

Computer spielen Computerspiele mit Hilfe von reinforcement learning am Beispiel Vier Gewinnt

Pablo Lubitz



BACHELORARBEIT

eingereicht am
Universitäts-Bachelorstudiengang

Informatik

in Bremen

im März 2020

Betreuung:

Holger Schultheis, Thomas Barkowsky

Erklärung

Ich erkläre eidesstattlich, dass ich die vorliegende Arbeit selbstständig und ohne fremde Hilfe verfasst, andere als die angegebenen Quellen nicht benutzt und die den benutzten Quellen entnommenen Stellen als solche gekennzeichnet habe. Die Arbeit wurde bisher in gleicher oder ähnlicher Form keiner anderen Prüfungsbehörde vorgelegt.

Bremen, am 9. März 2020

Pablo Lubitz

Inhaltsverzeichnis

Erklärung	iii
Kurzfassung	vi
Abstract	vii
1 Einleitung	1
1.1 Motivation	1
1.2 Problemstellung	1
1.3 Umsetzung	1
2 Grundlagen	2
2.1 Vier Gewinnt	2
2.2 Minimax Algorithmus	2
2.2.1 Negamax Algorithmus	3
2.2.2 Alpha-Beta Suche	3
2.2.3 Wahrscheinlichkeitsverteilung	3
2.3 Reinforcement Learning	3
2.3.1 Environment	4
2.3.2 Agent	5
2.3.3 Policy	5
2.3.4 State	5
2.3.5 Actions	5
2.3.6 Reward	5
2.3.7 Q-Funktion	6
2.3.8 Q-Values	6
2.3.9 Deep Q-Network	6
2.3.10 Markow-Entscheidungsproblem	6
2.3.11 credit assignment problem	6
2.3.12 Episode	6
2.3.13 Epoche	7
2.4 Software	7
2.4.1 Python	7
2.4.2 Tensorflow	7
2.4.3 Keras	7
2.4.4 OpenAI Gym	7

3	Stand der Technik	8
3.1	Arbeiten mit ähnlichem Inhalt	8
4	Konzept	9
4.1	Zielsetzung	9
4.2	Methode	9
4.3	Evaluation	9
5	Implementierung	10
5.0.1	Vier Gewinnt	10
5.1	Gegner mit Minimax Algorithmus	10
5.1.1	check_next_actions	11
5.1.2	find_best_move	11
5.1.3	Wahrscheinlichkeitsverteilung	11
5.2	Gegner mit schlechten Strategien	11
5.2.1	Gegner mit deterministischen Aktionen	11
5.2.2	Gegner mit zufälligen Aktionen	11
5.3	Reinforcement Learning	11
5.3.1	Environment	11
5.3.2	Agent	12
5.3.3	State	12
5.3.4	Actions	12
5.3.5	Reward	12
5.3.6	Policy	12
5.4	Lernen vom Gegner	14
5.5	Graphische Darstellung	14
5.6	Visualisierung des Spiels	14
6	Evaluierung	16
6.1	Reward-Funktion	16
6.2	Lernen vom Gegner	17
6.3	Minimax	19
6.4	Besonderheiten	21
7	Fazit	22
	Quellenverzeichnis	23
	Literatur	23
	Software	24
	Online-Quellen	24

Kurzfassung

In dieser Bachelorarbeit wurde sich mit dem Thema reinforcement learning in Spielen beschäftigt. Es wurde eine Implementierung dieses Verfahrens für das Spiel Vier gewinnt genutzt. Hieran wurde getestet ob dieses Verfahren durch das beobachten des gegnerischen Verhaltens optimiert werden kann. Hierbei wurde eine minimale Verbesserung festgestellt.

Abstract

In this bachelor thesis we are exploring reinforcement learning in the area of games. Therefore an implementation of this procedure was used for the game Connect Four. Furthermore a test was made to proof that the implementation could be optimized with the addition of a way to learn what your enemy did. This was shown to have a minimal positive impact.

Kapitel 1

Einleitung

1.1 Motivation

Spiele machen Spaß und können genutzt werden um neue Fähigkeiten zu lernen oder verbessern. Das Lernen beim Spielen funktioniert dann am besten, wenn das Spiel einen gewissen Schwierigkeitsgrad hat, welcher aber noch zu bewältigen sein muss. Ist das Spiel zu einfach, wird der Spieler gelangweilt, ist es zu schwer, wird er frustriert. In beiden Fällen sucht er sich wahrscheinlich sehr bald eine andere Beschäftigung. Hiergegen kann es helfen einen anderen Spieler zu haben, welcher ein ähnliches Level an Erfahrung mitbringt. Wenn nun die Rolle des Gegners nicht von einem Menschen besetzt werden kann ist es naheliegend einen Computergegner zu erschaffen. Solch ein Computergegner wird normalerweise mit Algorithmen umgesetzt, die einen optimalen Zug errechnen. Seit einigen Jahren ist es, durch den Vortschritt in Hardware und Wissenschaft, möglich Computergegner zu erschaffen, die ihr Können durch Neuronale Netze antrainiert haben.

1.2 Problemstellung

Computergegner die mit Hilfe von Neuronalen Netzen trainiert wurden, haben das Problem, dass sie nur gut werden können wenn es auch einen guten Gegner gibt oder eine enorme Menge an Spielaufzeichnungen existiert, an denen sie trainieren können. Da dieses Training erst richtigen Effekt zeigt, wenn sehr viele Spiele gespielt wurden, ist es nicht wirklich möglich einen Menschen als Trainer zu nutzen. Also wird hier normalerweise auf einen weiteren Computergegner zurückgegriffen, welcher das Spiel beherrscht ohne dies mit Neuronalen Netzen erlernt zu haben. Es wäre schön wenn die Anzahl der benötigten Spiele beim Training reduziert werden könnte, um schneller bessere Ergebnisse zu erzielen oder sogar nur mit einem Mensch lernen zu können.

1.3 Umsetzung

Meine Idee ist es, das Training durch das Mitlernen der Züge des Gegners zu verbessern. Dies erscheint sinnvoll, da der Gegner das Spiel schon beherrscht und daher tendenziell gute Züge machen wird.

Kapitel 2

Grundlagen

In dieser Arbeit werden einige Technologien genutzt um die Problemstellung zu lösen. Die wichtigsten dieser Technologien werden hier erklärt, um das weitere Verständniss des Lesers zu gewährleisten.

2.1 Vier Gewinnt

Vier Gewinnt ist ein Spiel für 2 Spieler in dem abwechselnd Spielsteine in ein vertikales Spielfeld der gröÙe 7×6 eingeworfen werden. Jeder Spieler kann sich in seinem Zug zwischen einem von sieben Einwurföchern entscheiden. Wird ein Spielstein eingeworfen so fällt dieser auf die niedrigste der sechs Positionen die noch nicht durch andere Spielsteine gefüllt ist. Hierdurch ergibt sich ein Spielfeld mit 42 Feldern. Sind schon sechs Spielsteine in ein bestimmtest Einwurfloch gesteckt wurden so darf hier kein weiterer Spielstein eingeworfen werden. Ein Spieler gewinnt das Spiel wenn er es geschafft hat, dass sich vier seiner Spielsteine in einer Reihe befinden. Eine Reihe kann horizontal, vertikal oder diagonal gebildet werden. Werden alle 42 Positionen mit Spielsteinen befüllt ohne eine Reihe von vier Steinen des selben Spielers zu bilden geht die Partie unentschieden aus.



