



Fakultät
Informatik und Medien

BACHELORARBEIT

zur Erlangung des akademischen Grades

Bachelor of Science

im Studiengang Informatik

der Fakultät Informatik und Medien

der Hochschule für Technik, Wirtschaft und Kultur Leipzig

Evaluierung von Trainingsstrategien im Reinforcement Learning anhand eines Würfelspiels in Unity

— —

vorgelegt von: Tony Lenz

Geburtsort- und datum: Leipzig, den 07.01.1993

Abgabe: Leipzig, den 15. März 2024

Erstgutachter: Prof. Dr.-Ing. Müller ERSTGUTACHTER, HTWK Leipzig

Zweitgutachter: Prof. Dr.-Ing. Bleymehl ZWEITGUTACHTER, HTWK Leipzig



Hochschule für Technik,
Wirtschaft und Kultur Leipzig

Eidesstattliche Erklärung

Hiermit erkläre ich, dass ich die an der Hochschule für Technik, Wirtschaft und Kultur Leipzig, konkret an der Fakultät Informatik und Medien (FIM), eingereichte Arbeit zum Thema Evaluierung von Trainingsstrategien im Reinforcement Learning anhand eines Würfelspiels in Unity selbstständig, ohne Hilfe Dritter und ohne Benutzung anderer als der angegebenen Quellen und Hilfsmittel verfasst habe.

Alle den benutzten Quellen wörtlich oder sinngemäß entnommenen Stellen sind als solche einzeln kenntlich gemacht. Die Abbildungen in dieser Arbeit wurden von mir selbst erstellt oder mit einem entsprechenden Hinweis auf die Quelle versehen.

Diese Arbeit ist bislang keiner anderen Prüfungsbehörde vorgelegt und auch nicht veröffentlicht worden.

Ich bin mir bewusst, dass eine falsche Erklärung rechtliche Folgen haben wird.

.....

Unterschrift Tony Lenz

Leipzig, den 15. März 2024

Kurzfassung

Die Kurzfassung - präziser: Abstrakt, sollte auf eine Seite beschränkt sein. Der Abstrakt sollte das Ziel der Arbeit explizit beschreiben, die angewandten *Methoden* aufführen, die wichtigsten *Ergebnisse* aufzählen und die *Hauptschlussfolgerungen* aufzeigen.

Einen Weg zu finden, hunderte von Seiten an Informationen in wenigen Sätzen zusammenzufassen ist eine Herausforderung. Jedes Wort muss sorgfältig abgewogen werden (Gibt es eine bessere, prägnantere Möglichkeit, die Hauptaussage auszudrücken?). Die einzelnen Sätze sollten mit größter sorgfalt kombiniert werden.

Das Verfassen des Abstraktes sollte bis zum Schluss aufgeschoben werden, aus dem selben Grund, warum der Titel der Arbeit erst dann endgültige Form erhalten sollte, wenn das entsprechende Arbeit ansonsten vollständig ist.

Die nachfolgenden Ausführungen sollen nochmal ins Gedächtnis rufen was die Anforderungen an einzelne Abschnitte einer Abschlussarbeit an der FIM der Hochschule für Technik, Wirtschaft und Kultur Leipzig (HTWK) sind.

Die Kurzfassung schließt mit der Nennung von 5 bis 10 Schlagwörtern (Keywords) ab, welche der Verschlagwortung bspw. für die Recherche nutzbar gemacht werden. Genauer gesagt, sind das inhaltliche Schlüsselwörter Ihrer Arbeit.

Die Struktur dieser Arbeit gliedert sich wie folgt:

- **Theoretische Grundlagen:** In diesem Abschnitt werden die grundlegenden

Konzepte des Reinforcement Learning erläutert, darunter Neuronale Netze Markov-Entscheidungsprozesse und Probleme welche beim dem RL auftreten können. Weiterhin wird eine detaillierte Beschreibung des Spiels 'Noch mal' gegeben, einschließlich der Regeln, Spielziele und möglicher Strategien.

- **Konzeption:** Dieser Abschnitt umfasst eine Anforderungsanalyse und beschreibt den Ansatz zur Implementierung des Reinforcement Learning-Agenten für 'Noch mal'.
- **Experimente und Ergebnisse:** Es werden die Ergebnisse der Experimente präsentiert, einschließlich der Leistung des RL-Agenten beim Spielen von 'Noch mal' und einer Analyse seiner Fähigkeiten und Schwächen.
- **Diskussion:** Eine Diskussion über die Ergebnisse, die Einschränkungen der Studie und mögliche Verbesserungen wird vorgenommen.
- **Fazit und Ausblick:** Abschließend werden die wichtigsten Erkenntnisse zusammengefasst und ein Ausblick auf mögliche zukünftige Forschungsrichtungen gegeben.

Keywords: IoT, SD-WAN, Machine Learning (ML), Fiber to the Home (FTTH)

Inhaltsverzeichnis

1	Einleitung	1
2	Theoretische Grundlagen	3
2.1	Neuronale Netze	3
2.2	Reinforcement Learning	5
2.2.1	Markov- Entscheidungsprozess	7
2.3	Das Würfelspiel: Noch Mal	8
2.3.1	Spielablauf Einspieler Variante	10
3	Konzeption	11
3.1	Anforderungsanalyse	11
3.2	Auswahl der Verwendeten Technologien	12
4	Implementierung	14
4.1	Umsetzung in Unity	14
4.2	Visualisierung	16
4.3	Implementierung des Agenten	16
4.4	Spezifikation des Alghorithmus	18
4.5	Auswertung der erreichten Punkte	21
5	Präsentation der Ergebnisse	22
5.1	Agent ohne Spielreglementierung	22
5.2	Agent mit zusätzlichen Belohnungen	24
5.3	Vergleich trainiert und untrainiert	24
5.4	Training auf Sonderfeldern	26
5.5	Training mit mehr Spielzügen	28
5.6	Überprüfung auf Überanpassung	29
5.7	Training auf Minifeld	30

5.8	Training blinder Agent	32
6	Auswertung und Ausblick	34
6.1	Bewertung der Ergebnisse	34
6.2	Schritte zur Verbesserung des Agenten	35
	Literaturverzeichnis	36
	Abbildungsverzeichnis	38
	Tabellenverzeichnis	40
	Quellcodeverzeichnis	41
	Abkürzungsverzeichnis	I
A	Anhang - Abbildungen	II
B	Anhang - Tabellen	IV
C	Anhang - Quelltexte	V

1 Einleitung

In der heutigen Zeit stehen die Menschheit an der Schwelle einer digitalen Revolution, in der künstliche Intelligenz eine immer bedeutendere Rolle spielt. Insbesondere das Gebiet des Reinforcement Learnings hat in den letzten Jahren enorme Fortschritte gemacht und findet Anwendung in einer Vielzahl von Bereichen. Diese Entwicklung bietet vielfältige Möglichkeiten, komplexe Probleme zu lösen und intelligente Systeme zu entwickeln, die in der Lage sind, eigenständig zu lernen und Entscheidungen zu treffen.

Reinforcement Learning (RL) hat beeindruckende Erfolge erzielt, indem es Algorithmen entwickelt hat, die komplexe Spiele wie 'Go', auf einem kompetitiven Niveau spielen können und sogar menschliche Fähigkeiten übersteigt. Diese Anwendung verdeutlicht die Vielseitigkeit und Leistungsfähigkeit von Reinforcement Learning bei der Bewältigung verschiedenster Herausforderungen und unterstreichen die Fähigkeiten von RL, komplexe Probleme zu lösen und neue Lösungswege zu finden.[1, 2]

In dieser Bachelorarbeit liegt der Fokus auf der Implementierung und dem Training eines Reinforcement Learning-Agenten für ein in Unity umgesetztes Würfelspiel. 'Noch mal' ist ein Würfelspiel, das Strategie und Glück erfordert. Durch die Implementierung eines RL-Agenten für dieses Spiel, lässt sich untersuchen, wie gut maschinelle Lernmodelle in der Lage sind, komplexe Entscheidungsprobleme zu lösen und Strategien zu entwickeln, um ein definiertes Ziel zu erreichen. Diese Arbeit zielt darauf ab, einen Beitrag zum Verständnis der Anwendung von Reinforcement Learning in der Spieleentwicklung zu leisten und Einblicke in die Leistungsfähigkeit dieser Techniken zu bieten.

Die Arbeit gliedert sich in folgende Abschnitte:

- **Theoretische Grundlagen:** Die theoretischen Grundlagen zum Reinforcement Learning und neuronalen Netzen werden hier dargestellt. Außerdem werden die Spielregeln des modellierten Spiels beschrieben.
- **Konzeption:** Dieses Kapitel beinhaltet eine Anforderungsanalyse. Auf Grundlage dieser die Auswahl der verwendeten Technologien begründet.
- **Implementierung:** In diesem Kapitel wird dargestellt, wie das Spiel und der Agent implementiert wurde.
- **Präsentation der Ergebnisse:** Für diesen Abschnitt wurden verschiedene Versuche durchgeführt. Die Ergebnisse dieser werden in diesem Kapitel dargestellt.
- **Auswertung und Ausblick:** In diesem Kapitel werden die Ergebnisse bewertet und weitere Schritte zur Verbesserung des Agenten diskutiert.

2 Theoretische Grundlagen

2.1 Neuronale Netze

Neuronale Netze (=NN) sind ein Modell für künstliche Intelligenz nach dem Vorbild des Gehirns.

Es besteht aus mehreren Schichten verknüpfter Neuronen. Diese verarbeiten numerische Informationen und wandeln diese in einen Output. In 2.1 ist der Aufbau eines künstlichen Neurons schematisch dargestellt. Jedes Neuron besitzt eine Aktivierungsfunktion, welche entscheidet ob es aktiv ist oder nicht. Neuronen geben Werte von 0 (passiv) bis 1 (aktiv) als Output weiter.

Um ein NN zu erzeugen, wird eine Vielzahl der Neuronen miteinander verknüpft. Die Ausgänge der vorderen Neuronen sind dabei immer mit den Eingängen der Neuronen der nächsten Schicht verknüpft. 2.2 zeigt den Aufbau eines solchen Netzes.

Wird ein NN initialisiert, werden Kantengewichte zwischen den Neuronen zufällig verteilt. Diese Kantengewichte sorgen dafür wie stark einzelne Neuronen in die nachfolgende Rechnung eingehen. Deshalb liefern untrainierte NN schlechte Ergebnisse und müssen lernen die Problemstellung richtig zu lösen. Dieser Prozess wird als Training bezeichnet. Während des Trainings eines NN werden Kantengewichte angepasst, um die Heuristik an ein optimales Ergebnis anzupassen. Richtig trainierte NN können gute Antworten für komplexe Problemstellungen geben, so liefern sie beispielsweise in der Mustererkennung gute Ergebnisse.[3] Für das Training von NN wird viel Rechenleistung gebraucht, da der Prozess mit vielen Rechenoperationen verbunden ist. [4]

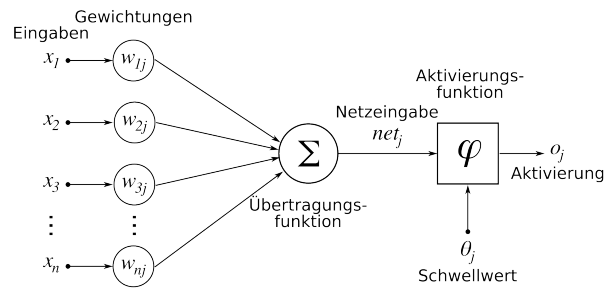


Abbildung 2.1: Aufbau eines künstlichen Neurons [5]

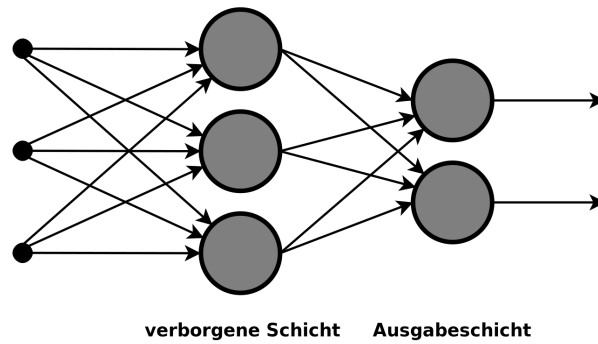


Abbildung 2.2: Schematische Darstellung eines neuronalen Netzes [5]

2.2 Reinforcement Learning

Reinforcement Learning ist ein Bereich des maschinellen Lernens, bei welchem ein Agent durch Interaktion mit seiner Umgebung (Environment) lernt, welche Aktionen in welchen Situationen am besten geeignet sind, um ein bestimmtes Ziel zu erreichen. Da ein Agent keine konkrete Vorgehensweise besitzt, versucht er über Trial-and-Error herauszufinden welche Aktionen zu Belohnungen führen.

