete and stochastic,

## 2.4. Model Based and Model Free Learning

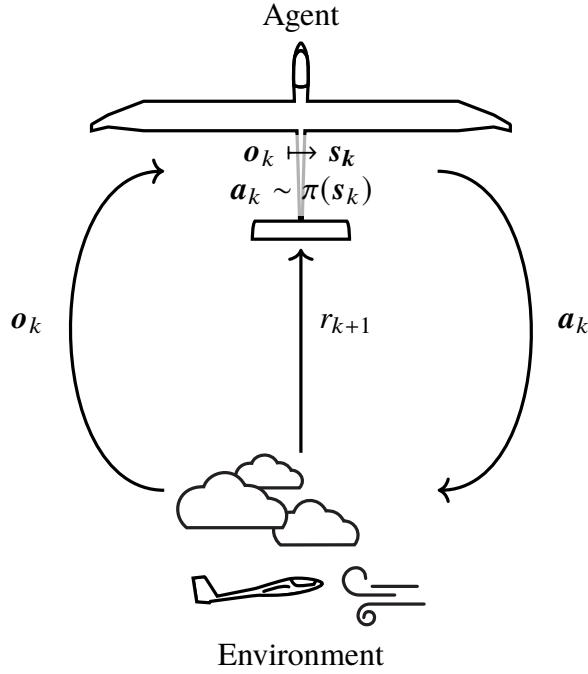


Figure 2.1.: The Agent-Environment-System [12]

there are multiple possible next states and rewards.<sup>1</sup> A *sample model* returns one of the possibilities, whereas a *distribution model* returns some representation of all possible next states and rewards. To distinguish experience from a model from real experience, the results from a model are referred to as *simulated experience* [17, section 8.1].

Regardless of the type, all models approximate what might happen to the agent in future time steps. Whether the experience gathered through a model is useful obviously depends on its quality. Every model adds complexity to an algorithm. The more complexity, the higher the cost of implementing and updating it. Model free methods are generally simpler to implement and more scalable.

In model free solution methods, there is no model of the environment. The agent is not explicitly aware of the environment behavior. Instead, the policy (and value function) are directly affected by the experience of the agent. Sutton and Barto refer to model free methods as some sort of "trial and error" (c.f. [17, section 1.5]). If the goal is to maximize the expected return, a model might not be needed at all. In TRPO, for example, the policy is optimized by maximizing the expected return with respect to the policy parameters  $\theta$  by sampling trajectories through the state space and utilizing a gradient based optimization method. This is simpler than learning a model of the environment and using it to plan.

A model of the environment is - in principle - not required for successful training. It

---

<sup>1</sup>In a continuous stochastic environment, there would be a distribution of successor states and rewards.

## 2. Solution Methods for Markov Decision Processes

can however be used to gather simulated experience (c.f. section 2.4) and therefore reduce training time on a real agent like a robot. Although imperfect, information from a model can be useful in such cases where training time is restricted which particularly holds true for robots flying outdoors.

For more information on how machine learning algorithms can be classified, please refer to chapter 2.11, which also contains one classification example in Fig. ??.

## 2.5. Agent

The agent in an MDP consists of at least one of the following parts.

- Policy:  
A mapping from states to actions that tells the agent what to do (c.f. section 2.7)
- Value Function:  
A mapping from states (or state-action pairs) to values (c.f. section 2.8)
- Model:  
A model of the environment that might be updated by the agent's experience

Most agents implement a policy. They also can have a value function, like in actor-critic algorithms, where the value function is used for so called bootstrapping <sup>2</sup>.

$$\mathbb{E}[G_t | s_t = s] = \mathbb{E}[r_{t+1} + \gamma G_{t+1} | s_t = s, a_t = a] \quad (2.3)$$

$$= \mathbb{E}[r_{t+1} + \gamma V(s_{t+1}) | s_t = s, a_t = a] \quad (2.4)$$

Instead of sampling a complete episode, bootstrapping makes it possible to learn from incomplete episodes.

## 2.6. Return

The return  $G_t$  at a time step  $t$  is the cumulative reward  $r_t$  at each step from the state  $s_t$  onwards until reaching a terminal state:

$$G_t = \sum_{k=0}^{T-t-1} \gamma^k r_{t+k+1} \quad (2.5)$$

---

<sup>2</sup>Bootstrapping is a general term, that applies to methods, that make use of previous results in order to obtain a new result. In this case, it means utilizing the Bellman Expectation Equation to estimate the return  $G_t$  by doing a one step lookahead. The state value  $V(s)$  is then updated utilizing  $V(s')$ , which can be seen as a "previous" result.

If  $\gamma \in (0, 1]$  is close to zero, immediate rewards are weighted stronger. If  $\gamma$  is close to one, rewards in the distant future are weighted likewise, making decisions more far sighted. In infinite horizon problems, where there is no terminal state,  $\gamma$  must be less than one in order to avoid infinite returns. In finite horizon (i.e. episodic) MDPs,  $\gamma$  is usually close to, or exactly one.

Each decision at a given state  $s_t$  also affects what the next state  $s_{t+1}$  is and therefore the next reward. More on the reward and what role it plays in Dynamic Programming is explained in section 4.1.3.

## 2.7. Policy

In reinforcement learning problems, a policy is usually the part that represents the agent. It contains all the information needed to decide what action to choose in each state the agent visits. A policy is a mapping from states to actions.

Generally, a policy is a mapping from states to actions. A policy can also be stochastic, in which case it maps from states and actions to a number between 0 and 1, which represents the probability to pick  $a$  when in  $s$ .

$$\pi : \mathcal{S} \times \mathcal{A} \mapsto [0, 1] \quad (2.6)$$

If the policy is deterministic, only one  $a_t$  out of  $\mathcal{A}$  has a probability of one with all other actions assigned to probability zero. In any policy, the sum of all possible action-probabilities must be one.

$$\pi(a|s) = \mathbb{P}[a_t = a | s_t = s] \quad (2.7)$$

This yields equations 2.8 and 2.9 for the state transition probability and the reward function under policy  $\pi$ . The reward function is explained in section 4.1.3.

$$\mathcal{P}_\pi(s, s') = \sum_a \pi(a|s) \mathcal{P}(s'|s, a) \quad (2.8)$$

$$\mathcal{R}_\pi(s) = \sum_a \pi(a|s) \mathcal{R}(s, a) \quad (2.9)$$

At each timestep, the agent receives a state observation  $o_{t-1}$  from the environment, passes it to the policy and the policy returns an action  $a_t$ . This action is taken and the agent gets a new observation  $o_t$  and a reward  $r_t$ . This process is repeated until the agent reaches a terminal state.

A policy can be implemented in many different ways. The simplest option is a table where each field contains the action for one state or - if the policy is stochastic - all possible actions and their respective probabilities. Other examples include linear combinations of the inputs and linear combinations of basis functions. The most popular way to implement a policy is - however - through an artificial neural network (ANN). More information about neural networks is given in chapter 3.

## 2.8. Value Functions

Some RL-algorithms try to find a policy directly from experience. Others, like Q-learning, additionally utilize a so called value function. A value function maps from states or state-action pairs to their values. In this context, the value of a state or an action is a measure of how useful it is for the agent to visit the respective state or choosing the respective action. There are two types of value functions, state value functions and action value functions. The state value  $V(s)$  is the expected total return from state  $s_t$  onwards until the episode terminates. The action value also takes into account what the agent chooses to do in state  $s_t$ . It is the expected total return from state  $s_t$  when choosing action  $a_t$ .

$$V(s) = \mathbb{E}[G_t | s_t = s] \quad (2.10)$$

$$Q(a|s) = \mathbb{E}[G_t | s_t = s, a_t = a] \quad (2.11)$$

In the following sections, the action space is regarded finite for simplicity. The state value function 2.10 can be expressed in terms of the action value function 2.11:

$$V_\pi(s) = \sum_a \pi(a|s) Q_\pi(s, a) \quad (2.12)$$

Similarly, the action value function can be expressed by means of the state value function:

$$Q_\pi(s, a) = \mathcal{R}(s, a) + \gamma \sum_{s'} \mathcal{P}(s'|s, a) V_\pi(s') \quad (2.13)$$

If an agent has found the optimal value function  $V^*$  or  $Q^*$  of an MDP, the optimal policy  $\pi^*$  can be derived directly by acting greedily with respect to the value function at each state.

Like a policy, a value function may be implemented as a table, a linear combination of the input features, a combination of basis functions, or an artificial neural network.

## 2.9. The Bellman Equation

One way to find an optimal trajectory is by utilizing a value function  $V(s)$ , which maps from a state (or a state-action pair) to a state (or action) value. It tells the agent how good it is to be in a specific state in terms of the expected return. The optimal value function tells the agent the maximum return he can collect from every state.

According to the Bellman Expectation Equation, every trajectory through the state space can be decomposed into two parts: the next step and the rest of the path. This principle also applies to value functions, which are introduced in the following sections.

The state value function (c.f. Eq. 2.10) can be written as a weighted sum of the action values, each multiplied with the probability to take the respective action according to

the policy. Replacing  $Q_\pi(s, a)$  in equation 2.12 with equation 2.13 yields equation 2.15, which is known as the *Bellman Expectation Equation*:

$$V_\pi(s) = \sum_a \pi(a|s) Q_\pi(s, a) \quad (2.14)$$

$$= \sum_a \pi(a|s) \sum_r \sum_{s'} \mathcal{P}(s', r|s, a) [r + \gamma V_\pi(s')] \quad (2.15)$$

In a deterministic environment, equation 2.15 becomes

$$V_\pi(s) = \sum_a \pi(a|s) [r + \gamma V_\pi(s')] \quad (2.16)$$

The Bellman Expectation Equation can be used to update a state value by looking at the state values of its successor states and weighting them by the probabilities to reach each successor state (c.f. equations 2.30 and 2.31).

Instead of a weighted sum of rewards and successor values, like in equation 2.15, one can take only the action that has the maximum action value. The resulting state value  $V_*(s)$  is called the *optimal state value*. It is the maximum expected return from state  $s$  onwards. This yields the following equations.

$$V_*(s) = \max_{a \in \mathcal{A}(s)} Q_{\pi_*}(s, a) \quad (2.17)$$

$$= \max_a \mathbb{E}_{\pi_*}[G_t | s_t = s, a_t = a] \quad (2.18)$$

$$= \max_a \mathbb{E}_{\pi_*}[r_{t+1} + \gamma G_{t+1} | s_t = s, a_t = a] \quad (2.19)$$

$$= \max_a \mathbb{E}[r_{t+1} + \gamma V_*(s_{t+1}) | s_t = s, a_t = a] \quad (2.20)$$

$$= \max_a \sum_{s', r} \mathcal{P}(s', r | s, a) [r + \gamma V_*(s')] \quad (2.21)$$

Again, if the environment is deterministic, equation 2.21 can be simplified:

$$V_*(s) = \max_a [r + \gamma V_*(s')] \quad (2.22)$$

All equations from 2.17 to 2.22 are alternative ways to express the *Bellman Optimality Equation*. It states that the value of a state under an optimal policy must equal the expected return for the best action in that state [17]. Intuitively it is obvious that the optimal policy must pick the best action in each state. In Value Iteration, the *Bellman Optimality Equation* is used as an update rule. See chapter 2.11.3 for more details.

## 2.10. Solving an MDP

A Markov Decision Process is a series of decision problems, where each decision has an impact on the total return the agent receives. In most RL-problems, the goal is to

## 2. Solution Methods for Markov Decision Processes

maximize the total return from a given initial state. This is done by trying to find the optimal policy  $\pi^*$  which yields the maximum expected return from each state.  $\pi^*$  can be achieved directly or by finding the optimal value function, i.e. the mapping from each state or state-action-pair to its true value  $V^*$ . The definition of a state value is given in section 2.8 If  $\pi^*$  is found for a specific MDP, the MDP is solved.

### 2.11. Dynamic Programming

This thesis deals with Dynamic Programming for policy optimization. More specifically, policy iteration and value iteration are applied to a trajectory optimization problem. In a prior work on this topic, TRPO and A3C were applied to a 2D trajectory optimization problem [20]. Both algorithms deal with the complete environment, i.e. the glider dynamics and the stochastic wind distribution, at the same time. Whereas the wind distribution is unknown a priori, the glider dynamics are well known. Therefore, a DP algorithm can be applied to any scenario with calm air. The resulting optimal policy  $\pi_{\text{calm}}^*$  then can be used as a starting point for TRPO in scenarios with unknown wind conditions.

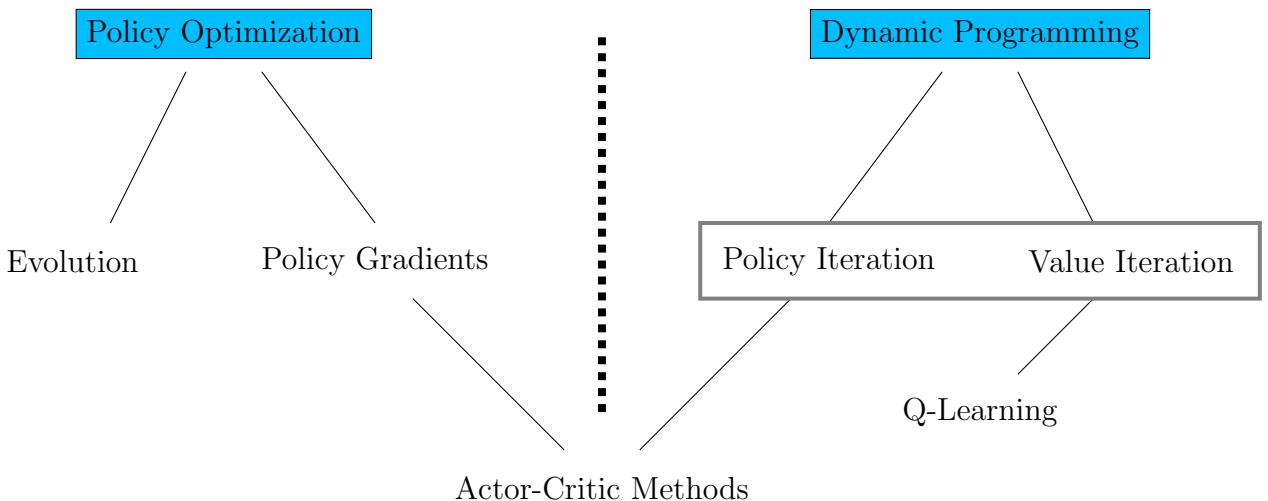


Figure 2.2.: MDP solution methods, adopted from [2]]

Dynamic programming is a general term, that is used for algorithms that solve a specific type of problem. Those problems can be decomposed into subproblems and solved by composing the solutions of those subproblems. If the subproblems are very similar, the solution of one subproblem can be used to calculate another one. In such cases, dynamic programming is applicable. In MDPs, these subproblems are similar in the sense that the Bellman Expectation Equation and the Bellman Optimality Equation hold true in all states, i.e. each state (and action) value is connected to its successor values by the Bellman Equations.

Also, the optimal action is found in all states the same way, as can be seen in section 2.11.3. If the algorithm is successful, the policy is optimal according to section 2.2. Unlike with TRPO, optimization is not done through sampling complete trajectories, but by iterating the state-values and policy for each state independently. After one sweep over the state space, the value function and policy are updated and the process is repeated.

In theory, starting TRPO or A3C with a policy that is already optimal with respect to the scenario in calm air should accelerate training with updrafts because the information about how to behave in a situation with no wind is already incorporated in the policy.

### 2.11.1. The Curses of Dimensionality

The following section explains the term *curses of dimensionality* in the context of discrete dynamic programming. It is adopted from [13, section 1.2].

As mentioned in the previous sections, dynamic programming aims at solving a complicated problem by decomposing it into simpler subproblems. Instead of solving a big problem at once, many small problems are solved one at a time. The time it takes to solve a specific number of subproblems obviously depends on the number of subproblems itself. In this work, estimating the value of a state or finding the optimal action in a state is such a subproblem. In the case of a discrete state vector  $s \in \mathbb{R}^N$ , where each component can have  $K$  different values, the number of subproblems is  $K^N$ . This number can quickly get very large depending on the number of dimensions, rendering application of DP to the underlying problem infeasible. This is referred to as the *first curse of dimensionality* by Warren B. Powell [13]<sup>3</sup>.

The *second curse of dimensionality* according to [13] can occur in problems, where randomness plays a role. An agent moving through a state space might encounter random winds or other events that influence his successor state  $s_{t+1}$ , if he is in the state  $s_t$  and chooses an action  $a_t$ . If the random variable depicting  $s_{t+1}$  has  $J$  dimensions and each of its components can take up to  $M$  values, the number of possible successor states is  $M^J$ .

In a discrete environment, where the agent can choose  $L$  actions and each can take  $I$  values, the number of possible combinations of these  $L$  actions is  $I^L$ . Even with only two actions (e.g. a glider controlling its bank angle  $\mu$  and its angle of attack  $\alpha$ ) and a discretized action space, where each action can take 20 different values, there are already 400 combinations. This is referred to in [13] as the *third curse of dimensionality*.

The second curse of dimensionality only occurs in stochastic environments, the two other curses can occur in any environment, where the number of actions and/or states is high.

---

<sup>3</sup>Commonly, the first curse of dimensionality is also referred to just as the *curse of dimensionality*.

## 2.11.2. Approximate Dynamic Programming

Dynamic programming has been utilized to solve problems in various disciplines since its invention. Most of the applications were in finance, computer science and control theory. The latter is also the topic of this thesis. Unfortunately, the descriptions of the relationship between RL and DP vary from discipline to discipline and seem inconsistent at first. This is because many of the principles of dynamic programming were discovered by people from various disciplines and each discipline developed its own terms to describe DP algorithms. This section contains some information about important keywords and methods in DP and RL and tries to bring the different viewpoints of control theory, computer science and finance together. For each discipline, the relevant nomenclature is shown and some sources containing further information are given. Then, the seemingly inconsistent viewpoints are brought together.

All DP applications have in common, that they deal with sequential decision processes and aim to find an optimal decision rule. Optimality is defined by either a cost function or a reward function which have to be minimized or maximized, respectively. The differences mainly arise in nomenclature and what problems are solved by DP in the respective field.