Abbildung 2.1: Das Spiel Vier Gewinnt [18]

2.2 Minimax Algorithmus

Ein Minimax Algorithmus beschreibt einen rekursiven Ansatz zum Finden einer optimalen Lösung für Probleme von zwei Parteien mit widersetzlich Zielen. Dies sind in der

Regel Nullsummenspiele mit perfekter Information, es gibt also für jeden Gewinn des einen Spielers genauso viel Verlust für den anderen und es gibt kein Spielelement, dass nur für einen Spieler einsehbar ist. Der Minimax Algorithmus berechnet sich hierfür den Suchbaum des Spiels, welcher alle möglichen Züge beider Spieler in nachvollziehbarer Reihenfolge enthält. Hierbei werden in jedem Rekursionsschritt alle möglichen Züge des derzeitigen Spielers betrachtet und bewertet. Wenn der derzeitige Zug nicht zum Gewinn führt, findet eine Bewertung durch das Betrachten des nächsten Rekursionsschrittes statt. Der Name Minimax ergibt sich durch das Betrachten der abwechselnden besten Züge von den Spielern. Da dies aus der Perspektive von einem Spieler betrachtet wird, ergibt sich hieraus ein abwechselnd minimaler und maximaler Zug. Dieses Verfahren führt bei einfachen Spielen wie Tic-Tac-Toe zu einem perfekten Spieler, da hier eine überschaubare Menge an möglichen Zügen betrachtet wird (9 Felder die von 2 Spielern gefüllt werden). Für komplexere Probleme gibt es die Möglichkeit eine Suchtiefe zu bestimmen, wodurch die Rekursion nur bis zu dieser Tiefe durchgeführt wird. Dies ist sinnvoll da die Berechnung einzelner Züge sonst sehr lange dauern kann. Der Minimax Algorithmus kann in bestimmten Spielen noch verbessert werden. Diese Verbesserungen sind der Negamax Algorithmus und die Alpha-Beta Suche welche im Folgenden einmal beschrieben werden. Der Einfachheit halber wird trotz dieser Erweiterungen in dieser Arbeit im Bezug auf diesen Algorithmus immer von Minimax geredet.

2.2.1 Negamax Algorithmus

Der Negamax Algorithmus basiert auf der Annahme, dass der maximal beste Zug für den einen Spieler der minimal beste Zug für den anderen bedeutet. Ist dies so in einem Spiel in dem beide Spieler in jedem Zug die selben Dinge tun können, kann hierdurch die Rekursion vereinfacht werden, indem für beide Spieler die selbe Formel benutzt wird. Hierfür muss dann in jedem Rekursionsschritt der Input negiert werden. Wegen dieser Negierung heißt dieser Algorithmus auch Negamax.

2.2.2 Alpha-Beta Suche

Alpha-Beta Suche ist eine Variante des Minimax Algorithmus, die die Berechnung des Suchbaumes beschleunigt, indem sie bestimmte Teile der Suche nicht ausführt wenn es schon einen klar besseren Pfad gibt.

2.2.3 Wahrscheinlichkeitsverteilung

Die Wahrscheinlichkeitsverteilung basiert auf der Annahme, dass in jedem State die Chancen der Menge aller möglichen Aktionen addiert 100% ergibt. Somit kann dann z.B. eine zufällige Aktion genommen werden die aber sehr stark Züge mit besseren Chancen bevorzugt.

2.3 Reinforcement Learning

Reinforcement learning (Bestärkendes Lernen) ist eine Maschine-Learning Methode bei der ein Agent mit Hilfe eines Environments, Policy, State, Actions und eines Rewards lernt, eine ihm gegebene Aufgabe zu lösen. Hierfür bieten sich Probleme an die als

Markow-Entscheidungsproblem[7] formuliert werden können, da diese alles abdecken was für ein Environment wichtig ist. Es wird dabei Grundsätzlich wie in Abbildung 2.2 beschrieben vorgegangen. Der Agent bekommt also einen State vom Environment, und durch diesen sucht er sich mittels seiner Policy eine Action aus. Diese Action beeinflusst dann das Environment und dieses gibt dem Agenten einen Reward für seine Action. Durch den Reward kann der Agent dann die Gewichtungen für seine Policy anpassen. Dazu bekommt der Agent dann den neuen State und kann dann, mit der angepassten Policy, wieder eine neue Action auswählen. Diese Teile des reinforcement learnings werden hier nun einmal genauer beschrieben.

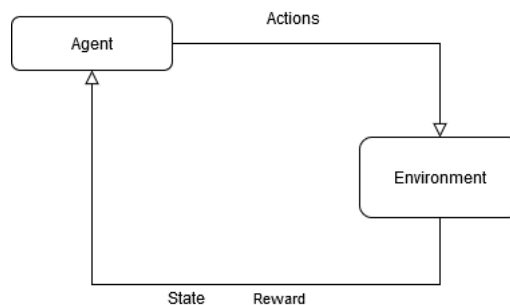


Abbildung 2.2: grundsätzliches Konzept von reinforcement learning

2.3.1 Environment

Das Environment beschreibt das gesamte Problem, dass der Agent versucht zu lösen. Es beinhaltet alles was für das Spielen wichtig ist außer dem Spieler, welcher von dem Agenten behandelt wird. Dies ist bei den meisten Spielen ein Spielfeld mit den Positionen aller Akteure, Spielsteine oder sonstige Elemente die zum Spielen genutzt werden. Dazu kümmert sich das Environment auch noch um die Spiellogik. Die Spiellogik beschreibt alle Abläufe die durch die Regeln des Spiels definiert wurden. In dem Beispiel Vier Gewinnt wäre das unter anderem die Funktion, dass die Spielsteine immer auf das unterste freie Feld fallen oder die Überprüfung ob ein Spieler gewonnen hat. Das letzte wichtige Element im Environment der meisten Spiele ist der Gegenspieler welcher im Falle von physikalischen Spielen normalerweise die selben Möglichkeiten besitzt wie der Spieler. Für das Verhalten des Gegenspielers werden im Normalfall Algorithmen genutzt die den bestmöglichen Zug errechnen, wie hier der Minimax Algorithmus. Für sehr komplexe Spiele bei denen sich kein optimaler Zug berechnen lässt wird oft auf menschliche Spieler zurückgegriffen. Da ein Computer aber viel schneller im spielen ist und reinforcement learning viele Durchläufe absolvieren muss um Wirkung zu zeigen wird versucht dies zu vermeiden oder bestehende Datensätze von Spielen mit Menschen genutzt.

2.3.2 Agent

Bei dem Agenten handelt es sich um den Akteur der die ihm gegebene Aufgabe meistern soll. Er sieht das Environment und wählt mittels der Policy die ihm am besten erscheinende Action aus um den Reward zu maximieren.

2.3.3 Policy

Mit der Policy wird das Verhalten beschrieben nach dem der Agent entscheidet welche der möglichen Aktionen die aktuell beste ist um den Reward zu maximieren.

2.3.4 State

Der State beschreibt den aktuellen Zustand des Environments welcher in jedem Zug an den Agenten weitergereicht wird, damit dieser seine Auswahl treffen kann.

2.3.5 Actions

Actions sind die möglichen Aktionen zwischen denen sich der Agent je nach State entscheiden muss. Diese sind alle möglichen Züge die ein Spieler zu einer bestimmten Zeitpunkt ausführen kann.

2.3.6 Reward

Der Reward ist die Belohnung die der Agent für seine Aktionen bekommt. Diese Belohnung beschreibt, wie gut es ist eine bestimmte Aktion auszuführen. Hieran passt der Agent dann während des Trainings seine Q-Values entsprechend seiner Policy an. Der Reward wird also benutzt um das Verhalten des Agenten zu beeinflussen. Möchte man ein bestimmtes Verhalten wird hierfür ein positiver Reward vergeben. Möchte man ein anderes Verhalten nicht wird ein negativer Reward vergeben. Je besser die Reward-Funktion gewählt wird desto eher wird der Agent gewolltes Verhalten zeigen. Bei der Wahl einer Reward-Funktion wird zwischen diskreten und kontinuierlichen Funktionen unterschieden, es ist aber auch eine Mischung dieser beiden möglich.