Bestandteile des RL sind der Agent und die Umwelt. Der Agent bekommt die Informationen der Umwelt übergeben und entscheidet welche Handlungen daraus folgen sollen. Das Environment setzt die Aktionen des Agenten in Handlungen (Aktionen) um und bewertet diese mit numerischen Belohnungen (Rewards). Anhand der erhaltenen Rewards, versucht der Agent seine Aktionen anzupassen, um die zukünftigen Belohnungen zu maximieren. Ein RL-Agent nimmt seine Umgebung als eine Menge an bestimmten Zuständen wahr. Jeder Zustand enthält eine Vielzahl von Merkmalen welche für die Entscheidungsfindung relevant sein können. Als Zustandsraum wird die Menge aller möglichen Zustände bezeichnet.

Auf jedem dieser Zustände ist es möglich eine bestimmte Menge an Aktionen auszuführen, welche das Umfeld in einen anderen Zustand überführen. Die Menge aller möglichen Aktionen wird als Aktionsraum bezeichnet.

Wird eine Aktion auf einem bestimmten Zustand ausgeführt, so bezeichnet man die Zustandsübergänge bei Ausführung der Aktion als Zustandsübergangsfunktion.

Die Verhaltensstrategie eines Agenten liefert zu jedem Zustand eine Aktion. Die Policy wird im Verlauf des Trainings erlernt.

Das Erlernen einer Policy erfolgt durch Training des Neuronalen Netzes. Zu Beginn des Trainings werden Kantengewichte des zur trainierenden Neuronalen Netzes zufällig verteilt. Um zu überprüfen ob seine Aktionen gut oder schlecht sind, muss der Agent durch einen numerischen Rückgabewert die Information erhalten. Diese numerischen Rückgabewerte werden als Belohnungen (Rewards) bezeichnet. Belohnungen werden vom Zustandsraum verteilt und können gutes Verhalten des Agenten durch positive Rewards und falsches Verhalten durch negative Rewards bestärken. Der Agent versucht den Wert der erhaltenen Belohnungen zu maximieren und erlernt auf diese Weise eine optimale Policy.

Probleme, welche beim RL auftreten können, sind unter anderem:

- Überanpassung
- Belohnungsumgehung,
- Sparse Rewards
- große Beobachtungsvektoren

Von **Überanpassung** spricht man, wenn der Agent ein Problem gut in der Lernumgebung in welcher er trainiert wurde lösen kann allerdings nicht auf abweichenden Umgebungen. Es tritt auf, wenn der Agent zu spezialisiert auf seine Trainingsumgebung ist und sich nicht an andere Situationen anpassen kann. Diese Problem lässt sich durch Generalisierung des Trainingsprozesses beheben. Ein Agent, welcher in vielen zufällig generierten Umgebungen trainiert, kann deutlich besser mit neuen Situationen umgehen, benötigt im Gegenzug dazu auch mehr Zeit zum Trainieren.

Belohnungsumgehung tritt auf, wenn der Agent lernt die Belohnungsfunktion zu manipulieren, um erhaltene Belohnungen zu maximieren, ohne das eigentliche Ziel der Aufgabe zu erreichen. Dies führt zu unerwünschtem Verhalten des Agenten. Verhindert werden kann dieses Problem, indem Belohnungen direkt mit dem Ziel verknüpft werden, also nur Belohnungen ausgelöst werden, wenn tatsächlich positive Aktionen ausgeführt werden.

Von **Sparse Rewards** spricht man, wenn die Lernumgebung dem Agenten nur selten oder unregelmäßig Belohnungen vergibt. Dies führt dazu, dass der Agent Schwierigkeiten hat, die ihm gegebene Aufgabe richtig zu erlernen. Um dieses Problem zu umgehen, kann man zusätzliche künstliche Belohnungen einführen, welche den Agenten auf den Weg zum richtigen Verhalten führen. Dies kann im Rückschluss jedoch wieder zur Belohnungsumgehung führen, wenn der Agent die Zwischenbelohnungen, nutzt um maximale Belohnungen zu erhalten.

Auch ein zu **großer Beobachtungsvektor** kann zu Komplikationen führen. Wird dem Agenten eine Vielzahl an Observationen zugeführt, müssen erheblich mehr Parameter im Modell verarbeitet werden, was zu größerem Berechnungsaufwand führt. Dadurch wird auch die Lerngeschwindigkeit verlangsamt und der Bedarf an Speicherplatz für die Modelle steigt. So ist es möglich, dass einfach erscheinende Aufgaben mehrere Tage an Rechenzeit benötigen, um ein gutes Modell zu generieren. [6]

2.2.1 Markov- Entscheidungsprozess

Jedes RL Problem lässt sich durch den Markov Entscheidungsprozess (= MDP) beschreiben. Der MDP ist ein mathematisches Modell für Entscheidungsprobleme, bei welchem der Agent abhängig von der Umwelt Entscheidungen trifft, um ein bestimmtes Ziel zu erreichen. MDP setzen die Einhaltung der Markov-Eigenschaft voraus. Diese ist erfüllt, wenn ein Zustandsübergang nur vom letzten Zustand und der letzten Aktion abhängig ist. Hauptkomponenten des MDP sind eine endliche Menge valider Zustände, eine endliche Menge valider Aktionen, Belohnungsfunktion und Verhaltensstrategie. Das Ziel eines MDP ist es eine optimale Policy, durch Maximierung der erhaltenen Belohnungen zu ermitteln.

2.3 stellt den Ablauf eines MDP dar. Zu Beginn einer Episode, befindet sich die Lernumgebung in einem gewissen Zustand. Dieser Zustand wird dem Agenten übergeben. Anhand der Beobachtungen führt der Agent eine Aktion aus, welche die Umgebung verändert. Die Veränderung der Umgebung erzeugt einen neuen Zustand, welcher im nächsten Schritt wiederum dem Agenten übergeben wird. Durch Auslösen von Aktionen in der Lernumgebung werden Belohnungen erzeugt, welche dem Agenten einen numerischen Wert liefern, ob die Aktion gut oder schlecht war. Der Agent versucht die erhaltenen Belohnungen zu maximieren. Eine hohe Belohnung impliziert das richtige Verhalten zum Lösen des Problems, welches der Agent bewältigen muss. Dieser Prozess wiederholt sich bis das Problem gelöst wurde. [6]

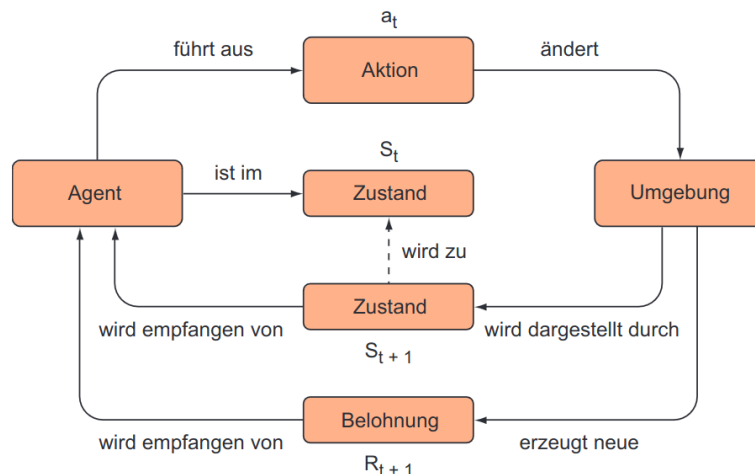


Abbildung 2.3: Ablauf eines MDP [6]

2.3 Das Würfelspiel: Noch Mal

Im modellierten Spiel 'Noch Mal!' geht es darum so viele Kästchen wie möglich anzukreuzen und damit viele Spalten und gleichfarbige Kästchen auszufüllen. Jeweils ein Farb- und ein Zahlenwürfel müssen kombiniert werden, um entsprechend verfügbare und zusammenhängende Felder der gewählten Farbe abzukreuzen. Dabei kann der Zahlenwürfel die Werte von eins bis fünf oder einem Joker ergeben. Der Farbwürfel besitzt die fünf Farben des Feldes und einem Joker. Gespielt wird nach folgenden Regeln:

1. Felder in der Spalte H sind von Beginn an verfügbar
2. Alle Kreuze müssen immer zusammenhängend in genau einem Farbblock der gewählten Farbe platziert werden
3. Kreuze müssen waagrecht oder senkrecht benachbart zu einem bereits abgekreuzten Feld oder Teil der Spalte H sein, um verfügbar zu werden
4. Es müssen genau so viele Felder angekreuzt werden, wie das Ergebnis des gewählten Zahlenwürfels
5. Es könnte nicht mehr als 5 Kästchen in einem Zug abgekreuzt werden
6. Wird ein Zahlenjoker gewählt, darf der Spieler eine Zahl von 1-5 bestimmen
7. Wird ein Farbjoker gewählt, darf der Spieler eine Farbe bestimmen

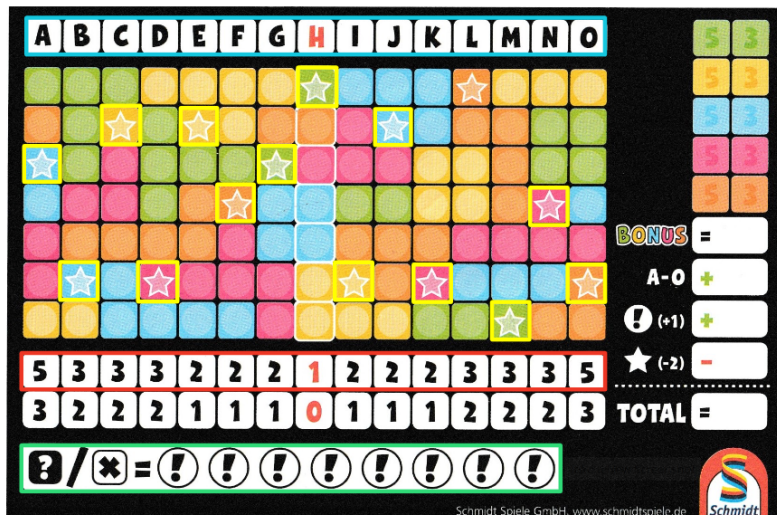
Um im Spiel 'Noch mal' eine möglichst hohe Anzahl an Punkten zu erhalten, ist es wichtig nach folgenden Strategien zu spielen:

Priorisierung äußerer Spalten: Äußere Spalten geben mehr Punkte, weshalb es wichtig ist diese komplett auszufüllen.

Beenden von Farben: Vollständig ausgefüllte Farben geben viele Extrapunkte. Im späten Spielverlauf kann es besser sein Farben komplett zu füllen anstatt Spalten zu werten.

Priorisieren von Sternfeldern: Jedes ausgefüllte Sternfeld gibt 2 Punkte, eine gute Spielweise ist es so viele Sternfelder wie möglich auszufüllen.

Strategische Nutzung von Jokern: Ungenutzte Joker geben zum Ende des Spiels Punkte. Es ist gut diese so wenig wie möglich zu nutzen, um Extrapunkte zu bekommen. Jedoch können mit Hilfe von Jokern einfach bestimmte Felder gewählt werden, welche benötigt werden um eine Wertung zu erzielen.