A problem, that occurs in finance, is the allocation of a budget in order to maximize (a measure of) profit. In [13, section 2.1.2 f.], there are examples of a discrete and continuous budgeting problem. In these examples, a resource  $x$  is allocated to a task  $t \in \mathcal{T}$ , where  $\mathcal{T}$  is the set of possible tasks. The state variable, which is called  $s_t$  in this thesis, is called  $R_t$  by Powell. It describes how much of the resources  $R$  (cf. Eq. 2.23) is left to assign to a task  $t$ . The goal is to maximize  $\sum_{t \in \mathcal{T}} C_t(x_t)$  with respect to  $x$ .  $C$  stands for contribution, which is equivalent to the *reward* in this work. Where, in this work, the goal is to maximize the sum of rewards, the example from [13] aims at maximizing a sum of contributions. The only difference is nomenclature. The value of having a specific amount of resources to allocate to task  $t$  is  $V_t(R_t)$ , which is consistent with this work.

$$R = \sum_{t \in \mathcal{T}} x_t \quad (2.23)$$

In control theory, a common problem is to minimize energy consumption, travel time or another measure of performance of a dynamic system. The state variable is mostly called  $s_t$ . Unlike in this work, it is often impossible to obtain  $s_t$  explicitly. In these cases, an observation  $o_t$ , which represents some information about  $s_t$ , is used as a substitute. The action at timestep  $t \in [0, T]$ , which is called  $a_t$  in this work, is often referred to as  $u_t$ . Instead of maximizing a sum of expected rewards  $\sum_t r_t$ , the goal is often to minimize a cost function  $J(u) = \sum_t c_t(s_t, u_t)$  with respect to  $u_t$ . For further information about trajectory optimization by optimal control, please refer to [8].

The nomenclature in computer science is mostly consistent with this work and therefore does not need to be explained here.

As for the relationship between the terms *reinforcement learning* and *dynamic programming*, there are differing views in literature. Some of the authors regard DP as a

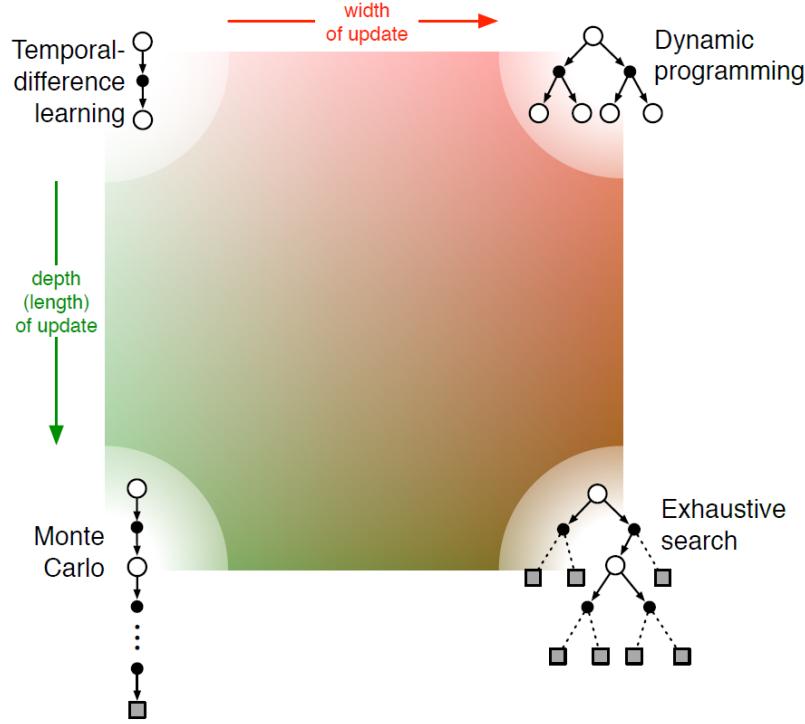


Figure 2.3.: Dynamic programming as a reinforcement learning method (taken from [17])

method from the pool of RL methods (cf. [16] and [17]). In their work, they use the term DP for RL methods that bootstrap<sup>4</sup> and sample multiple successor values instead of just one. Fig. 2.3 illustrates this viewpoint.

In [13], reinforcement learning is described as a so called forward dynamic programming method<sup>5</sup>. Powell considers RL to be a form of DP, which is, in some sense, contradictory to the view in the previous paragraph. But it makes sense inasmuch as many RL algorithms perform similar steps in all visited states, utilizing previous results to solve the current problem (e.g. when bootstrapping or picking a greedy action w.r.t. a value function). The term *forward* is also explained in [13]: Traditional DP methods typically start at a terminal state and step back until they reach a start state. This is usually more efficient than stepping forwards, but is very difficult, if the state space is multidimensional [13, section 1.5]. In many RL algorithms, trajectories are therefore either sampled completely from start to finish or one step lookaheads are performed in order to improve a policy or a value function estimate.

<sup>4</sup>Recall that this means, that they look one step into the future at each visited state and use an estimate of the successor state value  $V(s')$  to estimate  $V(s)$  instead of sampling complete trajectories like monte carlo methods.

<sup>5</sup>Other synonyms from [13] include *approximate dynamic programming*, *iterative dynamic programming*, and *adaptive dynamic programming*

## 2. Solution Methods for Markov Decision Processes

### 2.11.3. Types of Dynamic Programming Algorithms

#### Policy Evaluation

In Policy Evaluation, the goal is to find an approximation of the value function  $V_\pi$  that corresponds to a given policy  $\pi$ . The value of a state is defined as the expected total return from that state onwards until the episode terminates. Obviously, the expectation in equation 2.26 depends on the actions that are taken at each step and therefore on the policy that is evaluated. The state value that corresponds to a specific policy  $\pi$  is denoted by  $V_\pi(s)$ .

$$V_\pi(s_t) = \mathbb{E}[G_t | s_t = s] \quad (2.24)$$

$$= \mathbb{E} \left[ \sum_{k=0}^{T-t-1} \gamma^k r_{t+k+1} | s_t = s \right] \quad (2.25)$$

$$= \mathbb{E}[r_{t+1} + \gamma^1 r_{t+2} + \gamma^2 r_{t+3} + \gamma^3 r_{t+4} + \dots + \gamma^{T-t-1} r_T | s_t = s] \quad (2.26)$$

With a stochastic discrete policy, equation 2.26 becomes

$$V_\pi(s_t) = \sum_a \pi(a|s) (\mathcal{R}(s, a) + \gamma V_\pi(s')) \quad (2.27)$$

Recall the definition of a state value  $V(s) = \mathbb{E}[G_t | s_t = s]$  from Eq. 2.10. This equation can be used as a mapping to update the value of  $s_t$  with the expected return according to the current policy.

$$V_\pi(s_t) \leftarrow \mathbb{E} \left[ \sum_{k=0}^{T-t-1} \gamma^k r_{t+k+1} | s_t = s \right] \quad (2.28)$$

In Eq. 2.26, all terms except the first one can be replaced by the expected return from  $s_{t+1}$  onwards:

$$V_\pi(s_t) \leftarrow \mathbb{E}[r_{t+1} + \gamma G_{t+1} | s_t = s] \quad (2.29)$$

$$= r_{t+1} + \gamma V_\pi(s_{t+1}) \quad (2.30)$$

With a discrete stochastic policy, equation 2.30 becomes

$$V_\pi(s_t) \leftarrow \sum_a \pi(a|s) (r_t + \gamma V_\pi(s')) \quad (2.31)$$

In a state space with a finite set of states  $\mathcal{S}$ , writing down Eq. 2.31 for each state results in a system of  $|\mathcal{S}|$  linear equations that can be solved for  $V(s)$ . If the number of states is very large, this is not very efficient.

Instead, 2.26 is usually solved iteratively. The simplest way is to perform a one step lookahead from each state to get  $r_t$  and  $V(s_{t+t})$  from the current estimate of the state value function. Once all values are calculated, the values of all states are updated simultaneously<sup>6</sup>. Unlike other RL-algorithms, the state values are not updated along a trajectory. Instead, each state is updated independently.

Each set of state values approximates the true value function of the problem better than the previous set of estimates. Every state value converges to the true state value under the given policy. It is not efficient to wait until the values have fully converged, i.e. the true value function  $V_\pi$  is reached. Instead the process is repeated until the maximum absolute change in values lies beneath a certain threshold  $\epsilon_V$ . This maximum absolute change can be expressed by the supremum-norm  $\|\cdot\|_\infty$  (c.f. section 2.11.4). Algorithm 1 shows a possible implementation of iterative policy evaluation. In practice, there can be an upper bound for the number of evaluation steps.

$$\|V_k - V_{k-1}\|_\infty < \epsilon_V \quad (2.32)$$

---

**Algorithm 1** Iterative policy evaluation

---

```

function POLICY EVALUATION
    Initialize  $V(s) = 0 \forall s \in \mathcal{S}$  and load policy  $\pi$ .
    while true do
        for all  $s \in \mathcal{S}$  do
             $V_{\pi,new}(s) \leftarrow r + \gamma V_{\pi,old}(s')$ 
        end for
        if  $\|V_{new} - V_{old}\|_\infty \leq \epsilon_V$  then
            break
        end if
    end while
end function
```

---

## Policy Iteration

Policy iteration is an iterative way to calculate an estimate of the optimal policy of an MDP. Starting with an arbitrary policy and value function, one iteration of policy evaluation is performed to get an estimate of  $V_\pi$ . After that, after that, the policy is replaced by the greedy policy w.r.t. to the new value function estimate  $V_\pi$ . This alternating process of Policy Evaluation and Policy Improvement is repeated until the policy is satisfactory.

---

<sup>6</sup>There are variants of Policy Iteration, where the updates are performed asynchronously. This means that, for example, all neighbors of a terminal state are updated first. After that, the states, that lie next to these states, are updated, and so on. They usually converge faster, but at the expense of implementation complexity.

## 2. Solution Methods for Markov Decision Processes

$$a_{\text{greedy}}(s_t) = \underset{a}{\operatorname{argmax}}[Q(s_t, a_t)] \quad (2.33)$$

For a deterministic, greedy policy  $\pi_{\text{greedy}}$ , this yields:

$$\pi(a_{\text{greedy}}|s_t) = \mathbb{P}[a_t = a_{\text{greedy}}(s_t)|s_t] \quad (2.34)$$

$$= 1 \quad (2.35)$$

and for a stochastic policy  $\pi$ :

$$\mathbb{E}[\pi(a_t|s_t)] \leftarrow a_{\text{greedy}} \quad (2.36)$$

As mentioned in subsection 2.11.3, each policy evaluation step usually ends if the maximum difference between two subsequent values of the same state are sufficiently close.

There is a version of Policy Iteration where only one value update is done before updating the policy. This algorithm is called Optimistic Policy Iteration (OPI). Although the first estimate of the value function is only a rough approximation, OPI also converges to the optimal value function and policy. The following diagram illustrates the policy iteration scheme.

$$\pi_0 \xrightarrow{\text{evaluate}} V_1 \xrightarrow{\text{improve}} \pi_1 \xrightarrow{\text{evaluate}} V_2 \xrightarrow{\text{improve}} \pi_2 \xrightarrow{\text{evaluate}} \dots \xrightarrow{\text{evaluate}} V_* \xrightarrow{\text{improve}} \pi_*$$

Another way to picture policy iteration is shown in figure 2.4. Every PI-step brings the Value Function  $V(s)$  closer to  $V_*(s)$  and the policy  $\pi$  closer to  $\pi_*$ .

Policy iteration can be seen as two competing processes. On one hand, every policy evaluation step changes  $v$ , which leads to a change in the greedy policy. The previous greedy policy is therefore no longer greedy after the evaluation step. On the other hand, every policy improvement step makes the policy greedy w.r.t the current value function. This typically changes the value function, so the value function before the policy update is no longer the correct one (cf. [17, section 4.6]).

The black arrows in Fig. 2.4 refer to a policy iteration algorithm, that calculates  $v_\pi$  exactly in every evaluation step. This typically takes long for one iteration, but requires less iterations than optimistic policy iteration. The latter is what the blue arrows refer to. Each policy evaluation step only yields a rough estimate of  $v_\pi$  and takes less time than exact evaluation. OPI however usually takes more iterations than exact PI.

The performance of both algorithms on the trajectory optimization problem is compared in chapter 5.

In a discrete MDP, there is a finite number of states and actions. At each state  $s \in \mathcal{S}$ , the agent can choose the action  $a$  from a set  $\mathcal{A}$  of possible actions. In such a scenario, there exist  $|\mathcal{A}|^{|S|}$  different policies. It therefore takes  $|\mathcal{A}|^{|S|}$  iterations to find the optimal policy assuming no policy occurs twice. If that was the case, this policy would have

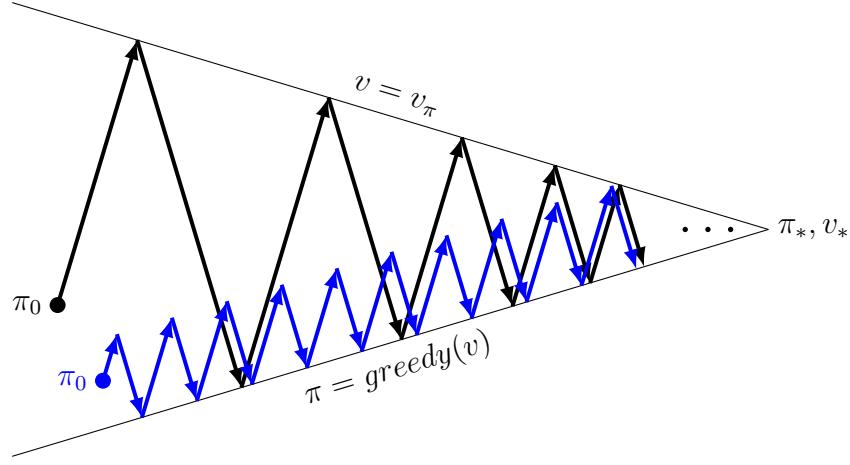


Figure 2.4.: The policy iteration algorithm (adopted from [17, section 4.6])

to occur in two subsequent iterations and therefore be the fixed point  $\pi_*$  of the policy sequence generated by PI.

As the results in chapter 5 show, PI converges much faster in practice than the upper bound  $|\mathcal{A}|^{|S|}$  may suggest. Algorithm 2 shows one way to implement generalized policy iteration. The algorithm alternates between PE and policy improvement until none of the actions change from one iteration to the next. The pseudo code for optimistic policy iteration can be found in algorithm 4 in appendix A.4.

## Value Iteration

Value iteration is similar to policy iteration, but instead of calculating a new policy at each iteration step, the state values  $V(s)$  are directly updated with the maximum possible successor value that is achievable starting from  $s$ . This yields a sequence of value functions, each one being a better estimate of  $V_*$  than its predecessor. Each of the intermediate value functions is abstract in so far as there does not have to exist an explicit policy corresponding to it (c.f. [16, lecture 3]). In a possibly discounted MDP with a deterministic environment, the update rule for each state value can be seen in equation 2.37.

$$V(s_t) \leftarrow \max_a [r_{t+1} + \gamma V(s_{t+1}) | s_t = s] \quad (2.37)$$

In principle, this is equivalent to synchronous OPI, where after one state value update at each state, the policy is replaced by the greedy policy with respect to the new value function. The only difference is, that value iteration does not output an intermediate policy at each iteration. Instead, it only iterates value functions.

$$V_1 \rightarrow V_2 \rightarrow V_3 \rightarrow \dots \rightarrow V_*$$

## 2. Solution Methods for Markov Decision Processes

In practice, a stopping criterion is introduced to limit calculation time. In this work, iteration is stopped if  $\|V_k - V_{k-1}\|_\infty < \epsilon_V$ , identical to equation 2.32. The last iterate  $V_n$  is considered an approximation of  $V_*$ . If  $V_n$  is close enough to the true optimal value function  $V_*$ , the (approximately) optimal policy can be calculated in the end by acting greedily with respect to  $V_n$ . Recall that, if Value Iteration is stopped too early, there might not even be a real policy that corresponds to the last value function iterate  $V_n$  [16, lecture 3]. A possible implementation of value iteration is shown in algorithm 3.

### 2.11.4. The Contraction Mapping Theorem

A contraction mapping  $f : M \rightarrow M$  on a metric space  $M$  has the following property:

$$|f(x_1) - f(x_2)| \leq \gamma|x_1 - x_2| \quad (2.38)$$

with  $x_1, x_2 \in M$  and  $\gamma \in [0, 1)$ . In this context,  $|(\cdot)|$  denotes an arbitrary metric on  $M$ .

The term contraction reflects the fact that, graphically speaking, a contraction mapping reduces the distance between  $x_1$  and  $x_2$ .

## Discounted MDPs

Policy Evaluation in discounted MDPs is a contraction mapping. This result can be obtained in few steps. The mapping in equation 2.31 and 2.39 is known as the *Bellman Operator*. It maps from a state value function to a state value function. If the Bellman Operator is applied multiple times, this can be expressed by  $T^n(v)$ .

$$T_\pi(v(s')) = \mathcal{R}_\pi + \gamma \mathcal{P}_\pi v(s') \quad (2.39)$$

If the Bellman Operator is applied to two state value functions  $V_1(s)$  and  $V_2(s)$ , it reduces the distance between the two in value function space. This distance can be expressed by the supremum norm, which is defined as follows:

$$\|V\|_\infty = \max_{s \in \mathcal{S}} |V(s)| \quad (2.40)$$

The supremum norm of  $V_1 - V_2$  is therefore

$$\|V_1(s) - V_2(s)\|_\infty = \max_{s \in \mathcal{S}} |V_1(s) - V_2(s)| \quad (2.41)$$

which is the biggest absolute difference of values of the same state in the state space. If the Bellman Operator is applied to both  $V_1$  and  $V_2$ , the contraction property proof is straightforward. For simplicity, the arguments of  $V_1(s)$  and  $V_2(s)$  are omitted in the

following equations.

$$\|T_\pi(V_1) - T_\pi(V_2)\|_\infty = \|(\mathcal{R}_\pi + \gamma \mathcal{P}_\pi V_1) - (\mathcal{R}_\pi + \gamma \mathcal{P}_\pi V_2)\|_\infty \quad (2.42)$$

$$= \|\gamma \mathcal{P}_\pi (V_1 - V_2)\|_\infty \quad (2.43)$$

$$\leq \|\gamma \mathcal{P}_\pi\| \|V_1 - V_2\|_\infty \quad (2.44)$$

$$\leq \|\gamma\| \|V_1 - V_2\|_\infty \quad (2.45)$$

$$= \gamma \|V_1 - V_2\|_\infty \quad (2.46)$$

Whenever  $s'$  happens to be the state where  $V_1$  and  $V_2$  differ the most, equality holds in line 2.44. If not, applying the supremum norm to  $V_1(s') - V_2(s')$  increases the value of the expression.  $\mathcal{P}_\pi$  must always be smaller than or equal to one. Similar to line 2.44, equality in line 2.45 holds only if  $\mathcal{P}_\pi = 1$ . Replacing  $\mathcal{P}_\pi$  by one in line 2.44 yields line 2.45. The left side of line 2.42 must thus be less than line 2.46 which proves the contraction property of the Bellman Operator in discounted MDPs.