Bei einer diskreten Reward-Funktion werden nur feste im vorhinein definierte Rewards gegeben die zu bestimmten Events verteilt werden. Solch ein Event kann zum Beispiel das Gewinnen und Verlieren am Ende eines Spiels sein oder das Betreten des Agenten in einen gewollten oder ungewollten Bereich.

Bei einer kontinuierlichen Reward-Funktion wird der Reward zu jedem Schritt durch verschiedene Faktoren aus dem Environment berechnet. Solche Faktoren können zum Beispiel die verstrichene Zeit oder die Anzahl der Züge oder die aktuelle Position des Agenten sein. Eine gute kontinuierliche Reward-Funktion hilft dabei den Agenten schneller dem gewünschten Verhalten anzunähern indem der Reward zu den gewünschten Zielstates hin immer weiter angehoben wird.

Da diese beiden Arten von Reward-Funktionen das Verhalten des Agenten an unterschiedlichen Stellen beeinflussen, ist es oft sinnvoll eine Mischung dieser zu nutzen.

Wird sich für eine Mischung entschieden ist es wichtig, dass die aus beiden Teilen kommenden Rewards in Relation zueinander stehen, denn sonst kann dies dazu führen, dass die Auswirkung von dem Teil mit kleineren Rewards von dem Teil mit größeren Rewards überschattet wird.

2.3.7 Q-Funktion

Bei Komplexeren Environments ist es schwer für jeden Zug vorzusehen wie Zielführend er ist. Um diesem credit assignment problem entgegen zu wirken wird dann eine Q-Funktion benutzt die anhand der gelernten Zustände versucht die ungenauen Q-Values an ihre potenziellen Werte anzunähern.

2.3.8 Q-Values

Die Q-Values beschreiben für jedes Paar von State und Action einen erwarteten Reward. Leider sind diese Q-Values für die meisten Probleme nicht bekannt und müssen daher mit einer Q-Funktion antrainiert werden.

2.3.9 Deep Q-Network

Da es, bei einem Spiel mit einer Komplexität wie Vier gewinnt, nicht möglich ist für alle Paare von State und Action einen erwarteten Reward zu bestimmen, wird ein neuronales Netz genutzt welches versucht eine Policy für diese erwarteten Rewards zu bestimmen. Dieses Verfahren nennt sich Deep Q-Learning. Hierbei wird ein Deep Q-Network erstellt, welches während des Trainings eine Policy erstellt. Hierfür werden alle gespielten Spiele abgespeichert. Von diesen Spielen werden dann zufällige Züge betrachtet und somit die aktuelle Policy durch die hier bestehenden Q-Values angepasst.

2.3.10 Markow-Entscheidungsproblem

Das Markow-Entscheidungsproblem beschreibt eine Menge von State, Action, Probability, Reward Tupeln. Es gibt also für jeden State, Actions die eine Chance (Probability) haben in einen anderen State zu wechseln. Jede dieser Transitionen hat einen Reward der dieser Transition einen Wert zuweist. Das Entscheidungsproblem bezieht sich hier auf die Transition von einem Start-State zu einem Ziel-State, wobei versucht wird den Reward zu maximieren.

2.3.11 credit assignment problem

Das credit assignment problem beschreibt ein häufig auftretendes Problem in Environments welches sich dadurch auszeichnet, dass ein Ablauf aus mehreren Schritten erst nach seinem Abschluss einen Reward zugewiesen bekommt. Es ist dann nicht klar zu beurteilen, welcher der Schritte, wie hilfreich zur Erfüllung des Ablaufs war.

2.3.12 Episode

Eine Episode beschreibt einen einzelnen Durchlauf des Agenten im Environment.

2.3.13 Epoche

Eine Epoche ist eine Sammlung von Episoden und beschreibt einen gesamten Lern-durchlauf.

2.4 Software

Um die Problemstellung zu bearbeiten wurde eine Implementierung des Spiels in Python genutzt welches als ein OpenAI Gym Environment fungiert. An diesem Environment trainiert dann ein Agenten der mit Hilfe von Keras erschaffen wurde.

2.4.1 Python

Als Programmiersprache wurde sich für Python[14] entschieden da ein gutes Vorwissen in dieser Sprache vorhanden war und alle für diese Arbeit wichtigen Packages in dieser Programmiersprache existieren. Der Schöpfer von Python beschreibt es selbst als einfach zu erlernende objektorientierte Programmiersprache, welche Leistung mit einer klaren Syntax verbindet.[4]

2.4.2 Tensorflow

Bei Tensorflow [15] handelt es sich um ein open source Framework für Maschine-Learning welches ursprünglich für den internen Bedarf bei Google, für zum Beispiel Spracherkennung oder Google Maps, entwickelt wurde.

2.4.3 Keras

Keras[12] ist eine open source Bibliothek die auf Tensorflow aufbaut welche die Möglichkeit bietet vereinfacht neuronale Netze zu erstellen um diese in reinforcement learning Algorithmen zu benutzen.

2.4.4 OpenAI Gym

Gym [13] ist ein Werkzeug zum Entwickeln und Vergleichen von reinforcement learning Algorithmen welches von OpenAI, einem Forschungslabor aus San Francisco, zur Verfügung gestellt wird. Es bietet einen Standard zwischen dem Environment und dem Agenten. Hierdurch können neue Agenten und Environments nach gewissen Vorgaben erstellt werden. Somit ist es mit Gym einfacher möglich verschiedene Agenten an einem Environment zu Trainieren und diese zu vergleichen oder den selben Agenten an verschiedenen Environments auf seine Anpassungsfähigkeit zu testen.

Kapitel 3

Stand der Technik

In den letzten Jahrzehnten wurde für viele Spiele, die bis dahin als zu komplex galten, Programme geschaffen, die diese spielen können und dies sogar auf dem Level der besten menschlichen Spieler oder besser. Dies liegt zu einem großen Teil an den Vortschritten in der Hardware welche es ermöglicht, große künstliche neuronale Netzwerke zu erstellen und trainieren. Einige der bekanntesten Erfolge sind zum Beispiel deep Blue für Schach und Alpha Go für das Spiel Go.[9]

Weitere Erfolge aus den letzten Jahren sind Agenten die Computerspiele meistern, wie zum Beispiel MarI/O [17] welcher Super Mario spielt oder AlphaStar [16] welcher Starcraft II spielen kann.

Die Veröffentlichung von OpenAI Gym[13], welche unter anderem eine Menge an Atari-Spielen zur Verfügung stellt, eröffnete das Thema reinforcement learning für Spiele für die gesamte Welt.

Wie in dem Buch von Sutton und Barto [9] gezeigt wurde lassen sich Markow-Entscheidungsprobleme mit Hilfe von reinforcement learning lösen. Dass das Lernen, durch das imitieren, von dem Verhalten eines Profis verbessert werden kann wurde unter anderem von Oliver Amantier gezeigt. [2; 5; 8]

3.1 Arbeiten mit ähnlichem Inhalt

Ich habe zwei Arbeiten gefunden die auch mit Hilfe von reinforcement learning versuchen, das Spiel Vier Gewinnt zu meistern. Einmal die Arbeit von Lukas Stephan [10], welcher das Problem mit einer Java Implementierung angegangen ist und den Lernalgorithmus mit einer Mustererkennung ausgestattet hat. Zweitens die Arbeit von Markus Thill, Patrick Koch und Wolfgang Koenen [6], welche N-Tupel-Systeme nutzt, um einen Agenten zu Trainieren. Anders als diese Arbeiten ist mein Ansatz nicht das Erkennen von Mustern, sondern das Lernen durch das Verhalten des Gegners.

Kapitel 4

Konzept

Im Folgenden wird nun beschrieben, was das inhaltliche Ziel dieser Bachelorarbeit ist und wie vorgegangen werden soll um dieses Ziel zu erreichen.

