

Experiment zur Verbesserung der Robustheit von Reinforcement Learning Policies anhand trainiertem Gegenspieler in Dronensimulationen

Bachelorarbeit

vorgelegt am 11. März 2023

Fakultät Wirtschaft

Studiengang Wirtschaftsinformatik

Kurs WWI2020F

von

LEON HENNE

Betreuerin in der Ausbildungsstätte: DHBW Stuttgart:

IBM Deutschland GmbH
Sophie Lang
Senior Data Scientist

Prof. Dr. Kai Holzweißig
Studiendekan Wirtschaftsinformatik

Unterschrift der Betreuerin

Vertraulichkeitsvermerk: Der Inhalt dieser Arbeit darf weder als Ganzes noch in Auszügen Personen außerhalb des Prüfungs- und Evaluationsverfahrens zugänglich gemacht werden, sofern keine anders lautende Genehmigung des Dualen Partners vorliegt.

Inhaltsverzeichnis

Abkürzungsverzeichnis	III
Abbildungsverzeichnis	IV
Tabellenverzeichnis	V
1 Einleitung	1
1.1 Problemstellung	1
1.2 Zielsetzung	2
1.3 Forschungsfrage	2
1.4 Forschungsmethodik	3
1.5 Aufbau der Arbeit	3
2 Diskussion des aktuellen Stands der Forschung und Praxis	4
2.1 Aufbau der Literaturrecherche	4
2.2 Verstärkendes Lernen	5
2.2.1 Wertebasierende Methoden	7
2.2.2 Strategiebasierende Methoden	7
2.2.3 Akteur-Kritiker Methoden	8
2.2.4 Abgrenzung zu Multi-Agent Reinforcement Learning (MARL) Algorithmen	9
2.2.5 Limitierungen und Herausforderungen von RL	9
2.3 Simulationsumgebungen für RL	10
2.3.1 Definitionen von Simulationsumgebungen	10
2.3.2 Entwicklung von Simulationsumgebungen für RL Anwendungen	11
2.3.3 aktuelle Physik-Engines und Simulationsanwendungen	13
3 Durchführung des Laborexperiments	14
4 Ergebnisse des Laborexperiments	15
5 Reflexion und Forschungsausblick	16
Anhang	17
Literaturverzeichnis	19

Abkürzungsverzeichnis

DHBW	Duale Hochschule Baden-Württemberg
RL	Reinforcement Learning
KPI	Key Performance Indicator
MARL	Multi-Agent Reinforcement Learning
DART	Dynamic Animation and Robotics Toolkit
ODE	Open Dynamics Engine

Abbildungsverzeichnis

1	vereinfachte Darstellung der Interaktion zwischen dem Agenten und seiner Umgebung	5
2	Klassifizierung von Algorithmen im Bereich des RL	6

Tabellenverzeichnis

1	Konzept Matrix für Artikel zu Simulationsumgebungen und zur Robustheit RL Algorithmen nach Webster/Watson 2002. Legende: RL (Reinforcement Learning), MARL (Multi-Agent Reinforcement Learning), ES (Entwicklung von Simulationsumgebungen), DS (Dronensimulation), KS (kompetitive Simulationsumgebungen), DR (Domain Randomization), RRLP (Robustheit von RL Policies), LE (Laborexperimente)	5
2	wichtigsten Kriterien zur Auswahl von Simulatoren ¹	12

¹Ivaldi/Padois/Nori 2014, S. 4

1 Einleitung

1.1 Problemstellung

Reinforcement Learning (RL) findet heutzutage bereits Anwendung in vielerlei Forschungsprojekten wie Deepmind AlphaStar oder OpenAI Five, aber auch in Produkten und Dienstleistungen wie AWSDeepRacer oder Metas Horizon open-source RL-Plattform.² RL ist im Bereich des maschinellen Lernens eine Herangehensweise zur Lösung von Entscheidungsproblemen.³ Ein Software-Agent leitet dabei durchzuführende Aktionen aus seiner Umgebung ab, mit dem Ziel die kumulierte erhaltene Belohnung zu maximieren, währenddessen sich seine Umgebung durch alle Aktionen verändert.⁴ Die Umgebungen beinhalten in ihrer einfachsten Form eine simulierte Welt, welche zu jedem Zeitschritt eine Aktion entgegennimmt, und den eigenen nächsten Zustand sowie einen Belohnungswert zurückgibt.⁵ Da ein Problem beim Einsatz von RL Algorithmen die Limitierungen sein können, Daten in der echten Welt zu sammeln und fürs Training zu verwenden, werden häufig hierfür Simulationsumgebungen eingesetzt.⁶ Eine Limitierung können bspw. Sicherheitsaspekte sein, welche beim Training von Roboterarmen, oder sich autonom bewegenden Systemen auftreten, da die einzelnen physischen Bewegungen nicht vorhersehbar abschätzbar sind.⁷ Simulationen nehmen damit als Testumgebung eine wichtige Rolle ein in der Entwicklung von Kontrollalgorithmen.⁸ Insgesamt bedarf die erfolgreiche Anwendung von Reinforcement Learning demnach nicht nur effiziente Algorithmen, sondern auch geeignete Simulationsumgebungen.⁹ Besonders schwierig, und daher sehr wichtig zu erforschen, ist es die Trainingsumgebung bestmöglich an die echte Welt anzupassen, sodass bspw. die Agenten für Roboter und autonome Fahrzeuge, nach dem Training mit generalisierten Policies in der Realität eingesetzt werden können.¹⁰ In der Forschungsliteratur wird diese beschriebene Problematik als „Sim to real“-Transfer beschrieben.¹¹ Eine Domäne der echten Welt wird dabei eher selten ausschließlich von veränderten dynamischen Parametern und nur einer Person oder nur einer Organisation geprägt. Oftmals beeinflussen mehrere Parteien teilweise kooperierend aber auch teilweise konkurrierend den eigenen Erfolg, wie bspw. einen dem Wettbewerb unterliegenden Markt. Stellt man sich ein solches Szenario vor, ist es naheliegend, dass auch jene Einflüsse möglichst präzise in die Simulationsumgebung integriert sein müssen, um ein generalisierendes Modell erlernen zu können. Während bereits in Produkten wie Powertac nach Collins/Ketter 2022 die Simulation von Märkten entwickelt wurde, scheint der Einfluss des Gegenspielers in kompetitiven

²Vgl. Li 2019, S. 4

³Vgl. Schuderer/Bromuri/van Eekelen 2021, S. 3

⁴Vgl. Schuderer/Bromuri/van Eekelen 2021, S. 3

⁵Vgl. Reda/Tao/van de Panne 2020, S. 1

⁶Vgl. Zhao/Queralta/Westerlund 2020, S. 737

⁷Vgl. Zhao/Queralta/Westerlund 2020, S. 738

⁸Vgl. Cutler/Walsh/How 2014, S. 2

⁹Vgl. Reda/Tao/van de Panne 2020, S. 8

¹⁰Vgl. Slaoui u. a. 2019, S. 1

¹¹Vgl. Zhao/Queralta/Westerlund 2020, S. 738

Simulationen auf die Robustheit von RL Algorithmen und demnach auf die Lösung des „Sim to real“-Transfers unerforscht.