Another proof of convergence of discounted optimistic policy iteration is shown in [18].

## Undiscounted MDPs

An MDP, where the discount factor  $\gamma$  is one, is called an undiscounted MDP. In such an MDP, the Bellman Operator looks like in equation 2.47.

$$T_\pi(v(s)) = \mathcal{R}_\pi + \mathcal{P}_\pi V(s') \quad (2.47)$$

If  $\gamma$  is one, the above proof of the contraction property does not hold anymore. It can however be proven, that the Bellman Operator still contracts. In the undiscounted case, this takes more than one step. Equation 2.48 shows the n-step contraction property. For simplicity, the index  $\pi$  and the argument  $s$  are omitted.

$$\|T^{(n)}(V_1) - T^{(n)}(V_2)\|_\infty \leq \|V_1 - V_2\|_\infty \quad (2.48)$$

Equation 2.48 means, that applying the Bellman Operator  $n$  times to two value functions brings them closer together in value function space. As a result, the Bellman Operator also leads to  $V_1 \approx V_2$  eventually also in undiscounted MDPs.

More information about under which circumstances policy iteration and value iteration also converge in undiscounted MDPs can be found in [9] and [6]. In [5], there is some general information about discretization procedures in DP and how they affect the convergence properties of the discretized MDP.

## 2.12. Exploiting the Problem Structure

According to [13], there is no such thing as a universal dynamic programming algorithm that is capable of solving all possible problems efficiently. Instead, making use of the

## 2. Solution Methods for Markov Decision Processes

structure of a problem is often vital in order to reduce computation time or make DP feasible in the first place. The following paragraph illustrates how this can be done in the case of path planning through value iteration.

In numerics, making use of new information as soon as it is available often reduces calculation time significantly. In value iteration, this can mean to update the value of a state as soon as a new estimate of it has been calculated. This way, subsequent calculations can make use of more recent data than in the case of synchronous updates<sup>7</sup>. The most efficient way to exploit the problem structure in the context of this work (cf. chapter 4) is to go backwards through the state space. The states directly next to the terminal states are updated first and simultaneously. After that, the layer<sup>8</sup> next to these states is updated and so on. This way, the final reward can propagate to the start state in one iteration, i.e. one backward pass through the state space. As the value iteration algorithm in this work also considers movement away from the goal, it typically takes more than one iteration to converge, but only about a third of the number of iterations in the synchronous case (cf. section 5.1.1).

---

<sup>7</sup>Recall that, in synchronous DP, the state values are updated at the same time after one complete sweep over the state space. If the start state  $s_0$  is 50 steps away from the target state  $s_T$ , it takes 50 iterations (i.e. complete sweeps over the state space) to find a trajectory from  $s_0$  to  $s_T$ .

<sup>8</sup>In this context, one layer is one set of points that have the same x-coordinate. This set forms a vertical "layer" of states that have the same distance to the goal.

---

**Algorithm 2** Generalized policy iteration

---

**function** GENERALIZED POLICY ITERATION

    Initialize  $V(s) = 0 \forall s \in \mathcal{S}$

    Load arbitrary initial policy  $\pi_0$ .

    Initialize  $m, \epsilon_V$

**while** true **do**

**function** POLICY EVALUATION

**while** true **do**

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

$V_{\pi,new}(s) \leftarrow r + \gamma V_{\pi,old}(s')$

**end for**

**if**  $\|V_{new} - V_{old}\|_\infty \leq \epsilon_V$  **then**

                    break

**end if**

**end while**

**end function**

**function** POLICY IMPROVEMENT

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

                sample  $m$  actions  $a_n(s)$

**for all**  $a_n$  **do**

$Q(s, a_n) = r + \gamma V(s')$

**end for**

$a_{greedy}(s) = \underset{a_n}{\operatorname{argmax}}[Q(s, a_n)]$

**end for**

**end function**

    train policy on  $a_{greedy,new}$  to obtain  $\pi_{new}$

**if**  $\|a_{greedy,new} - a_{greedy,old}\|_\infty = 0$  **then**

        break

**end if**

**end while**

---

**end function**

---

---

**Algorithm 3** Value iteration

---

```

function VALUE ITERATION
    Initialize  $V(s) = 0 \forall s \in \mathcal{S}$ 
    Initialize  $m, \epsilon_V$ 
    while true do
        for all  $s \in \mathcal{S}$  do
            sample  $m$  actions  $a_n(s)$ 
            for all  $a_n$  do
                 $Q(s, a_n) = r + \gamma V(s')$ 
            end for
             $V_{new}(s) = \max_{a_n} [Q(s, a_n)]$ 
        end for
        if  $\|V_{new} - V_{old}\|_\infty \leq \epsilon_V$  then
            break
        end if
    end while
function POLICY IMPROVEMENT
    for all  $s \in \mathcal{S}$  do
        sample  $m$  actions  $a_n(s)$ 
        for all  $a_n$  do
             $Q(s, a_n) = r + \gamma V(s')$ 
        end for
         $a_{greedy}(s) = \operatorname{argmax}_{a_n} [Q(s, a_n)]$ 
    end for
end function
    train policy on the latest  $a_{greedy,new}$  to obtain an estimate of  $\pi_*$ 
end function

```

---

# 3. Function Approximation

A common problem in engineering is developing a model from measurement data. Such a model makes it possible to generalize from the given data to new data, for which a real measurement does not exist. The process of developing a model is often combined with function approximation. The goal here is to find a mapping as a linear combination of elements, which constitute a basis from functions, that approximates a given behavior, also called a target function, as close as possible.

Mathematically, function approximation is the process of finding a mapping from inputs  $x$  to outputs  $y$ . Usually, a set  $\mathcal{X} = \{x_{fit,1}, x_{fit,2}, \dots\}$  of inputs and a set  $\mathcal{Y} = \{y_{fit,1}, y_{fit,2}, \dots\}$  corresponding desired outputs is given and the goal is to find a function that returns a value close to  $y_{fit,i}$  for each  $x_{fit,i}$ .

There are two types of function approximation. In some cases, knowledge about the specific problem can help to find suitable base functions. If, for example, a mapping is known to be linear ( $y = a_0 + a_1x$ ), the function approximation process degenerates to finding the two coefficients  $a_0$  and  $a_1$ . The same applies to quadratic, trigonometric or exponential target functions or linear combinations of them. If, however, the target function is not known, function approximation is more elaborate. One can try multiple target functions or combine them in order to find one, that fits best, or use a *universal function approximator* like an artificial neural network.

## 3.1. Tables

The simplest way to express a discrete mapping is by a table where each input/output pair is represented by a cell of the table. The number of inputs and corresponding outputs in a table is inherently finite. The set of states and actions in a continuous MDP is however infinite, so the policy and value functions must also be continuous. If a discrete mapping like a table is used as a representation of a policy or value function in a continuous MDP, it will, at some point, be confronted with unknown inputs. The inability to deal with unknown inputs is a big disadvantage of tables. There are, however, ways to overcome this weakness.

### Nearest Neighbor

The most straightforward way to obtain an output for an unknown input is to look for the table entry that is closest to the unknown input. This yields a piecewise constant mapping from inputs to outputs and enables the table to generalize. Whether this suffices

### 3. Function Approximation

as an approximation for a continuous mapping  $f$  heavily depends on the properties of  $f$ . If there are large differences between the output of two subsequent inputs, the error can be very high.

#### Linear Interpolation

Another way to obtain outputs for unknown inputs is to interpolate linearly. In a first step, the two nearest known inputs  $x_0$  and  $x_1$  and their respective outputs  $y_0$  and  $y_1$  are looked up (note that  $x_0 \leq x \leq x_1$  in this case). In the second step, the output is calculated by the following formula:

$$y_{\text{interp}} = y_0 (1 - \gamma) + y_1 \gamma, \quad (3.1)$$

where  $\gamma = \frac{x-x_0}{x_1-x_0}$ .

One advantage of linear interpolation of table values is that the output is at least continuous, though not differentiable at the known input points.

## 3.2. Artificial Neural Networks

The structure of artificial neural networks resembles the structure of a human brain according to the current state of research [11, section 1.1]. The human brain consists of a very high number of neurons, each receiving signals from other neurons, processing them and itself sending a signal to other neurons if certain conditions are met. Each neuron is very simple, it is the high number of them and the fact that they all work in parallel that makes the human brain so powerful.

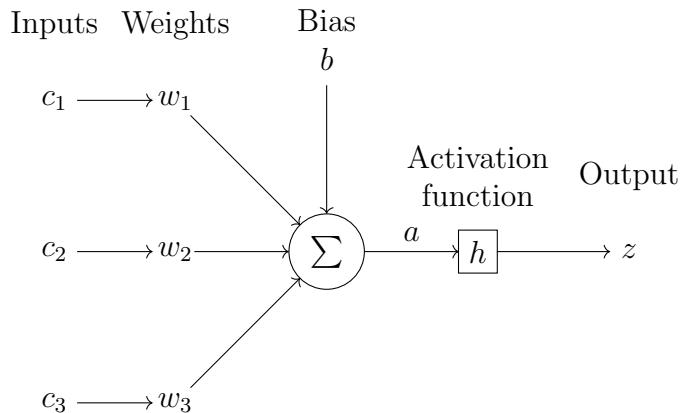


Figure 3.1.: A neuron in an artificial neural network

A neuron in an artificial neural network looks like what is depicted by Fig. 3.1. It takes an arbitrary number of inputs  $c_i$ , multiplies each input by the corresponding input weight

### 3.2. Artificial Neural Networks

$w_i$ , adds a bias and applies to the result  $a$  a non-linear function, the so called activation function  $h$ , yielding  $z$ , the neuron output. This output can either serve as an input to another neuron or as the network output. Multiple activation functions are common, most of which are somewhat s-shaped. These activation functions map the set of real numbers to the range between  $-1$  and  $1$  or sometimes between  $0$  and  $1$ , limiting the neuron output and therefore its influence on the output of the ANN network. Examples include the step function, ( $z = 1$  if  $x \geq x_0$ ,  $z = 0$  otherwise), the rectifier nonlinearity ( $h(x) = \max(0, x)$ ), the sigmoid function ( $h(x) = 1/(1 + e^{-x})$ ) and the  $\tanh(\cdot)$ -activation function ( $h(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$ ). Nonlinearity is essential here, because if linear activation functions are used, each neuron output is a linear function of its inputs. Any combination of linear mappings is itself linear. So if a deep neural network has only linear activation functions, the ANN output is always a linear function of the inputs, regardless of the number of hidden layers. Therefore, a deep neural network, where all activation functions are linear, can be replaced with a network, that has only one hidden layer (or something that is even simpler). The computational expense of dealing with a deep ANN is - however - very high compared to one with a single hidden layer. A deep ANN with only linear activation functions therefore does not make sense.

Figure 3.2 shows the two activation functions that were used in the policy neural network. The  $\tanh(x)$ -function is drawn in blue, the red line represents a linear activation function.

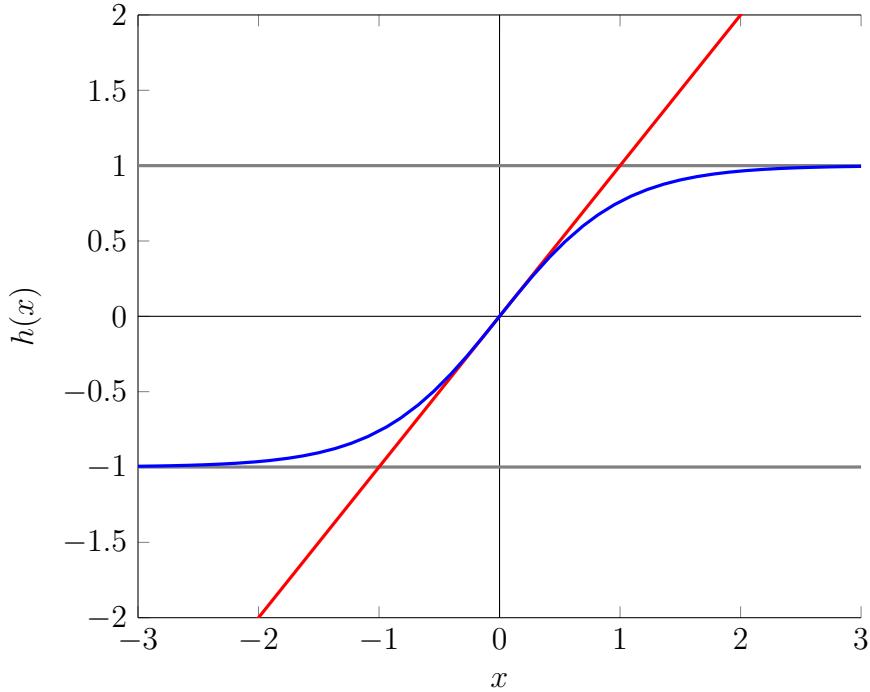


Figure 3.2.: Activation functions (adopted from [20])

### 3. Function Approximation

In an artificial neural network, the neurons are arranged in layers. The ANNs in this work are all fully connected feed forward networks, also known as Multi Layer Perceptrons (MLPs). Each neuron in a given layer  $k$  gets an input from each neuron in the previous layer  $k - 1$  and outputs a signal to each neuron in the next layer  $k + 1$  (c.f. Fig. 3.1).

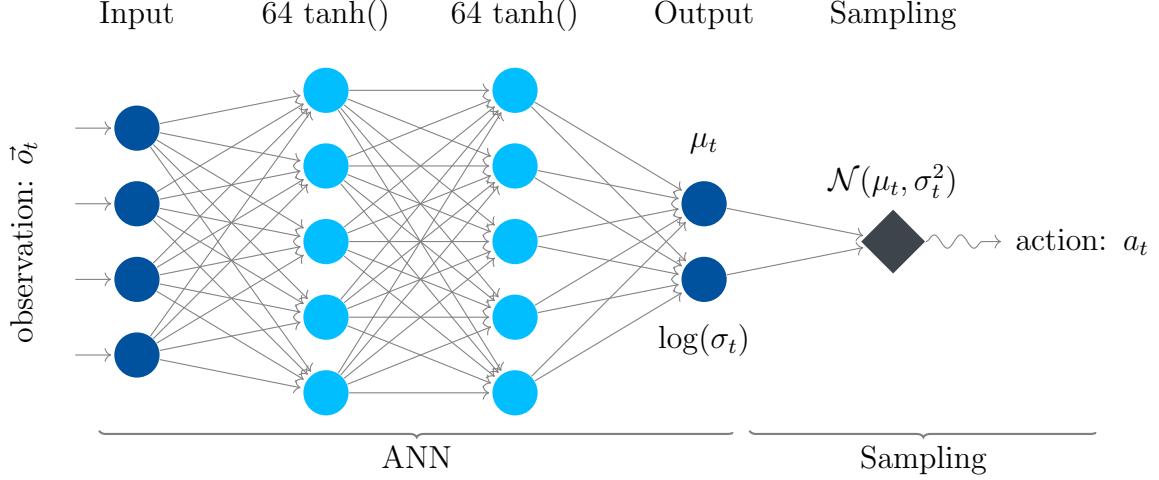


Figure 3.3.: The policy MLP [20]

If a policy or value function is represented by an artificial neural network (ANN), the output for each state (or state-action-pair) is determined by the weights and biases of the network and the activation functions. As the latter are usually given or predetermined (and not trained), only the weights and biases are usually regarded as parameters. One set of weights and biases forms the *parameter vector*  $\theta$  of the neural network. If, for example, a stochastic policy is parametrized by  $\theta$  and represented by a neural net with  $k$  layers, the notation is as follows:

$$\pi_\theta(a|s) = \mathbb{P}[a|s, \theta] \quad (3.2)$$

with

$$\theta = \{\mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{W}^{(2)}, \mathbf{b}^{(2)}, \dots, \mathbf{W}^{(k)}, \mathbf{b}^{(k)}\} \quad (3.3)$$

The policy ANN does not necessarily need to output the action directly. Instead, most RL-related policies are stochastic, i.e. the neural net outputs the mean and standard deviation of a normal distribution:

$$\pi_\theta(\mu_\theta, \sigma_\theta|s) = \frac{1}{\sqrt{2\pi}\sigma_\theta} e^{-\frac{(a-\mu_\theta)^2}{2\sigma_\theta^2}} \quad (3.4)$$

The action  $a$  is then drawn from  $\mathcal{N}(\mu_\theta(s), \sigma_\theta^2(s))$ . In this work, the policy is represented with a neural network that has two outputs, the mean and standard deviation. Dynamic

Programming is used only to find the optimal mean values. The standard deviation is set to 0.5. This serves as an initial solution for the policy net to be trained by an RL method.

## Input Standardization

All training algorithms that are used in this work rely on the partial derivative of a loss function with respect to the network weights and biases. As can be seen in Eq. 3.8 and the following equations, they perform a step in (an approximation of) the current gradient direction. This gradient heavily depends on the network inputs. If some of the inputs have much larger values than the rest, this can result in some weights being updated faster than others, which is not desirable. Therefore, the inputs to a neural network should have the same scale. In this work, the inputs are scaled to a normal distribution with  $\mu = 0$  and  $\sigma = 1$ . Note that  $\mu_t$  and  $\sigma_t$  are not the same as  $\mu$  and  $\sigma$ .  $\mu_t$  and  $\log(\sigma_t)$  are the outputs of the policy MLP at timestep  $t$ .

$$\hat{x} = \frac{x - \mu}{\sigma} \quad (3.5)$$

The specific values for  $\mu$  and  $\sigma$  for each of the scenarios are presented in table A.1.