- Spaltennamen
- Mögliche Punkte pro Spalte
- Anzeiger für restliche Joker
- Sternfelder

Abbildung 2.4: Vorlage des Spielfedes mit Indikation der Punktwertung [7]

Die Grafik 2.4 ist ein Spielfeld aus 'Noch mal!'. Nach Vorlage dieses Spielfelds wurde die Lernumgebung für den Agenten modelliert. Die blau markierten Felder geben den Spaltennamen der Spalten an. Die rot markierten Felder zeigen an, wie viele Punkte beim ausfüllen der jeweiligen Spalte erzielt werden. Das grün markierte Feld zeigt die Anzahl der verbleibenden Joker an, wird ein Joker benutzt, muss eines der Felder abgestrichen werden. Zum Ende des Spiels erhält der Spieler Extrapunkte für die verbleibenden Joker. In jeder Spalte befindet sich ein gelb markiertes Sternfeld. Diese geben zum Ende des Spiels Minuspunkte, weshalb es wichtig ist alle Sternfelder auszufüllen.

2.3.1 Spielablauf Einspieler Variante

Der Spieler hat 30 Züge Zeit, maximale Punkte zu erreichen und würfelt jede Spielrunde alle 4 Würfel, bestehend aus 2 Farb- und 2 Zahlenwürfeln. Anschließend wählt er ein Paar aus Farb- und Zahlenwürfel aus und kreuzt entsprechend des gewürfelten Paares verfügbare Felder auf dem Spielfeld ab. Ein Spieler darf immer entscheiden, ob er Würfelwürfe zum Ankreuzen verwenden möchte oder nicht. Um Kästchen anzukreuzen, wählt der Spieler eine Kombination aus Zahlen- und Farbwürfel aus. Wählt er beispielsweise 'Grün' und '2' so müssen 2 zusammenhängende grüne Felder angekreuzt werden.

Gelingt es dem Spieler eine Spalte komplett auszufüllen, erhält er je nach Spalte Punkte dafür. Für äußere Spalten werden mehr Punkte vergeben als für die inneren Spalten. Für das komplette Ausfüllen einer Farbe erhält der Spieler fünf Punkte pro ausgefüllter Farbe. Jedes nicht angekreuzte Sternfeld gibt zum Spielende zwei Minuspunkte. Für jeden übrig gebliebenen Joker erhält der Spieler zum Ende des Spiels einen Punkt.

2.5 ist aus der Spielanleitung des Spiels. Sie zeigt an, wie viele Punkte erreichbar sind und bewertet das Ergebnis. Anhand dieser erreichten Punkte wird die Güte des Trainingsfortschritts bewertet.

PUNKTE	LEVEL
> 40	★★★★★ Es gibt also doch Superhelden!
37-40	★★★★★ Wirst du „Glückspilz“ oder „The Brain“ genannt?
33-36	★★★★★ Du könntest auch professioneller „ NOCH MAL! “-Spieler sein.
29-32	★★★★☆ Super! Welch grandioses Ergebnis!
25-28	★★★★☆ Hoffentlich ohne Schummeln geschafft!
21-24	★★★★☆ Klasse! Das lief ja gut.
17-20	★★★★☆ Das war wohl nicht dein erstes Mal...
13-16	★★★★☆ Gut, aber das geht noch besser.
9-12	★★★☆☆ Na, wird doch langsam.
5-8	★★☆☆☆ Nicht ganz schlecht.
1-4	★☆☆☆☆ Da muss wohl noch etwas geübt werden.
0	☆☆☆☆☆ Dabei sein ist alles.
< 0	✂☆☆☆☆ Das grenzt ja schon an Arbeitsverweigerung.

Abbildung 2.5: Grafik zur Bewertung der gesammelten Punkte der Einspieler Variante [7]

3 Konzeption

3.1 Anforderungsanalyse

In dieser Arbeit soll das Spiel 'Noch mal!' implementiert werden und von einem RL Agenten gespielt werden. Das Spiel muss nicht vom Nutzer selbst spielbar sein, sondern dem Agenten lediglich eine Lernumgebung bereitstellen, in welcher er trainieren kann. Es müssen alle Funktionalitäten des Spiels abgedeckt werden. Weiterhin muss die Lernumgebung das Durchführen von illegalen Zügen unterbinden. Spielparameter müssen konfigurierbar sein, die Flexibilität des Trainings zu erhöhen.

Die Visualisierung des Spiels steht nicht im Vordergrund. Trotzdem soll sie vorhanden sein, da sie ermöglicht Verhaltensweisen des Agenten besser überwachen zu können. Um das Spiel zu programmieren, wird eine leistungsfähige Gameengine vorausgesetzt, welche die Umsetzung vereinfacht. Da der Lernprozess des Agenten im Vordergrund steht, wird eine benutzerfreundliche Schnittstelle vorausgesetzt, welche ein RL Framework bereit stellt und die Verwaltung von Modellen vereinfacht. Um das Training zu überwachen und verschiedene Modelle miteinander zu vergleichen, wird ein Framework benötigt, welche den Trainingsprozess grafisch darstellt. Dieses soll ohne großen Mehraufwand nutzbar sein. Das Training von RL Modellen, benötigt viel Rechenzeit. Deshalb ist es nötig Trainingsprozesse zu parallelisieren um die Dauer des Trainings zu minimieren. Die Parallelisierung sollte von der Entwicklungsumgebung bereitgestellt werden. Um verschiedene Trainingsszenarien zu erstellen und situativ einsetzen zu können, ist es notwendig mehrere Lernumgebungen konfigurieren und speichern zu können. In 3.1 werden die Anforderungen zusammengefasst und in funktionale und nicht funktionale Anforderungen kategorisiert.

Funktional	Nicht-funktional
Implementierung des Spiels	Überwachen des Trainings
Bereitstellen der Lernumgebung	Parallelisierung des Trainingprozesses
Abdecken der Spielfunktionalitäten	Konfigurierbare und speicherbare Lernumgebungen
Unterbinden von illegalen Spielzügen	Visualisierung des Spiels
Schnittstelle für RL-Framework	Verwendung einer leistungsfähigen Gameengine
	Benutzerfreundliche Entwicklungsumgebung
	Konfigurierbare Spielparameter

Tabelle 3.1: Anforderungen an die Implementierung des Spiels 'Noch mal!' für einen RL-Agenten

3.2 Auswahl der Verwendeten Technologien

Basierend auf der Anforderungsanalyse ergeben sich Anforderungen an die Technologien, welche verwendet werden. Um das Spiel 'Noch mal!' zu programmieren, wurde Unity als bevorzugte Gameengine gewählt. Unity stellt eine Vielzahl von Bibliotheken zur Verfügung und ermöglicht die unkomplizierte Umsetzung der Visualisierung des Spiels. C# ist die gängige Programmiersprache in Unity, deshalb wird das Projekt in C# umgesetzt. Mit der Auswahl von Unity als Gameengine ist bereits die Mehrheit der funktionalen Anforderungen aus 3.1 teilweise erfüllt, müssen jedoch programmiertechnisch umgesetzt werden.

Für die Erstellung des RL-Agenten wurde das ML-Agents Framework verwendet. Es bietet alle Funktionalitäten zum Übergeben von Beobachtungen an den Agenten und Schnittstellen zum Ausführen von Aktionen im Lernumfeld. Weiterhin ist es möglich Aktionen mit Belohnungen zu bewerten. Durch die Integration von ML Agents wird sichergestellt, dass der Agent den aktuellen Zustand des Spiels richtig übergeben bekommt und darauf reagieren kann. Durch die Nutzung des Frameworks, ist auch die Schnittstelle für RL geschaffen.

Zur grafischen Darstellung des Trainingsprozesses wurde Tensorboard genutzt. Die Integration erfolgt ohne großen Mehraufwand und bietet die automatisierte Erstellung von Grafiken des Trainingsprozesses. Dies erleichtert die Evaluierung von verschiedenen Modellen.

Zur Visualisierung der erreichten Punkte der Agenten, wurde Python mit Matplotlib verwendet. Mit Matplotlib ist es möglich grafische Darstellungen von Daten zu erstellen. Mithilfe dieser Grafiken ist es möglich Rückschlüsse auf das Verhalten der Agenten zu

führen, beziehungsweise deren Trainingsfortschritte zu bewerten.

Mit der Verwendung jener genannten Technologien und Frameworks ist es möglich den Rahmen dieser Arbeit zu bearbeiten und zu bewerten. In 3.2 wird ersichtlich welche Anforderung über welche Technologie erfüllt wird.

Anforderung	Technologie
Implementierung des Spiels	Unity
Bereitstellen der Lernumgebung	ML-Agents
Abdecken der Spielfunktionalitäten	Unity
Unterbinden von illegalen Spielzügen	Unity
Schnittstelle für RL-Framework	ML-Agents
Überwachen des Trainings	TensorBoard, Matplotlib
Parallelisierung des Trainingsprozesses	ML-Agents
Konfigurierbare Lernumgebungen	ML-Agents
Visualisierung des Spiels	Unity, Matplotlib
Benutzerfreundliche Entwicklungsumgebung	Unity
Konfigurierbare Spielparameter	Unity

Tabelle 3.2: Anforderungen an die Implementierung des Spiels 'Noch mal!' für einen RL-Agenten und die entsprechenden Technologien

4 Implementierung

4.1 Umsetzung in Unity

Der **Controller** stellt alle Funktionalitäten bereit, welche gebraucht werden, um die Lernumgebung zu initialisieren. Er besitzt Prefabs des Agents, des GameFields und der Würfel und initialisiert diese zum Start. Weiterhin implementiert der Controller die Funktionalität der Punktevergabe, welche für eine Mehrspielervariante genutzt werden kann. Der Controller ist das Parent aller anderen Elemente und so ist er das zentrale Element der Steuerung. Auch das wiederholte Rollen der Würfel wird im Controller angestoßen. Der Controller war hilfreich beim Erstellen paralleler Trainings, da dieser mehrfach in die Szene aufgenommen werden musste um mehrere Spielfelder, welche gleichzeitig bespielt werden zu initialisieren.

Der **NumberDice** implementiert die Logik, welche für das Würfeln und Visualisieren der Zahlenwürfel benötigt wird. Die Visualisierung funktioniert mit selbst angefertigten Sprites, welche in einem sortierten Array liegen und je nach gewürfelter Zahl initialisiert und gerendert werden. Beim wiederholten Würfeln, wird das aktuelle Sprite zerstört und ein neues initialisiert. Damit ist gewährleistet, dass immer das aktuelle Würfelergebnis angezeigt wird. Die Zahl des Würfels wird als Integer gespeichert, wobei er die Zahlen 1-6 annehmen kann. Die Zahl '6' entspricht dem Zahlenjoker. A.2 zeigt die verwendeten Sprites für den Zahlenwürfel.

Wie der Zahlenwürfel implementiert der **ColorDice** die Funktionalität des Würfels der Farben. Diese werden als String dargestellt und können folgende Werte annehmen: {'blue', 'green', 'red', 'yellow', 'orange', 'joker'}. Zur Visualisierung wird ein Sprite erstellt, was in der gewürfelten Farbe eingefärbt wird. Ein schwarzes Feld entspricht dem gewürfelten

Farbjoker.

Das **GameField** stellt das tatsächliche Spielfeld dar. Es implementiert die benötigten Methoden um die SquareFields zu verwalten und rückzusetzen. Außerdem wird die Anzahl der Joker in ihm gehalten. A.1 zeigt neun dieser Spielfelder, welche nach Vorlage von 2.4 modelliert wurden.

Funktionalitäten:

- Visualisierung des Spielfeldes
- Aktualisieren der Gruppen aller Felder
- Abkreuzen der Felder
- Berechnen der validen Nachbarn der Felder
- Berechnen der verbleibenden Felder einer bestimmten Farbe
- Rückgabe der validen Felder für die aktuell gewählten Würfel.
- Reduzieren der verbleibenden Joker
- Rücksetzen der Felder um ein neues Spiel zu Starten

Die **FieldSquares** stellen die einzelnen Teilfelder des Spielfeldes dar. In Tabelle 4.1 wird dargestellt, welche Informationen gehalten werden.