1.2 Zielsetzung

Daher soll im Rahmen dieser Arbeit untersucht werden, ob die Integrierung eines RL basierten Gegenspielers in einer Simulation die Umgebung so beeinflussen kann, dass die erlernten Verhaltensmodelle, welche im Kontext von RL oftmals als Policies referenziert werden, robuster agieren unter den veränderten dynamischen Bedingungen und alternativen deterministischen Gegenspielern im Testszenario.

Dazu soll eine kompetitive Simulationsumgebung entwickelt werden, in welcher sich zwei konkurrierender Spieler in Form von Flugobjekten spielerisch gegenseitig bekämpfen. In der Simulation werden folgend Policies in drei verschiedenen Szenarien trainiert.

- Training mit regelbasiertem Gegenspieler unter gleichbleibenden Dynamikparametern
- Training mit RL basiertem Gegenspieler unter gleichbleibenden Dynamikparametern
- Training mit regelbasiertem Gegenspieler unter sich verändernden Dynamikparametern

Anschließend werden alle trainierten Policies in einer Reihe von Testszenarien untersucht. Jedes Testszenario verfügt dabei über festgelegte sich vom Training unterscheidende Dynamikparameter und jeweils leicht unterschiedliche Handlungspräferenzen des deterministischen Gegenspielers. Bei der Untersuchung werden jeweils die folgenden Variablen als Key Performance Indicator (KPI) betrachtet.

- durchschnittlich erzielte Belohnung
- Varianz der Belohnungen
- Anzahl an unbeabsichtigten Abstürzen

Aus der Auswertung der Testszenarien kann der Effekt des RL basierten Gegenspielers auf die Robustheit mittels des Vergleichs mit dem regelbasierten Gegenspieler und der Domain Randomization evaluiert werden.

1.3 Forschungsfrage

Aus der beschriebenen Problemstellung und der für den Rahmen dieser Arbeit festgelegten Zielsetzung ergibt sich folgende Forschungsfrage:

Kann durch den Einsatz eines mittels RL trainierten Gegenspielers die Robustheit der gelernten Policy verbessert werden?

Zur Beantwortung der Forschungsfrage werden folgende Hypothesen aufgestellt und im Rahmen der Arbeit untersucht:

Hypothese 1: *Die in den Testszenarien durchschnittlich erzielte Belohnung ist unter Verwendung der Policy aus dem Training mit RL basiertem Gegenspieler signifikant und zuverlässig höher als die Policy aus dem Training mit regelbasiertem Gegenspieler.*

Hypothese 2: *Die Varianz der in den Testszenarien erzielten Belohnung ist unter Verwendung der Policy aus dem Training mit RL basiertem Gegenspieler signifikant und zuverlässig geringer als die Policy aus dem Training mit regelbasiertem Gegenspieler.*

Hypothese 3: *Die in den Testszenarien erreichte Anzahl von unbeabsichtigten Abstürzen ist unter Verwendung der Policy aus dem Training mit RL basiertem Gegenspieler signifikant und zuverlässig geringer als die Policy aus dem Training mit regelbasiertem Gegenspieler.*

1.4 Forschungsmethodik

Als Forschungsmethodik soll im Rahmen dieser Arbeit ein quantitatives Laborexperiment nach Recker 2021 durchgeführt werden. Hierbei wird häufig nach dem hypothetisch-deduktives Modell vorgegangen, in welchem Hypothesen formuliert, empirische Studien entwickelt, Daten gesammelt, Hypothesen anhand dessen evaluiert und gewonnene Erkenntnisse berichtet werden.¹² Eine Möglichkeit der Untersuchung der Ursache- und Wirkungsbeziehung stellt das Laborexperiment dar.¹³ Dabei wird die kontrollierte Umgebung der Simulation erschaffen, deren Aufbau die unabhängige Variable darstellt. Die Metriken anhand welcher die Performance und die Robustheit der trainierten Policies gemessen werden, bilden im Experiment die abhängigen Variablen.

1.5 Aufbau der Arbeit

Insgesamt gliedert sich die Arbeit nach einem Schema von Holzweißig 2022. Die Arbeit beginnt mit einem einleitenden Kapitel in welchem Motivation, Problemstellung, Zielsetzung und Forschungsmethodik erläutert sind. Anschließend wird im zweiten Kapitel der aktuelle Stand der Forschung zu den relevanten Konzepten der Problemstellung wiedergegeben. Im dritten Kapitel wird die Forschungsmethodik dargestellt, indem die Simulationsumgebung als Messinstrument entwickelt wird sowie verschiedene Messszenarien erläutert und entsprechende Daten gesammelt werden. Daraufhin sind im folgenden vierten Kapitel die Messdaten auszuwerten und aufgestellte Hypothesen zu überprüfen. Im Zuge dessen kann ebenso die Forschungsfrage anhand der Annahme oder Ablehnung der Hypothesen beantwortet werden. Abschließend wird im letzten Kapitel ein Fazit zu den erzielten Forschungsergebnissen dargelegt und ein Ausblick auf weitere Forschung gegeben.

¹²Vgl. Recker 2021, S. S.89f.