4.1 Zielsetzung

Es soll am Beispiel Vier Gewinnt gezeigt werden wie ein Spieler funktioniert, welcher die Spielstrategie mit Hilfe von reinforcement learning erlernt hat und ob diese Technik durch das imitieren des Gegners verbessert werden kann.

4.2 Methode

Es wird eine digitale Version des Spiels Vier Gewinnt genutzt und eine Schnittstelle erschafft die es ermöglicht auszuwählen ob ein Mensch oder der Computer die Rolle der Spieler übernimmt. Dies ist wichtig da ein selbstlernender Algorithmus viele Spiele durchlaufen muss um einen Lernfortschritt aufzuzeigen. Somit kann in der Simulation dann trainiert werden ohne jedes einzelne Spiel gegen einen Menschen spielen zu müssen. Spielt der Computer soll zwischen verschiedenen effektiven Strategien ausgewählt werden können. Die Strategien sind: komplett zufällige Züge wählen, immer so weit wie möglich links ins Spielfeld werfen und ein Minimax Algorithmus welcher verschiedenen starke Varianten unterstützt. An diesen Strategien soll der Agent mit verschiedenen Voraussetzungen trainieren um dann vergleichen zu können was besser und schlechter beim lernen hilft.

4.3 Evaluation

Je nachdem wie gut der selbstlernende Agent nach dem Trainieren gegen die unterschiedlich effektiven Strategien ist, zeigt dann wie gut das reinforcement learning mit den verschiedenen Voraussetzungen funktioniert hat. Somit kann dann gezeigt werden welche Veränderungen besser sind um einen guten reinforcement learning Agenten zu erschaffen.

Kapitel 5

Implementierung

Es wurde sich für eine bestehende grundlegende Implementierung entschieden welche schon die Spiellogik von Vier Gewinnt sowie einen sehr einfachen reinforcement learning Ansatz enthält. In dieser Version lernt der Agent gegen einen Gegner zu gewinnen der zufällige Züge macht.[11]

Diese Implementierung wurde um einen Gegner erweitert der seine Züge nach dem Minimax Algorithmus auswählt.

Weitergehend wurde die reinforcement learning Implementierung überarbeitet um gegen den Minimax Algorithmus gewinnen zu können.

Diese Implementierungen werden hier einmal genauer beschrieben:

5.0.1 Vier Gewinnt

Die Implementierung des Spiels Vier Gewinnt deckt alle, für das Spiel wichtigen, Regeln ab. Es wird sich darum gekümmert welcher Spieler aktuell dran ist. Hat dieser sich dann für eine Action entschieden übernimmt das Spiel die korrekte Platzierung der Spielsteine, sowie die Überprüfung auf den Abschluss eines Spiels.

Es wurde sich dafür entschieden, dass der reinforcement learning Agent immer mit einer Partie beginnt, um somit die Menge an möglichen States zu halbieren, wodurch das gesamte Problem vereinfacht wird.

5.1 Gegner mit Minimax Algorithmus

Der Minimax Algorithmus wurde als Teil des Environments implementiert und besteht im Groben aus zwei Teilen. Einmal die Funktion *check_next_actions* die den Suchbaum für die nächsten Züge aufbaut. Und die Funktion *find_best_move* die den Suchbaum nach dem besten Zug durchsucht. Hierdurch kann sich dieser Gegner zu jedem State den besten nächsten Zug errechnen. Da dies im Optimalfall ein unschlagbarer Gegner ist wird hier noch mit einer Wahrscheinlichkeitsverteilung gearbeitet um den Gegner schlagbarer zu machen.

5.1.1 check_next_actions

Das Erstellen des Suchbaums funktioniert nach dem Negamax Algorithmus, es wird also eine Funktion für beide Spieler genutzt und immer abwechselnd das Maximum und das Minimum gesucht.

5.1.2 find_best_move

Das Suchen des besten Zuges wird durch die Alpha-Beta Suche beschleunigt indem bestimmte Äste des Suchbaumes, die sicher nicht zum gewünschten besten Zug führen, nicht weiter betrachtet werden. Dazu wurde noch eine Funktion zur Wahrscheinlichkeitsverteilung implementiert, die dafür sorgt, dass der Minimax Algorithmus in bestimmten Situationen mit einer gewissen Chance nicht optimal spielt.

5.1.3 Wahrscheinlichkeitsverteilung

Durch die Wahrscheinlichkeitsverteilung wird es ermöglicht einen zufälligen Zug aus den vom Minimax Algorithmus errechneten Zügen zu wählen. Hierbei macht die Gewinnchance eines Zuges aus, wie wahrscheinlich es ist, dass dieser ausgewählt wird. Diese Wahrscheinlichkeit kann mit einer *Variable* angepasst werden, wodurch die Wahrscheinlichkeit steigt dass einer der suboptimalen Züge, die vom Minimax Algorithmus berechnet wurden, ausgewählt wird.

5.2 Gegner mit schlechten Strategien

Um das Lernverhalten zu beobachten wurden zwei Gegner mit schlechten Strategien geschaffen.

5.2.1 Gegner mit deterministischen Aktionen

Es wurde ein Gegner geschaffen der zu jedem Zug immer so weit links wie möglich in das Spielfeld einwirft.

5.2.2 Gegner mit zufälligen Aktionen

Es wurde ein Gegner geschaffen der zu jeden Zug in ein zufälliges Loch des Spielfeldes einwirft, das noch frei ist.

5.3 Reinforcement Learning

Für das reinforcement learning wird ein Deep Q-Network genutzt, welches einen Agenten an dem Environment trainiert. Im Folgenden wird die Implementierung der hierfür notwendigen Komponenten beschrieben.

5.3.1 Environment

In Vier Gewinnt ist dies das Spielfeld und der Gegner. Das Spielfeld ist hier ein zweidimensionales Array der Länge $6 * 7$. Der Gegner ist während des Lernens ein Compu-

tergegner der nach unterschiedlichen Strategien wie zum Beispiel dem Minimax Algorithmus handelt.

5.3.2 Agent

Der Agent ist einer der beiden Spieler von Vier Gewinnt, welcher lernt seine Züge durch die ihm gegebene Policy auszuwählen, um einen möglichst hohen Reward zu erhalten.

5.3.3 State

Für Vier Gewinnt beschreibt der State welche der 42 Felder mit welchen Steinen befüllt sind. Ein Beispiel für so einen State kann man in dem Abschnitt Visualisierung des Spiels finden.

5.3.4 Actions

Die Actions in Vier Gewinnt beschreiben die sieben Löcher in die der Agent Steine werfen kann.

5.3.5 Reward

Zuerst wurde eine einfache diskrete Implementierung erstellt, welche den Reward immer zum Ende einer Partie ausgibt, wobei dieser beim Gewinn 1 beim Verlieren -1 und wenn die Partie unentschieden ausgeht 0 beträgt. Dies wurde durch die Anzahl der Züge um eine kontinuierliche Komponente erweitert. Hierfür wird der Reward, der durch die diskrete Implementierung ermittelt wurde durch die Anzahl der Züge geteilt. Da es nicht möglich ist, vor dem vierten Zug zu gewinnen, wird erst ab diesem Zug gezählt. Es ergibt sich also die Formel:

$$\frac{Reward}{AnzahlZüge - 3}$$

Wie man in Abbildung 5.1 gut sehen kann wird somit das schnelle Gewinnen mehr belohnt (obere Linie) und das schnelle Verlieren mehr bestraft (untere Linie). Da nur wenn das Spielfeld komplett voll ist unentschieden gespielt wird, wird hier ein Reward von 0 gegeben (einzelner Punkt). Diese Reward-Funktion wird im weiteren Text als zugabhängige Reward-Funktion benannt.