## Action Scaling

In an artificial neural network, every neuron can output a value between -1 and 1. Therefore, the output of a neural network is limited by the number of neurons in the last hidden layer. Obtaining a large output requires a lot of neurons. On the other hand, lots of neurons favor overfitting and increase training time. If the desired outputs are too small, on the other hand, then so are the neuron outputs. The  $\tanh(\cdot)$  activation function that is used for the policy is approximately linear in this area. Therefore, the whole network behaves almost linearly if it is trained to output values close to zero. With linear mapping being possible to approximate by an ANN without any hidden layer, there would not be any benefit of applying a deep net. The most important advantage of deep neural networks is however their ability to represent highly complex and a priori unknown behavior.

Both problems can be avoided if outputs lie between -1 and 1. This way, it is unlikely that any of the activation functions reach their saturation limit, nor is the network restricted to linear behavior. Therefore, the neural network is usually trained to output values between -1 and 1 and these values are then scaled to the desired magnitude by multiplying them with a scaling factor. Another way to mitigate this problem is choosing activation functions that are also nonlinear if the input and/or output is close to zero.

### 3. Function Approximation

## Weight initialization

Training an artificial neural network is an iterative process. Therefore it requires initial values for the weights and biases. As training a neural network through supervised learning comes down to solving an optimization problem iteratively, training success highly depends on the initial weights.

The two most popular ways to initialize neural network weights are the Glorot uniform initialization and the Glorot normal distribution. According to the Glorot uniform initialization, the weights in one layer must be distributed uniformly within  $\pm \sqrt{\frac{6}{n_{in}+n_{out}}}$  where  $n_{in}$  and  $n_{out}$  are the numbers of neurons in the prior and current layer, respectively.

The Glorot normal initialization states that the weights must be initialized according to a normal distribution with zero mean and a variance of  $\sigma = \sqrt{\frac{2}{n_{in}+n_{out}}}$ .

In this work, all neural networks are initialized according to the Glorot uniform distribution.

## 3.3. Supervised Learning

All machine learning algorithms can be divided into two categories: supervised learning and unsupervised learning. In supervised learning, the agent is given the desired action for each state explicitly, so it knows what it is supposed to do. In unsupervised learning, however, the agent does not know what the best action is, so it has to find out itself. One way of doing that is from moving in the state space and learning from experience, for example by applying a reinforcement learning algorithm.

In the supervised learning case, training is done via minimizing a loss function that depends on the desired action and the actual action taken by the agent. A common loss-function is the mean squared error (MSE) :

$$\bar{e} = \frac{1}{N} \sum_{n=1}^N e_n \quad (3.6)$$

$$= \frac{1}{N} \sum_{n=1}^N \left[ \frac{1}{2} (a_n(s) - a_{fit,n})^2 \right] \quad (3.7)$$

The goal of supervised learning is to minimize the difference between  $a_{n,fit}$  and  $a_n(s)$  in each state. If  $a_n(s)$  is the output of a function approximator, minimizing  $L_{MSE}$  is equivalent to fitting the function approximator to the desired outputs.

One advantage of the application of a loss function is that numerous iterative optimization algorithms have been developed to solve such a problem.

## 3.4. Optimization Techniques

As mentioned before, supervised learning is usually done by defining a loss function that depends on the network weights and biases and finding a local or global minimum numerically.

In this work, the Adam algorithm is used for training the policy to output the desired action at each point in the state space. As most optimization algorithms, Adam is gradient-based. Therefore, some insight into the general ideas behind gradient based numeric optimization is useful for understanding Adam. Sections 3.4.1 and 3.4.2 explain briefly how gradient based optimization works in general. In section 3.4.3, Adam is described and explained.

All presented algorithms produce a series of points in the space that the network weights and biases establish. The point at the  $k$ th optimization step is defined by a distinct set of weights and biases  $\theta$ , similar to coordinates in the state space (c.f. Eq. 3.3).

### 3.4.1. Gradient Descent

In gradient descent, the next point  $\theta_{k+1}$ , i.e. the next set of weights and biases, is reached by computing the gradient of the loss function with respect to the network weights,  $\nabla_{\theta_k} L_k$ , and performing one step in the direction of steepest descent.  $\theta_{k+1}$  is calculated by the following equation:

$$\theta_{k+1} = \theta_k - \nabla_{\theta_k} \bar{e}_k \cdot \alpha \quad (3.8)$$

with

$$\nabla_{\theta_k} \bar{e}_k = \frac{1}{N} \cdot \sum_{n=1}^N \nabla_{\theta_k} e_{n,k} \quad (3.9)$$

$$= \frac{1}{N} \cdot \sum_{n=1}^N [(a_n(s) - a_{n,fit}) \cdot \nabla_{\theta_k} a_n(s)] \quad (3.10)$$

and an arbitrary step size  $\alpha$ .

The choice of  $\alpha$  has a big effect on the convergence of the optimization algorithm. If  $\alpha$  is too small, each  $\theta_{k+1}$  is close to the respective  $\theta_k$  so it takes many steps to find a minimum. If  $\alpha$  is too big, the algorithm might not converge at all.

Ordinary gradient descent methods are not suitable for optimizing the output of large artificial neural networks because computing the gradient  $\frac{\partial L_k}{\partial \theta_k}$  for all the weights and biases is very computation-intensive and time-consuming.

### 3.4.2. Stochastic Gradient Descent

As mentioned before, the output of a neural network cannot be optimized efficiently with a gradient descent algorithm. One way to bypass this weakness is Stochastic Gradient

### 3. Function Approximation

Descent (SGD). If the loss function is a sum like in Eq. 3.7, it is sufficient to calculate the gradient of only one summand of the loss function, i.e. only consider one random training sample at each step  $k$ , and perform a small step in the corresponding direction:

$$\nabla_{\theta_k} \bar{e}_k \approx \nabla_{\theta_k} e_{n,k} \quad (3.11)$$

$$\theta_{k+1} = \theta_k - \nabla_{\theta_k} e_{n,k} \cdot \alpha \quad (3.12)$$

If this is done for all samples successively, it also converges to the loss from Eq. 3.7. Thanks to the faster computation of  $\nabla_{\theta_k} L_{n,k}$ , SGD converges faster than ordinary gradient descent although Eq. 3.11 is only a rough estimate of the true gradient. (todo: Beleg und Erklärung)

#### 3.4.3. The Adam Algorithm

Basic Stochastic Gradient Descent has a fixed stepsize  $\alpha$ . Most advanced optimization algorithms have a variable stepsize and an update direction that depends on the previous updates to accelerate convergence. One example is to add a momentum term to the gradient from Eq. 3.12. This affects the step size and the direction of the update.

It takes into account the previous update and changes the direction of the next update. Like a ball rolling down a slope, the algorithm does not immediately change direction if the gradient of the slope changes direction or stop if the gradient is zero. If there is a plateau in the loss function or a weak local minimum, ordinary algorithms tend to get stuck there. With momentum, the risk of getting stuck is reduced.

The Adam-algorithm utilizes a modified momentum term for better convergence. Equation 3.13 shows one update step, the parameters  $\hat{m}_t$  and  $\hat{v}_t$  are calculated with equations 3.14 and 3.15.

$$\theta_{k+1} = \theta_k - \alpha \cdot \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon} \quad (3.13)$$

with

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t} = \frac{\beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t}{1 - \beta_1^t}, \quad (3.14)$$

$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t} = \frac{\beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2}{1 - \beta_2^t} \quad (3.15)$$

The algorithm utilizes the first and second moments of the previous gradients to calculate the update. Consider a pair of  $m_t$  and  $v_t$  that follow recursively from  $m_{t-1}$  and  $v_{t-1}$  according to equations 3.14 and 3.15. For the calculation of each  $m_t$  and  $v_t$ , the previous value  $m_{t-1}$  and the current gradient  $g_t$  are used.  $m_t$  and  $v_t$  are biased estimates of the first and second moments, there is a correction term in each of the denominators of equations 3.14 and 3.15 to account for that. See [10] for more details.

## 3.5. Overfitting

It is seldom desirable to minimize the loss function to its absolute minimum value on the training data. At first, this might sound counterintuitive. But apart from the additional time it takes to get there, it increases the chances of running into problems.

When dealing with a continuous state space, it is common practice to discretize it and deal with the discretized state space instead. While a continuous state space consists of an infinite number of states, the number of states in the discretized state space is finite. This reduces computation time and makes some algorithms applicable in the first place. If an appropriate grid is used (i.e. if it is dense enough in relevant areas of the state space), the solution is likely to be close to the solution of the continuous problem. The more grid points, the closer the solution is to the continuous solution.

If, however, a function approximator like an ANN is optimized on the discretized state set, it is possible that, although it might fit the training data very closely, it may not generalize well, i.e. give unexpected outputs when fed with states that have not been used in training. This problem is called *overfitting*. It especially occurs if at least one of the following conditions is met:

- Noisy data (every data point differs slightly from the expected value)
- Limited dataset (only a few points to fit to)
- High number of parameters (many degrees of freedom)<sup>1</sup>

Overfitting is an issue in many function- and curve-fitting applications. Fig. 3.4 shows overfitting in case of a high order polynomial. Although it is much simpler than a neural network, the problems are similar.

The plot shows two polynomials that have been fitted to a set of seven points. The blue line represents a quadratic polynomial, the red line belongs to a 10th order polynomial. A  $n$ th order polynomial has  $n+1$  parameters and can therefore be fit exactly to  $n+1$  points. The quadratic polynomial only has three parameters and can therefore not pass through all seven points exactly. In such cases, fitting is done by minimizing the mean quadratic fitting error:

$$\bar{e}_{fit} = \frac{1}{m} \sum_{i=1}^m [y_i - y_{i,fit}]^2 \quad (3.16)$$

$$= \frac{1}{m} \sum_{i=1}^m [y_i - \sum_n (a_n x^n)]^2 \quad (3.17)$$

where  $m$  is the number of given points.

---

<sup>1</sup>i.e. lots of neurons in the case of an ANN.

### 3. Function Approximation

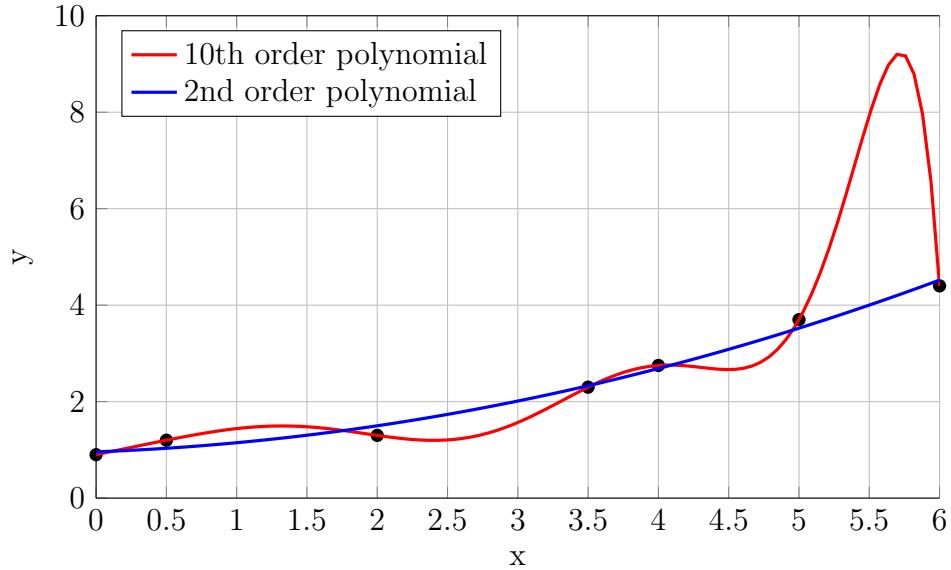


Figure 3.4.: Overfitting in case of a high order polynomial and noisy data

The 10th order polynomial has eleven parameters. It is therefore possible to fit the function exactly to the seven input points. This can be done by solving a system of linear equations where each point contributes one equation. In the case of the 10th order polynomial, there are eleven parameters. If this polynomial is fitted to seven points, there are seven equations and eleven parameters to be determined. The corresponding system of equations is called *underdetermined*. Four of the parameters have to be set prior to solving the linear system. Their choice is arbitrary, which can lead to a shape like the red plot. Although the seven conditions at the points are met exactly by the red line, the blue plot intuitively seems like a more appropriate approximation.

Of course, any high order polynomial can also be fitted to a set of points by minimizing 3.17. Unfortunately, this can also lead to overfitting, because it uses the same information (input point coordinates). Therefore, it will eventually find a set of parameters that fits the points exactly and is likely to have the same oscillations as the red line shown in Fig. 3.4.

In any application, where there are more parameters than points to fit, overfitting can be an issue. Any function approximator should therefore only have as many parameters as necessary. Any additional parameter can increase the risk of overfitting. There are several other precautions that can be taken in order to avoid overfitting. Some of them are presented in the following paragraphs.

One way to deal with overfitting is splitting the state set into a set of training data and a set of validation data. After training on the training dataset, the value of the loss function on the training data is compared to the value of the loss function on the validation data. If - from one training step to the next - the training loss decreases while

### 3.5. Overfitting

the validation loss does not, overfitting is likely. This can be used as a criterion to stop the training. Fig. 3.5 shows how training loss and validation loss typically evolve during training. After iteration 63, only the training loss decreases. At this point, training should be stopped in order to avoid overfitting.

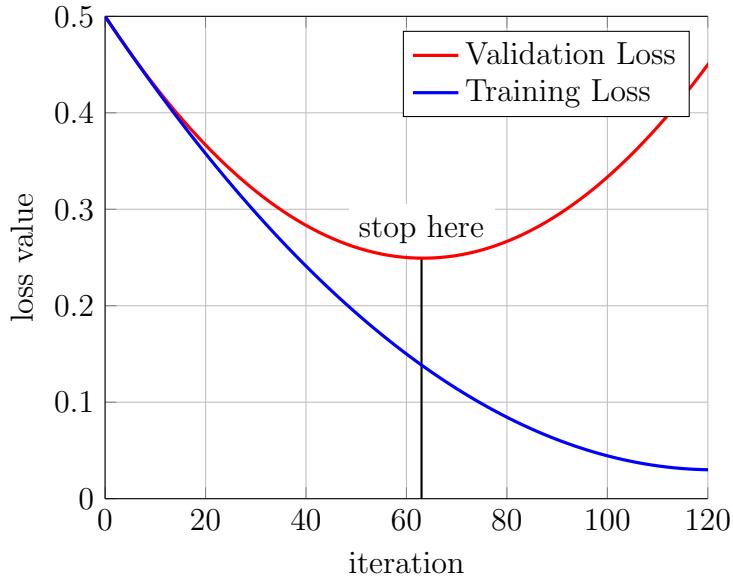


Figure 3.5.: Stop criterion to avoid overfitting

The training set can also be split into three datasets, one for training, one for validation during training, and the third for post-training overfitting detection. Typically, the third dataset contains approximately 20% of the samples, the remaining 80% are divided into training and validation data at the same ratio, so the parameter updates are done with 64% of the total sample count, and 16% are used to check for overfitting during training.

Another way of preventing overfitting is to penalize large weights and biases in the hidden layers. If the weights and biases have high values, they can cause similar behavior to what is shown in Fig. 3.4. Penalizing large weights and biases typically smoothes out the neural network output, thus averting oscillations.

In this work, there are two datasets. At the beginning of training, 80% of the data are chosen randomly and used for training. The remaining 20% are used to detect overfitting during training according to what is depicted in Fig. 3.5. Note that it typically takes more iterations until overfitting occurs than in the previous example. Large weights and biases are also penalized. This is done by adding a term to the loss function from Eq. 3.7, yielding Eq. 3.18. Note that, in Eq. 3.18, the set of parameters  $\theta$  is interpreted as a vector that just contains all weights and biases of the hidden layers as a list to simplify the notation of the sum in the last term.

### 3. Function Approximation

$$\bar{e} = \frac{1}{N} \sum_{n=1}^N \left[ \frac{1}{2} (a_n(s) - a_{fit,n})^2 \right] + a_\theta \sum_{i=1}^I \theta_i^2, \quad (3.18)$$

where  $I$  is the number of scalar weights and biases in the hidden layers and  $\theta_i$  is the  $i$ th parameter of the ANN. A value of 0.06 for  $a_\theta$  has proven to suppress overfitting without affecting the neural network output too much.

# 4. Trajectory Optimization with Dynamic Programming

## 4.1. 2D Environment

In the context of this work, the aircraft is treated as a mass point. Rotational dynamics are neglected. This is sufficient for trajectory optimization algorithms. See [7] for more details.

In this scenario, the glider moves in the geodetic vertical plane. Its control is the angle of attack  $\alpha$ . When  $\alpha$  is increased, the lift  $L$  increases. This results in an increase of  $\gamma$ , as can be seen in equation 4.4. This way the agent can control its velocity vector and therefore its position. All parameters of the glider are shown in appendix A.1. Any 3D-Scenario where the horizontal trajectory is restricted by obstacles or the optimal ground trace is obvious can be regarded as a 2D problem. The only control variable to be optimized is then the angle of attack or an equivalent quantity (e.g. vertical acceleration).

### 4.1.1. Policy Representation

As the policy iteration algorithm is running on a rectangular grid and the states that are updated do not change, a table that stores the greedy action for every state is sufficient for the policy iteration algorithm. Once the optimal policy  $\pi_{calm}^*$  is found, a neural network is trained to approximate the optimal action for each of the states on the grid. This policy is then used to do TRPO.

Value iteration utilizes the same grid as the policy iteration algorithm, i.e. the states do not change during calculation. A table is therefore also sufficient for VI. Once (an approximation of) the optimal value function  $v_{calm}^*$  is found, the greedy policy with respect to it is derived by one step of policy improvement and stored in a table. The resulting (approximation of the) optimal policy  $\pi_{calm}^*$  is used to train a neural network and do TRPO.