¹³Vgl. Recker 2021, S. 106

2 Diskussion des aktuellen Stands der Forschung und Praxis

2.1 Aufbau der Literaturrecherche

In Anlehnung der Literaturrecherche nach Webster/Watson 2002 wurden alle voraussichtlich benötigten Konzepte für die Durchführung der beschriebenen Forschungsmethodik in Tabelle 1 festgehalten. Alle angeführten Konzepte wurden mittels verschiedener Suchbegriffe in Suchmaschinen, Datenbanken und Bibliotheken wie *Google Scholar* 2/28/2023, *IEEE Xplore* 2/28/2023 oder die digitale Bibliothek der Association for Computing Machinery (ACM) ACM Digital Library 2/28/2023 recherchiert. In der daraus gefundenen Literatur wurden zitierte Werke ebenfalls nach den beschriebenen Konzepten durchsucht und insgesamt jede Literaturquelle in Tabelle 1 den in ihnen enthaltenen Konzepten zugeordnet.

Artikel	Konzepte							
	RL	MARL	ES	DS	KS	DR	RRLP	LE
Sutton/Barto 2018	X							
Li 2019	X							
Zhao/Queralta/Westerlund 2020	X		X			X		
Wang/Hong 2020	X							
Zhang/Wu/Pineau 2018	X		X				X	
Cutler/Walsh/How 2014	X							
Canese u. a. 2021	X	X						
Reda/Tao/van de Panne 2020	X		X			X		
Ningombam 2022	X							
Wong u. a. 2022	X	X						
Schuderer/Bromuri/van Eekelen 2021	X	X	X					
Körber u. a. 2021			X					
Bharadhwaj u. a. 2019			X					
Foronda 2021			X					
Maria 1997			X					
Brockman u. a. 2016	X		X					
Yan Duan u. a. 2016	X		X				X	
Ivaldi/Padois/Nori 2014			X					
Ayala u. a. 2020			X					
Todorov/Erez/Tassa 2012			X					

Artikel	Konzepte
---------	----------

Tab. 1: Konzept Matrix für Artikel zu Simulationsumgebungen und zur Robustheit RL Algorithmen nach Webster/Watson 2002. Legende: RL (Reinforcement Learning), MARL (Multi-Agent Reinforcement Learning), ES (Entwicklung von Simulationsumgebungen), DS (Dronensimulation), KS (kompetitive Simulationsumgebungen), DR (Domain Randomization), RRLP (Robustheit von RL Policies), LE (Laborexperimente)

2.2 Verstärkendes Lernen

Verstärkendes Lernen oder als RL in der Fachsprache bezeichnet, definiert einen konzeptionellen Ansatz zielorientiertes Lernen von Entscheidungen zu verstehen und zu automatisieren.¹⁴ Dabei besteht der Fokus darauf, dass ein Agent aus der direkten Interaktion mit seiner Umgebung lernt, ohne dass explizite Überwachung notwendig ist.¹⁵ Der Agent lernt über die Zeit eine optimale Strategie zur Lösung des Entscheidungsproblems aus dem Ausprobieren und Scheitern mittels verschiedener Aktionen die gewünschte Veränderung in seiner Umwelt herzustellen.¹⁶ Notwendig dabei ist es, dass der Agent den Zustand seiner Umgebung wahrnehmen, und auch durch entsprechende Aktionen beeinflussen kann, sodass die Erreichung des Zielzustandes möglich ist.¹⁷ Zur Erreichung dieses Zielzustandes muss der Agent alle Aktionen entdecken, welche ihm die größtmögliche kumulierte Belohnung liefern, wobei Aktionen nicht nur die unmittelbare sondern auch zukünftige Belohnungen beeinflussen.¹⁸ Zusammengefasst lässt sich die beschriebene Interaktion des Agenten mit seiner Umgebung wie folgt in Abbildung 1 darstellen.

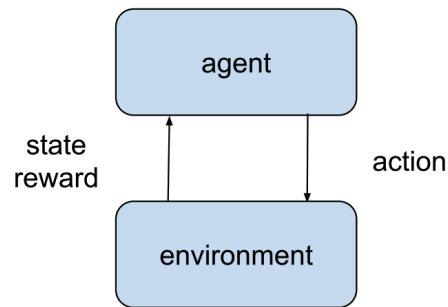


Abb. 1: vereinfachte Darstellung der Interaktion zwischen dem Agenten und seiner Umgebung¹⁹

Ein Standardaufbau einer Aufgabe für verstärkendes Lernen kann demnach verstanden werden, als sequentielles Entscheidungsproblem zu dessen Lösung ein Agent zu jedem diskreten Zeitschritt eine Aktion ausführt, welche den Zustand der Umgebung verändert.²⁰ Betrachtet man

¹⁴Vgl. Sutton/Barto 2018, S. 13

¹⁵Vgl. Sutton/Barto 2018, S. 13

¹⁶Vgl. Li 2019, S. 4

¹⁷Vgl. Sutton/Barto 2018, S. 2

¹⁸Vgl. Sutton/Barto 2018, S. 1

¹⁹Enthalten in: Li 2019, S. 5

²⁰Vgl. Zhao/Queralt/Westerlund 2020, S. 2

die technische Umsetzung einer solchen Interaktion zwischen dem Agenten und dessen Umgebung, wird häufig zur Modellierung ein Markov Entscheidungsprozess verwendet. Im Kontext von RL ist der Entscheidungsprozess definiert nach einem Tupel aus folgenden Elementen:²¹

- Alle Zustände S
- Alle Aktionen A
- initiale Zustandsverteilung $p_0(S)$
- Übergangswahrscheinlichkeit $T(S_{t+1}|S_t, A_t)$
- Belohnungswahrscheinlichkeit $R(r_{t+1}|S_t, A_t)$

Zum Finden der optimalen Strategie existieren modellbasierende und modellfreie Algorithmen des verstärkenden Lernens.²² Bei modellbasierenden Algorithmen wird das Umgebungsverhalten, also die Übergangs- und Belohnungswahrscheinlichkeiten als bekannt vorausgesetzt.²³ Unter modellbasierenden Algorithmen wird dynamische Programmierung eingesetzt, um mittels Strategieevaluation und Strategieiteration die optimale Strategie zu finden.²⁴ Unter modellfreien Algorithmen werden die drei verschiedenen Ansätze Wertebasierend, Strategiebasierend und Akteur-Kritiker basierend unterschieden²⁵ Der Agent im Kontext von modellfreien RL Methoden kennt nur die Zustände S und die Aktionen A , jedoch nicht die Umgebungsverhalten T und die Belohnungswahrscheinlichkeit R .²⁶ Fasst man die Klassifizierung der Algorithmen und Methoden von RL zusammen, lässt sie sich wie folgt darstellen:

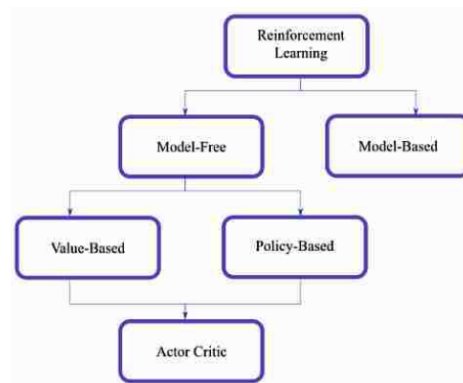


Abb. 2: Klassifizierung von Algorithmen im Bereich des RL²⁷

²¹Vgl. Zhang/Wu/Pineau 2018, S. 2

²²Vgl. Wang/Hong 2020, S. 3

²³Vgl. Wang/Hong 2020, S. 3

²⁴Vgl. Li 2019, S. 5

²⁵Vgl. Li 2019, S. 5

²⁶Vgl. Cutler/Walsh/How 2014, S. 2

²⁷Enthalten in: Canese u. a. 2021, S. 6

2.2.1 Wertebasierende Methoden

Der Agent sucht in diesem Kontext die optimale Strategie π^* , welche allen Zuständen S die jeweilige Aktion $A(S)$ zuordnet, sodass die kummulierte Belohnungswahrscheinlichkeit $R(r_{t+1}|S_t, A_t)$ über alle Zeitschritte t maximal ist.²⁸ Neben dieser kurzfristigen direkten Belohnung müssen auch die langfristigen zukünftigen Belohnungen aus den neuen Zuständen betrachtet werden, wofür das Konzept der Wertigkeit eingeführt wird.²⁹ Über eine Zustands- oder Aktionswertigkeitsfunktion, oftmals als Q-Funktion referenziert, wird eine Vorhersage über die zu erwartende kumulierte abgezinste zukünftige Belohnung berechnet.³⁰ Durch den Abzinsungsfaktor $\gamma \in [0, 1)$ wird der Einfluss zukünftiger Belohnungen nach ihrer zeitlichen Reihenfolge priorisiert.³¹ Mit der Wertigkeitsfunktion kann evaluiert werden, welche Strategie langfristig am erfolgreichsten ist, da bspw. manche Aktionen trotz geringer sofortiger Belohnung einen hohen Wert aufweisen können, wenn aus dem zukünftigen Zustand eine hohe Belohnung zu erwarten ist.³² Die Wertigkeitsfunktion und die daraus berechneten Wertigkeiten von Aktionen oder Zuständen werden über alle Zeitschritte neu geschätzt und stellen mit die wichtigste Komponenten in Algorithmen des verstärkenden Lernens dar.³³ Methoden basierend auf diesem Wertigkeitswert lernen eine Schätzfunktion der Wertigkeit für alle Zustände ($V_\pi(s) \forall S$) und alle Zustandsaktions-Paare ($Q_\pi(s_t, a_t) \forall s, a \in (S, A)$) der optimalen Strategie π^* durch aktualisieren der folgenden Funktionen eins und zwei.³⁴

$$(1) \quad Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left[r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t) \right]$$

$$(2) \quad V(s_t) = \max_a Q(s_t, a | \omega)$$

Aus den geschätzten Wertigkeit jedes Zustandsaktions Paares kann die optimale Strategie $\pi^*(s)$ durch $\arg \max_a Q(s, a)$ bestimmt werden.³⁵

2.2.2 Strategiebasierende Methoden

Methoden, welche die Strategie durch deren direkte Parametrisierung anstelle einer Bewertung aller Handlungsalternativen mittels Wertigkeitsfunktion optimieren, werden als strategiebasierend bezeichnet.³⁶ Diese Methodik kann beim Trainieren deterministischer Strategien zu unerwarteten Aktionen führen, weshalb häufig das Optimieren einer Wahrscheinlichkeitsverteilung für alle Aktionen bevorzugt wird.³⁷ Als Subklasse der RL Methoden wird der statistische Gradientenabstieg verwendet um die parametrisierte Strategie π_θ hinsichtlich der maximalen langfristigen

²⁸Vgl. Reda/Tao/van de Panne 2020, S. 2

²⁹Vgl. Wang/Hong 2020, S. 3

³⁰Vgl. Li 2019, S. 5

³¹Vgl. Li 2019, S. 5

³²Vgl. Sutton/Barto 2018, S. 6

³³Vgl. Sutton/Barto 2018, S. 6f.

³⁴Vgl. Zhang/Wu/Pineau 2018, S. 2

³⁵Vgl. Zhang/Wu/Pineau 2018, S. 2

³⁶Vgl. Zhang/Wu/Pineau 2018, S. 2

³⁷Vgl. Ningombam 2022, S. 3

kumulierten Belohnung zu optimieren.³⁸ Die Strategie π_θ oder auch $\pi(a|s, \theta)$ beschreibt dabei die Wahrscheinlichkeit Aktion a im Zustand s auszuwählen unter dem Parametervektor θ .³⁹ Zur Optimierung der Strategie wird die Funktion der kumulierten Belohnungen J nach dem Parameter der Gewichte θ wie folgt in Formel drei abgeleitet und der optimierte Parametervektor anhand Formel vier aktualisiert.⁴⁰

$$(3) \nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left[\left(\sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_t | s_t) \right) \left(\sum_{t=1}^T r(s_t, a_t) \right) \right]$$

$$(4) \theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$$

Zusammengefasst kann Formel vier und fünf dabei so interpretiert werden, dass die logarithmierte Wahrscheinlichkeit Aktion a_t im Zustand s_t auszuwählen erhöht werden soll, wenn a_t in einer höheren kumulierten Belohnung resultiert.⁴¹

2.2.3 Akteur-Kritiker Methoden

Unter Akteur-Kritiker Methoden werden hybride wertebasierende und strategiebasierende Methoden verstanden, welche zugleich die Strategie optimieren und eine Wertefunktion approximieren.⁴² Die strategiebasierende Methodik mit der lernenden Strategie agiert dabei als Akteur, wohingegen die Wertefunktion, welche jeder Aktion und jedem Zustand einen Belohnungswert zuweist, als Kritiker handelt.⁴³ Der Akteur wählt somit aus seiner Wahrscheinlichkeitsverteilung die auszuführende Aktionen aus, während der Kritiker diese anhand seiner Wertigkeit bewertet.⁴⁴ Betrachtet man den Trainingsprozess von Akteur-Kritiker basierten Methoden ist dieser wie folgt aufgebaut:⁴⁵

