

Technische Hochschule Ostwestfalen-Lippe
University of Applied Sciences and Arts

Wintersemester 2022/2023

Vergleich verschiedener Reinforcement Learning Algorithmen und deren Hyperparametern

Modularbeit

Vorgelegt im Kontext des Moduls „Anwendung des maschinellen Lernens“

am Fachbereich Technische Informatik und Elektrotechnik
im Studiengang Data Science

Veranstalter:	Prof. Dr. Burkhard Wrenger
Vorgelegt von:	Bjarne Seen Liebigstraße 130 32657 Lemgo bjarne.seen@stud.th-owl.de
Matr. Nr.:	15467085
Vorgelegt von:	Joshua Henjes Hanseweg 11 32657 Lemgo joshua.henjes@stud.th-owl.de
Matr. Nr.:	15467024
Abgabetermin:	01.03.2023

1. März 2023

Inhaltsverzeichnis

Abkürzungen	I
1 Einleitung	1
2 Grundlagen	2
2.1 Definitionen	2
2.2 Algorithmen	6
2.2.1 Q-Learning	6
2.2.2 SARSA	6
2.3 Hyperparameter	6
3 Implementierung	8
3.1 Probleme	8
3.2 Optimierung der Hyperparameter	10
3.3 Vergleichen von Algorithmen	11
4 Ergebnisse	12
4.1 Optimierung der Hyperparameter	12
4.1.1 Anzahl an Episoden	12
4.1.2 Maximale Anzahl an Steps pro Episode	14
4.1.3 Learning Rate	15
4.1.4 Discount Rate	15
4.1.5 Exploration Rate	16
4.2 Vergleichen von Algorithmen	17
4.2.1 Taxi	17
4.2.2 Cliff	18
4.2.3 Frozen Lake	19
5 Zusammenfassung und Ausblick	20
A Selbstständigkeitserklärung Joshua Henjes	21
B Selbstständigkeitserklärung Bjarne Seen	22

Abkürzungen

MDP	Markov Decision Process
ML	Machine Learning
RL	Reinforcement Learning
SARSA	State-Action-Reward-State-Action
TD	Temporal Difference

1 Einleitung

Reinforcement Learning ist neben Supervised und Unsupervised Learning eines der elementaren Felder des maschinellen Lernens. Im Gegensatz zu den anderen Feldern benötigt Reinforcement Learning keine Trainingsdaten, denn der Algorithmus lernt durch wiederholtes Interagieren mit einer dynamischen Umgebung eine Strategie, um eine Belohnungsmetrik zu maximieren. Es wird daher auch als bestärktes Lernen oder verstärktes Lernen bezeichnet.

Bekannt wurde das Reinforcement Learning vor allem durch das Meistern von bekannten Brett- und Computerspielen, so ist Googles „AlphaGo“ in der Lage, die besten Go Spieler der Welt zu schlagen [v7labsReinforcementLearning]. Trotz dieser beeindruckenden Erfolge findet RL in der Industrie bisher nur geringe Anwendung.

Immer kürzer werdende Produktzyklen und steigende Produktvielfalt stellen für die heutigen Produktionsprozesse eine große Herausforderung dar. Zukünftige Produktionen müssen immer anpassungsfähiger werden. Zeitgleich soll der Personalaufwand aufgrund des anhaltenden Fachkräftemangels möglichst gering ausfallen. Maschinelles Lernen, insbesondere das Reinforcement Learning, kann bei der Bewältigung dieser Herausforderungen eine relevante Rolle übernehmen.

Auch bei der Bekämpfung des Klimawandels kann Reinforcement Learning unterstützen. Um unsere Klimaziele zu erreichen, ohne unseren Lebensstandard signifikant zu senken, ist eine Optimierung des Ressourcenbedarfs nötig. Mit ausreichender Trainingszeit sind RL-Algorithmen sehr gut in der Optimierung von Prozessen und somit auch in dessen Ressourcenverbrauch. Google konnte als einer der Vorreiter im Gebiet des maschinellen Lernens durch ML-Algorithmen den Energieverbrauch der Kühlung ihrer Rechenzentren um bis zu 40 Prozent reduzieren [v7labsReinforcementLearning].

Mittlerweile existiert eine Vielzahl an unterschiedlichen RL-Algorithmen. Während die mathematischen und strukturellen Unterschiede meist gut dokumentiert und einsehbar sind, ist ein direkter Vergleich der Leistungsfähigkeit der Algorithmen in verschiedenen Umgebungen nur schwer zu finden. Aus diesem Grund beschäftigt sich diese Ausarbeitung mit dem Vergleich von beliebten RL-Algorithmen anhand von Umgebungen mit geringer Komplexität.