5.3.6 Policy

Da in Vier Gewinnt nur am Ende eines Spiels ein Reward ausgegeben werden kann, wird hier ein Deep Q-Network genutzt um das credit assignment problem zu lösen. Dieses Deep Q-Network wird nach einer Epsilon-Greedy Policy trainiert.

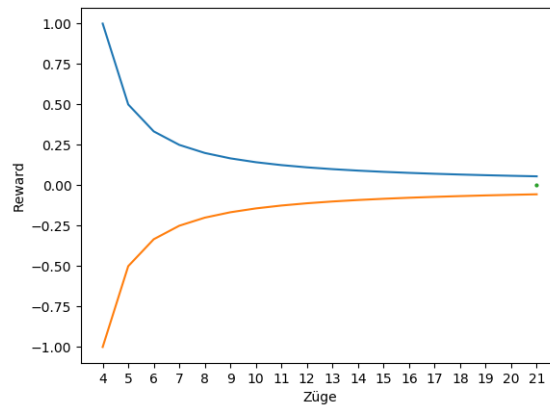


Abbildung 5.1: Reward Funktion

Deep Q-Network

Das Deep Q-Network speichert zu jedem gespielten Zug den derzeitigen State, die ausgeführte Aktion, den erhaltenen Reward, der daraus folgende State und ob der neue State die Episode beendet hat in sein Memory. Aus diesem Memory wird dann, mit der Funktion *experience_replay*, eine gewisse Menge an zufälligen Proben genommen. Es ist sehr sinnvoll hier eine zufällige Probe zu nehmen, da sonst durch die in Reihe betrachteten Züge ein ungewollter Bias entstehen kann, welcher die Konvergenz des Trainingsalgorithmus stark verlangsamen kann [3, Seite 469]. Für diese Proben werden dann die Q-Values vorhergesagt. Mit diesen Q-Values wird dann das Q-Network trainiert. Dieser Prozess wird so häufig wiederholt, dass hierdurch die Q-Values immer genauer werden und dadurch ergibt sich dann eine ziehlführende Policy.

Epsilon-Greedy

Als Unterstützung des Deep Q-Network wurde sich für die Epsilon-Greedy Policy entschieden, da diese sehr gut für Deep Q-Networks funktioniert [1]. Epsilon-Greedy beschreibt das Verhalten des Agenten welches entweder explorativ oder ausbeutend ist. Ist das Verhalten explorativ entscheidet sich der Agent für einen Zug, den er vorher noch nicht oder selten gemacht hat um diesen zu erlernen. Ist das Verhalten ausbeutend so wird der Zug genommen welcher derzeit die höchste Gewinnchance bietet. Zum Anfang des Lernens macht es natürlich Sinn, dass der Agent ganz viel erkundet, da er noch kein Wissen über das Spiel besitzt. Hat der Agent dann einige Spiele hinter sich kann er auf die ausbeutende Strategie umschalten. Um dieses Verhalten des erstigen Erkundens und späteren Ausbeutens zu ermöglichen nutzt die Epsilon-Greedy Policy die Formel: $np.random.rand() < 1 - \epsilon$. Gibt diese Formel *True* zurück so wird erkundet, gibt sie *False* aus wird ausgebeutet. Der Wert von Epsilon ist am Anfang etwas mehr als Null und wird während der Lernphase langsam erhöht bis er etwas weniger als 1 erreicht hat. Somit ist das Ergebniss der Formel am Anfang meistens *True* und später meistens *False*. Somit kann das DQN am Anfang viele Daten bekommen um hieraus

eine Policy zu erarbeiten und dann diese später nach und nach verbessern, wenn durch die Ausbeutung dann nach dieser Policy vorgegangen wird.

5.4 Lernen vom Gegner

Um das Lernen des Agenten zu verbessern wurde sich dazu entschieden, dass der Agent auch die Züge seines Gegners zum Lernen nutzen kann. Dies ist für Vier Gewinnt angemessen, da es sich um ein Spiel mit perfekter Information handelt. Es sind also jedem Spieler zum Zeitpunkt einer Entscheidung alle Informationen über das Spiel ersichtlich. Da ein menschlicher Anfänger in diesem Spiel, bei einer Niederlage an den Spielzügen seines Gegners lernen kann, wurde dies auch als sinnvoll für den Agenten angesehen. Hierfür mussten einige Anpassungen vorgenommen werden. Da der reinforcement learning Agent immer der erste der beiden Spieler ist, der seinen Zug macht, fehlen dem Deep Q-Network Informationen über den Zug nach seinem Gegner. Um diese Informationen zu erhalten wird das gesamte Spiel einen halben Zug in der Vergangenheit betrachtet. Somit kann dann der alte Zug des Gegners als der eigene betrachtet werden und der eigene Zug als Antwort auf diesen. Damit der reinforcement learning Agent dies dann verarbeiten kann muss noch das ganze Spielfeld negiert werden. Es werden also die Steine beider Spieler getauscht damit es wieder wie eine ganz normale Zugabfolge aussieht.

5.5 Graphische Darstellung

Um die Effektivität der verschiedenen Versionen zu vergleichen wird der Reward in einem Graphen geplottet. Dies wird mit dem Paket pyplot von matplotlib gemacht. Beispiele hierfür findet man im Kapitel 6 Evaluierung.

5.6 Visualisierung des Spiels

Da das Spielfeld im Code nur als zweidimensionales Array, welches mit 1, 0 und -1 gefüllt ist, existiert ist es für den Menschen nicht so einfach zu erkennen was gerade im Spiel passiert ist. Hierfür wurde mit Hilfe des Python Pakets colorama eine Visualisierung erschaffen. Diese wandelt alle Felder in den Unicode Character 'BLACK CIRCLE' (U+25CF) um wobei die 1 Gelb und die -1 Rot eingefärbt werden. Hierzu wird noch der Hintergrund Blau eingefärbt um dem Design des Spiels möglichst nah zu kommen. Abbildung 5.3 ist ein durch diese Weise erstelltes Spielfeld für das Interne Array 5.2

$$\begin{bmatrix} [& 0 & 0 & 0 & 0 & 0 & 0 &] \\ [& 0 & 0 & 0 & 0 & 0 & 0 &] \\ [& 0 & 0 & 0 & 0 & 0 & -1 & 0] \\ [& 0 & 0 & 0 & 0 & 0 & -1 & 0] \\ [& 1 & 0 & 1 & -1 & 0 & -1 & 0] \\ [& 1 & 1 & -1 & 1 & 0 & -1 & 1] \end{bmatrix}$$

Abbildung 5.2: Array-Representation

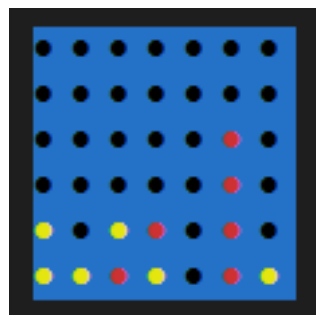


Abbildung 5.3: Das Spielfeld auf der Konsole

Kapitel 6

Evaluierung

In diesem Kapitel wird nun der implementierte reinforcement learning Agent evaluiert. Hierfür wird zuerst, mit Hilfe von einem deterministischen Gegner und einem Gegner der in jedem Zug eine zufällige Aktion auswählt, gezeigt, ob die neu gewählte Reward-Funktion einen Vorteil beim Lernen gibt. Danach wird dann an diesen beiden Gegnern gezeigt, ob das Lernen mit Einbezug der gegnerischen Züge eine Veränderung im Lernverhalten zeigt. Hiernach wird dann gegen den Minimax Algorithmus trainiert, um zu zeigen wie gut der reinforcement learning Agent werden kann.