Beschreibung	Typ	Wertebereich
Feld ist ein Sternfeld	Boolean	True / False
Farbe des Feldes	String	-
Feld ist ausgefüllt	Boolean	True / False
Feld ist verfügbar	Boolean	True / False
Clustergröße	Integer	1-6
X-Koordinate des Feldes	Integer	0-14
Y-Koordinate des Feldes	Integer	0-6

Tabelle 4.1: Übersicht Informationen der Fieldsquares

4.2 Visualisierung

Die Visualisierung des Spielfeldes erfolgt über ein angefertigtes Prefab. In diesem wurden die 105 Kästchen in einem Raster von 15x7 instanziiert und manuell mit den Informationen versehen. Dieses manuell angefertigte Spielfeld wurde als Prefab gespeichert und dient als Umgebung für den Agenten. Zu Beginn des Spiels werden die Felder instanziiert in die Farben der hinterlegten Information in den richtigen Farben eingefärbt. Ausgefüllte Kästchen werden grau eingefärbt, diese Funktionalität wird im Fieldsquare Prefab ausgeführt. A.1 zeigt den Aufbau der Trainingsumgebung. Jedes der Spielfelder besitzt einen eigenen Agenten und ist somit die Lernumgebung für genau einen Agenten. Alle Agenten trainieren gemeinsam ein Modell. Es können beliebig viele weitere Lernumgebungen hinzugefügt werden. Mehrere gleichzeitige Trainingsvorgänge bedeuten jedoch wiederum einen höheren Rechenaufwand.

4.3 Implementierung des Agenten

Der Agent ist die Schnittstelle zwischen dem Environment und dem RL. Dem Agent werden alle nötigen Informationen des Spielfeldes übergeben. Diese werden in ein neuronales Netz übertragen, welches wiederum die Ausgabewerte in einem Vektor zurück an den Agenten leitet. Anschließend wird der Vektor verarbeitet und die gewählten Aktionen werden ausgeführt. Für gute Aktionen erhält der Agent positive Rewards, bei schlechten Aktionen wird der Zug übersprungen.

Zu Beginn jeder Episode, welche einem Spielzug entspricht, muss dem Agenten der aktuelle Zustand des Feldes übermittelt werden, aus welchem er die bestmögliche Option für einen Zug berechnet. In der ML Agents Bibliothek gibt es hierfür eine vorgefertigte Methode mit dem Namen `CollectObservations`. Diese erzeugt einen Observationsvektor 4.2, zu welchem die Informationen der aktuellen Zustands hinzugefügt werden. Während des Trainings eines Neuronalen Netzes, muss die Größe des Vektors und die Art der Werte gleich bleiben. Deswegen ist nicht möglich ein Modell auf unterschiedlichen Spielfeldern zu trainieren, da sich so die Anzahl der Beobachtungen unterscheiden würden. ?? erstellt anhand der Daten den Vektor aus beobachtungen und übergibt ihn an den Agenten.

Aufbau der Beobachtungen:

Index	Beschreibung	Type	Wertebereich
0	Anzahl der verbleibenden Joker	Float	$[0 - 1]$
1	Anzahl der gespielten Runden	Float	$[0 - 1]$
2	Ergebnis des ersten Zahlenwürfels	Float	$[0 - 1]$
3	Ergebnis des zweiten Zahlenwürfels	Float	$[0 - 1]$
4-9	Ergebnis des ersten Farbwürfels	Vector6 (Binary)	$(0, 1)^6$
10-15	Ergebnis des zweiten Farbwürfels	Vector6 (Binary)	$(0, 1)^6$
16-24	Informationen für Feld 1	FeldVektor	-
25-33	Informationen für Feld 2	FeldVektor	-
...
953-961	Informationen für Feld 105	FeldVektor	-

Tabelle 4.2: Zusammenfassung der Observations und Feldinformationen

Stelle im Vektor	Beschreibung	Type	Wertebereich
$k * 9 + 16 - k + 21$	Farbe des Feldes k	Vector6 (Binary)	$(0, 1)^6$
$k * 9 + 22$	Ist Feld k verfügbar	Boolean	True / False
$k * 9 + 23$	Ist Feld k abgestrichen	Boolean	True / False
$k * 9 + 24$	Ist Feld k ein Sternfeld	Boolean	True / False

Tabelle 4.3: Aufbau FeldVektor

Anhand der Observations berechnet das neuronale Netz einen Ausgabevektor. Tabelle 4.4 zeigt den Aufbau der hier verwendeten Beobachtungen. Mit diesem führt der Agent nun bestimmte Aktionen aus und versucht sein Ergebnis (Rewards) zu maximieren. Anhand der gesammelten Rewards wird das neuronale Netz nun angepasst, um das bestmögliche Ergebnis zu erreichen.

Das Vergeben von Rewards lehnt sich an die Punktevergabe im Spiel wie nach Vorlage 2.4 an. Der Agent erhält Belohnungen, wenn er auch Punkte im gespielten Spiel erlangen würde. In 4.5 ist eine Übersicht der Belohnungen. Aus dieser lässt sich ablesen, welche Aktion welchen Reward auslöst.

Tabelle 4.4: Index und Beschreibung der Variablen

Index	Beschreibung	Typ	Wertebereich
1	Index des gewählten ZahlenWürfels	Integer	0-1
2	Index des gewählten Farbwürfels	Integer	0-1
3	Jokerzahl	Integer	0-4
4	X-Koordinate des gewählten Feldes	Integer	0-14
5	Y-Koordinate des gewählten Feldes	Integer	0-7
6	Action 1 für die Auswahl der Nachbarn	Continuous	-
7	Action 2 für die Auswahl der Nachbarn	Continuous	-
8	Action 3 für die Auswahl der Nachbarn	Continuous	-
9	Action 4 für die Auswahl der Nachbarn	Continuous	-

Aktion	Erhaltene Belohnung
Abkreuzen eines Sternfeldes	2
Ausfüllen einer kompletten Spalte	1-5 abhängig der Spalte
Ausfüllen einer kompletten Farbe	5
Verbleibende Joker zum Ende des Spiels	1 / verbleibendem Joker

Tabelle 4.5: Belohnungen für bestimmte Aktionen

4.4 Spezifikation des Algorithmus

Im folgenden wird der Ablauf zum Wählen der Felder erläutert. Im Beispiel wird der Ausgabevektor $(1, 1, 0, 3, 4, 0.6, 0.5, 0.4, 0.8)$ verwendet.

Die ersten beiden Stellen des Ausgabevektors entsprechen den gewählten Würfeln. Im ersten Schritt werden alle Felder des Spielfelds untersucht, ob sie ein valides Ziel für das gewürfelte Ergebnis bilden. Dies ergibt sich aus der Gruppe der Spielfelder, der Farbe und ob das Kästchen verfügbar ist. Valide Felder werden in eine Liste (availableFields) aus Verfügbaren Feldern geschrieben. Abb. 4.1 zeigt den beschriebenen Zustand des Feldes. Abb. 4.2 bildet den nächsten Schritt im Algorithmus ab. Es wird geprüft, ob die gewählten Koordinaten in availableFields vorhanden sind. Wenn kein Feld verfügbar ist, wird die Episode abgebrochen und der Agent überspringt seinen Zug. Sofern der Agent ein valides Feld gewählt hat, wird dieses in eine weitere Liste (pickedFields) geschrieben und benachbarte Felder der selben Gruppe werden zurückgegeben.

Im nächsten Schritt wird jedem der verfügbaren Nachbarn abhängig der Gesamtanzahl ein Wertebereich zwischen 0 und 1 zugewiesen, dies wird in 4.3 und 4.6 verdeutlicht. Anhand des diskreten Wertes des Ausgabevektors wird das zugehörige Feld in eine Liste

aus den gewählten Feldern geschrieben.

FeldKoordinaten	von	bis
(3,5)	0	0.33
(4,4)	0.33	0.66
(3,3)	0.66	0.99

Tabelle 4.6: Bereiche für bestimmte Felder

Für alle gewählten Felder werden die benachbarten Felder zurückgegeben und der vorherige Schritt wiederholt. Wenn so viele Felder gewählt, wie erwürfelt wurden, werden die Felder anschließend ausgefüllt und auf Rewards überprüft. Die wird in Abbildung 4.4 dargestellt.

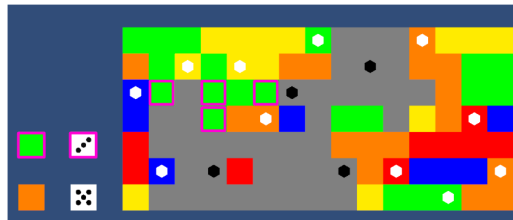


Abbildung 4.1: Valide Felder für das gewählte Würfelergebnis wurden markiert

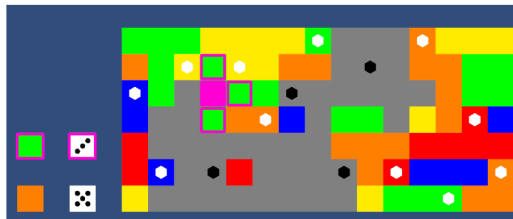


Abbildung 4.2: Feld(3,4) wird in die pickedField Liste aufgenommen und benachbarte Felder werden zurückgegeben.

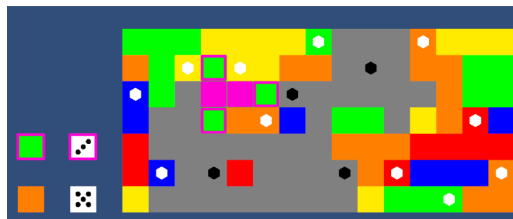


Abbildung 4.3: Feld(4,4) ist das nächste gewählte Feld und wird in pickedFields aufgenommen

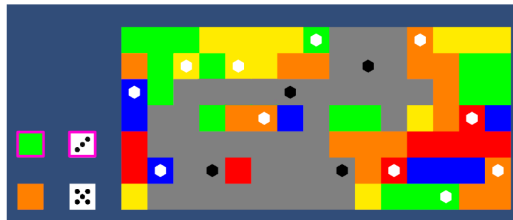


Abbildung 4.4: Felder wurden gewählt und ausgefüllt

4.5 Auswertung der erreichten Punkte

Um die verschiedenen Modelle bewerten zu können, wurden die erreichten Punkte aufgezeichnet und anschließend ausgewertet. Im Verlauf eines Spiels werden alle erreichten Punkte im Spielfeld gespeichert. Nach 30 gespielten Zügen wird die erreichte Punktzahl in eine Logdatei geschrieben. In 4.5 lässt sich ableiten, dass Belohnungen nur dann verteilt werden, wenn es auch zu einer Punktwertung kommt. Die beiden Werte sind lediglich um 30 Punkte verschoben, da der Spieler das Spiel mit -30 Punkten beginnt.

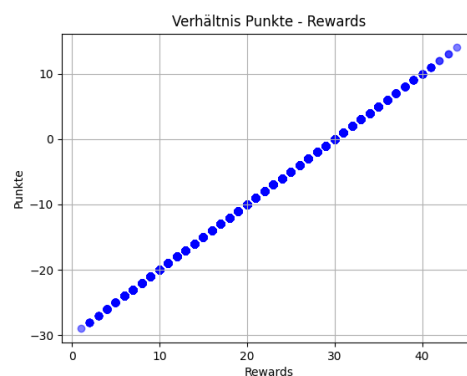


Abbildung 4.5: Verhältnis Rewards zu Punkten

Die Logdateien wurden anschließend mit Python ausgewertet. In C.4 werden die angegebenen Punkteaufzeichnungen importiert und jedem Datensatz einen Titel zur Beschriftung der Graphen hinzugefügt. Im Anschluss wird hier der gleitende Durchschnitt abhängig der gewählten Anzahl an Spielen berechnet.