Während die Zeit, welche ein RL-Algorithmus zum Lernen benötigt, bei der Anwendung auf Brett- und Computerspielen eher eine untergeordnete Rolle spielt, ist sie in der Anwendung in industriellen Umgebungen deutlich relevanter. Zum einen verlangsamen hohe Trainingszeiten den Entwicklungsprozess, was wiederum zu höheren Lohn- und Entwicklungskosten führt. Zum anderen ist es in vielen Anwendungsfällen nötig, dass die Umgebung während des Trainingsprozesses dem Algorithmus zur Verfügung steht. Im Fall von Produktionsanlagen ist Trainingszeit somit sehr kostspielig. Aus diesem Grund wird neben der Leistungsfähigkeit auch die Lerngeschwindigkeit der Algorithmen im Folgenden untersucht.

2 Grundlagen

2.1 Definitionen

Um Reinforcement Learning im Folgenden besser beschreiben zu können, ist zunächst die Klärung einiger Grundbegriffe nötig. Diese sind aus der englischen Sprache entstanden. Auf eine Übersetzung dieser Begriffe in das Deutsche wurde verzichtet, um eine Vergleichbarkeit zu anderen Werken in diesem Themenbereich zu gewährleisten.

1. Markov Decision Process

Ein Markov Decision Process (MDP) [**builtinUnderstandingMarkov**] ist ein formales Modell, das zur Beschreibung von Entscheidungsproblemen verwendet wird, bei denen ein Entscheidungsträger (oft als Agent bezeichnet) in einer Umgebung handelt und dabei versucht, eine bestimmte Zielsetzung zu erreichen. Ein MDP basiert auf dem Konzept eines Markov-Prozesses, der ein stochastischer Prozess ist, bei dem der zukünftige Zustand nur vom gegenwärtigen Zustand abhängt und nicht von früheren Zuständen.

Ein MDP besteht aus den folgenden Komponenten:

(a) *Agent*

Der Agent [**mediumBeginnersGuide**] ist der Entscheidungsträger, welcher Aktionen in einem Szenario/Umfeld ausführt und dafür eine Belohnung bekommt.

(b) *Environment*

Das Environment [**datasolutReinforcementLearning**] ist, wie die deutsche Übersetzung schon vermuten lässt, die Umgebung in dem sich der Agent befindet. Das Environment legt dabei die grundlegenden Regeln fest und definiert, welche Aktionen möglich sind. Das Environment trägt somit ausschlaggebend zur Komplexität der zu lösenden Aufgabe bei. In vielen Fällen ist das Environment eine Simulation. Dies ermöglicht einen deutlich schnelleren Lernprozess, da jegliche Interaktion ohne nennenswerte Verzögerung ausgeführt werden kann. Bei komplexen Aufgabestellungen ist es so zudem möglich, mehrere Agenten parallel zu trainieren.

(c) *Action*

Als Action **A** wird eine Interaktion des Agent mit dem Environment beschrieben. Die Lösung eines Problems kann somit als Abfolge bestimmter Actions angesehen werden. Welche Actions der Agent ausführen kann, hängt dabei von den Grundregeln des Environments ab.

(d) *State*

Der State **S** ist die eindeutige und vollständige Beschreibung des Zustands, in welchem sich das Environment befindet. Aus technischer Sicht ist der State meist ein Vektor, eine Matrix oder ein Tensor, welcher alle relevanten Information des aktuellen Zustands enthält.

(e) *Reward*

Der Reward **R** [**datasolutReinforcementLearning**] ist die unmittelbare Belohnung, welche der Agent als Feedback zu einer Action erhält. In der Praxis ist dies ein numerischer Wert, welche entweder erhöht oder reduziert werden kann. Der Agent kann so für eine Action belohnt oder bestraft werden, dabei versucht er sein Handeln so auszurichten, dass er die größtmögliche Belohnung erzielt. Die Art und Weise, wie der Reward vergeben wird, bestimmt somit das Verhalten des Agents.