- Aktueller Zustand der Umgebung als Eingabe dem Akteur und Kritiker übergeben
- Akteur liefert eine auszuführende Aktion basierend auf dem Umgebungszustand
- Der Kritiker bekommt die Aktion als Eingabe und berechnet dessen Wertigkeit mittels Q-Funktion
- Durch die Wertigkeit seiner Aktion kann der Akteur seine Strategie anpassen
- Mit der neuen Strategie führt der Akteur die nächste Aktion im folgenden Zustand aus
- Die Q-Funktion des Kritikers wird mit den neuen Informationen aus der erhaltenen Belohnung angepasst

³⁸Vgl. Ningombam 2022, S. 3

³⁹Vgl. Sutton/Barto 2018, S. 321

⁴⁰Vgl. Wang/Hong 2020, S. 6

⁴¹Vgl. Wang/Hong 2020, S. 6

⁴²Vgl. Zhang/Wu/Pineau 2018, S. 2f.

⁴³Vgl. Sutton/Barto 2018, S. 321

⁴⁴Vgl. Ningombam 2022, S. 3

⁴⁵Vgl. Ningombam 2022, S. 4

2.2.4 Abgrenzung zu Multi-Agent Reinforcement Learning (MARL) Algorithmen

Innerhalb dieses Unterkapitels soll der beschriebene Aufbau von RL Algorithmen und deren Optimierungsproblem zu den von MARL Systemen abgegrenzt werden. Bei MARL Systemen wird anstatt einem Agenten eine Menge von Agenten eingesetzt welche alle mit ihrer Umgebung interagieren um den Weg der Zielerreichung zu lernen.⁴⁶ Dieser Ansatz dient dazu Problemstellungen welche nicht vollständig durch einen Agenten lösbar sind zu bearbeiten.⁴⁷ Einsatzgebiete von MARL sind dabei unter anderem das Routing von Netzwerkpaketen, Wirtschaftsmodellierung oder zusammenhängende Robotersysteme.⁴⁸ Je nach Ziel und der demnach definierter Belohnungsfunktion können die Agenten auf die drei unterschiedlichen Arten vollständig kooperativ, vollständig kompetitiv und der Mischung aus beiden miteinander interagieren.⁴⁹ Aus der unterschiedlichen Interaktion jedes Agenten mit der selben Umgebung ergibt sich der Unterschied, dass die Umgebungsdynamik aus der Kombination aller Aktionen der Agenten beeinflusst wird anstatt aus der Aktion des einzelnen Agenten.⁵⁰ Da dieser Effekt auch die Annahme der Stationarität von Markov Entscheidungsprozessen verletzt, bedarf die Umgebung auch einer anderen Representation.⁵¹ Ein Konzept was dafür häufig verwendet ist das Markov Spiel, welches sich anders als der Entscheidungsprozess durch einen mehrdimensionalen Aktions- und Belohnungsraum aus der Kombination aller N Agenten auszeichnet.⁵² Betrachtet man die Limitierungen von MARL erkennt man aus den beschriebenen Punkten die Herausforderungen der nicht vorhandenen Stationarität und der Skalierbarkeit, welcher sich die Herausforderung der teilweisen Beobachtbarkeit der Umgebung anschließt.⁵³

2.2.5 Limitierungen und Herausforderungen von RL

Trotz signifikanter Errungenschaften birgt der Einsatz von den besprochenen RL Algorithmen weiterhin Limitierungen und Risiken für ungewolltes Verhalten.⁵⁴

Eine der Herausforderungen zeigt sich bei der Representation der Agentenumwelt, da RL stark auf diesem Konzept basiert.⁵⁵ Daraus ergibt sich die Aufgabe, die Umwelt und dessen Verhalten sowie die Wahrnehmung durch den Agenten realitätsgetreu und präzise zu gestalten.⁵⁶ Neben der Definition und Wahrnehmung des Umweltverhaltens ist die Spezifikation des Ziels des Agenten ein ebenso kritischer Teil, da unerwartete Intentionen aus der Zielstellung abgeleitet werden könnten.⁵⁷ Zusätzlich teilen RL Algorithmen auch Herausforderungen aus anderen Gebieten des

⁴⁶Vgl. Wong u. a. 2022, S. 6

⁴⁷Vgl. Canese u. a. 2021, S. 1

⁴⁸Vgl. Canese u. a. 2021, S. 1

⁴⁹Vgl. Canese u. a. 2021, S. 8f.

⁵⁰Vgl. Wong u. a. 2022, S. 2

⁵¹Vgl. Wong u. a. 2022, S. 6

⁵²Vgl. Canese u. a. 2021, S. 4

⁵³Vgl. Canese u. a. 2021, S. 9ff.

⁵⁴Vgl. Li 2019, S. 7

⁵⁵Vgl. Sutton/Barto 2018, S. 8

⁵⁶Vgl. Sutton/Barto 2018, S. 7

⁵⁷Vgl. Li 2019, S. 7

maschinellen Lernens wie Genauigkeit, Interpretierbarkeit und die im Rahmen dieser Arbeit untersuchte Robustheit von Modellen.⁵⁸

Eine weitere Limitierung stellt der große Suchraum an Aktionen und das unbekannte Verhalten der Umgebung dar. Dies sorgt dafür, dass häufig die Effizienz einzelner Daten sehr gering ist und die Abwägung zwischen Exploration neuer Strategie und der Optimierung bekannter Verhaltensmuster ein wichtiger Bestandteil ist.⁵⁹ Aufgrund der geringen Effizienz der Daten aber des dennoch hohen Bedarfs an bewerteter Agentenerfahrung wird häufig auf simulierte Daten zurückgegriffen.⁶⁰ Simulierte Daten werden dabei häufig von möglichst hoch qualitativen Simulationsumgebungen bereitgestellt, da zu dem hohen Bedarf der Methodik häufig Limitierungen in der Sammlung von Daten in der echten Welt bestehen.⁶¹

Aufgrund dieser Bedeutung der Simulationsumgebung für RL Algorithmen und deren Transfer in die echte Welt werden im nachfolgenden Kapitel die Merkmale und Entwicklungen von Simulationen genauer betrachtet.