Dabei werden erst alle Werte addiert und der Mittelwert gebildet. Danach wird der herausrotierende Wert abgezogen und um den Neuen ergänzt. Auf diese Weise muss jeder Wert lediglich 2 mal gelesen werden.

In der Methode C.5 können nun variabel Graphen gewählt und anschließend gerendert werden. Vor Methodenaufruf, muss die Liste an anzuzeigenden Graphen bearbeitet werden und in den Methodenaufruf übergeben werden.

5 Präsentation der Ergebnisse

Um den Lernprozess evaluieren zu können, wurden verschiedene Experimente durchgeführt. Diese Experimente werden in diesem Kapitel vorgestellt und ausgewertet.

5.1 Agent ohne Spielreglementierung

Im ersten Schritt überprüfte ich, ob der Agent die Spielregeln selbstständig über Belohnungen erlernen kann. In diesem Versuch sollte der Agent Felder via Index auswählen. Jedes Kästchen im Spielfeld besaß einen Index zwischen 0-104. Je nach Höhe der gewürfelten Zahl, wurden ebenso viele Feldindizes überprüft. In ?? ist der Aufbau des Actionbuffers dargestellt.

Index	Bezeichnung	Datentyp	Wertebereich
1	Index des zu wählenden Farbwürfels	int	0-1
2	Index des zu wählenden Zahlenwürfels	int	0-1
3	Feldindex	int	0-104
4	Feldindex	int	0-104
5	Feldindex	int	0-104
6	Feldindex	int	0-104
7	Feldindex	int	0-104

Tabelle 5.1: Aufbau Actionbuffer

Nach der Auswahl aller potenziell abzukreuzenden Feldern wurde überprüft, ob der Zug legal ist. Die Überprüfung belief sich auf folgende Kriterien: Mindestens ein Feld musste verfügbar sein, alle Felder mussten benachbart und der selben Farbe sein. Illegale Züge zogen negative Rewards mit sich, legal ausgeführte Züge dagegen positive.

Diese Herangehensweise führte nicht zum gewünschten Ergebnis. Auch nach einigen

Stunden des Trainings, konnte der Agent nur sehr selten legale Züge durchführen. Dem Agent war es nicht möglich in der begrenzten Trainingszeit den Zusammenhang der Beobachtungen zu den gegebenen Rewards festzustellen. Dennoch ist der Ansatz nicht gänzlich falsch. Mit einem hohen Rechenaufwand, könnte der Agent auch die Spielregeln erlernen, es würde nur sehr viel Zeit kosten.

5.2 Agent mit zusätzlichen Belohnungen

In diesem Versuch bekam der Agent zusätzlich zu den ursprünglichen Rewards, zusätzliche Belohnungen für bestimmte Aktionen. Diese Belohnungen sollten den Agenten zu komplexeren Spielzielen leiten, um so schneller zu einer optimalen Policy zu gelangen. In reftab:rewards2 sind alle zusätzlichen Rewards ersichtlich. Der Versuch wurde 2.4Mio Spielzüge ausgeführt. Dieser Zeitraum ist, wie sich im Verlauf der folgenden Experimente herausstellt, relativ kurz. Da sich jedoch ein negatives Ergebnis abzeichnete, wurde auf längeres Training verzichtet.

Aktion	Erhaltene Belohnung
Abkreuzen von Feldern	0.02f pro Feld
Abkreuzen eines gesamten Clusters	0.04f pro Feld
Wahl eines Würfelpaars ohne legale Züge	-50.0f
Wahl eines Jokers ohne verfügbare Joker	-50.0f

Tabelle 5.2: Zusatzbelohnungen für verschiedene Aktionen

Die Grafiken 5.1 und 5.2 zeigen deutlich, dass je länger das Training voranschritt, die durchschnittlich erreichten Punkte pro Spiel sanken. Die zusätzlichen Belohnungen führten dazu, dass der Agent nicht versuchte Spalten abzukreuzen oder Farben komplett auszufüllen, da er das Abkreuzen von Clustern für deutlich effizienter hielt, um seine Belohnungen zu maximieren. Obwohl die Belohnungen für erreichte Spalten oder komplett ausgefüllte Farben mehr Punkte ergaben.

5.3 Vergleich trainiert und untrainiert

In diesem Experiment, werden die erreichten Punkte und Rewards eines untrainierten Agenten gegenüber den erzielten Ergebnissen eines trainierten Agenten gestellt. Der trainierte Agent hat bereits 25 Mio. Spielzüge absolviert, was ungefähr 830k gespielten Spielen entspricht. Der untrainierte Agent trainiert ein neu initialisiertes NN, welches zufällig gewählte Kantengewichte zwischen den Neuronen erhält. Wie an den Grafiken 5.4 und 5.3 zu erkennen ist, hat der trainierte Agent tatsächlich einen höheren Durchschnitt an erzielten Punkten pro Spiel. Auch die gesammelten Rewards sind bei dem trainierten

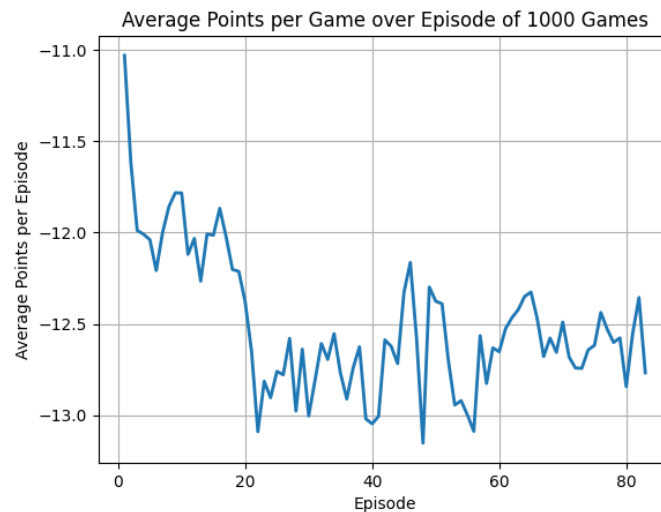


Abbildung 5.1: Durchschnittliche Punkte des Agenten mit Zusatzbelohnungen

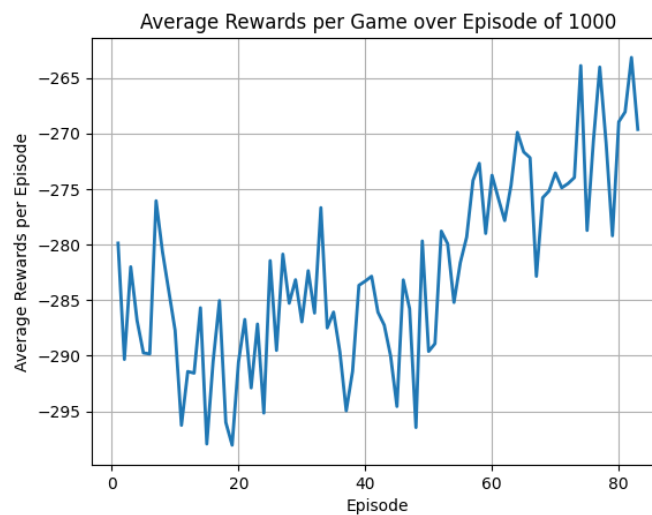


Abbildung 5.2: Durchschnittliche Belohnungen des Agenten mit Zusatzbelohnungen

Agenten höher. Dies liegt daran, dass die Rewards so festgelegt sind, dass der Agent sie nur erhält, wenn er auch im Spiel punktet. Schon während des Trainings war ein merklicher Unterschied festzustellen, deshalb war das Ergebnis dieses Experiments zu erwarten.

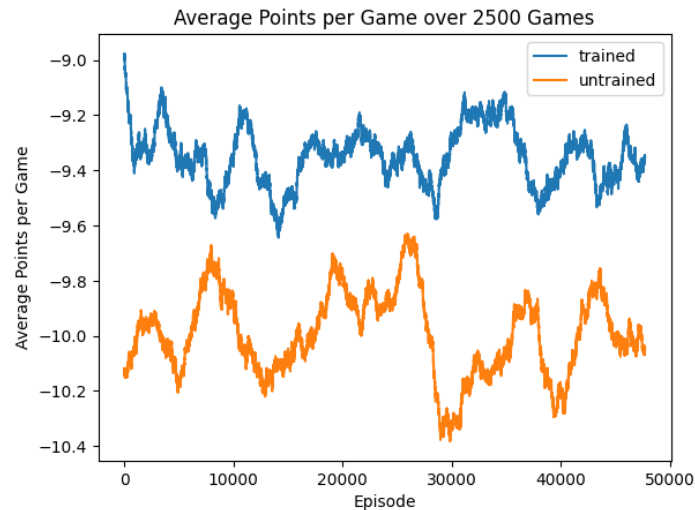


Abbildung 5.3: Durchschnitt der erreichten Punkte pro Spiel

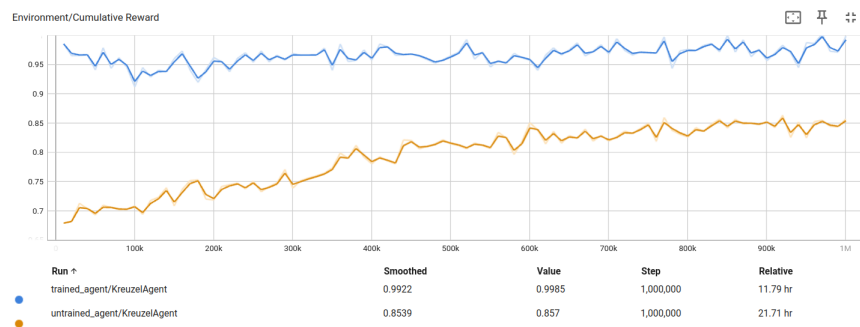


Abbildung 5.4: Übersicht der gesammelten Belohnungen

5.4 Training auf Sonderfeldern

In diesem Experiment wurden die erreichten Punkte und Rewards des trainierten Agenten gegen einen Agenten, welcher mit Sonderfeldern trainiert wurde verglichen. Beide Agenten wurden mit dem selben vortrainierten Modell instanziiert. 5.5 stellt den Aufbau der Trainingsumgebung dar. Dieser Versuch sollte überprüfen, ob ein Agent im Verlauf besser erlernen kann welche Stellen im Beobachtungsvektorvektor für welche Information zuständig sind.

In der Grafik 5.5 ist das Training mit den speziellen Feldern dargestellt. Jedes der Feld hat gewisse Besonderheiten, welche sich zu den normalen Spielfeldern abgrenzen. Fünf

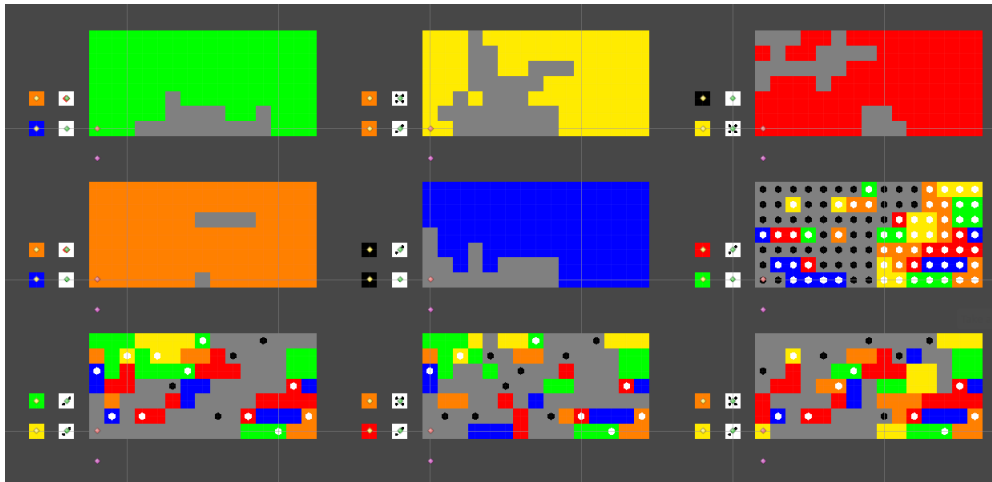


Abbildung 5.5: Übersicht der speziellen Felder

Felder sind in einer kompletten Farbidentität eingefärbt. In C.2 wurde der Zahlenwürfel manipuliert um häufiger die entsprechende Farbe zu werfen. Diese Felder sollten dem Agenten besser den Zusammenhang des Farbwürfels und des gewählten Farbidentität der Kästchen näherbringen. In dem anderen Spielfeld ist jedes Feld als Sternfeld markiert. Dies sollte dem Agenten zeigen, dass jedes Feld mit markierten Sternen mehr Punkte bringt. Bei den anderen drei Feldern ist jedes Feld von vorn herein als verfügbar markiert. Dies sollte zum einen das Konzept des verfügbaren Feldes vermitteln, zum anderen dem Agenten ermöglichen das Feld weiter als normal zu explorieren, um die komplexen Ziele des Spiels leichter zu erreichen.