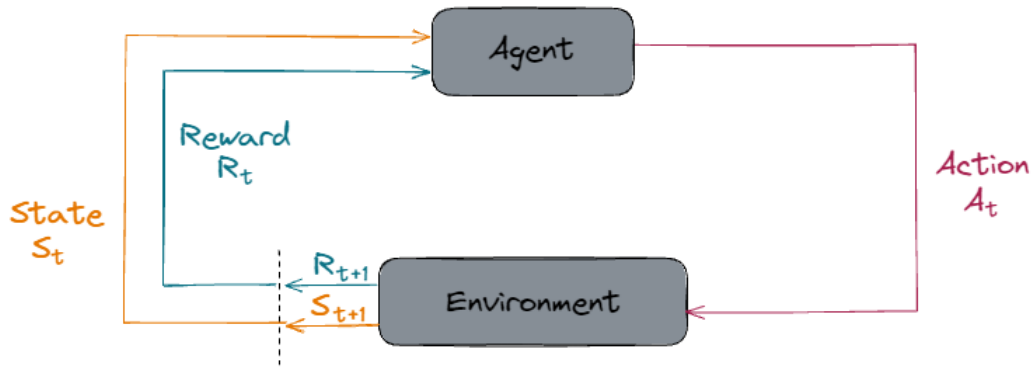


Abbildung 1: Markov Decision Process

Das Ziel des Agenten in einem MDP besteht darin, eine Strategie zu entwickeln, die ihm dabei hilft, die maximale kumulierte Belohnung im Laufe der Zeit zu erhalten. Eine Strategie ist eine Abbildung von Zuständen auf Actions, die angibt, welche Action der Agent in jedem Zustand ausführen sollte. Die optimale Strategie maximiert die erwartete zukünftige Belohnung über alle Zustände und Aktionen.

Der MDP besitzt, wie eben beschrieben, eine Menge an States \mathcal{S} , eine Menge an Actions \mathcal{A} und eine Menge an Rewards \mathcal{R} . In dem Prozess werden die Schritte $t = 0, 1, 2, \dots$ durchlaufen und der Agent befindet sich jeweils in einem State $\mathbf{S}_t \in \mathcal{S}$. Basierend auf diesem State kann der Agent eine Action $\mathbf{A}_t \in \mathcal{A}$ wählen. Dies ergibt dann das State-Action Paar $(\mathbf{S}_t, \mathbf{A}_t)$.

In dem nächsten Schritt $t + 1$ wird das Environment in den State $\mathbf{S}_{t+1} \in \mathcal{S}$ überführt. Hier bekommt der Agent nun den entsprechenden Reward $\mathbf{R}_{t+1} \in \mathcal{R}$ für die Action \mathbf{A}_t , welcher er zuvor in State \mathbf{S}_t genommen hat. Dieser Prozess ist in der Abbildung 1 abgebildet.

2. Episode

Als Episode [lesswrongWhatTraining] wird ein vollständiger Durchlauf während des Trainings bezeichnet. Jede Episode startet mit dem Anfangszustand des Environments und kann auf mehreren Wegen enden. Im besten Fall wird die Episode beendet, weil die gestellte Aufgabe vom Agent gelöst worden ist. In vielen Fällen wird eine Episode jedoch abgebrochen, weil die maximale Anzahl an Actions überschritten wurde. Eine solche Grenze wird implementiert, um sicherzustellen, dass das Training effektiv und effizient abläuft. Wenn keine Obergrenze festgelegt wird, kann der Agent endlos versuchen, das Ziel zu erreichen, ohne jemals erfolgreich zu sein. Zudem kann so verhindert werden, dass der Agent in einer Schleife von kleinen Belohnungen feststeckt. Die letzte Möglichkeit, wie eine Episode enden kann, ist vom Environment definiert. In vielen Fällen kann der Agent durch bestimmte Fehlentscheidungen die Episode beenden.

3. Policy

Im Reinforcement Learning bezeichnet eine Policy π [mediumReinforcementLearningPolicyValue] eine Funktion, die Entscheidungen trifft, um eine bestimmte Aufgabe zu lösen. Eine Policy entscheidet, welche Action ein Agent in einem bestimmten State ausführen soll, basierend auf den Informationen, die der Agent in der Vergangenheit gesammelt hat.

Die Policy wird durch das Optimierungsproblem des Reinforcement Learning bestimmt, das darin besteht, die optimale Strategie zu finden, um die Belohnung des Agents zu maximieren. Die optimale Policy ist diejenige, die in jedem State die Action empfiehlt, die die höchste erwartete Belohnung ergibt.

4. Value Function

Bei dem Begriff Value Function [**mediumReinforcementLearningPolicyValue**] handelt es sich um eine Funktion, die den erwarteten Wert einer bestimmten State-Action-Kombination oder eines States wiedergibt. Der Wert gibt an, wie nützlich es ist, sich in diesem State oder dieser Kombination zu befinden, um das Gesamtziel zu erreichen, also die maximale Belohnung zu erhalten.

Die Value Function kann verwendet werden, um die optimale Policy zu finden. Im Reinforcement Learning gibt es zwei Arten von Value Functions: Die State-Value Function und die Action-Value Function.