### 4.1.2. Equations of Motion

The physical state of the glider in the geodetic vertical plane consists of the position  ${}_g\mathbf{r} = [x, z]^\top$  and velocity  ${}_g\mathbf{v}_K = [{}_gu_K, {}gw_K]^\top$ . Apart from the geodetic reference frame (index g), a body fixed frame (index f), a trajectory fixed frame (index k) and an aerodynamic frame (index a) are used to describe the glider dynamics. In the vertical plane, coordinates

#### 4. Trajectory Optimization with Dynamic Programming

can be transformed between the frames by a rotation around the y-axis. Figure 4.1 illustrates the transformation angles that connect the reference frames. Note that, in this work, there is no wind.  $\vec{v}_W$  and  $\alpha_W$  are therefore both zero and the aerodynamic and trajectory fixed frames are the same.

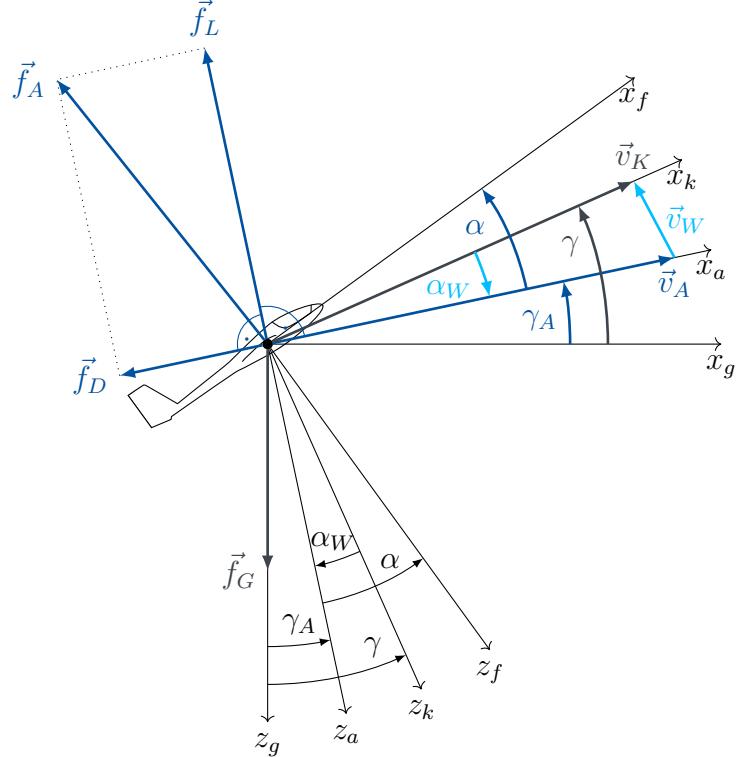


Figure 4.1.: Reference frame transformation [12]

In the vertical plane, there are 4 state variables,  $x$ ,  $z$ ,  $u = \dot{x}$  and  $w = \dot{z}$ , each with respect to a geodetic reference frame that has its origin at a point on the earth surface. The aerodynamic forces are often expressed in an aerodynamic reference frame. The coordinate transformations  $u$  and  $w$  can also be expressed by  $V$  and  $\gamma$ .

The dynamics in the vertical plane are given by the following equations. They are taken from [7].

$$\dot{x} = V \cos(\gamma) \quad (4.1)$$

$$\dot{z} = -V \sin(\gamma) \quad (4.2)$$

$$\dot{V} = -\frac{D + g \cos(\gamma)}{m} \quad (4.3)$$

$$\dot{\gamma} = \frac{L}{V m} - \frac{g \cos(\gamma)}{V} \quad (4.4)$$

where  $L = q c_L S$ ,  $D = q c_D S$  with  $q = \frac{\rho}{2} V^2$  and  $V = \sqrt{u^2 + w^2}$ .

### 4.1.3. Reward Function

As shown in Fig. 2.1, the agent interacts with his environment by, at each time step, picking an action  $a_t$  and receiving an observation  $o_t$  and a reward  $r_t$ . In RL, the reward function is the only means to give the agent information about what he should or should not do.

The reward function in this work is as follows:

$$\mathcal{R}(s, a) : \mathcal{S} \mapsto \mathbb{R} : r_t = \begin{cases} -\Delta t & \text{if } s_{t+1} \neq s_{T+} \\ -\Delta t + \frac{d}{10} & \text{if } s_{t+1} = s_{T+} \end{cases} \quad (4.5)$$

In Fig. 4.2, all state transitions, where the agent crosses a thick line (i.e. leaves the admissible part of the state space), lead to termination of the episode. If the agent hits one of the black thick lines, he gets a final reward of 0. Only passing the green line earns him a final reward of  $\frac{d}{10}$ .

where  $s_{T+}$  is a terminal state that is also a target state. There are terminal states  $s_{T-}$  where the agent is not meant to go. These terminal states are those where the agent touches the ground before reaching his target range. If the agent reaches his goal, he receives a final reward  $r_T$  according to the second line of Eq. 4.5.

It is possible to penalize undesired behavior of the agent by defining a large negative reward if, for example, the agent touches the ground before reaching the target distance or flies too far in the wrong direction. This is, however, not necessary, if the final reward is sufficiently high. The important thing here is, that reaching a target state must always yield a higher return than ending an episode on the ground from any state  $s \in \mathcal{S}$ .<sup>1</sup>

### 4.1.4. Discretization of the State- and Action Space

The state space is discretized with a rectangular grid, i.e. it is replaced by a set of points  $s_i = [x_k, z_l, u_m, w_n]^\top$  evenly distributed across all dimensions.

Fig. 4.2 illustrates the discretization of the state space. The agent starts its flight at the top left corner (at the glider symbol) with zero velocity (hence the green dot on the speed grid). For the grid cell at  $x_4$  and  $z_3$ , the grid for the speed vector is drawn. All of the cells in the position grid contain a speed grid, like the small one in Fig. 4.2.

---

<sup>1</sup>This is true in the scenarios in chapter 5. In scenario 1, for example, the agent has to cover a horizontal distance of 500m. This takes him between 19.7s (cf. table 5.2) and approx. 40s and earns him a time penalty of up to  $-40$ . With a final reward of  $\frac{d}{10} = +50$ , the resulting total reward is still positive and therefore sufficiently high for the agent to find the goal.

#### 4. Trajectory Optimization with Dynamic Programming

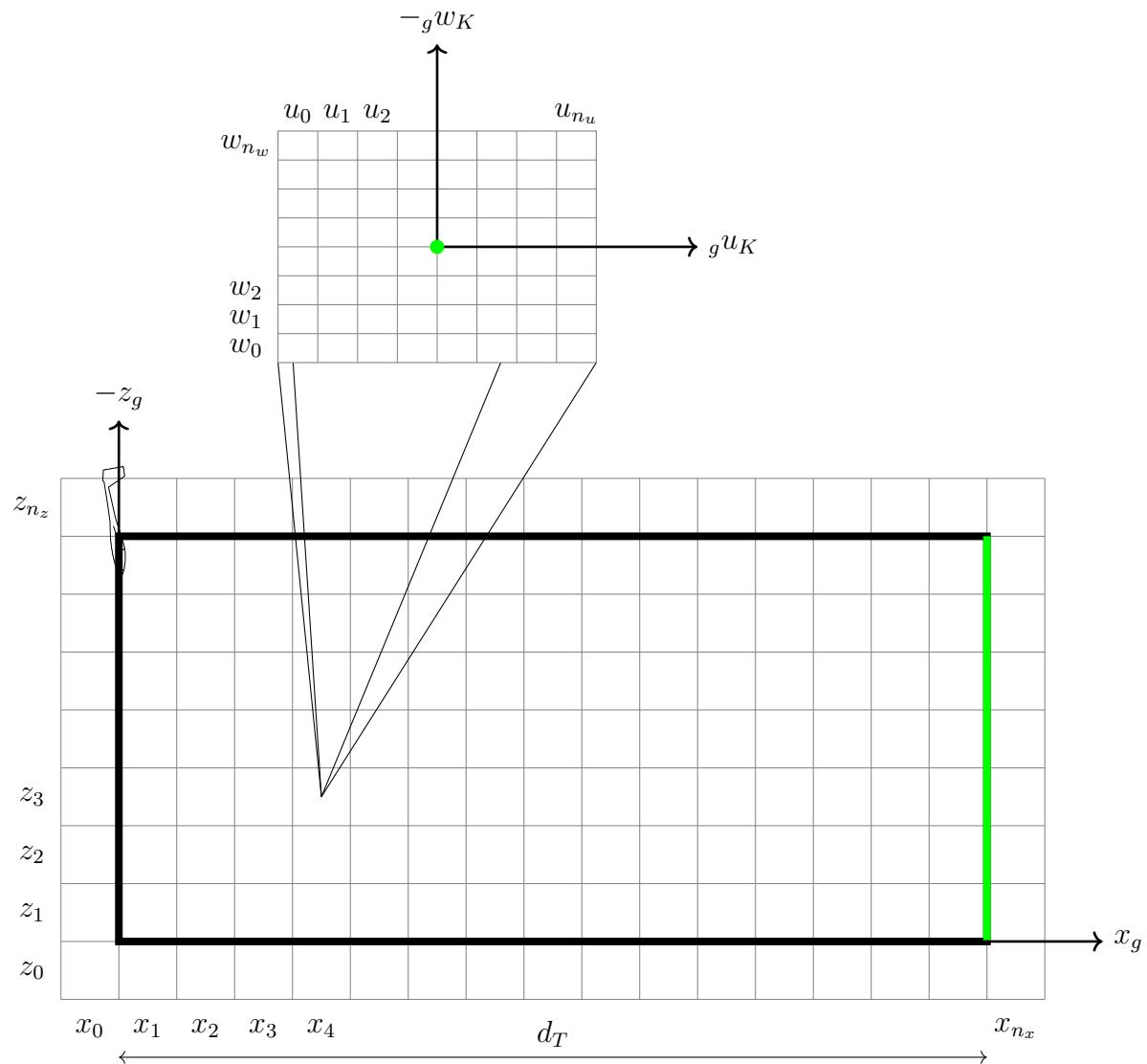


Figure 4.2.: The discretized state space in all 2D scenarios.

### 4.1. 2D Environment

In the 2D-scenario, the action space is one-dimensional, with the only action being the angle of attack  $\alpha$ . For the policy improvement step, a finite number of evenly spaced actions is sampled from the infinite set of possible actions between  $\alpha_{min} = 0$  and  $\alpha_{max} = 0.2$ . Assuming that the mapping from actions to returns is continuous, the action that yields the maximum expected return out of the sampled actions is an approximation of the true greedy action.

By discretization, the continuous trajectory optimization problem is made time- and space-discrete. For the purpose of keeping calculation time low, discretization should be as coarse as possible without affecting the result too much. Therefore, two discretization methods are used.

For a first comparison of policy iteration and value iteration, a coarse state grid is used in combination with a sample time of  $\Delta t = 2s$  in order to keep computation time low (cf. table 4.1, column 1 (left)). The goal is to determine which algorithm obtains a reasonable result in the least time<sup>2</sup>. The algorithm that performs best on a coarse grid can be expected to also do so on a denser state grid. The process of training a neural network with greedy actions is stochastic and can therefore distort the result of the algorithm comparison. Therefore, no artificial neural network is used to represent the policy in the first comparison. Instead, the actions are taken from a table and interpolated linearly between states. The goal is to find

For further evaluation of the best algorithm from the first comparison, the grid is refined and the sample time is reduced to  $\Delta t = 1s$  (cf. table 4.1, columns 1 (right), 2 and 3). The grid density is matched with an appropriate sample time as follows: With a horizontal speed of  $15 - 25 \frac{m}{s}$ , the agent covers about  $15 - 25$  m within one second. The horizontal separation of grid points should therefore be about 10m. This way, the agent can skip one grid cell, if its velocity is sufficient. Grid cells, that are closer to the goal, usually have a higher state value. With the possibility to skip a cell, the agent is thus encouraged to fly at higher speeds. For the different scenarios, this yields the grid resolutions shown in table 4.1. In theory, more horizontal grid layers can yield a better solution. But because the tabular VI data is used to train a neural network, which "smoothes" out the given data, more than 52 layers (i.e. 2m between cells vertically) have proven not to improve the result significantly.

---

<sup>2</sup>As the number of grid points, i.e. number of states, directly affects calculation time, it is essential to keep the number of grid points as low as possible. See section 2.11.1 for further details.

#### 4. Trajectory Optimization with Dynamic Programming

scenario	1	2
distance	500 m	1000 m
$n_x$	27	52
$n_z$	27	52
$n_u$		8
$n_w$	8	8

Table 4.1.: Grid parameters for trajectory optimization

# 5. Results

## 5.1. 2D Dynamic Programming

### 5.1.1. Optimal Control Benchmark

In the 2D scenario, three Dynamic Programming algorithms are compared:

- Generalized policy iteration (cf. section 2.11.3)
- Optimistic policy iteration (cf. section 2.11.3)
- Value iteration (cf. section 2.11.3)

The agent faces three different scenarios, each with calm air. Table 5.1 shows the distance to cover, start state, and the mean and standard deviation for state normalization (cf. section 3.2) for each scenario. The agent starts at a height of 100m and with zero velocity in all scenarios. For each scenario, there is a separate mean  $\mu$  and standard deviation  $\sigma$  for each coordinate. Therefore, each column contains a vector of means and standard deviations.

scenario	1	2
distance $d_T$	500m	1000m
start state $s_0$	$s_0 = [0, 100, 0, 0]^\top$	$s_0 = [0, 100, 0, 0]^\top$
normalization $\mu$	[250, 50, 10, 5]	[500, 50, 10, 5]
normalization $\sigma$	[150, 30, 5, 2]	[300, 30, 5, 2]

Table 5.1.: Some information about the investigated scenarios

Recall that policy iteration algorithms produce a sequence of value-functions and policies that approximate  $V_*$  and  $\pi_*$ . Each iterate  $V_{k+1}$  and  $\pi_{k+1}$  is closer to  $V_*$  and  $\pi_*$  than its predecessor  $V_k$  and  $\pi_k$ . In the GPI algorithm, the iterative evaluation stops, if the values of all states change by less than the reward for one time step  $\Delta t$  within one evaluation step. When this point is reached, only the time penalty propagates through the network. At states  $s_t$ , where the agent only moves slowly, chances are that the successor state  $s_{t+1}$  is (roughly) equal to  $s_t$  (and therefore  $V(s_t) = V(s_{t+1})$ ). In such cases, the state value  $V(s_t)$  is overwritten with  $r_t + V(s_t)$  at each iteration. This goes on indefinitely if PE is

## 5. Results

not stopped manually. Every policy iteration where the evaluation step is repeated multiple times obviously takes longer than an optimistic policy iteration. The whole policy iteration algorithm has converged if - at the last policy improvement step - none of the actions have changed. This means that all actions were already optimal with respect to the current value function as well as the previous one.

Value iteration only iterates state value functions. It combines one step of policy evaluation and policy improvement in every iteration. One value iteration therefore usually takes less time than one step of policy iteration. If a close approximation of  $V_*$  is found, a policy can be derived by acting greedily with respect to the last value function iterate.

### Scenario 1 - Coarse Discretization

In scenario 1, the agent has to cover 500 m while flying through calm air. Table 5.2 shows the flight times from the start-state to the goal after policy optimization with synchronous<sup>1</sup> GPI, OPI and VI on a coarse grid (cf. Tab. 4.1, column 1 (left)). This comparison has two goals: Investigating how the algorithms deal with a coarse grid and comparing calculation times. As a benchmark, the optimal control is shown in the first column. Calculation time is the time it took to obtain the given trajectories and control sequences. Note that lower calculation and flight times are better.

Although all three DP algorithms theoretically take up to  $|\mathcal{A}|^{|S|}$  iterations to converge (c.f. section 2.11.3), they are known to typically converge after only few iterations. This is also the case in this scenario. All algorithms yield similar flight times that are significantly higher than the OC result. The main reason for this is the coarse state grid, i.e. poor approximation of the true (optimal) value function. As can be seen in table 5.2, VI takes the least calculation time. As all DP algorithms converge to a reasonable solution, only optimistic policy iteration and value iteration are considered for further investigation. Calculation time gets more critical at higher distances, i.e. more grid points (cf. section 2.11.1). Fig. 5.1 shows the trajectory obtained by synchronous VI in scenario 1. The thick line is the glider trajectory  $[x(t), z(t)]^T$ , the thin line is the control sequence  $[x(t), \alpha(t)]^T$ . Figures 5.1, 5.2 and 5.3 contain a few of the grid cells in the lower left corner to illustrate the respective discretization method. The other algorithms yield similar trajectories, which can be found in appendix A.6.

To reduce computing time, OPI and VI have also been implemented as asynchronous algorithms (cf. section 2.12). As can be seen in table 5.3, the number of required iterations is significantly lower than in the synchronous case. Each iteration takes longer than in synchronous PI and VI, because splitting the state-set into layers and distributing the data to the CPU cores requires more effort here than in the synchronous case. The trajectories obtained by the asynchronous implementations are identical to their synchronous

---

<sup>1</sup>Recall from chapter 2, that synchronous algorithms update all state values simultaneously after a sweep over the complete state space. They usually are easier to implement and do not exploit the problem structure, but take longer to converge.

### 5.1. 2D Dynamic Programming

	OC	GPI	OPI	VI
iterations (-)	-	15	18	17
calculation time (h:min)	0:15	3:15	1:54	1:42
flight time (s)	19.7	21.2	21.2	21.2

Table 5.2.: Scenario 1, comparison of synchronous algorithms on a coarse grid

counterparts (cf. Fig. 5.1, A.3 and A.4). They are therefore not presented here.

	OC	OPI	VI
iterations (-)	-	6	4
calculation time (h:min)	0:15	0:37	0:29
flight time (s)	19.7	21.2	21.2

Table 5.3.: Scenario 1, comparison of asynchronous algorithms on a coarse grid

	OC	VI
iterations (-)	-	7
calculation time (h:min)	0:15	2:04
flight time (s)	19.7	20.0

Table 5.4.: Scenario 2, performance of asynchronous VI on a fine grid

Figures A.3 and A.4 in the appendix show the results of value iteration, generalized policy iteration and optimistic policy iteration, respectively.

### Scenario 1 - Fine Discretization

In this section, the result of asynchronous value iteration on a fine state grid is shown. The number of points in the x- and z-directions are now 52 instead of 27. The speed vector discretization is like in the previous section.

The asynchronous VI algorithm works itself through the state space backwards, updating the state values close to the target first and moving backwards from there (cf. 2.12). This increases convergence speed significantly. The number of required iterations is also independent of the number of states with this method. This is not the case with synchronous updates (i.e. updating all state values at once after a complete state space sweep).