Das Training mit speziellen Feldern führte zu einer Verschlechterung des Ergebnisses. Aus der Grafik 5.6 lässt sich ableiten, dass das Modell, welches mit speziellen Feldern trainiert wurde, im Durchschnitt weniger Punkte erhalten erreichte, als der normal trainierte Agent. Dieses Training führte nicht zu einer Verbesserung des Modells. Ursache hierfür liegt sicher im Spielfeld, in welchem alle Felder als Sternfelder markiert wurden. Hier konnte der Agent willkürlich Züge ausführen und bekam überdurchschnittlich viele Punkte. Deshalb priorisierte der Agent nicht mehr die eigentlichen Ziele, was wiederum zur Folge hatte, dass die Leistung des Agenten auf dem eigentlichen Feld schlechter wurde.

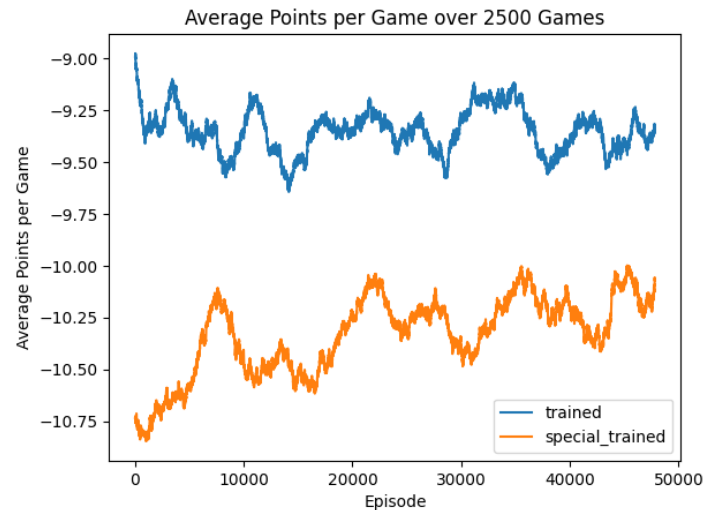


Abbildung 5.6: Durchschnitt der erreichten Punkte beider Agenten

5.5 Training mit mehr Spielzügen

In diesem Experiment trainierte der Agent mit mehr zur Verfügung stehenden Spielzügen. Dadurch konnte der Agent das Feld besser explorieren und insgesamt mehr Aktionen auslösen welche zu Belohnungen führten. Dies hat zur Folge, dass auch schwierig erreichbare Rewards ausgelöst wurden, welche somit vom Agent erlernt werden konnten. Die Grafik 5.7 zeigt, dass das Experiment eine Verbesserung des Models zur Folge hatte, da im Durchschnitt etwas mehr Punkte erreicht wurden.

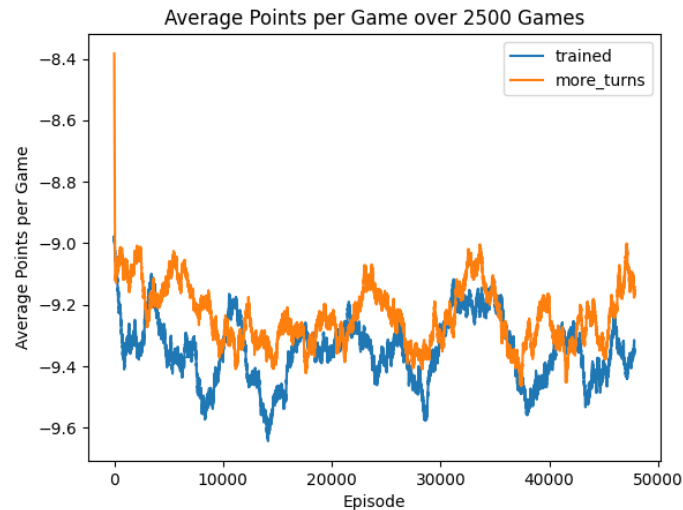


Abbildung 5.7: Vergleich 'Mehr Züge' und 'trainierter Agent'

5.6 Überprüfung auf Überanpassung

In diesem Experiment sollte der Agent auf Überanpassung überprüft werden. Trainierter und untrainierter Agent spielten das Spiel nach normalen Spielregeln auf einem anderen Spielfeld. In den Grafiken 5.8 und 5.9 ist erkennbar, dass beide Agenten ungefähr die selben Rewards gesammelt haben. Der untrainierte Agent konnte im Durchschnitt jedoch etwas mehr Punkte sammeln. Dies schließt darauf, dass der Agent tatsächlich nur auf dem im Training verwendeten Spielfeld gut performen kann und neue Spielfelder erst erlernen muss. Interessant ist weiterhin, dass der Durchschnitt aller Punkte etwa 2 Punkte über dem des anderen Spielfeldes liegt, was auf eine höhere Schwierigkeit des schwarzen Spielfeldes hinweist.

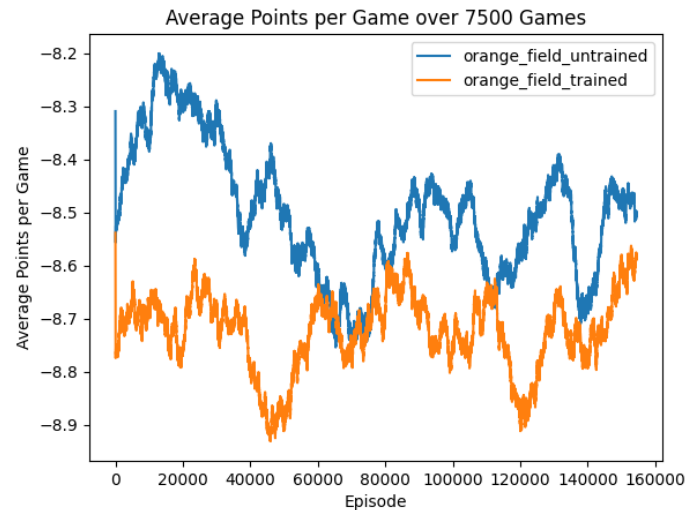


Abbildung 5.8: Vergleich Punkte 'trainiert' und 'untrainiert' auf orangen Spielfeld

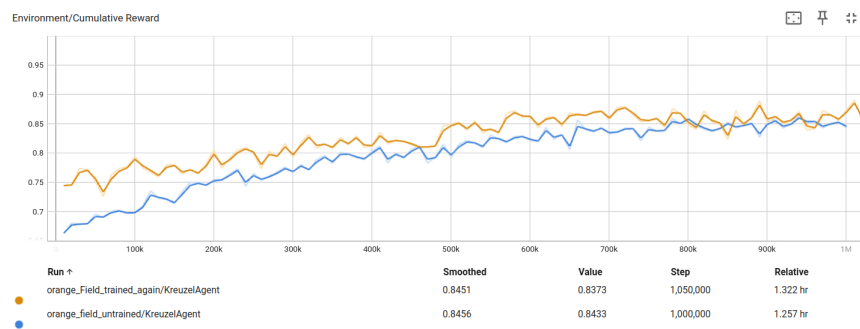


Abbildung 5.9: Übersicht gesammelte Rewards auf orangen Spielfeld

5.7 Training auf Minifeld

Um zu überprüfen, ob ein kleineres Feld einen positiven Effekt auf das Training hat, habe ich bei einem Spielfeld 3 Zeilen abgeschnitten und ließ einen Agenten darauf trainieren. A.3 zeigt die Lernumgebung für dieses Experiment für einen Agenten. Im Anschluss überprüfte ich die Leistung des Agenten auf dem normalen Spielfeld gegenüber einem untrainierten Agenten. Da die Observations von Modellen gleich bleiben müssen, entschied ich mich nicht mehr vorhandene Kästchen mit Nullen im Vektor zu präsentieren. 5.1 zeigt den Quellcode zum übergeben eines leeren Feldes.

Quellcode 5.1: Befüllen des Beobachtungsvektor mit leerem Feld

```

1 private void PushEmptyField(VectorSensor sensor){
2     sensor.AddObservation(GetColorIndexOneHotFromColor(""));

```

```

3   sensor.AddObservation(0);
4   sensor.AddObservation(0);
5   sensor.AddObservation(0);
6   }

```

Die Grafiken 5.10 und 5.11 zeigen den Durchschnitt der erreichten Punkte pro Spiel und die gesammelten Blohnungen während des Trainings. Es wird ersichtlich, dass der mit einem kleinen Feld vortrainierte Agent schlechter performt, als die beiden anderen. Belohnungen wurden auch hier nur verteilt, wenn es zur Punktwertung kommt. Deshalb ist es interessant, dass der untrainierte Agent mehr Punkte erreicht, als der auf dem kleinen Feld vortrainierte Agent, obwohl dieser wiederum einen höheren Durchschnitt an Belohnungen erhält. Ursache für das schlechtere Ergebnis, ist auf die Überanpassung zurückzuführen. So wurde beim Training auf dem kleinen Feld ausschließlich dieses Feld erlernt und nicht eine gute Spielweise.

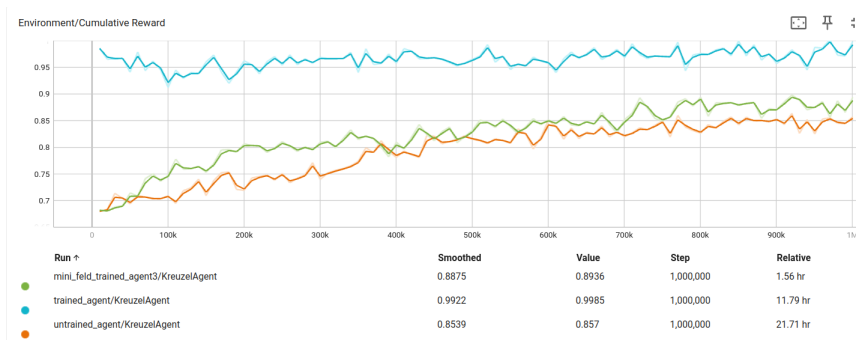


Abbildung 5.10: Gesammelte Belohnungen mit Minifeld

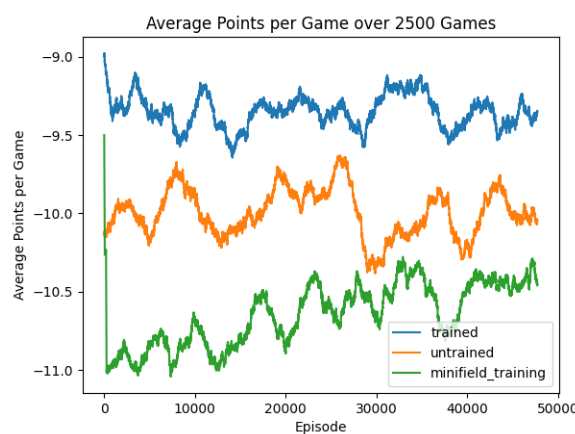


Abbildung 5.11: Vergleich erreichte Punkte

5.8 Training blinder Agent

In diesem Experiment bekam der Agent lediglich die Würfel, die Anzahl der verbleibenden Joker und die aktuelle Runde des Spiels übergeben. Dieser Versuch sollte überprüfen, wie gut ein Agent, der die aktuellen Informationen vom Spielfeld nicht erhält, performt. Dies sollte überprüfen, ob das trainierte Modell tatsächlich besser ist, als eine rein zufällige Spielweise. Da die Auswahl der Felder durch zufällige Interpolation aller möglichen Felder abläuft, kann der blinde Agent normal spielen. Durch den Versuchsaufbau verringert sich der Beobachtungsvektor von 916 auf eine Größe von 15 Informationen. Dies führte dazu, dass der Agent schnell zu seiner optimalen Policy gelangen konnte. Auch wenn der Agent das Spielfeld nicht sieht, kann er dieses implizit erlernen. Dieser Prozess wäre allerdings nicht sonderlich robust, wäre sehr lernintensiv und würde nicht auf anderen Feldern funktionieren. Der Agent konnte bereits nach sehr kurzer Zeit von etwa 400k Lernschritten gegen sein Maximum konvergieren, wie in 5.12 ersichtlich ist. Dort schneiden sich die beiden grünen Graphen und verbleiben auf ungefähr dem selben Niveau. Dieses Modell wurde insgesamt 5 Mio. Episoden trainiert. Der trainierte Agent konnte sich dagegen kontinuierlich minimal verbessern und erreichte auch nach 25Mio Lernschritten noch keinen Maximalwert.