Die State-Value Function V gibt den erwarteten Wert der Gesamtbelohnungen an, den ein Agent in einem bestimmten State erzielen kann. Mit anderen Worten, sie gibt an, wie nützlich es ist, sich in einem bestimmten State zu befinden, um das Ziel der maximalen Belohnung zu erreichen. Die State-Value Function wird oft als $V(s)$ bezeichnet, wobei s der State ist.

Die Action-Value Function Q gibt den erwarteten Wert der Gesamtbelohnungen an, den ein Agent in einem bestimmten State erreichen kann, wenn er eine bestimmte Action ausführt. Die Action-Value Function wird oft als $Q(s,a)$ bezeichnet, wobei s der State und a die Action sind.

Value Functions können auf verschiedene Weise geschätzt werden, wie z.B. mit Hilfe von Monte-Carlo-Methoden und Temporal Difference Learning.

Im weiteren Verlauf der Arbeit wird sich mit der Action-Value Function Q und dem Temporal Difference Learning auseinandergesetzt.

5. Temporal Difference Learning

Temporal Difference (TD) Learning [**mediumTemporalDifference**] ist eine Methode des Reinforcement Learnings, die es einem Agent ermöglicht, aus Erfahrungen zu lernen, indem er die erzielte Belohnung mit der erwarteten Belohnung vergleicht. Im Gegensatz zu Monte-Carlo-Methoden, die die Gesamtbelohnungen aus der Erfahrung berechnen, werden bei TD-Learning die Value Functions schrittweise durch den Vergleich von aufeinanderfolgenden Schätzungen aktualisiert.

Die TD-Learning Methode verwendet dabei die Bellman Equation, um die Value Functions zu aktualisieren. Diese Reihe von Gleichungen beschreibt die Beziehung zwischen den Zustands- und Aktionswerten und der optimalen Policy.

Wenn der Agent eine Action ausführt und in einen neuen Zustand gelangt, wird die erhaltene Belohnung zusammen mit dem geschätzten zukünftigen Wert des nächsten Zustands verwendet, um eine neue Schätzung der Value Function zu berechnen. Diese Schätzung wird dann mit der Vorherigen kombiniert, um die Value Function schrittweise zu aktualisieren.

Im TD-Learning gibt es zwei wichtige Methoden: SARSA (State-Action-Reward-State-Action) und Q-Learning.

6. Optimalität

(a) *Optimale Policy*

Das Ziel beim Reinforcement Learning ist es, eine Policy π zu finden, welche die Belohnung des Agents maximiert. Dafür muss die Policy π gefunden werden, welche dem Agenten mehr Belohnung liefert, als jede andere Policy π' .

Mathematisch formuliert muss folgendes gelten:

$$\pi \geq \pi' \text{ wenn } q_\pi(s, a) \geq q_{\pi'}(s, a) \text{ für alle } s \in S \text{ und } a \in A \quad (1)$$

Hierbei ist $q_\pi(s, a)$ die Action-Value Function, welche der optimalen Policy π folgt.

(b) *Optimale Action-Value Function*

Für die optimale Action-Value Function gilt folgende Gleichung q_* :

$$q_*(s, a) = \max_{\pi} q_\pi(s, a) \quad (2)$$

Somit liefert q_* die größten erwarteten Belohnungen für jede Policy π für jedes mögliche State-Action Paar.

(c) *Bellman Optimality Equation*

Die Bellman Optimality Equation [floydhubIntroductionQLearning] für die Action-Value Function beschreibt den optimalen Wert eines State-Action Paares (s, a) unter der Annahme, dass der Agent die optimale Policy verfolgt.

$$q_*(s, a) = E[R_{t+1} + \gamma \max_{a'} q_*(s', a')] \quad (3)$$

Sie kann genutzt werden um in einem iterativen Prozess die optimale Action-Value Function q_* zu finden.

7. Exploration vs Exploitation

Exploration [towardsdatascienceExplorationReinforcement] bezieht sich auf die Strategie des Agents, neue Entscheidungen zu erkunden und zu lernen, indem er Actions ausprobiert, die er noch nicht ausprobiert hat. Der Agent versucht, neue Möglichkeiten zu erforschen, um mehr über das Environment zu erfahren.

Exploitation hingegen bezieht sich auf die Strategie des Agents, Entscheidungen auf der Grundlage der bisherigen Erfahrungen zu treffen, um den maximalen Reward zu erzielen. Der Agent nutzt die bisherigen Erfahrungen, um die Action auszuwählen, die ihm die höchste Belohnung in der Vergangenheit gegeben hat.