2.3 Simulationsumgebungen für RL

Anders als im klassischen Bereich des maschinellen Lernens wie überwachtes- und unüberwachtes Lernen, werden beim verstärkenden Lernen viele der Testdatensätze nicht aus der echten Welt akquiriert.⁶² Um entsprechend realistische Daten für das Training bereitzustellen, werden Simulationsumgebungen in Abhängigkeit von ihrer RL Anwendung ausgewählt.⁶³ Dennoch bleibt nahezu immer eine gewisse Diskrepanz zwischen der Dynamik in der Simulation und der Dynamik in der echten Welt.⁶⁴ Daher ist es kaum garantiert, dass erlernte Strategien der Agenten sich auch auf nur leicht veränderte Umgebungen übertragen lassen.⁶⁵

2.3.1 Definitionen von Simulationsumgebungen

Ausgehend von der Literaturrecherche zeigte sich, dass in der Forschungsliteratur die allgemeine Definition von Simulationen kaum aufgegriffen wird. Eine mögliche Definition nach Maria 1997 wird wie folgt dargelegt:

Eine Simulation eines existierenden Systems stellt die Anwendung eines Modells dar, welches konfigurierbar zu experimentellen Zwecken das eigentliche System vertritt, um wirtschaftliche oder systematische Herausforderungen des existierenden Systems zu umgehen. Das Model wird

⁵⁸Vgl. Li 2019, S. 7

⁵⁹Vgl. Li 2019, S. 7

⁶⁰Vgl. Zhao/Queralta/Westerlund 2020, S. 7

⁶¹Vgl. Li 2019, S. 8

⁶²Vgl. Zhang/Wu/Pineau 2018, S. 1

⁶³Vgl. Körber u. a. 2021, S. 7

⁶⁴Vgl. Bharadhwaj u. a. 2019, S. 1

⁶⁵Vgl. Bharadhwaj u. a. 2019, S. 1

in diesem Kontext definiert als Repräsentation des Aufbaus und der Verhaltensweise des existierenden Systems.

Innerhalb bestimmter Anwendungsgebiete wie der Medizin und der Pflege werden zusätzlich virtuelle Simulationen wie nachstehend definiert.

Unter virtuellen Simulationen versteht man eine digitale Lernumgebung, welche durch teilweiser Immersion eine wahrnehmbare Erfahrung bereitstellt.⁶⁶

2.3.2 Entwicklung von Simulationsumgebungen für RL Anwendungen

Im weiteren Teil dieses Kapitels wird aufbauend auf den Definitionen, die Entwicklung von Simulationen betrachtet. Allgemein lässt sich dieser Entwicklungsprozess in die folgenden Teilschritte gliedern:⁶⁷

1. Identifikation der Herausforderungen im existierenden System und Ableitung von Anforderungen für die Simulation.
2. Zielgruppe, Funktionsrahmen und quantitative Bewertungskriterien der Simulation definieren.
3. Analyse des zu simulierenden Systemverhaltens durch Sammeln und Verarbeiten von realen Daten des existierenden Systems.
4. Entwicklung einer schematischen Darstellung des Modells und dessen Überführung in nutzbare Software.
5. Validierung des Modells durch bspw. den Vergleich mit dem existierenden System.
6. Dokumentierung des Modells, dessen Variablen, Metriken und getroffene Annahmen.

Die Entwicklung von Simulationen wurde in der Forschungsliteratur besonders durch den Fortschritt im Bereich des verstärkenden Lernens vorangetrieben, da der Vergleich von RL-Algorithmen zuverlässige Benchmarks in Form von Simulationsumgebungen benötigt.⁶⁸ Aus dieser Motivation wurde 2016 durch die OpenAI der OpenAI Gym Werkzeugkasten entwickelt, welcher eine Sammlung an Benchmarksimulationen mit einheitlicher Schnittstelle für RL Algorithmen enthält.⁶⁹ Seither wurde diese definierte Schnittstelle vielfach verwendet um RL Umgebungen mit dem Ziel zu entwickeln, diese zu publizieren und dessen Wiederverwendung zu ermöglichen.⁷⁰ Die Schnittstelle ist definiert als Python Klasse *gym.Env*, von welcher weitere Klassen erben und die vorgeschriebenen Funktionen zum Zeitschritt und zum Zurücksetzen der Simulation implementieren.⁷¹ Der Werkzeugkasten von OpenAI fokussiert sich auf einen Episoden ähnlichen Rahmen,

⁶⁶Vgl. Foronda 2021, S. 1

⁶⁷Vgl. Maria 1997, S. 8f.

⁶⁸Vgl. Brockman u. a. 2016, S. 1

⁶⁹Vgl. Brockman u. a. 2016, S. 1

⁷⁰Vgl. Schuderer/Bromuri/van Eekelen 2021, S. 4

⁷¹Vgl. Schuderer/Bromuri/van Eekelen 2021, S. 4

in welchem der Agent durch zunächst zufälliges Auswählen von Interaktionen lernt.⁷² Weitere Entwicklungsentscheidungen des OpenAI Gym Werkzeugkastens umfassen z.B. die bewusst fehlende Schnittstelle des Agenten, die strikte Versionierung der Umgebung oder die standardmäßige Simulationsüberwachung.⁷³

Werden Lernumgebungen nach der Gym Schnittstelle oder nach eigener Definition für RL Anwendungen eingesetzt, kann sich deren Gestaltung unterschiedlich auf die Leistung der Anwendung auswirken.⁷⁴ Eine enge initiale Wahrscheinlichkeitsverteilung des Umgebungszustandes kann die Lerneffizienz erhöhen, wohingegen eine weite Wahrscheinlichkeitsverteilung positiv die Robustheit der erlernten Strategie beeinflusst.⁷⁵ Die Robustheit kann zusätzlich durch die Einbindung von Fehlverhalten in der Wahrnehmung der Umgebung beeinflusst werden, da auch in realen Szenarien ein Risiko für Fehlverhalten besteht.⁷⁶ Weiterer Gestaltungsspielraum ergibt sich aus der Wahl des Handlungsbereichs, welcher häufig Drehmoment basierend gewählt wird, obwohl durch PID-Regler ein effektiverer und effizienterer Lernprozess stattfinden kann.⁷⁷