Figure 5.2 shows in blue the trajectory after only 7 value iterations. To smoothen the control sequence and trajectory, the results of VI have been used to train a neural

## 5. Results

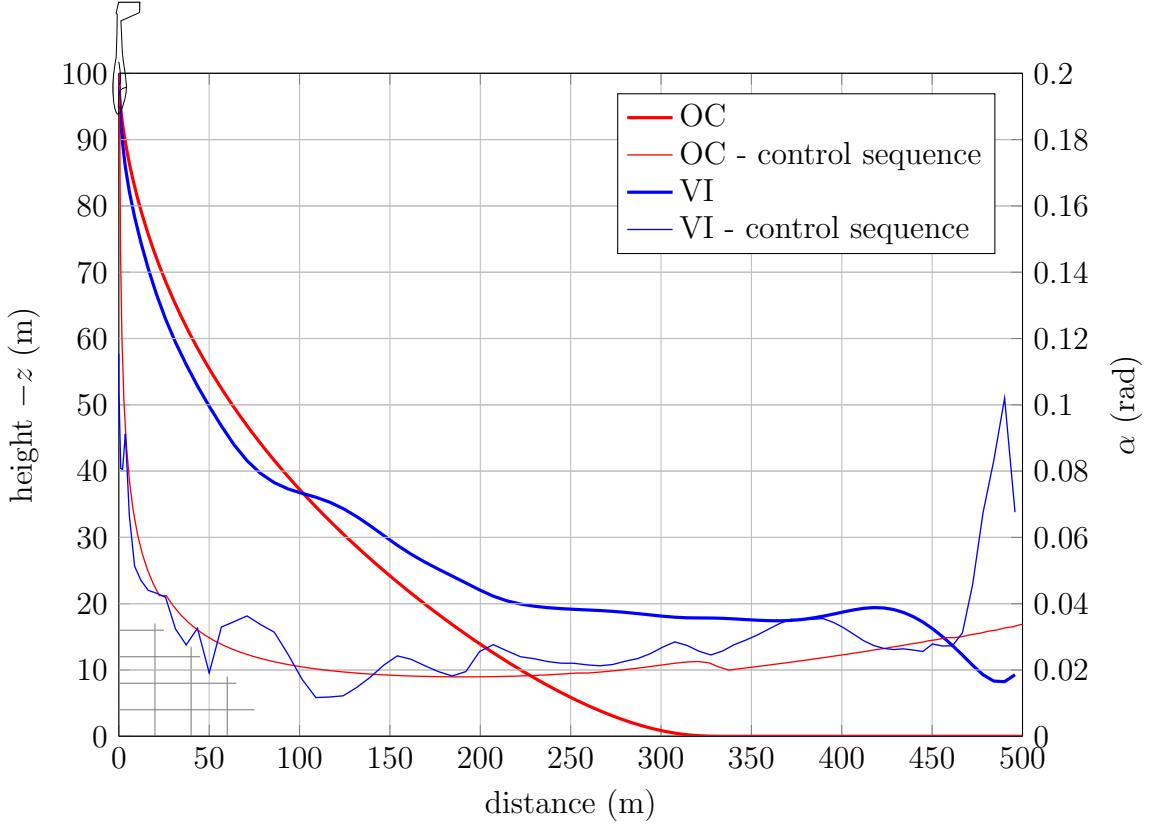


Figure 5.1.: Results of synchronous value iteration and optimal control in scenario 1

network. The actions during the simulation are ANN outputs. As can be seen, both the policy from VI and the optimal control direct the agent towards the target. The differences between the control sequences are only minor, which leads to a difference in flight times of only 0.3 s. Although the trajectories differ more, the result is satisfactory.

algorithm	OC	VI
iterations (-)	-	8
calculation time (h:min)	1:31	5:00
flight time (s)	47.0	49.8

Table 5.5.: Flight- and computation-times for OC and asynchronous VI in scenario 2

### 5.1. 2D Dynamic Programming

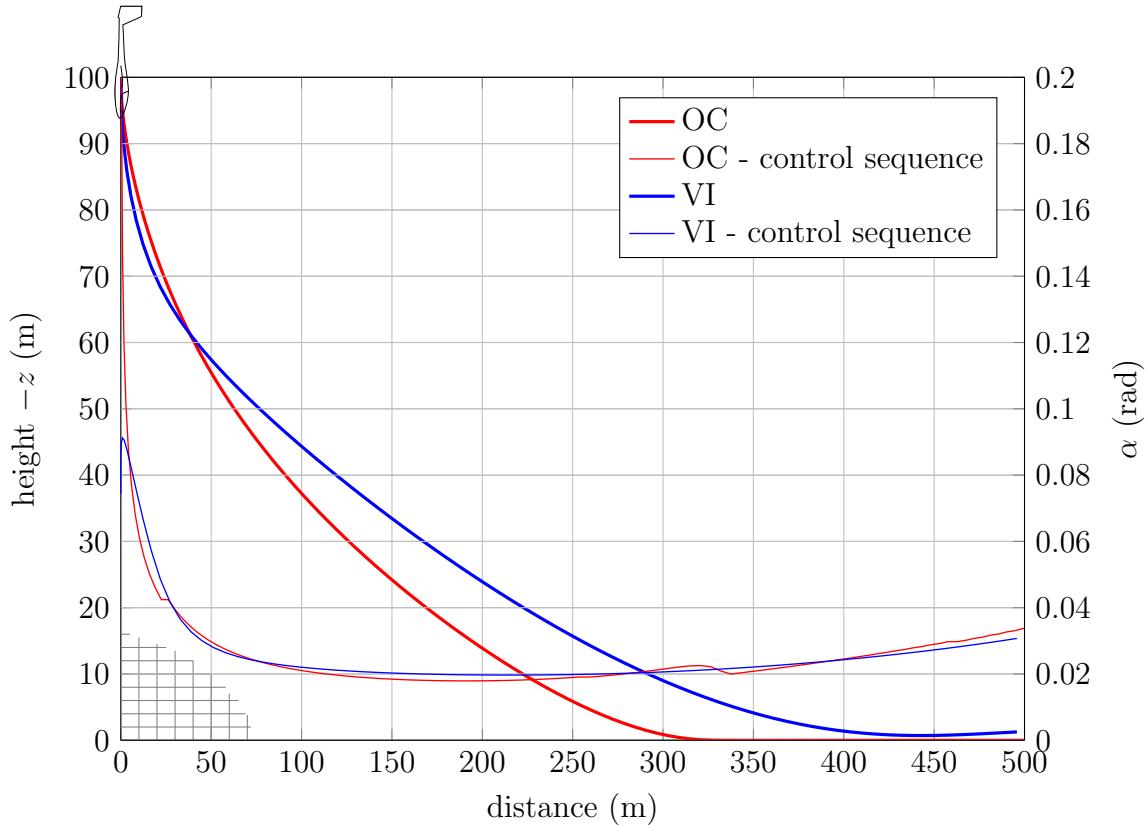


Figure 5.2.: Results of value iteration and optimal control in scenario 1 with a fine grid

### Scenario 2

In scenario 2, the agent has to cover a distance of 1000m, the start state and wind conditions are the same as in scenario 1. Fig. 5.3 shows the trajectory after 8 value iterations. The flight and calculation times can be found in table 5.5.

## 5. Results

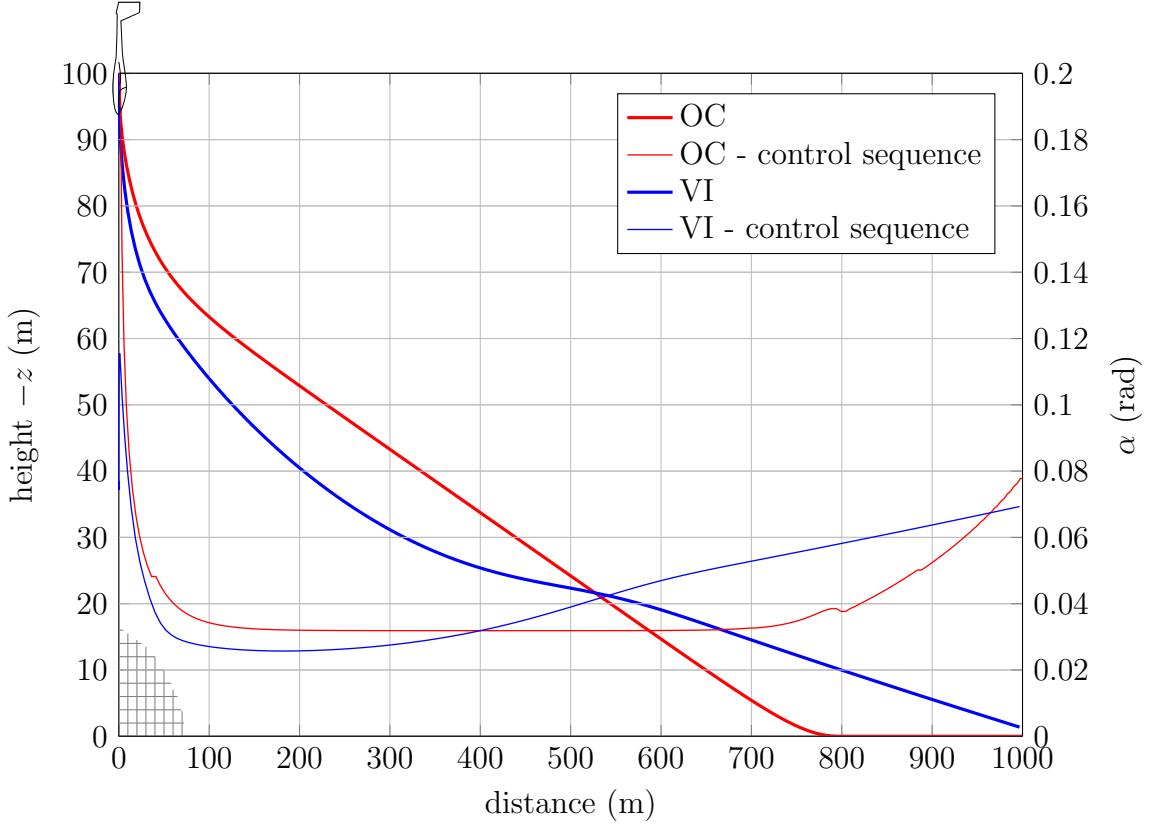


Figure 5.3.: Results of value iteration and optimal control in scenario 2

### 5.1.2. Policy Initialization for TRPO

In the previous section, it was shown that value iteration is a viable alternative for optimal control in certain path planning scenarios. The goal of this section is to find out whether value iteration can also serve to accelerate convergence of a TRPO algorithm.

The parameters of a neural network are usually initialized randomly according to a specific probability distribution (cf. section 3.2). The idea behind policy initialization by DP is to use a set of weights, that have already been optimized in a similar scenario, in order to speed up convergence of a subsequent training procedure.

The TRPO training procedure differs from the one in dynamic programming. Instead of sweeping over the complete state space and performing a one step lookahead at each point, TRPO samples complete trajectories.

The performance of TRPO can either be measured by flight times (like in section 5.1.1) or by the return per iteration. Figures 5.4 and 5.5 show the average and maximal return per iteration of TRPO, respectively. In this TRPO algorithm, an iteration consists of 25 episodes, i.e. 25 flights through the state space. After each iteration, the policy parameters are changed in order to increase the expected return. More information about

### 5.1. 2D Dynamic Programming

TRPO can be found in [15]. A detailed description of glider trajectory optimization through TRPO is contained in [20].

The benchmark for DP pre-training is scenario 2. The agent has to cover 1000m through calm air and starts at a height of 100m with zero velocity.

The plots of the average and maximum return per iteration typically have an exponential shape. The average return has numerous downward peaks, which is due to the exploration aspect of TRPO. This is the desired behavior and is not affected by pre-training, as Fig. 5.4 shows. The pre-trained policy performs better at early iterations, but the effect fades away after about 270 iterations.

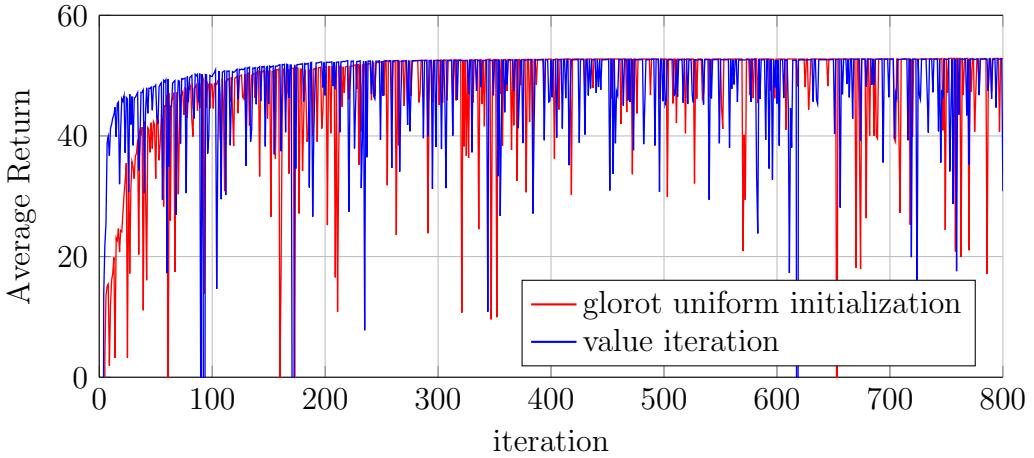


Figure 5.4.: Average return per iteration with and without policy initialization

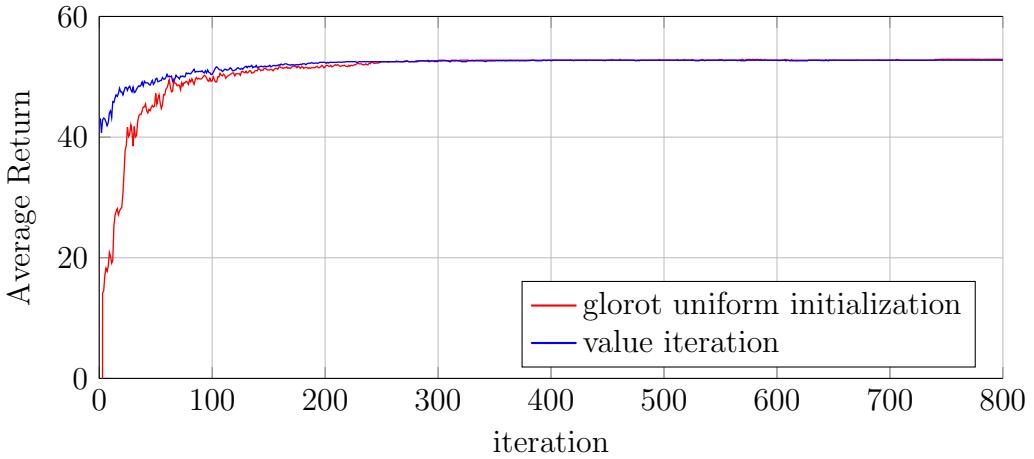


Figure 5.5.: Maximum return per iteration with and without policy initialization

## 5. Results

To see if TRPO spoils the policy's ability to find the goal from initial states that differ from table 5.1, the agent is launched from other states than before. Three policies are compared:

- $\pi_1$ : A policy obtained by VI
- $\pi_2$ : A policy obtained by TRPO<sup>2</sup>
- $\pi_3$ : A policy initialized by VI and optimized by TRPO

Table 5.6 contains the flight times obtained with the three policies from varying initial states. In the first three rows, the agent's initial state is  $s_0 = (x_0, 100\text{m}, 0\frac{\text{m}}{\text{s}}, 0\frac{\text{m}}{\text{s}})^\top$ . Rows four to six contain flight times of flights starting at  $s_0 = (x_0, 80\text{m}, -20\frac{\text{m}}{\text{s}}, 0\frac{\text{m}}{\text{s}})^\top$ . Red table cells mark flights that ended on the ground and did not reach a goal state. Green table cells mark the shortest flight time from the given initial state, i.e. the policy that did best in that case.

	$x_0 \text{ (m)}$	0	100	200	300	400	500	600	700	800	900	999
	$u_0 \text{ (\frac{m}{s})}$	0	100	200	300	400	500	600	700	800	900	999
$\pi_1$	0	49.8	44.0	38.2	33.0	27.7	23.0	18.5	14.0	10.0	6.2	-
$\pi_2$	0	47.5	43.7	35.9	31.3	8.6	14.5	8	7.8	7.4	6.1	-
$\pi_3$	0	47.3	40.9	34.9	29.3	18.2	14.2	12.3	11.2	9.4	6.2	-
$\pi_1$	-20	-	44.7	38.7	33.2	27.7	22.7	18.0	13.5	9.2	5.2	1.7
$\pi_2$	-20	-	41.8	35.3	31.3	16.4	6.1	6.0	5.6	5.4	4.9	1.4
$\pi_3$	-20	-	41.6	35.3	29.4	23.9	18.9	10.2	8.8	7.9	5.3	1.8

Table 5.6.: Flight times of agents with three different policy initialization techniques and different start states

The first policy  $\pi_1$  is the same that yielded the results in Fig. 5.3. Figures 5.6 and 5.7 show how the agent with  $\pi_1$  behaves when confronted with varying initial states. He is able to reach the goal from all tested initial states.

The trajectories obtained by  $\pi_2$  can be seen in Fig. 5.8 and 5.9. The agent manages to reach the goal only from states that lie close to  $s_{0,\text{TRPO}}$  or very close to the goal. This is not surprising, since TRPO only optimizes one trajectory, namely the one from  $s_{0,\text{TRPO}}$  to the goal. Trajectories, that lead to the goal, take less time than those obtained with  $\pi_1$  (cf. table 5.6). Trajectories, that do not start in  $s_{0,\text{TRPO}}$ , seem to try to follow the optimal trajectory from  $s_{0,\text{TRPO}}$ , but, depending on the specific initial state, make a more or less significant "detour" and are not as smooth as the VI results in Fig. 5.6 and 5.7. However, the trajectories, that reach the goal, perform better than those obtained with  $\pi_1$ .

<sup>2</sup>The start state for TRPO is  $s_{0,\text{TRPO}} = (0, 100, 0, 0)^\top$ .