Ein Gleichgewicht zwischen Exploration und Exploitation ist wichtig, um optimale Entscheidungen zu treffen. Wenn der Agent zu sehr auf Exploration setzt, wird er immer wieder neue Entscheidungen ausprobieren und keine optimale Entscheidung treffen. Wenn der Agent hingegen zu sehr auf Exploitation setzt, wird er sich auf die Entscheidungen konzentrieren, die ihm in der Vergangenheit den höchsten Reward gebracht haben. Dementsprechend wird er nicht in der Lage sein, neue Lösungswege zu erkunden, die möglicherweise noch höhere Rewards erzielen können.

Daher verwenden Reinforcement Learning-Algorithmen oft eine Strategie namens epsilon-greedy-Strategie. Bei dieser wählt der Agent mit einer Wahrscheinlichkeit von epsilon eine zufällige Action aus, um Exploration durchzuführen, und mit einer Wahrscheinlichkeit von $1 - \epsilon$ eine Action aus, die bisher die höchste Belohnung erzielt hat. Auf diese Weise kann der Agent gleichzeitig erkunden und von seinen bisherigen Erfahrungen profitieren, um optimale Entscheidungen zu treffen.

2.2 Algorithmen

2.2.1 Q-Learning

Q-Learning ist ein iteratives, modellfreies Reinforcement Learning Verfahren aus dem Bereich der Temporal Difference Learning Verfahren. Dieses wird verwendet, um die optimale Action-Value Function $Q_*(s, a)$ zu lernen, indem es Erfahrungen sammelt und die Action-Values iterativ aktualisiert. Das Verfahren basiert auf der Bellman Optimality Equation für die Action-Value Function $Q_*(s, a)$, die besagt, dass der optimale Wert eines State-Action Paares (s, a) die maximale erwartete zukünftige Belohnung ist, die durch die Wahl der optimalen Actions in diesem State und danach unter Verwendung der optimalen Policy erreicht wird.

Q-Learning nutzt eine Tabelle, die die Werte für jedes State-Action Paar (s, a) speichert. Das Verfahren sammelt Erfahrungen, indem es den Agenten in der Umgebung handeln lässt, und verwendet diese Erfahrungen, um die Werte zu aktualisieren. Die Aktualisierung erfolgt durch die Formel:

$$Q(s, a) \leftarrow Q(s, a) + \alpha[R + \gamma \max_{a'} Q(s', a') - Q(s, a)] \quad (4)$$

wobei $Q(s, a)$ der aktuelle Wert für das State-Action Paar (s, a) ist, α die Learning Rate ist, R die unmittelbare Belohnung ist, γ die Discount Rate ist, die die zukünftige Gesamtbelohnung gewichtet, und $\max_{a'} Q(s', a')$ die erwartete zukünftige Gesamtbelohnung ist, die durch Wahl der optimalen Action a' im nächsten State s' erreicht wird.

2.2.2 SARSA

Ähnlich wie bei Q-Learning basiert SARSA auf der Bellman Optimality Equation. Im Gegensatz zu Q-Learning lernt SARSA direkt die Action-Values der Policy, die tatsächlich ausgeführt wird. Dies bedeutet, dass SARSA bei der Aktualisierung der Action-Values die nächste Action a' basierend auf der aktuellen Policy wählt, anstatt die Action mit dem höchsten Action-Value zu wählen, wie dies bei Q-Learning der Fall ist. Daher spricht man bei Q-Learning auch von einem off-Policy Algorithmus, während SARSA ein on-Policy Algorithmus ist.

Die Aktualisierungsregel von SARSA ist wie folgt:

$$Q(s, a) \leftarrow Q(s, a) + \alpha[R + \gamma Q(s', a') - Q(s, a)] \quad (5)$$

wobei der Unterschied in dem Term $\gamma Q(s', a')$ liegt. Denn hierbei wird wieder das epsilon-greedy Verfahren benutzt um den Action-Value zu bestimmen, im Vergleich zu dem maximalen Wert, welcher bei Q-Learning gewählt wurde.

2.3 Hyperparameter

Unter Hyperparametern versteht man Einstellungen oder Konfigurationen für einen Machine-Learning-Algorithmus, um das Verhalten und die Leistung zu steuern. Diese werden vor Beginn des Trainings festgelegt. Da die Wahl der Hyperparameter einen großen Einfluss auf die Fähigkeit der Modelle hat, ist es für den späteren Vergleich wichtig diese Parameter für jeden Algorithmus zu optimieren, um einen fairen Vergleich zu gewährleisten.