Es folgen nun die Graphen dieser Evaluierung, wobei immer ein Paar von Graphen die selbe Lernepoche beschreibt. Der eine Graph beschreibt den absoluten Reward, also den genauen Wert, zu jeder Episode. Der andere Graph beschreibt den durchschnittlichen Reward zu jeder Episode im Bezug auf alle vorherigen Episoden dieser Epoche. Es wurde sich für eine Epoche von 1000 Episoden entschieden, da sich hier schon das Verhalten des Agenten abzeichnet und das Berechnen von längeren Epochen einen sehr hohen Zeitaufwand bedeuten würde.

Es sei darauf hin zu weisen, dass der durchschnittliche Reward (immer der rechte Graph) meist nicht das komplette Spektrum von 1 bis -1 abdeckt. Dies wurde gemacht um eine bessere Erkennbarkeit des Ausschnittes zu gewähren.

6.1 Reward-Funktion

Als erstes wird die Veränderung durch die zugabhängige Reward-Funktion evaluiert. Hierfür wird Abbildung 6.1 mit Abbildung 6.2 sowie Abbildung 6.3 mit Abbildung 6.4 verglichen.

Abbildung 6.1 (b) und Abbildung 6.2 (b) verhalten sich fundamental sehr ähnlich. Beide erreichen einen durchschnittlichen Reward von 0.5 nach etwa 250 Episoden. Wobei die einfache Reward-Funktion aber einen abfallenden, und die zugabhängige Reward-Funktion einen konstanteren Trend aufweist.

Abbildung 6.3(b) zeigt einen starken Abfall in gewonnenen Spielen nach etwa der 300. Episode. Dieses Verhalten ist auch sehr stark in Abbildung 6.3(a) zu erkennen. Im Vergleich zu der einfachen Reward-Funktion wird in Abbildung 6.4(b) eine Annähe-

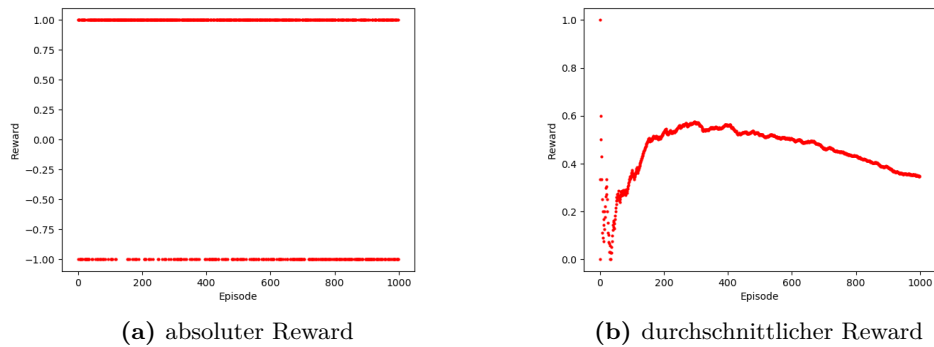


Abbildung 6.1: Reinforcement learning an zufällig werfendem Gegner mit einfacher Reward-Funktion und ohne die Züge des Gegners mitzulernen.

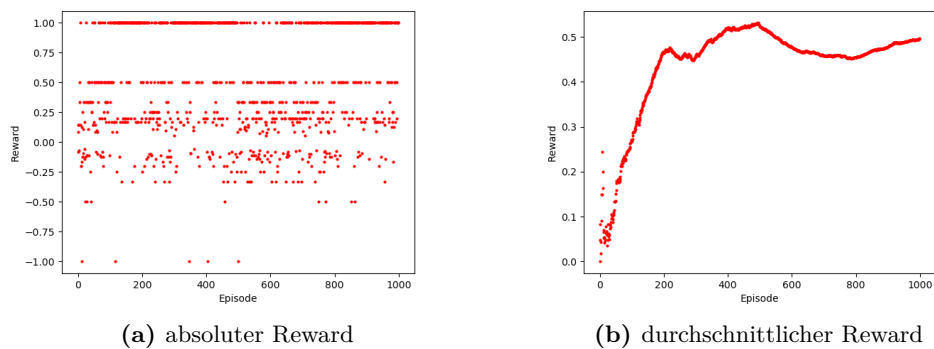


Abbildung 6.2: Reinforcement learning an zufällig werfendem Gegner mit verbesserter Reward-Funktion und ohne die Züge des Gegners mitzulernen.

rung an 0 aufzeigt was auf eine 50% Gewinnchance hinweist.

Es wurde gezeigt, dass durch die zugabhängige Reward-Funktion, besser gegen den zufälligen und den deterministischen Gegner gespielt wurde. Hieraus schließe ich, dass die zugabhängige Reward-Funktion einen positiven Effekt auf das Lernverhalten des reinforcement learning Agenten hat. Aus diesem Grund wird in allen fortlaufenden Evaluierungen mit der verbesserten Reward-Funktion gearbeitet.

6.2 Lernen vom Gegner

Nun wird evaluiert, ob das Lernen von den Zügen des Gegners einen positiven Einfluss auf das Lernverhalten des Agenten hat. Hierfür werden wieder der deterministisch und der zufällig ziehende Gegner genutzt.

Wie in Abbildung 6.5 (b) zu sehen ist, ist das Ergebnis hier im Vergleich zu Abbildung

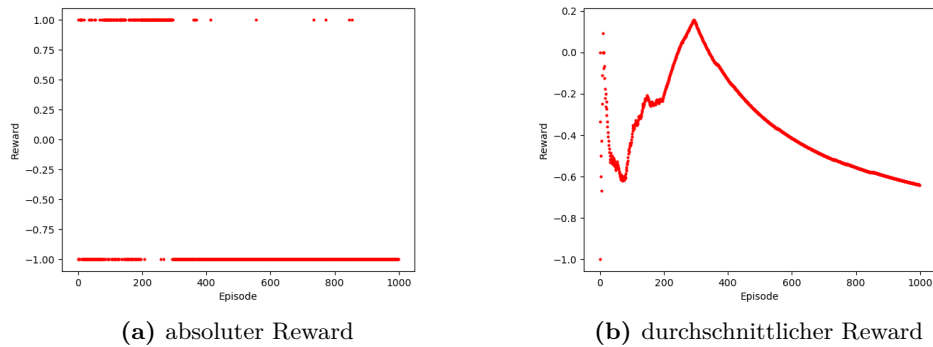


Abbildung 6.3: Reinforcement learning an deterministisch werfendem Gegner mit einfacher Reward-Funktion und ohne die Züge des Gegners mitzulernen.

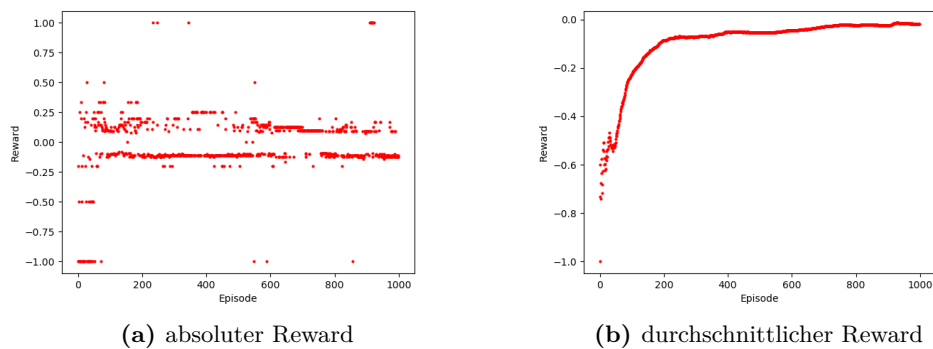


Abbildung 6.4: Reinforcement learning an deterministisch werfendem Gegner mit verbesserter Reward-Funktion und ohne die Züge des Gegners mitzulernen.