### 5.1. 2D Dynamic Programming

In Fig. 5.10, the agent faces the same start states as in Fig. 5.6. The trajectory starting at  $x = 0m$  is almost identical to the optimal control result in Fig. 5.3). This is the expected result, because  $s_0 = (0, 100, 0, 0)^\top$  is also the initial state of the TRPO calculation, which, according to [20], obtains nearly optimal trajectories from a predefined start state in the given scenario. All trajectories, that start at different states, try to reach the optimal trajectory from  $s_0$  and follow it to the goal. Those, that start close to  $s_0$ , intercept the optimal trajectory at a moderate angle and succeed in reaching the goal. The trajectories starting between  $x = 400m$  and  $x = 800m$  do not end in a goal state, but on the ground. They pass through the optimal trajectory from above and touch the ground between  $x = 750m$  and  $x = 920m$ . The respective flight times are marked red in table 5.6.

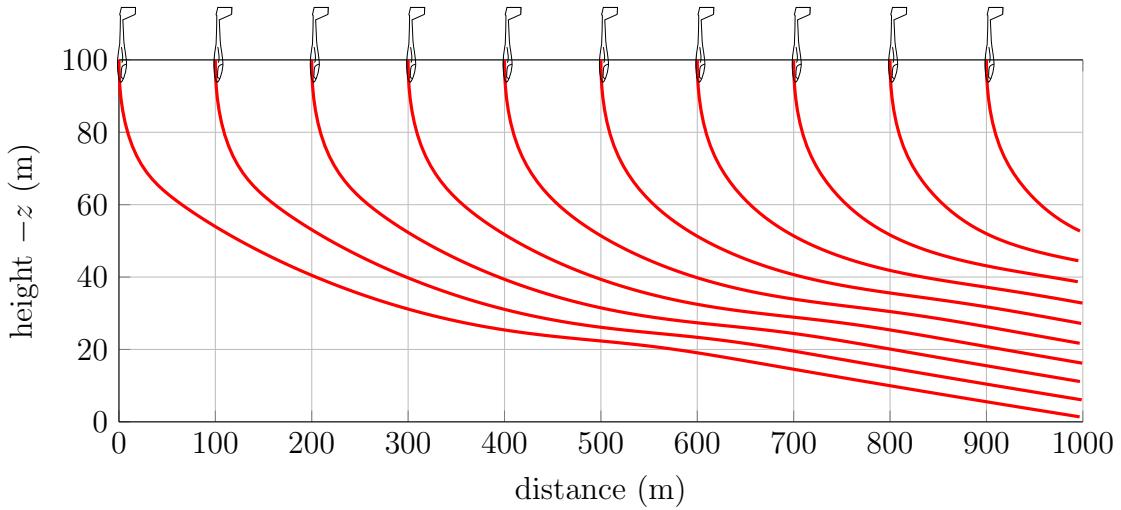


Figure 5.6.: VI - trajectories from various initial states with zero velocity

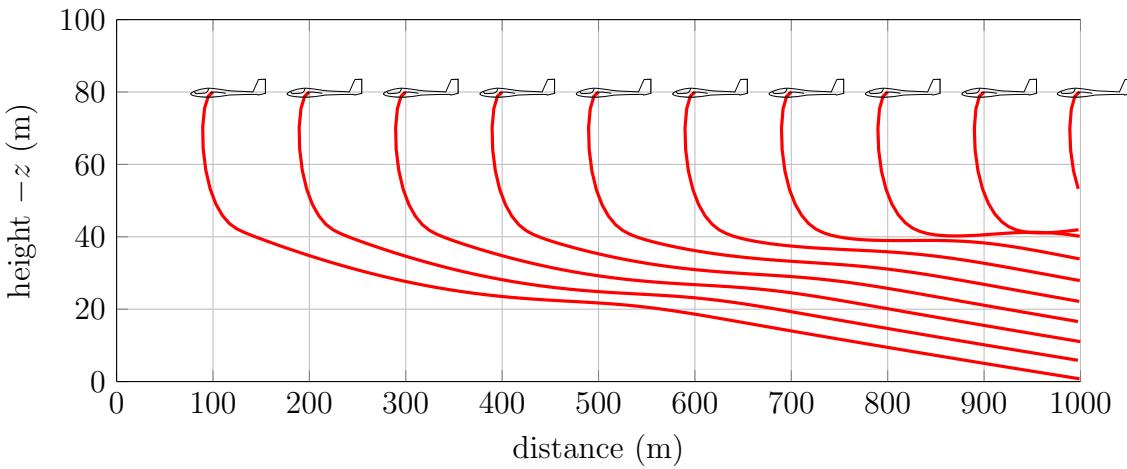


Figure 5.7.: VI - trajectories from various initial states with  $g u_K = -20 \frac{m}{s}$

## 5. Results

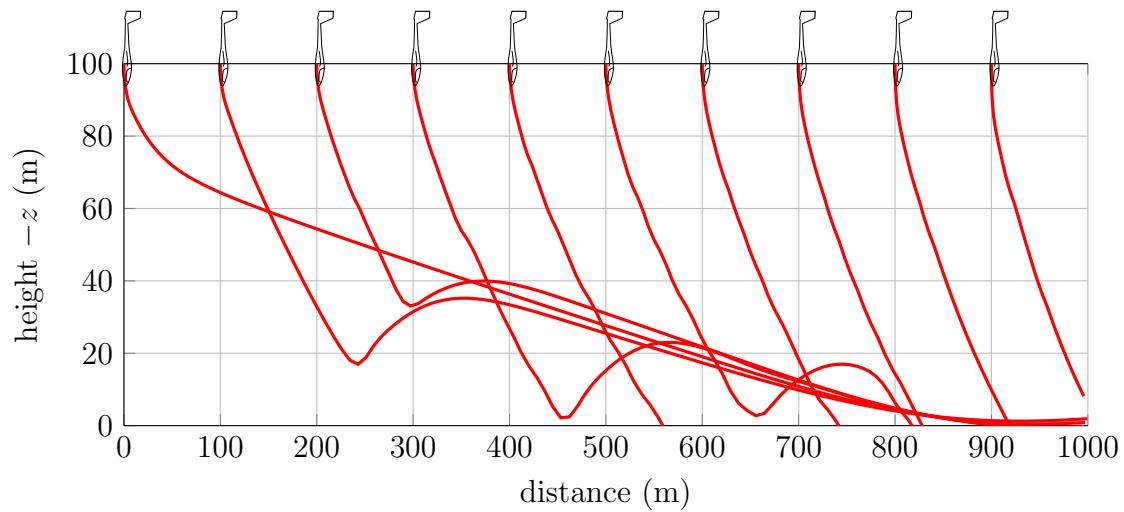


Figure 5.8.: TRPO - trajectories from various initial states with zero velocity

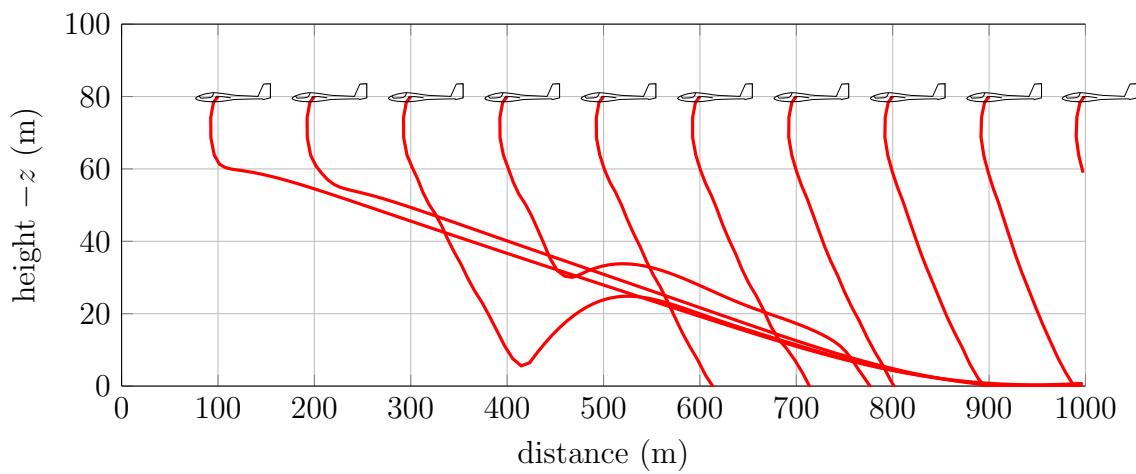


Figure 5.9.: TRPO - trajectories from various initial states with  ${}_g u_K = -20 \frac{m}{s}$

### 5.1. 2D Dynamic Programming

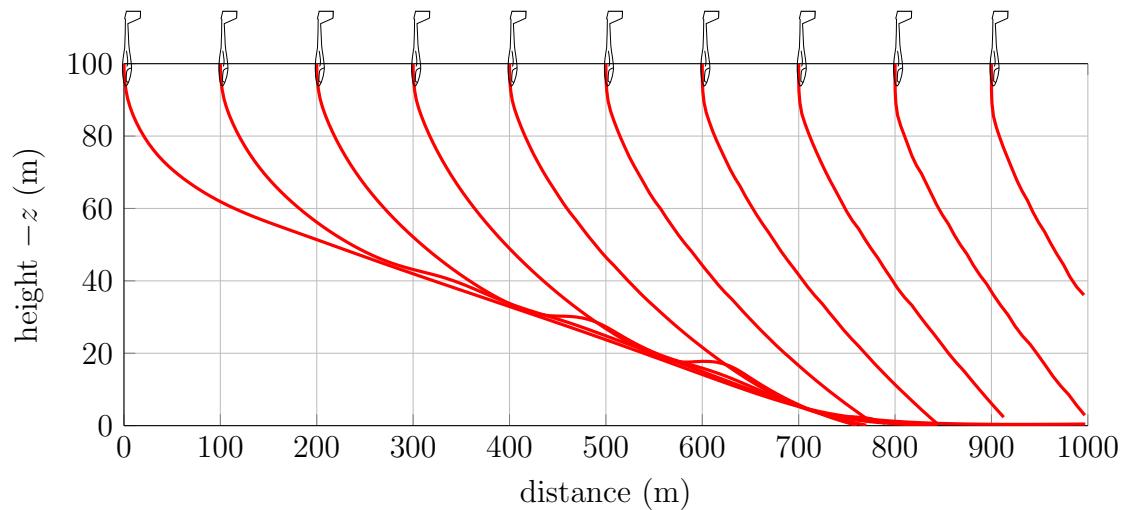


Figure 5.10.: VI and TRPO - trajectories from various initial states with zero velocity

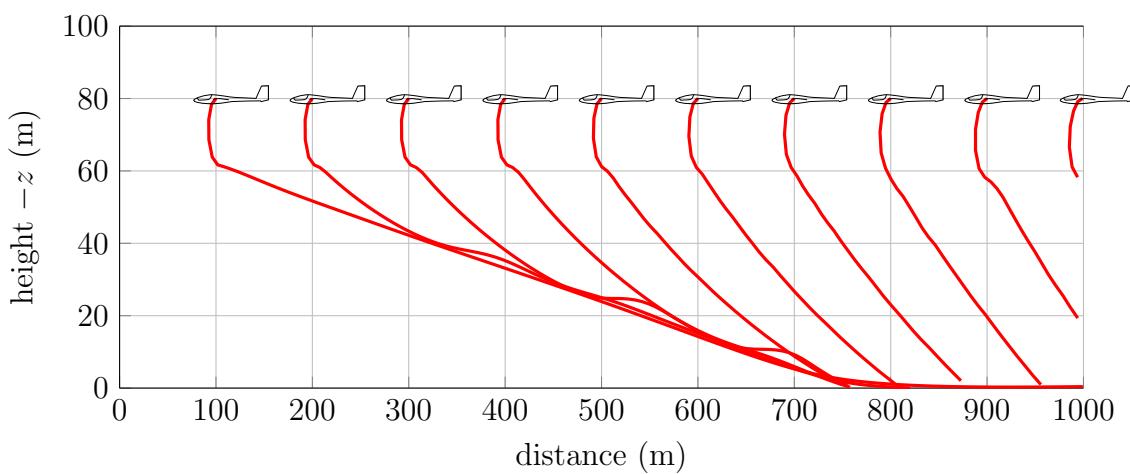


Figure 5.11.: VI and TRPO - trajectories from various initial states with  $g u_K = -20 \frac{m}{s}$



## 6. Discussion

In this work, three dynamic programming algorithms have been implemented and used for glider trajectory optimization in a calm environment. Their results were compared to an optimal control result and it has been investigated, if pre-training with DP helps convergence of a TRPO algorithm. The three algorithms all converged to a near-optimal solution in a scenario where 500m were to cover and the initial state was known a priori. Value iteration took the least time and was therefore chosen for studies in scenarios with distances of 1000m and 2000m. Results showed that TRPO performance is improved by DP. But if the pre-training time had been used for more TRPO iterations instead, the total calculation time would have been even lower.

As for the performance of PI and VI, the latter has several advantages. It is easier to implement, because it takes less steps than PI. This also makes it easier to optimize VI for a specific problem structure. In this case, performing VI layer by layer only took minor changes to the code, whereas PI took more adjustment in order to be performed this way. Because VI takes less steps for a state value update, each value iteration also takes less time in total (cf. appendix A.2).

Based on the results in chapter 5, one could argue, that dynamic programming does not make sense in this context. With a sufficiently fine state grid, even value iteration takes eight times as long as the optimal control benchmark in scenario 1 and only yields a suboptimal result in all scenarios. If compared to TRPO, a policy optimization algorithm (cf. [20]), it performs worse with respect to calculation time as well as flight time. It also cannot deal with unknown upwinds, which TRPO can.

A big advantage of DP is however, that it inherently covers the whole state space. If the result of one of the DP algorithms is used to train an artificial neural network, an agent utilizing it can get to the goal from every point in the state space<sup>1</sup>. This is not true for TRPO and OC, which both optimize a specific trajectory from a specific start state to the goal. As a tool for policy initialization, DP is especially helpful, if the start state (and therefore the approximate trajectory) is not known a priori.

Even if the start state is known, DP has an impact on the results of TRPO. As the results in Fig. 5.5 show, pre training with a DP algorithm increases the initial average and maximum return per iteration for a TRPO algorihm. If the stopping criterion of TRPO is reaching a specific average return (i.e. a specific average flight time), Fig. 5.4 shows that this goal can be reached earlier than with standard TRPO and a Glorot-initialized policy ANN. For example, a pre-trained policy net reaches an average return of 50 after only 80

---

<sup>1</sup>This is only true, if reaching the goal is physically possible at all. A glider can, for example, not reach a point that is 1000m away, if its altitude is insufficient and there are no upwinds.

## 6. Discussion

iterations, while a Glorot-initialized policy net takes about 130 iterations. Whether this higher return in early TRPO iterations justifies the additional expense of performing a DP pre-training, depends on the performance of the TRPO algorithm and can therefore not be answered generally. A DP result can also be calculated only once and then stored and used for multiple TRPO calculations afterwards (e.g. with random updrafts or horizontal winds). This way, the time DP takes carries less weight if compared to TRPO without pre-training.

Another topic in section 5.1.2 is, whether policy optimization by TRPO distorts the outputs of a policy obtained by a DP algorithm. As can be seen in table 5.6, this is the case to some extent. An agent with a pre-trained policy ( $\pi_1$ ) still finds the goal after policy optimization through TRPO ( $\pi_3$ ), if the start state is sufficiently close to the start state of TRPO. More concisely, all trajectories, that start at states between  $x_0 = 0\text{m}$  and  $x_0 = 300\text{m}$ , are improved by TRPO after DP. In fact, the results are even better than the results obtained by only performing TRPO ( $\pi_2$ ). Between  $x_0 = 400\text{m}$  and  $x_0 = 800\text{m}$ ,  $\pi_1$  is the only policy that is able to direct the agent to the target reliably. As these results suggest, a combination of DP and TRPO can retain some of the advantages of both algorithms, making DP a viable method to enhance TRPO performance and results. All in all, a combination of  $\pi_1$  and  $\pi_3$  would probably make most sense. At points, where  $\pi_3$  fails, the agent can still fall back on  $\pi_1$  and reach the goal without the need of a new TRPO.

# A. Appendix

## A.1. Glider Parameters

mass m	3.366kg
wing area S	0.568m <sup>2</sup>
aspect ratio $\Lambda$	10.2
oswald efficiency factor e	0.9
$c_{D0}$	0.015

Table A.1.: glider data

## A. Appendix

## A.2. Computer Configuration and Implementation

### A.2.1. Computer Configuration

Table A.2 contains the technical data of the computer used for the calculations.

CPU	Intel i7-7700HQ @ 2.80 GHz (4 physical cores, 8 threads)
RAM	16 GB DDR3 RAM @ 3200 MHz
GPU	Nvidia Geforce GTX 1050 Ti (4 GB VRAM, 768 CUDA cores)
OS	Ubuntu 16.04 LTS

Table A.2.: computer specifications

All algorithms were implemented in Python (Version 3.5), utilizing the rllab framework [1]. The policy ANN used in the simulations was handled by the GPU using CUDNN, Theano and Lasagne. All other code including the tables representing the greedy actions and state values was executed with eight parallel processes on the CPU.

The rllab framework already contains a step-method that lets the agent make a step through the state space. The original step method takes the current state  $s_t$  of the agent, which is stored as a class attribute, and an action  $a_t$ , which is passed in the function call. It returns the next state  $s_{t+1}$ , the reward  $r_{t+1}$  and a flag telling the agent whether he has reached a terminal state. This step-method has been modified so it also takes the current state  $s_t$  as an argument. This way, a one step lookahead can be performed from any state regardless of the state that is currently stored in the class attribute.

### A.2.2. Computational Subtleties of PI and VI

As mentioned in chapter 6, a value iteration takes less time than a policy iteration. This is due to a few subtle differences between VI and PI. In every Policy Iteration step, the following actions are performed for each state:

- Update the state value by performing a one step lookahead with the current greedy action (cf. Eq. 2.30).
- Sample a discrete set of actions  $\mathcal{A}_D$  and perform a one step lookahead for each action. This yields an estimate of  $Q(s, a)$  for each  $a \in \mathcal{A}_D$ .
- Pick the action that yields the maximum expected return out of  $\mathcal{A}_D$  according to Eq. 2.33.