- **Anzahl an Episoden**

Wie bereits bei 2.1 [Definitionen](#) beschrieben, stellt eine Episode einen vollständigen Durchlauf des Environments dar. Die Anzahl dieser Durchläufe während des Trainingsprozesses wird als

Hyperparameter festgelegt. Umso mehr Durchläufe der Agent zur Verfügung hat, um so häufiger kann er dazulernen und sich anpassen. Somit hat dieser Hyperparameter einen großen Einfluss auf die Leistung des Agent. Die Anzahl der Episoden steht zudem im direkten Zusammenhang zu der für das Training benötigten Zeit. Umso mehr Episoden durchlaufen werden, desto länger dauert der gesamte Trainingsprozess. Die optimale Anzahl der Episoden hängt von verschiedenen Faktoren wie der Komplexität der Umgebung, der Anzahl der States sowie dem verwendeten Lernalgorithmus ab. Bei zu wenigen Episoden ist der Algorithmus nicht in der Lage sein volles Potenzial zu entfalten. Zu viele Episoden können, besonders bei komplexen Modellen, zu Overfitting führen, dabei lernt der Agent die Umgebung auswendig und kann so nicht auf neue Situationen reagieren.

- **Discount Rate**

Die Discount Rate bestimmt, wie stark der Agent eine Belohnung in der Zukunft gewichtet im Vergleich zu einer sofortigen Belohnung. Ein höhere Discount Rate bedeutet, dass der Agent die zukünftigen Belohnungen stärker gewichtet und daher eher langfristige Entscheidungen trifft. Eine niedrigere Discount Rate führt dagegen zu kurzfristigeren Entscheidungen, da zukünftige Belohnungen weniger stark gewichtet werden. Der Wertebereich für diesen Hyperparameter liegt zwischen null und eins. Im ersten Moment könnte man jetzt denken, dass die höchstmögliche Discount Rate die beste sei, da so der komplette Fokus auf dem Endergebnis der Episode liegt. Dies ist aber in den meisten Fällen nicht korrekt. Denn nicht jede Entscheidung hat einen Effekt, der bis zum Ende der Episode reicht. Der Agent wäre somit nicht in der Lage solche Entscheidungen zu treffen. Zudem würden bei jeder Entscheidung eine enorme Menge an irrelevanten Informationen mitberücksichtigt, welche den Lernprozess enorm erschweren. Wie auch die maximale Anzahl an Episoden, hängt auch der optimale Wert für diesen Hyperparameter von der Komplexität der Problemstellung und dem verwendeten Lernalgorithmus ab. Zusätzlich ist die zeitliche Entfernung zwischen den Actions und den Rewards relevant.

- **Learning Rate**

Die Learning Rate gibt an, wie stark das Verhalten des Agents nach jeder Episode aufgrund der Geschehnisse in der Episode angepasst wird. Bei einer hohen Learning Rate werden die Gewichte des Agent stark verändert, und somit das Verhalten stark angepasst. Auf der einen Seite ist das gut, um das Training zu beschleunigen, und so schneller zum Ergebnis zu kommen. Gleichzeitig führt eine zu hohe Learning Rate dazu, dass das Modell instabil wird und keine Konvergenz erreicht und somit nie das Optimum erzielt. Eine zu niedrige Learning Rate führt dagegen dazu, dass das Modell zu langsam lernt und möglicherweise in einem lokalen Minimum stecken bleibt. Daraus folgt, dass das Modell eine sehr lange Zeit zum Konvergieren benötigt oder möglicherweise nie die optimale Leistung erreicht.

- **Exploration Rate**

Die Exploration Rate wurde bereits in dem Abschnitt [Definitionen](#) erwähnt und entspricht dem dort beschriebenen epsilon. In der Praxis wird die Exploration Rate häufig über den Zeitraum des Trainings verringert. Dies ist sinnvoll, da der Agent zu Beginn des Trainings das Environment noch gar nicht kennt und somit eine hohe Explorationsrate notwendig ist, damit er die vielen unterschiedlichen States kennenlernen kann. Im späteren Verlauf des Trainings kennt der Agent bereits die meisten States und wählt nur zu einer geringen Wahrscheinlichkeit andere Aktionen, welche womöglich besser sein könnten.

3 Implementierung

3.1 Probleme

1. Taxi

Für das Taxiproblem [faramaGymnasiumDocumentation] wurde ein 5 x 5 großes Gitternetz implementiert. Der Agent (Taxifahrer) hat dabei die Möglichkeit, sich in alle vier Himmelsrichtungen zu bewegen. Die Aufgabe besteht darin, den Passagier, welcher sich in einer zufälligen Ecke des Gitternetzes befindet, abzuholen und ihn dann zu seinem gewünschten Ziel zu bringen. Dieses ist immer einer der verbleibenden Ecken des Spielfeldes. Um diese Aufgabe noch etwas zu erschweren, wurden zusätzlich noch Wände in das Gitternetz integriert. Diese befinden sich an den gleichen Positionen und verhindern Bewegungen in bestimmte Richtungen.