6.2 (b) wesentlich schlechter. Dies ergibt aber Sinn, da der Gegner ja komplett zufällig spielt und somit das Lernen seiner Züge eher als Nachteil betrachtet werden muss. Abbildung 6.6 zeigt hingegen eine sehr starke Verbesserung im Vergleich zu Abbildung 6.4. Wie in Abbildung 6.6(a) zu sehen ist sind die meisten Rewards oberhalb von 0, dies war in Abbildung 6.4(a) nicht der Fall.

Somit wurde gezeigt, dass das Nutzen der Züge eines Gegners mit zielführender Strategie zu einer Verbesserung des Lernverhaltens des reinforcement learning Agenten führt. Da in den folgenden Tests nur noch gegen Gegner mit einer zielführenden Strategie gelernt wird, wird das Lernen von Gegnerischen Zügen für alle folgenden Evaluationen genutzt.

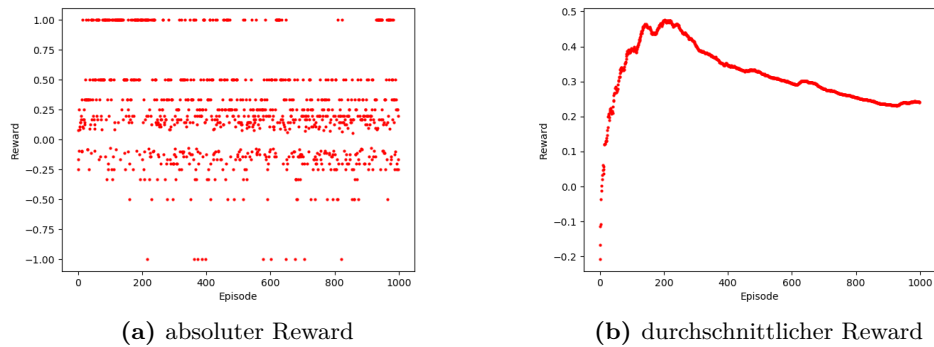


Abbildung 6.5: Reinforcement learning an zufällig werfendem Gegner mit verbesserter Reward-Funktion, wobei die Züge des Gegners mitgelernt werden.

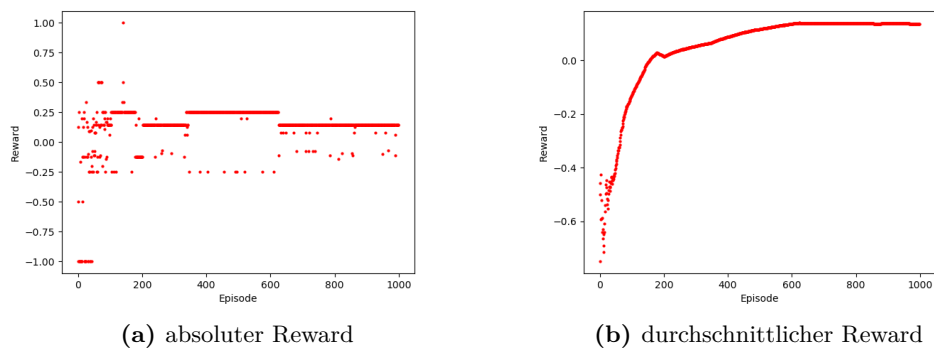


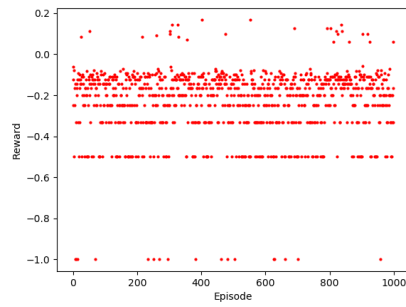
Abbildung 6.6: Reinforcement learning an zufällig werfendem Gegner mit verbesserter Reward-Funktion, wobei die Züge des Gegners mitgelernt werden.

6.3 Minimax

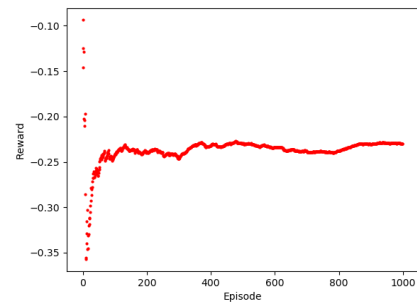
Die letzte Evaluierung sollte zeigen, dass ein reinforcement learning Agent mit den zuvor gezeigten Verbesserungen unterschiedlich schnell zu einer guten Strategie finden kann. Hierfür wurde der reinforcement learning Agent gegen den Minimax-Algorithmus mit verschiedenen starker Wahrscheinlichkeitsverteilung trainiert.

Die in Abbildung 6.7 und Abbildung 6.8 gezeigten Graphen zeigen das Trainingsverhalten gegen den Minimax Algorithmus. Diese sind nach steigender Wahrscheinlichkeit, den besten Zug zu wählen, geordnet. Das erste Paar in Abbildung 6.7 (a und b) hat somit die niedrigste Wahrscheinlichkeit den besten Zug zu wählen und das letzte Paar in Abbildung 6.8 (c und d) den höchsten.

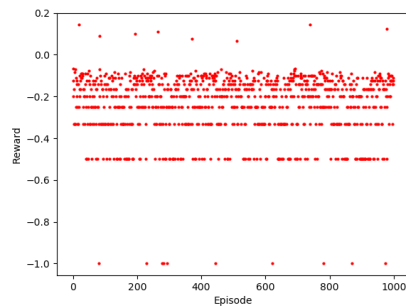
Leider ließ sich hier kein Lernvortschritt feststellen nicht mal bei der niedrigsten Wahrscheinlichkeit.