Neben den beschriebenen Eigenschaften von Umgebungen für verstärkendes Lernen unterliegen auch die verwendeten Simulationen bestimmten Merkmalen welche in der Entwicklung zu berücksichtigen sind. Laut einer Umfrage nach Ivaldi/Padois/Nori 2014 sind diese wichtigsten Eigenschaften die Stabilität, Geschwindigkeit, Präzision, Genauigkeit, Bedienbarkeit und der Ressourcenverbrauch. Die Entwicklung des Modells, welches das existierende System ersetzt, sollte sich demnach möglichst positiv auf die beschriebenen Eigenschaften auswirken. Neben den beschriebenen Leistungsbezogenen Merkmalen sind die folgenden weiteren Kriterien mitunter die wichtigsten zur Auswahl einer Simulation:

Rank	Most important criteria
1	Simulation very close to reality
2	Open-source
3	Same code for both real and simulated robot
4	Light and fast
5	Customization
6	No interpenetration between bodies

Tab. 2: wichtigsten Kriterien zur Auswahl von Simulatoren⁷⁸

Aus Tabelle zwei lässt sich entnehmen, das besonders die Nähe zur Realität ein wichtiges Auswahlkriterium ist. Im Kontext von verstärkendem Lernen im Robotik Bereich ist ein wichtiger Baustein die Physik-Engine zur Modellierung von Dynamiken.⁷⁹

⁷²Vgl. Brockman u. a. 2016, S. 1

⁷³Vgl. Brockman u. a. 2016, S. 2f.

⁷⁴Vgl. Reda/Tao/van de Panne 2020, S. 1

⁷⁵Vgl. Reda/Tao/van de Panne 2020, S. 3

⁷⁶Vgl. Yan Duan u. a. 2016, S. 2

⁷⁷Vgl. Reda/Tao/van de Panne 2020, S. 7

⁷⁸Enthalten in: Ivaldi/Padois/Nori 2014, s. 4

⁷⁹Vgl. Ayala u. a. 2020, S. 2

2.3.3 aktuelle Physik-Engines und Simulationsanwendungen

Innerhalb dieses Abschnittes wird aufgrund der Bedeutung der Physik-Engine für den Grad der Simulationsrealität, eine Auswahl der aktuellen Physik-Engines und deren Simulationsanwendung betrachtet.

Gazebo

Gazebo ist eine durch die Open Source Robotics Foundation entwickelte Simulationsanwendung, welche mehrere Physik-Engines unterstützt.⁸⁰ Mittels Gazebo lassen sich Interaktionen zwischen Robotern in Innen- und Außenbereichen unter realistischer Sensorik simulieren.⁸¹ Die unterstützten Physik-Engines umfassen Bullet, Dynamic Animation and Robotics Toolkit (DART), Open Dynamics Engine (ODE) und Simbody.⁸²

MuJoCo

MuJoCo stellt eine Physik-Engine für modellbasierte Steuerung dar, dessen Objekte durch C++ oder XML definiert und Gelenkzustände im Koordinatensystem beschrieben werden.⁸³ Diese beschriebene Eigenschaft lässt sich auch aus dem Namen als Abkürzung für **M**ulti-**J**oint dynamics with **C**ontact ableiten.⁸⁴ Die MuJoCo Anwendung ist lizenziert, was diesen Besitz einer Lizenz für die Installation oder die Virtualisierung innerhalb eines Containers voraussetzt.⁸⁵

PyBullet

PyBullet basiert als Simulationssoftware auf der Bullet-Engine und fokussiert funktional auf die Anwendung von RL im Robotikbereich.⁸⁶ Die Bullet-Engine ist hingegen eine offene Softwarebibliothek welche neben verstärkenden Lernen auch bei Computeranimationen angewendet wird.⁸⁷ Die Handhabung von PyBullet profitiert von der ausgiebigen Dokumentation, der großen Entwicklergemeinschaft und der Unterstützung von verschiedenen Dateiformaten wie SDA, URDF und MJCF zur Einbindung von Objekten.⁸⁸

⁸⁰Vgl. Ivaldi/Padois/Nori 2014, S. 7

⁸¹Vgl. Ayala u. a. 2020, S. 4

⁸²Vgl. Körber u. a. 2021, S. 3

⁸³Vgl. Todorov/Erez/Tassa 2012, S. 1

⁸⁴Vgl. Todorov/Erez/Tassa 2012, S. 2

⁸⁵Vgl. Körber u. a. 2021, S. 3

⁸⁶Vgl. Körber u. a. 2021, S. 3

⁸⁷Vgl. Ivaldi/Padois/Nori 2014, S. 7

⁸⁸Vgl. Körber u. a. 2021, S. 6

3 Durchführung des Laborexperiments

4 Ergebnisse des Laborexperiments

5 Reflexion und Forschungsausblick

Anhang

Anhangverzeichnis

Anhang 1	Interview Transkripte	18
Anhang 1/1	Interview Transkript: Mitarbeiter eines Unternehmens	18