Für das Bewältigen dieser Aufgabe kann der Agent zu jedem Zeitpunkt zwischen sechs Actions entscheiden; Bewegung in jeweils einer der vier Himmelsrichtungen, Aufnehmen des Passagiers und das Absetzen des Passagiers. Natürlich sind nicht alle Actions zu jedem Zeitpunkt sinnvoll.

Der Zustand des Gitternetzes kann durch 500 diskrete States beschrieben werden. Diese Anzahl ergibt sich aus der Multiplikation der 25 möglichen Position des Taxis, der fünf möglichen Positionen des Passagiers (beinhaltet den Fall, dass sich der Passagier im Taxi befindet), mit den vier möglichen Zielorten.

Da das Abliefern des Passagiers auf dem richtigen Feld das Ziel der Aufgabe ist, bekommt der Agent dafür auch die höchste Belohnung. Damit die Erfüllung der Aufgabe so effizient wie möglich durchgeführt wird, bekommt der Agent für jede Action, die er durchführt und die nicht zu einer Belohnung führt, eine kleine Strafe. Zudem wird der Agent stark bestraft, wenn er den Passagier an einem falschen Ort absetzt oder versucht, einen Passagier an einer falschen Position aufzunehmen.

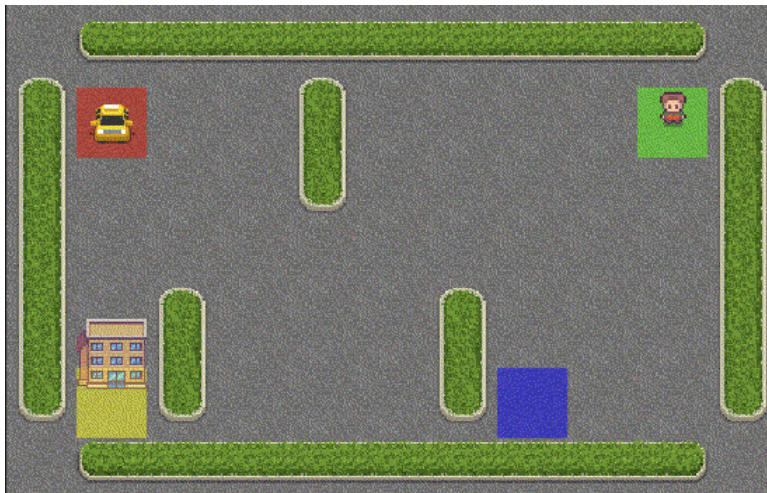


Abbildung 2: Taxi Environment

2. Cliff

Wie für das Taxiproblem wurde auch für das Cliff Problem [faramaGymnasiumDocumentation] ein Gitternetz implementiert. Dieses besitzt eine Größe von 4×12 . Die Aufgabe des Agent besteht darin, von einer Seite des Gitternetzes zur anderen zu gelangen, ohne dabei in die auf dem Feld befindliche Klippe zu fallen. Als Klippe wurden alle Felder definiert, welche sich auf dem direkten Weg zwischen Agent und Ziel befinden. Der Agent hat somit die Aufgabe, um diese Klippe herum zu navigieren und so das Ziel zu erreichen. Um die Herangehensweise der verschiedenen Algorithmen besser vergleichen zu können, befinden sich bei diesem Experiment alle Objekte (Agent, Klippe, Zielpunkt) zu Beginn jeder Episode an derselben Position.

Damit sich der Agent auf dem Feld bewegen kann, stehen ihm vier verschiedene Actions zur Verfügung, welche den Agent jeweils um ein Feld in einer der Himmelsrichtungen verschiebt.

Um den Zustand des Environments zu beschreiben, ist die aktuelle Position des Agent ausreichend. Insgesamt gibt es 48 (4×12) verschiedene Positionen auf dem Feld. Da das Betreten der Klippenfelder jedoch zum Ende der Episode führt, sind diese Positionen kein gültiger State. Gleiches gilt auch für die Zielposition. Bei zehn Klippen und einem Ziel ergeben sich so 37 States.