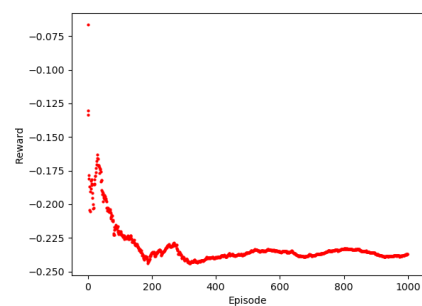
(a) absoluter Reward



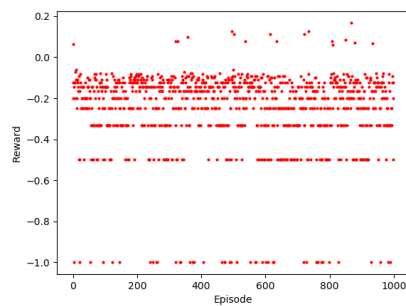
(b) durchschnittlicher Reward



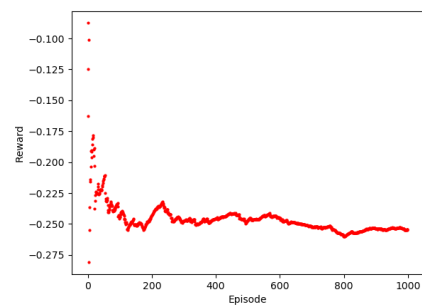
(c) absoluter Reward



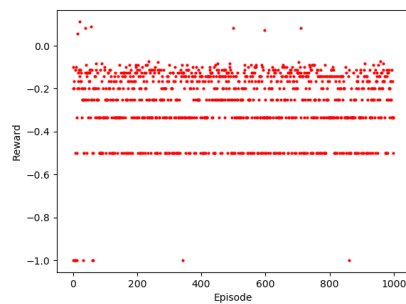
(d) durchschnittlicher Reward



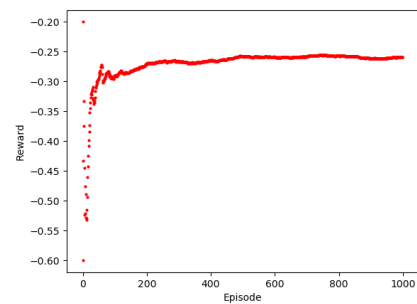
(e) absoluter Reward



(f) durchschnittlicher Reward

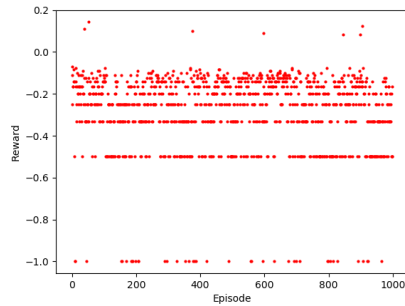


(g) absoluter Reward

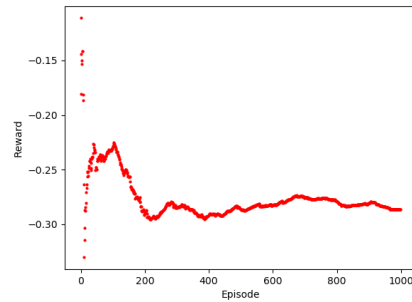


(h) durchschnittlicher Reward

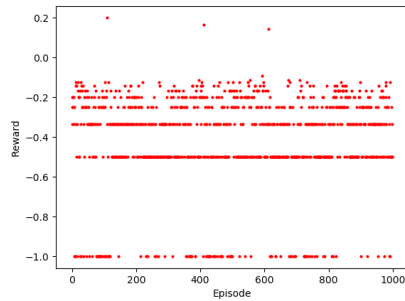
Abbildung 6.7: Reinforcement learning an Minimax Gegner mit verbesserter Reward-Funktion, wobei die Züge des Gegners mitgelernen werden. Teil 1



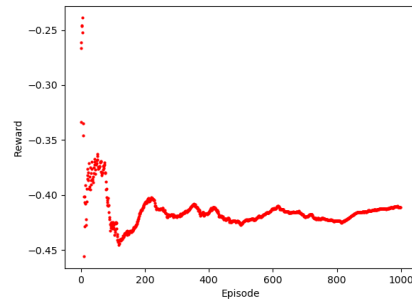
(a) absoluter Reward



(b) durchschnittlicher Reward



(c) absoluter Reward



(d) durchschnittlicher Reward

Abbildung 6.8: Reinforcement learning an Minimax Gegner mit verbesserter Reward-Funktion, wobei die Züge des Gegners mitgelernen werden. Teil 2

6.4 Besonderheiten

In den Graphen der Evaluation existiert eine Auffälligkeit die bei ungefähr 250 bis 300 Episoden auftritt. Hiernach wird das Lernverhalten des Agenten oft schlagartig schlechter. Dies ist am deutlichsten in der Abbildung 6.3(b) zu sehen. Ich vermute, dass dies mit der Epsilon-Greedy Policy zusammenhängt die in dem Raum etwa den Maximalen wert für ϵ erreicht hat.

Kapitel 7

Fazit

Es wurde festgestellt, dass das Lernen der Züge des Gegners einen positiven Effekt auf das Lernverhalten eines reinforcement learning Agenten haben kann. Des weiteren wurde gezeigt, dass die gewählte Reward-Funktion angemessen war und somit zur Verbesserung des Lernverhalten beigetragen hat. Dies konnte leider nur am Beispiel von Gegnern mit schlechten und sehr einfach zielführenden Strategien, und nicht am Gegner mit sehr guter Strategie gezeigt werden. Ob sich dieses Verhalten auch gegen bessere Gegner zeigt, müsste in weiteren Untersuchungen durch Vereinfachung der Problemstellung oder durch Verlängerung der Lernepoche getestet werden. Diese Veränderungen waren auf Grund des Rahmens der Bachelorarbeit leider nicht möglich.

Quellenverzeichnis

Literatur

- [1] Meire Fortunato et al. „Noisy Networks for Exploration“. Machine Learning (cs.LG). Ithaca, New York: Cornell University, 2017. URL: https://arxiv.org/pdf/1706.10295.pdf?source=post_page----- (siehe S. 13).
- [2] Olivier Armantier. „Does observation influence learning?“ *Games and Economic Behavior* 46.2 (2004), S. 221–239. URL: <http://www.sciencedirect.com/science/article/pii/S0899825603001246> (siehe S. 8).
- [3] Aurélien Géron. *Hands-On Machine Learning with Scikit-Learn and Tensorflow*. 2. Aufl. Gravenstein Highway North, Sebastopol, CA95472: O'Reilly Media Inc., 2017 (siehe S. 13).
- [4] Fred L. Drake Guido Van Rossum. *Python 3 Reference Manual*. 1. Aufl. Scotts Valley: CreateSpace, 2009 (siehe S. 7).
- [5] Jonathan Ho und Stefano Ermon. „Generative Adversarial Imitation Learning“. In: *Advances in Neural Information Processing Systems 29*. Hrsg. von D. D. Lee u. a. Curran Associates, Inc., 2016, S. 4565–4573. URL: <http://papers.nips.cc/paper/6391-generative-adversarial-imitation-learning.pdf> (siehe S. 8).
- [6] Patrick Koch und Wolfgang Konen Markus Thill. „Reinforcement Learning with N-tuples on theGame Connect-4“. Part of the Lecture Notes in Computer Science book series (LNCS, volume 7491). Gummersbach, Germany: Department of Computer Science, Cologne University of Applied Sciences, 2012. URL: https://www.researchgate.net/profile/Wolfgang_Konen/publication/235219697_Reinforcement_Learning_with_N-tuples_on_the_Game_Connect-4/links/0912f510811c0942f6000000.pdf (siehe S. 8).
- [7] Marco Wiering Martijn van Otterlo. „Reinforcement Learning and Markov Decision Processes“. Reinforcement Learning. Adaptation, Learning, and Optimization, vol 12. Berlin, Heidelberg: Springer, 2012. URL: https://www.ai.rug.nl/~mwiering/Intro_RLBOOK.pdf (siehe S. 4).
- [8] Bob Price und Craig Boutilier. „Implicit imitation in multiagent reinforcement learning“. In: *ICML*. Citeseer. 1999, S. 325–334 (siehe S. 8).
- [9] Andrew G. Barto Richard S. Sutton. *Reinforcement Learning An Introduction*. 2. Aufl. Cambridge, Massachusetts: The MIT Press, 2018 (siehe S. 8).

- [10] Lukas Stephan. „Implementierung und Evaluation eines intelligenten Spielagenten mit Reinforcement Learning“. Bachelorarbeit. Bremen, Deutschland: Universität Bremen, Informatik, Juni 2018. URL: unpublished (siehe S. 8).

Software

- [11] *connect-four-tensorflow by alexisjanvier*. URL: <https://github.com/marmelab/connect-four-tensorflow> (siehe S. 10).
- [12] *Keras*. URL: <https://keras.io/> (siehe S. 7).
- [13] *OpenAI Gym*. URL: <https://gym.openai.com/> (siehe S. 7, 8).
- [14] *Python*. URL: <https://www.python.org/> (siehe S. 7).
- [15] *Tensorflow*. URL: <https://www.tensorflow.org/> (siehe S. 7).

Online-Quellen

- [16] *AlphaStar: Mastering the Real-Time Strategy Game StarCraft II*. Jan. 2019. URL: <https://deepmind.com/blog/article/alphastar-mastering-real-time-strategy-game-starcraft-ii> (besucht am 05.03.2020) (siehe S. 8).
- [17] *MarI/O - Machine Learning for Video Games*. Juni 2015. URL: <https://youtu.be/qv6UVOQ0F44> (besucht am 05.03.2020) (siehe S. 8).
- [18] *Vier Gewinnt Spiel*. März 2020. URL: <https://cdn.anyfinder.eu/assets/5a6903838df2f5833889cbf45eb7599de9b858bb1821dc321540fb0bf4dca517> (besucht am 05.03.2020) (siehe S. 2).