Anhang 1: Interview Transkripte

Anhang 1/1: Interview Transkript: Mitarbeiter eines Unternehmens

Literaturverzeichnis

- ACM Digital Library (2/28/2023): ACM Digital Library. URL: <https://dl.acm.org/>.
- Ayala, A./Cruz, F./Campos, D./Rubio, R./Fernandes, B./Dazeley, R. (2020): A Comparison of Humanoid Robot Simulators: A Quantitative Approach. In: *2020 Joint IEEE 10th International Conference on Development and Learning and Epigenetic Robotics (ICDL-EpiRob)*, S. 1–6. DOI: 10.1109/ICDL-EpiRob48136.2020.9278116.
- Bharadhwaj, H./Wang, Z./Bengio, Y./Paull, L. (2019): A Data-Efficient Framework for Training and Sim-to-Real Transfer of Navigation Policies. In: *2019 International Conference on Robotics and Automation (ICRA)*. IEEE. DOI: 10.1109/icra.2019.8794310.
- Brockman, G./Cheung, V./Pettersson, L./Schneider, J./Schulman, J./Tang, J./Zaremba, W. (2016): OpenAI Gym. In: *CoRR* abs/1606.01540. arXiv: 1606.01540. URL: <http://arxiv.org/abs/1606.01540>.
- Canese, L./Cardarilli, G. C./Di Nunzio, L./Fazzolari, R./Giardino, D./Re, M./Spanò, S. (2021): Multi-Agent Reinforcement Learning: A Review of Challenges and Applications. In: *Applied Sciences* 11.11, S. 4948. DOI: 10.3390/app11114948. URL: <https://www.mdpi.com/2076-3417/11/11/4948>.
- Collins, J./Ketter, W. (2022): Power TAC: Software architecture for a competitive simulation of sustainable smart energy markets. In: *SoftwareX* 20, S. 101217. ISSN: 2352-7110. DOI: <https://doi.org/10.1016/j.softx.2022.101217>. URL: <https://www.sciencedirect.com/science/article/pii/S2352711022001352>.
- Cutler, M./Walsh, T. J./How, J. P. (2014): Reinforcement learning with multi-fidelity simulators. In: *2014 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE. DOI: 10.1109/icra.2014.6907423.
- Foronda, C. L. (2021): What Is Virtual Simulation? In: *Clinical Simulation in Nursing* 52, S. 8. ISSN: 18761399. DOI: 10.1016/j.ecns.2020.12.004.
- Google Scholar (2/28/2023). URL: <https://scholar.google.de/>.
- Holzweißig, K. (2022): Wissenschaftliches Arbeiten. In: 6.6.
- IEEE Xplore (2/28/2023). URL: <https://ieeexplore.ieee.org/Xplore/home.jsp>.
- Ivaldi, S./Padois, V./Nori, F. (2014): Tools for dynamics simulation of robots: a survey based on user feedback. URL: <https://arxiv.org/pdf/1402.7050>.
- Körber, M./Lange, J./Rediske, S./Steinmann, S./Glück, R. (2021): Comparing Popular Simulation Environments in the Scope of Robotics and Reinforcement Learning. URL: <https://arxiv.org/pdf/2103.04616>.
- Li, Y. (2019): Reinforcement Learning Applications. DOI: 10.48550/ARXIV.1908.06973. URL: <https://arxiv.org/abs/1908.06973>.
- Maria, A. (1997): Introduction to modeling and simulation. In: *Proceedings of the 29th conference on Winter simulation - WSC '97*. New York, New York, USA: ACM Press. DOI: 10.1145/268437.268440.
- Ningombam, D. D. (2022): Deep Reinforcement Learning Algorithms for Machine-to-Machine Communications: A Review. In: *2022 13th International Conference on Computing Commu-*

- nication and Networking Technologies (ICCCNT)*. IEEE. DOI: 10.1109/icccnt54827.2022.9984457.
- Recker, J. (2021)**: Scientific research in information systems: A beginner's guide. Second Edition. Progress in IS. Cham: Springer International Publishing. ISBN: 9783030854362. URL: <https://ebookcentral.proquest.com/lib/kxp/detail.action?docID=6789173>.
- Reda, D./Tao, T./van de Panne, M. (2020)**: Learning to Locomote: Understanding How Environment Design Matters for Deep Reinforcement Learning. In: *Motion, Interaction and Games*. Hrsg. von Daniele Reda/Tianxin Tao/Michiel van de Panne. New York, NY, USA: ACM, S. 1–10. DOI: 10.1145/3424636.3426907.
- Schuderer, A./Bromuri, S./van Eekelen, M. (2021)**: Sim-Env: Decoupling OpenAI Gym Environments from Simulation Models. In: *International Conference on Practical Applications of Agents and Multi-Agent Systems*. Springer, Cham, S. 390–393. DOI: 10.1007/978-3-030-85739-4{\textunderscore}39. URL: https://link.springer.com/chapter/10.1007/978-3-030-85739-4_39.
- Slaoui, R. B./Clements, W. R./Foerster, J. N./Toth, S. (2019)**: Robust Domain Randomization for Reinforcement Learning. In: *CoRR* abs/1910.10537. arXiv: 1910.10537. URL: <http://arxiv.org/abs/1910.10537>.
- Sutton, R. S./Barto, A. G. (2018)**: Reinforcement Learning, second edition: An Introduction. MIT Press. ISBN: 9780262352703.
- Todorov, E./Erez, T./Tassa, Y. (2012)**: MuJoCo: A physics engine for model-based control. In: *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, S. 5026–5033. DOI: 10.1109/IRoS.2012.6386109.
- Wang, Z./Hong, T. (2020)**: Reinforcement learning for building controls: The opportunities and challenges. In: *Applied Energy* 269, S. 115036. ISSN: 0306-2619. DOI: 10.1016/j.apenergy.2020.115036.
- Webster, J./Watson, R. T. (2002)**: Analyzing the Past to Prepare for the Future: Writing a Literature Review. In: *MIS Q.* 26.2, S. xiii–xxiii. ISSN: 0276-7783.
- Wong, A./Bäck, T./Kononova, A. V./Plaat, A. (2022)**: Deep multiagent reinforcement learning: challenges and directions. In: *Artificial Intelligence Review*. ISSN: 0269-2821. DOI: 10.1007/s10462-022-10299-x.
- Yan Duan/Xi Chen/Rein Houthoofd/John Schulman/Pieter Abbeel (2016)**: Benchmarking Deep Reinforcement Learning for Continuous Control. In: *International Conference on Machine Learning*, S. 1329–1338. ISSN: 1938-7228. URL: <https://proceedings.mlr.press/v48/duan16.html>.
- Zhang, A./Wu, Y./Pineau, J. (2018)**: Natural Environment Benchmarks for Reinforcement Learning. DOI: 10.48550/ARXIV.1811.06032. URL: <https://arxiv.org/abs/1811.06032>.
- Zhao, W./Queralta, J. P./Westerlund, T. (2020)**: Sim-to-Real Transfer in Deep Reinforcement Learning for Robotics: a Survey. In: *2020 IEEE Symposium Series on Computational Intelligence (SSCI)*. IEEE. DOI: 10.1109/ssci47803.2020.9308468.

Erklärung

Ich versichere hiermit, dass ich meine Bachelorarbeit mit dem Thema: *Experiment zur Verbesserung der Robustheit von Reinforcement Learning Policies anhand trainiertem Gegenspieler in Dronensimulationen* selbstständig verfasst und keine anderen als die angegebenen Quellen und Hilfsmittel benutzt habe. Ich versichere zudem, dass die eingereichte elektronische Fassung mit der gedruckten Fassung übereinstimmt.

(Ort, Datum)

(Unterschrift)