Abbildung 5.13 zeigt, dass auch der blinde Agent im Durchschnitt weniger Punkte erreichen konnte. Dies beweist, dass das Training, trotz der geringen erreichten Punkte, positiv verlaufen ist.

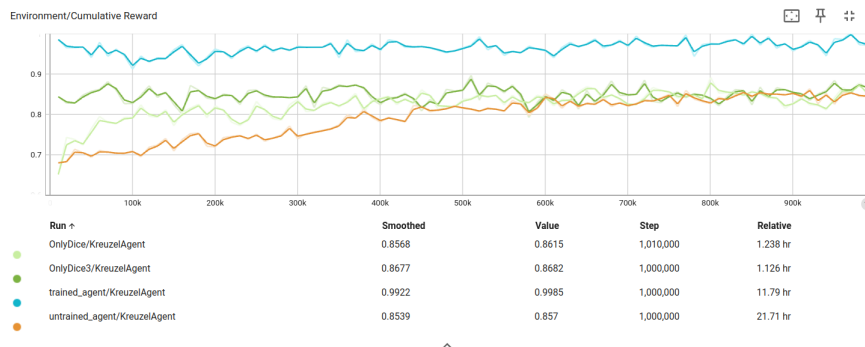


Abbildung 5.12: Gesammelte Belohnungen blinder Agent

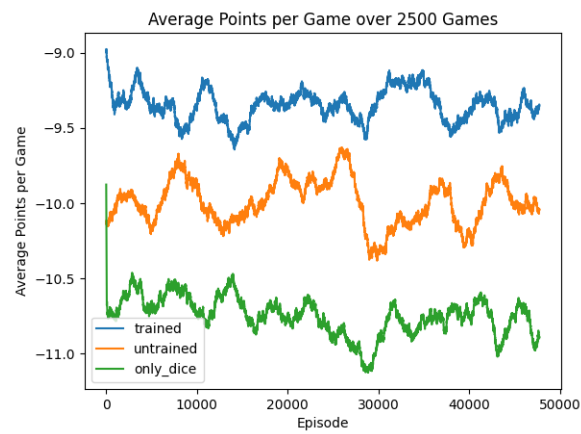


Abbildung 5.13: Vergleich erreichte Punkte blinder Agent

6 Auswertung und Ausblick

6.1 Bewertung der Ergebnisse

Die Experimente aus dem vorigen Kapitel zeigen deutlich, dass der Trainingsfortschritt des Agenten vorhanden war. Das trainierte Modell erhielt deutlich mehr Belohnungen als die untrainierten oder mit speziellen Lernumgebungen trainierten Agenten. Leider lässt sich dieser Fortschritt nur bedingt auf die Metrik der durchschnittlich erreichten Punkte anwenden. Wünschenswert wäre ein Modell gewesen, welches im Durchschnitt positive Spielergebnisse sammelt. Dem Agenten gelang es nicht kontinuierlich nach einer Strategie zu spielen, welche viele Punkte nach sich zieht. So konnte er fast niemals eine Farbe komplett ausfüllen, was viele Punkte bringen würde. Die liegt vor allem daran, dass er dieses für RL Agenten schwere Ziel nur sehr selten erreicht wurde und somit auch nicht erlernt werden konnte.

Das Ergebnis auf einem anderen Spielfeld war für den trainierten Agenten überraschend schlecht. Dies spricht dafür, dass das Training auf einem einzelnen Spielfeld nicht förderlich ist für ein Modell, der auf verschiedene Situationen gut agieren soll.

Während des Verlaufs dieser Arbeit musste ich feststellen, dass das Problem, ein wenig kompliziertes Spiel wie 'Noch mal!' zu erlernen, für einen RL Agenten deutlich komplexer ist als erwartet. Der Agent musste eine Vielzahl von Daten verarbeiten und diese semantisch zuordnen. Dies führte zu einer immensen Dauer des Trainingsprozesses. Der Agent, welcher 25 Mio Trainingsschritte durchlaufen hatte, brauchte für diesen Prozess etwa 30 Stunden. Während dieser Zeit konnte das Modell nur minimale Verbesserungen aufweisen. Weiterhin kommt hinzu, dass das modellierte Spiel zum großen Teil ein Glücksspiel ist. Es kann vorkommen, dass der Spieler Pech im Würfeln hat und so trotz perfekter Spielweise ein schlechtes Spielergebnis erreicht. Deshalb ist es auch nicht möglich ein Modell zu erzeugen, was immer ein gutes Ergebnis liefert. Allerdings wäre

mit weiterer Rechenzeit und vermehrter Nutzung von verschiedenen Trainingsarten noch eine weitere Verbesserung des Modells möglich.

6.2 Schritte zur Verbesserung des Agenten

Mit mehr Rechenleistung und Trainingszeit könnte der Agent durchaus noch zu besseren Ergebnissen kommen. Ein Flaschenhals war unter anderem, dass jede Trainingsiteration von 1 Mio. Schritten ungefähr 1,5 Stunden dauerte. Falls es zu Fehlern während des Trainings kam, musste die aufgewendete Zeit erneut investiert werden. Weiterhin konnte sich der Agent, wenn auch nur langsam, stetig verbessern. Mit einer deutlich erhöhten Trainingsdauer und mehr Rechenleistung kann das Modell deutlich verbessert werden. Ein kompletter Verzicht der Visualisierung während des Trainings kann zu einer Verbesserung der Trainingsdauer führen. Es würden mehr Ressourcen zur Berechnung des neuronalen Netzes zur Verfügung stehen und damit die Dauer verkürzen. Weiterhin kam es insbesondere während langen Trainingsepisoden zu Crashes, welche durch Unity hervorgerufen wurden. Dies könnte durch eine performantere Lernumgebung substituiert werden, um einen robusteren Ablauf zu erschaffen.

Ein weiteres Problem des trainierten Modells war die Überanpassung. Um dies zu umgehen und Modelle zu trainieren, welche auf verschiedene Zustände des Spielfeldes optimal reagieren können, wäre es hilfreich, das Spielfeld für jede Spielrunde zufällig generisch zu erzeugen. Dies würde dazu führen, dass der Agent nicht ein einzelnes Spielfeld, sondern das Spielen des Spiels erlernt. Generisch erzeugte Trainingsumgebungen würden die Trainingsdauer erhöhen, dafür aber in ein robusteres Modell resultieren.

Anpassen der Hyperparameters des Agenten können zu einer Optimierung des Trainingsprozesses führen. Leider ist das Anpassen der Hyperparameter nicht intuitiv nachvollziehbar, weshalb Erfahrung im Umgang mit RL erforderlich ist. Durch geschicktes Verändern dieser Parameter lassen sich lokale Maxima überspringen oder die Exploration des Agenten beschleunigen.

Eine weitere Verbesserung des Trainingsprozessen könnte auch eine ausgedehnte Selektion nach Vorbild der Evolutionären Algorithmen liefern. Dies hätte einen höheren Aufwand der Trainings zur Folge, könnte aber dazu führen, dass sich gute Modelle durchsetzen die schneller eine optimale Spielstrategie etablieren.

Literaturverzeichnis

- [1] Jonathan Hui. *AlphaGo: How it works technically?* en. Mai 2018. URL: <https://jonathan-hui.medium.com/alphago-how-it-works-technically-26ddcc085319> (besucht am 07.02.2024).
- [2] *AlphaGo*. en. Dez. 2020. URL: <https://deepmind.google/technologies/alphago/> (besucht am 04.03.2024).
- [3] Uwe Lorenz. *Reinforcement Learning: Aktuelle Ansätze verstehen - mit Beispielen in Java und Greenfoot*. de. Berlin, Heidelberg: Springer Berlin Heidelberg, 2020. ISBN: 978-3-662-61650-5 978-3-662-61651-2. DOI: 10.1007/978-3-662-61651-2. URL: <http://link.springer.com/10.1007/978-3-662-61651-2> (besucht am 08.03.2024).
- [4] Wolfgang Ertel. *Grundkurs Künstliche Intelligenz: Eine praxisorientierte Einführung*. de. Computational Intelligence. Wiesbaden: Springer Fachmedien Wiesbaden, 2021. ISBN: 978-3-658-32074-4 978-3-658-32075-1. DOI: 10.1007/978-3-658-32075-1. URL: <https://link.springer.com/10.1007/978-3-658-32075-1> (besucht am 09.02.2024).
- [5] *Künstliche Neuronale Netze / EF Informatik 2023*. de-CH. URL: <https://informatik.mygymer.ch/ef2023/07-ki/08-knn.html#kunstliches-neuron> (besucht am 06.03.2024).
- [6] Alex Zai und Brandon Brown. *Einstieg in Deep Reinforcement Learning: KI-Agenten mit Python und PyTorch programmieren*. en. München: Hanser, 2020. ISBN: 978-3-446-45900-7.
- [7] Schmidt Spiele GmbH. *Spielregeln 'Noch mal!'* URL: https://www.schmidtspiele.de/files/Produkte/4/49327%20-%20Noch%20mal!/49327_Noch_Mal_DE.pdf.
- [8] Richard S. Sutton und Andrew Barto. *Reinforcement learning: an introduction*. en. Nachdruck. Adaptive computation and machine learning. Cambridge, Massachusetts: The MIT Press, 2014. ISBN: 978-0-262-19398-6.

- [9] John Schulman u. a. *Proximal Policy Optimization Algorithms*. en. arXiv:1707.06347 [cs]. Aug. 2017. URL: <http://arxiv.org/abs/1707.06347> (besucht am 26. 01. 2024).
- [10] Manfred Broy. *Logische und Methodische Grundlagen der Programm- und Systementwicklung: Datenstrukturen, funktionale, sequenzielle und objektorientierte Programmierung - Unter Mitarbeit von Alexander Malkis*. de. Wiesbaden: Springer Fachmedien Wiesbaden, 2019. ISBN: 978-3-658-26301-0 978-3-658-26302-7. DOI: 10.1007/978-3-658-26302-7. URL: <http://link.springer.com/10.1007/978-3-658-26302-7> (besucht am 08. 02. 2024).
- [11] Oliver Kramer. *Computational Intelligence: Eine Einführung*. de. Informatik im Fokus. Berlin, Heidelberg: Springer Berlin Heidelberg, 2009. ISBN: 978-3-540-79738-8 978-3-540-79739-5. DOI: 10.1007/978-3-540-79739-5. URL: <https://link.springer.com/10.1007/978-3-540-79739-5> (besucht am 09. 02. 2024).
- [12] *Using TensorBoard to Observe Training - Unity ML-Agents Toolkit*. URL: <https://unity-technologies.github.io/ml-agents/Using-Tensorboard/> (besucht am 09. 02. 2024).
- [13] Yinlong Yuan u. a. „A novel multi-step reinforcement learning method for solving reward hacking“. en. In: *Applied Intelligence* 49.8 (Aug. 2019), S. 2874–2888. ISSN: 0924-669X, 1573-7497. DOI: 10.1007/s10489-019-01417-4. URL: <http://link.springer.com/10.1007/s10489-019-01417-4> (besucht am 11. 02. 2024).
- [14] L. P. Kaelbling, M. L. Littman und A. W. Moore. „Reinforcement Learning: A Survey“. en. In: *Journal of Artificial Intelligence Research* 4 (Mai 1996), S. 237–285. ISSN: 1076-9757. DOI: 10.1613/jair.301. URL: <https://www.jair.org/index.php/jair/article/view/10166> (besucht am 14. 03. 2024).