The most computationally expensive part here is the third item, which involves performing multiple one step lookaheads, picking the largest action value from a list, finding the respective list index and picking the greedy action from  $\mathcal{A}_D$ . Another disadvantage is that PI requires more memory, because the greedy actions have to be stored explicitly for each iteration  $i$ , as they represent  $\pi_i$  which is required for iteration  $i + 1$ .

In value iteration, every iteration step involves the following actions:

- Sample a discrete set of actions  $\mathcal{A}_D$  and perform a one step lookahead for each action. This yields an estimate of  $Q(s, a)$  for each  $a \in \mathcal{A}_D$ .
- Pick the maximum action value out of all  $Q(s, a)$  from step one and save it as the new state value of  $s_t$  (cf. Eq. 2.37<sup>1</sup>)

Both algorithms perform multiple one step lookaheads. Value iteration, however, takes less time. Unlike policy iteration, it does not calculate each greedy action explicitly. Instead, it only picks the maximum achievable state value. Only after convergence of VI, one step of policy improvement is performed with the last state value function iterate to achieve a policy (cf. section 2.11.3). In total, value iteration is about 20 % faster than policy iteration.

---

<sup>1</sup>Note that, in this case,  $\max_a[r_{t+1} + \gamma V(s_{t+1})|s_t = s]$  is equivalent to  $\max_a[Q(s, a)]$ .

## A. Appendix

### A.3. Illustration of a One Step Lookahead

Fig. A.1 illustrates what happens within a one step lookahead in a 2D-scenario. In this simplified case, the velocity vector does not influence the state value and five actions are sampled. One or more actions are picked, the glider dynamics are integrated over one timestep to find the successor state  $s'(s, \alpha)$  for each action. After that, the successor value is taken from the value function table. The green numbers in the grid cells are the corresponding state values from the table. Of course, the state value function can be represented by an arbitrary function approximator instead.

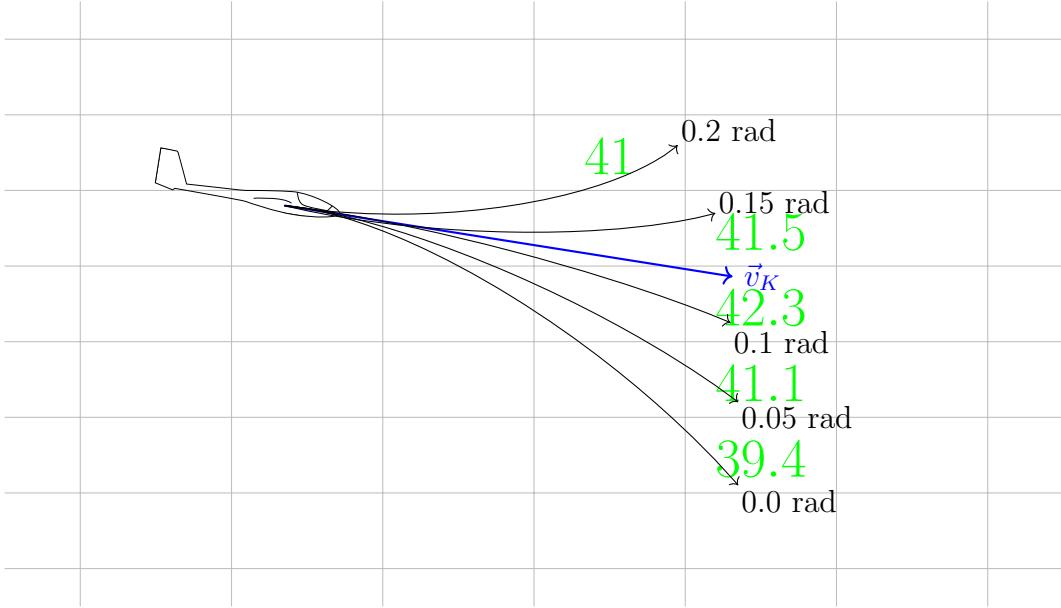


Figure A.1.: A one step lookahead in the 2D environment.

In this case, the maximum achievable successor state value is 42.3 and the corresponding greedy action is  $\alpha = 0.1$  rad.

## A.4. Optimistic Policy Iteration

---

**Algorithm 4** Optimistic Policy Iteration

---

```

function OPTIMISTIC POLICY ITERATION

    Initialize  $V(s) = 0 \forall s \in \mathcal{S}$ 
    Load arbitrary initial policy  $\pi_0$ .
    Initialize  $m, \epsilon_V$ 
    while true do

        function POLICY EVALUATION
            for all  $s \in \mathcal{S}$  do
                 $V_{\pi,new}(s) \leftarrow r + \gamma V_{\pi,old}(s')$ 
            end for
        end function

        function POLICY IMPROVEMENT
            for all  $s \in \mathcal{S}$  do
                sample  $m$  actions  $a_n(s)$ 
                for all  $a_n$  do
                     $Q(s, a_n) = r + \gamma V(s')$ 
                end for
                 $a_{greedy}(s) = \underset{a_n}{\operatorname{argmax}}[Q(s, a_n)]$ 
            end for
        end function
        Train policy on  $a_{greedy,new}$  to obtain  $\pi_{new}$ 

        if  $\|a_{greedy,new} - a_{greedy,old}\|_\infty = 0$  then
            break
        end if
    end while

end function

```

---

## A. Appendix

### A.5. Table Representation of the Value Function and Policy

The state value function and the greedy action for each state on the grid is saved in a four-dimensional table, which has been implemented in the context of this work. The table is implemented as a python class and contains methods to read data from a table cell and update cell contents. Because it also contains information about the state grid itself, it can transform a physical state into a set of table indices and back. This makes its usage in the DP algorithms simple. Every cell can either be accessed by an array of four physical coordinates  $[x, z, u, w]$  or four indices  $[i_x, i_z, i_u, i_w]$ . Unlike two-dimensional tables, a 4D table cannot be visualized easily. Fig. A.2 visualizes how a four-dimensional table works by means of multiple 2D tables.

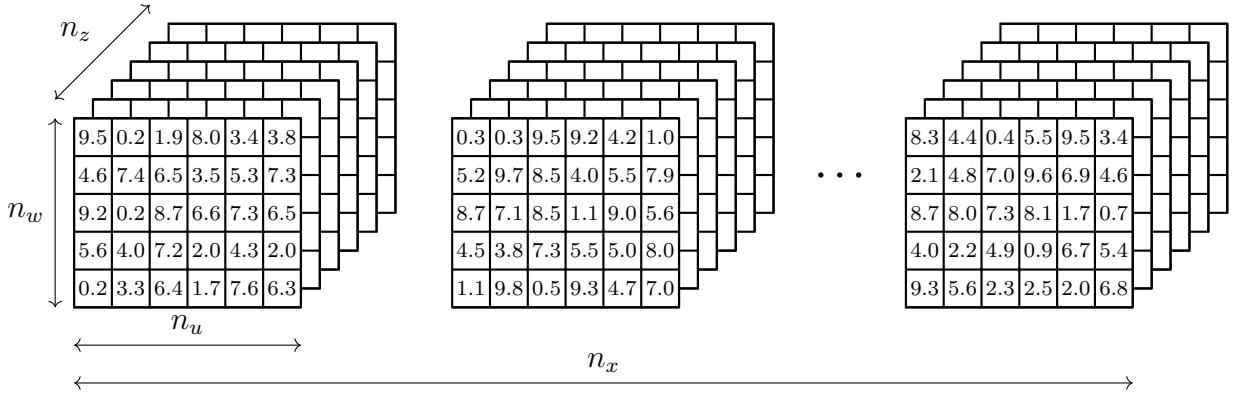


Figure A.2.: A four-dimensional table.

A 3D table can be considered as a stack of 2D tables. Likewise, a 4D table can be considered as a stack of 3D tables.

## A.6. Scenario 1: Additional Results

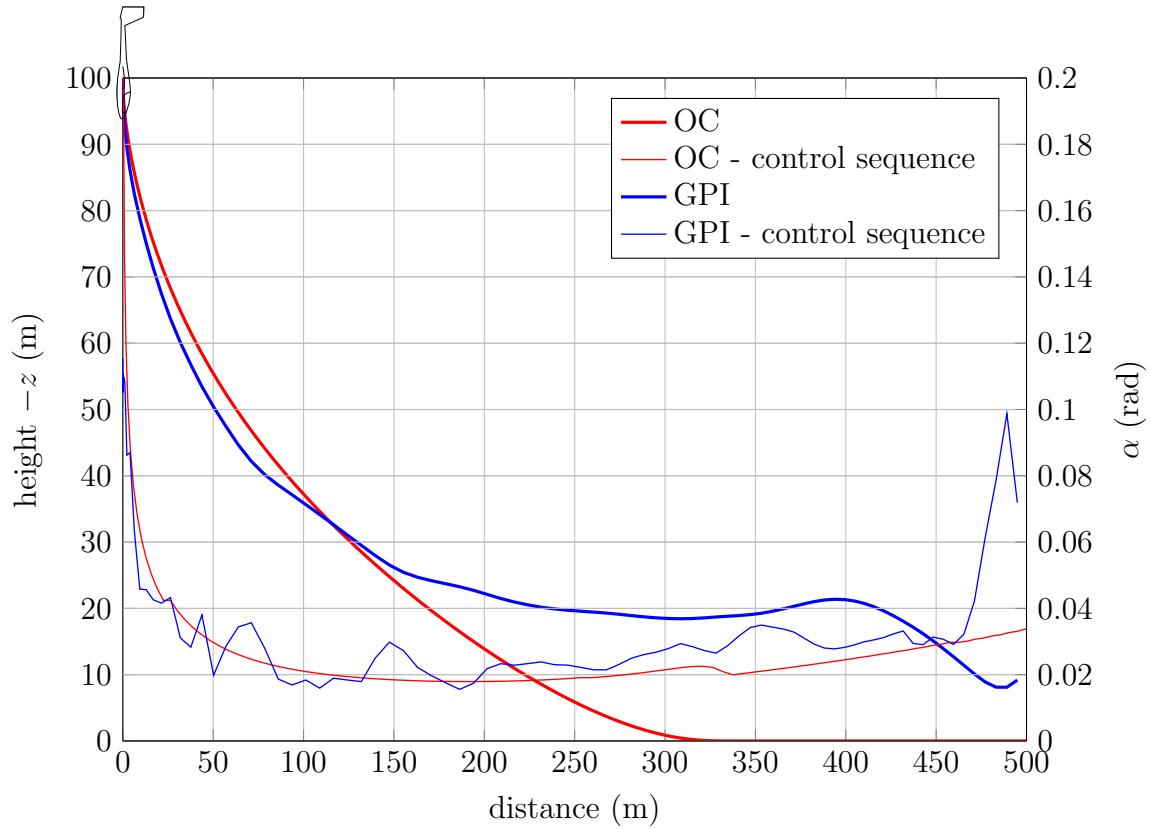


Figure A.3.: Results of synchronous generalized policy iteration and optimal control in scenario 1 with a coarse grid

## A. Appendix

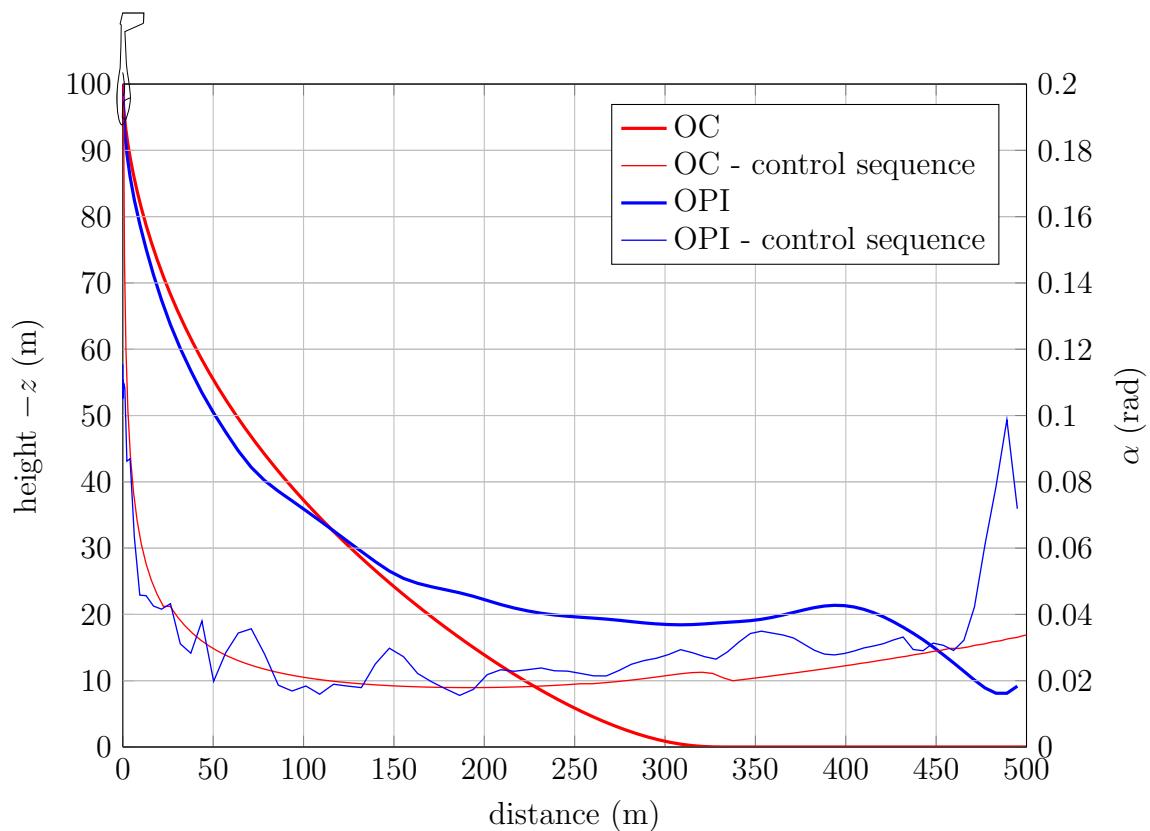


Figure A.4.: Results of synchronous optimistic policy iteration and optimal control in scenario 1 with a coarse grid

# Bibliography

- [1] *rllab-website*. <https://rllab.readthedocs.io/en/latest/>. Version: 2018. – visited on August 17, 2018
- [2] ABBEEL, P. ; SCHULMANN, J. : *Deep Reinforcement Learning through Policy Optimization*. <https://people.eecs.berkeley.edu/~pabbeel/nips-tutorial-policy-optimization-Schulman-Abbeel.pdf>. Version: 2016
- [3] BELLMAN, R. E.: The Theory of Dynamic Programming. (1954)
- [4] BELLMAN, R. E.: Applied Dynamic Programming. (1962)
- [5] BERTSEKAS, D. : Convergence of discretization procedures in dynamic programming. In: *IEEE Transactions on Automatic Control* 20 (1975), June, Nr. 3, S. 415–419. <http://dx.doi.org/10.1109/TAC.1975.1100984>. – DOI 10.1109/TAC.1975.1100984. – ISSN 0018–9286
- [6] CAVAZOS-CADENA, R. : Undiscounted Value Iteration in Stable Markov Decision Chains with Bounded Rewards. (1998), 04
- [7] FICHTER, W. ; GRIMM, W. : *Flugmechanik*. Shaker Verlag, 2009
- [8] GRIMM, W. : *Bahnoptimierung für Luft- und Raumfahrzeuge*. 2004. – Skriptum zur Vorlesung
- [9] HORDIJK, A. ; PUTERMAN, M. L.: On the Convergence of Policy Iteration in Finite State Undiscounted Markov Decision Processes: The Unichain Case. In: *Mathematics of Operations Research* 12 (1987), Nr. 1, 163–176. <http://www.jstor.org/stable/3689678>. – ISSN 0364765X, 15265471
- [10] KINGMA, D. P. ; BA, J. : Adam: A Method for Stochastic Optimization. In: *CoRR* abs/1412.6980 (2014). <http://arxiv.org/abs/1412.6980>
- [11] KRIESEL, D. : *A Brief Introduction to Neural Networks* available at <http://www.dkriesel.com>
- [12] NOTTER, S. ; ZUERN, M. ; GROSS, P. ; FICHTER, W. : Reinforced Learning to Cross-Country Soar in the Vertical Plane of Motion. (2018)

## Bibliography

- [13] POWELL, W. B.: *Approximate dynamic programming: solving the curses of dimensionality*. Wiley-Interscience (Wiley series in probability and statistics). – XVI, 460 Seiten S. <http://swbplus.bsz-bw.de/bsz275179400cov.htm>. – ISBN 978-0-470-17155-4. – UB Vaihingen
- [14] REDDY, G. ; CELANI, A. ; SEJNOWSKI, T. J. ; VERGASSOLA, M. : Learning to soar in turbulent environments. In: *Proceedings of the National Academy of Sciences* 113 (2016), Nr. 33, E4877–E4884. <http://dx.doi.org/10.1073/pnas.1606075113>. – DOI 10.1073/pnas.1606075113. – ISSN 0027-8424
- [15] SCHULMAN, J. ; LEVINE, S. ; MORITZ, P. ; JORDAN, M. I. ; ABBEEL, P. : Trust Region Policy Optimization. In: *CoRR* abs/1502.05477 (2015). <http://arxiv.org/abs/1502.05477>
- [16] SILVER, D. : *UCL Course on Reinforcement Learning*. 2015
- [17] SUTTON, R. S. ; BARTO, A. G.: *Reinforcement Learning - An Introduction*. The MIT Press, 2018
- [18] TSITSIKLIS, J. N.: On the Convergence of Optimistic Policy Iteration. In: *J. Mach. Learn. Res.* 3 (2003), März, 59–72. <http://dx.doi.org/10.1162/153244303768966102>. – DOI 10.1162/153244303768966102. – ISSN 1532-4435
- [19] WIRTH, L. ; OETTERSHAGEN, P. ; AMBUHL, J. ; SIEGWART, R. : Meteorological path planning using dynamic programming for a solar-powered UAV. In: *2015 IEEE Aerospace Conference* (2015), S. 1–11
- [20] ZUERN, M. : *Autonomous Soaring through Unknown Atmosphere, Master Thesis*. Institute for Flight Mechanics and Control, Stuttgart University, 2017