Abbildungsverzeichnis

2.1	Aufbau eines künstlichen Neurons [5]	4
2.2	Schematische Darstellung eines neuronalen Netzes [5]	4
2.3	Ablauf eines MDP [6]	7
2.4	Vorlage des Spielfedes mit Indikation der Punktwertung [7]	9
2.5	Grafik zur Bewertung der gesammelten Punkte der Einspieler Variante [7]	10
4.1	Valide Felder für das gewählte Würfelerggebnis wurden markiert	19
4.2	Feld(3,4) wird in die pickedField Liste aufgenommen und benachbarte Felder werden zurückgegeben.	19
4.3	Feld(4,4) ist das nächste gewählte Feld und wird in pickedFields aufgenommen	19
4.4	Felder wurden gewählt und ausgefüllt	20
4.5	Verhältnis Rewards zu Punkten	21
5.1	Durchschnittliche Punkte des Agenten mit Zusatzbelohnungen	25
5.2	Durchschnittliche Belohnungen des Agenten mit Zusatzbelohnungen	25
5.3	Durchschnitt der erreichten Punkte pro Spiel	26
5.4	Übersicht der gesammelten Belohnungen	26
5.5	Übersicht der speziellen Felder	27
5.6	Durchschnitt der erreichten Punkte beider Agenten	28
5.7	Vergleich 'Mehr Züge' und 'trainierter Agent'	29
5.8	Vergleich Punkte 'trainiert' und 'untrainiert' auf orangen Spielfeld	30
5.9	Übersicht gesammelte Rewards auf orangen Spielfeld	30
5.10	Gesammelte Belohnungen mit Minifeld	31
5.11	Vergleich erreichte Punkte	31
5.12	Gesammelte Belohnungen blinder Agent	32
5.13	Vergleich erreichte Punkte blinder Agent	33

A.1	Aufbau der Trainingsumgebung	II
A.2	Verwendete Sprites für den Zahlenwürfel	II
A.3	Lernumgebung für einen Agenten auf dem 'MiniFeld'	III

Tabellenverzeichnis

3.1	Anforderungen an die Implementierung des Spiels 'Noch mal!' für einen RL-Agenten	12
3.2	Anforderungen an die Implementierung des Spiels 'Noch mal!' für einen RL-Agenten und die entsprechenden Technologien	13
4.1	Übersicht Informationen der Fieldsquares	15
4.2	Zusammenfassung der Observations und Feldinformationen	17
4.3	Aufbau FeldVektor	17
4.4	Index und Beschreibung der Variablen	18
4.5	Belohnungen für bestimmte Aktionen	18
4.6	Bereiche für bestimmte Felder	19
5.1	Aufbau Actionbuffer	22
5.2	Zusatzbelohnungen für verschiedene Aktionen	24

Quellcodeverzeichnis

5.1	Befüllen des Beobachtungsvektor mit leerem Feld	30
C.1	Befüllen des Beobachtungsvektor mit validem Feld	V
C.2	Erhöhen der Warscheinlichkeit eine bestimmte Farbe zu würfeln	V
C.3	Erstellen des Beobachtungsvektors	VI
C.4	Einlesen aller angegebenen Logdateien	VII
C.5	Methode zum Anzeigen des Graphen	VIII

Abkürzungsverzeichnis

FIM Fakultät Informatik und Medien

FTTH Fiber to the Home

HTWK Hochschule für Technik, Wirtschaft und Kultur Leipzig

ML Machine Learning

A Anhang - Abbildungen

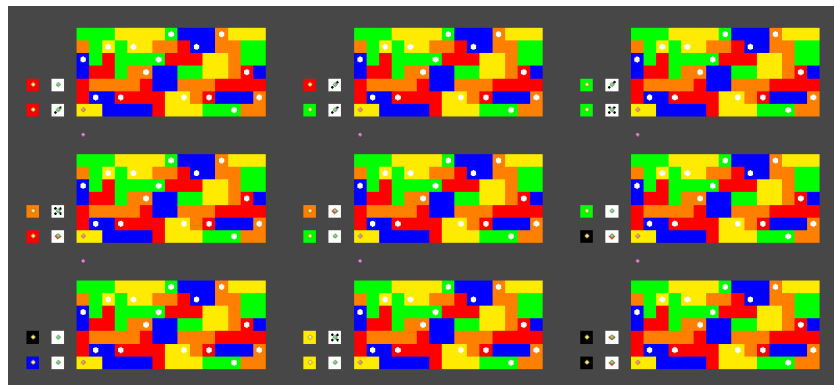


Abbildung A.1: Aufbau der Trainingsumgebung

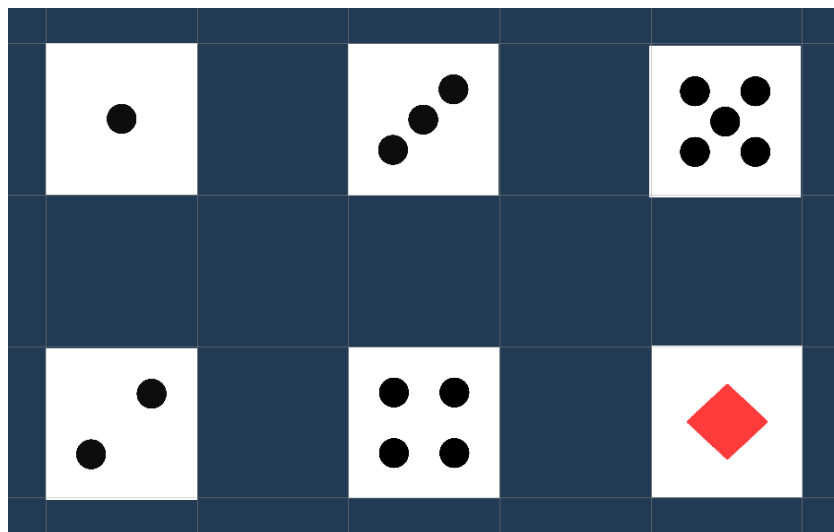


Abbildung A.2: Verwendete Sprites für den Zahlenwürfel

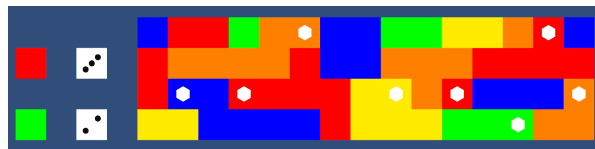


Abbildung A.3: Lernumgebung für einen Agenten auf dem 'MiniFeld'

B Anhang - Tabellen

C Anhang - Quelltexte

Quellcode C.1: Befüllen des Beobachtungsvektor mit validem Feld

```
1 private void PushValidField(Transform squareField, VectorSensor sensor) {  
2     sensor.AddObservation(GetColorIndexOneHotFromColor(squareField.GetComponent<FieldSquare>().color));  
3     sensor.AddObservation(squareField.GetComponent<FieldSquare>().available);  
4     sensor.AddObservation(squareField.GetComponent<FieldSquare>().starField);  
5     sensor.AddObservation(squareField.GetComponent<FieldSquare>().crossed);  
6 }
```

Quellcode C.2: Erhöhen der Warscheinlichkeit eine bestimmte Farbe zu würfeln

```
1 public void Roll() {  
2     List<string> colors = new List<string> { "red", "blue", "green", "yellow", "orange", "joker" };  
3     if(colorShift != ""){  
4         for (int i=0;i<10;i++){  
5             colors.Add(colorShift);  
6         }  
7     }  
8     int randomIndex = Random.Range(0, colors.Count);  
9     string colorResult = colors[randomIndex];  
10    this.color = colorResult;  
11    SetColor(color); // show dice  
12 }
```

Quellcode C.3: Erstellen des Beobachtungsvektors

```
1 public override void CollectObservations(VectorSensor sensor) {
2     sensor.AddObservation(GetJokersObservation());
3     sensor.AddObservation(GetRoundCountObservation());
4
5
6     foreach (Transform childTransform in controller.transform) {
7         if (childTransform.CompareTag("NumberDice")) {
8             GameObject child = childTransform.gameObject;
9             sensor.AddObservation(GetNumberDiceObservation(child));
10        }
11    }
12
13    foreach (Transform child in controller.transform) {
14        if (child.CompareTag("ColorDice")) {
15            sensor.AddObservation(GetColorIndexOneHotFromColor(child.GetComponent<ColorDice>().color));
16        }
17    }
18
19    int counter = 0;
20    foreach (Transform squareField in GameField.transform) {
21        if (squareField.CompareTag("Square")) {
22            PushValidField(squareField, sensor);
23            counter += 1;
24        }
25    }
26    do {
27        if (counter < 105){
28            PushEmptyField(sensor);
29            counter += 1;
30        }
31    }
32    while (counter < 105);
33 }
```

Quellcode C.4: Einlesen aller angegebenen Logdateien

```

1 import matplotlib.pyplot as plt
2 import numpy as np
3
4 def format_float(value):
5     return "{:.3f}".format(value)
6
7 file_paths = [
8     'LogPoints_trained.txt',
9     'LogPoints_trained_with_more_turns_2.txt',
10    'LogPoints_orange_field_untrained.txt',
11    'LogPoints_orange_field_trained2.txt',
12    'LogPoints_special_trained.txt',
13    'LogPoints_untrained.txt',
14    'miniField_trained.txt',
15    'OnlyDice.txt'
16 ]
17
18 filenames = [
19     'trained',
20     'more_turns',
21     'orange_field_untrained',
22     'orange_field_trained',
23     'special_trained',
24     'untrained',
25     'minifield_training',
26     'only_dice'
27 ]
28
29
30 pointLists = []
31 for file_path in file_paths:
32     points_list = []
33     with open(file_path, 'r') as file:
34         points = []
35         amount_of_games = 2500
36         for line in file:
37             if line.strip():
38                 points.append(float(line.strip()))
39                 if len(points) > amount_of_games:
40                     moving_average = sum(points[-amount_of_games:]) / amount_of_games
41                     points_list.append(moving_average)
42         pointLists.append(points_list)
43
44 min_length = min(len(points_list) for points_list in pointLists)
45 pointLists_clipped = [points_list[:min_length] for points_list in pointLists]

```

Quellcode C.5: Methode zum Anzeigen des Graphen

```
1
2 def display_selected_graphs(pointLists, filenames, selected_indices):
3     selected_pointLists = []
4     for index in selected_indices:
5         selected_pointLists.append(pointLists[index])
6
7
8     min_length = min(len(points_list) for points_list in selected_pointLists)
9     selected_pointLists_clipped = [points_list[:min_length] for points_list in selected_pointLists]
10
11     fig, ax = plt.subplots()
12
13     for i, points_list in zip(selected_indices, selected_pointLists_clipped):
14         x = np.arange(1, len(points_list) + 1)
15         ax.plot(x, points_list, label=f'{filenames[i]}')
16
17     ax.set_xlabel('Episode')
18     ax.set_ylabel('Average Points per Game')
19     ax.set_title(f'Average Points per Game over {amount_of_games} Games')
20     ax.legend()
21
22     plt.show()
23
24
25 selected_indices = [0,5] # Beispiel Trained vs Untrained
26 display_selected_graphs(pointLists, filenames, selected_indices)
```