Neben dem Erreichen des Ziels ist es besonders wichtig, dass der Agent nicht in die Klippe fällt. Daher ist dies mit einer hohen Bestrafung für den Agent versehen. Zudem soll das Ziel so schnell wie möglich erreicht werden. Aus diesem Grund ist, wie auch bei dem Taxi Problem, jede Action mit einer kleinen Strafe belegt. Eine explizite Belohnung des Agent ist für diesen Anwendungsfall nicht nötig, da das Ziel zum Ende der Episode führt und das beste Ergebnis somit jenes ist, welches zur geringsten Bestrafung führt.

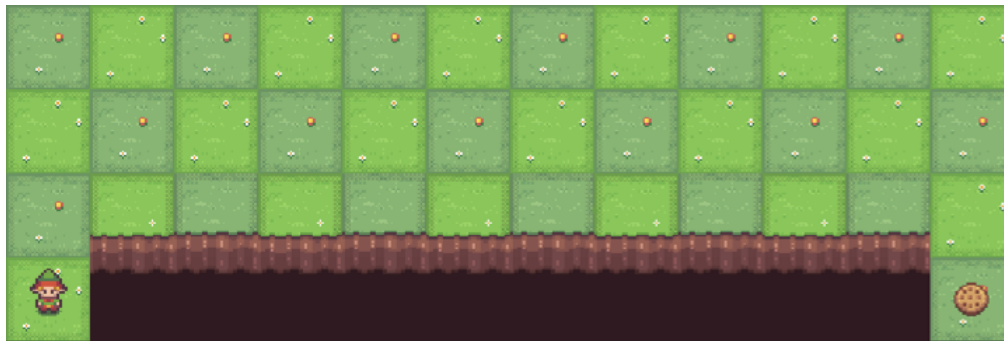


Abbildung 3: Cliff Environment

3. Frozen Lake

Auch das Frozen Lake Problem [faramaGymnasiumDocumentation] ist auf einem Gitternetz aus Feldern implementiert. Einige Felder sind Eisplatten, die der Agent sicher betreten kann, während andere Felder Löcher darstellen, in die der Agent fallen kann und damit das Spiel verliert. Das Ziel des Agent besteht darin, sicher zum Ziel zu gelangen, welches sich auf der anderen Seite des Sees (Gitternetz) befindet. Die besondere Herausforderung bei diesem Problem sind die Eisplatten. Bewegt sich der Agent auf einer dieser Platten in eine bestimmte Richtung besteht eine Chance, dass er ausrutscht und sich in eine andere Richtung bewegt.

Wie auch beim Cliff Problem kann der Agent nur durch Bewegung mit dem Environment interagieren und hat somit vier Actions zur Verfügung.

Ein 4 x 4 großes Environment mit vier Löchern hat somit 11 mögliche States (16 Felder minus die Löcher und des Ziels)

Die Struktur des Rewards ist für dieses Problem simpel. Der Agent wird belohnt, wenn er das Ziel erreicht und das Spiel wird beendet, wenn er in ein Loch fällt.



Abbildung 4: Frozen Lake Environment

3.2 Optimierung der Hyperparameter

Die theoretischen Grundlagen der Hyperparameter wurden bereits in 2.3 [Hyperparameter](#) erläutert. Um konkrete Werte für die im Rahmen dieser Arbeit untersuchten Probleme zu finden, wurden einige Versuche durchgeführt. Als Ausgangslage wurde die in Tabelle 1 dargestellten Werte angenommen. Im Folgenden werden diese initialen Werte anhand von Experimenten für die verschiedenen Algorithmen und Probleme optimiert.

Tabelle 1: Hyperparameter Ausgangslage

Paramater	Wert
<i>num_episodes</i>	1000
<i>max_steps_per_episode</i>	1000
<i>learning_rate</i>	0.1
<i>discount_rate</i>	0.99
<i>exploration_rate</i>	1
<i>max_exploration_rate</i>	1
<i>min_exploration_rate</i>	0.01
<i>exploration_decay_rate</i>	0.005

- Anzahl an Episoden

Als optimale Anzahl an Episoden wird der Punkt gewählt, ab wann sich der Agent nicht mehr verbessert. So wird ein zu langes Training vermieden. Dieser Punkt lässt sich sehr gut in einer grafischen Darstellung der erreichten Belohnung des Agent im Vergleich zu den während des Trainings durchlaufenen Episoden ablesen. Die Rohdaten aus dem Training sind aufgrund der Exploration sehr verrauscht und der eigentliche Trainingsfortschritt ist nur schwer zu erkennen. Um die Grafik lesbarer zu machen wurde ein gleitender Durchschnitt auf die Daten angewandt. Der Effekt dieser Operation wurde in der Abbildung 5 veranschaulicht.

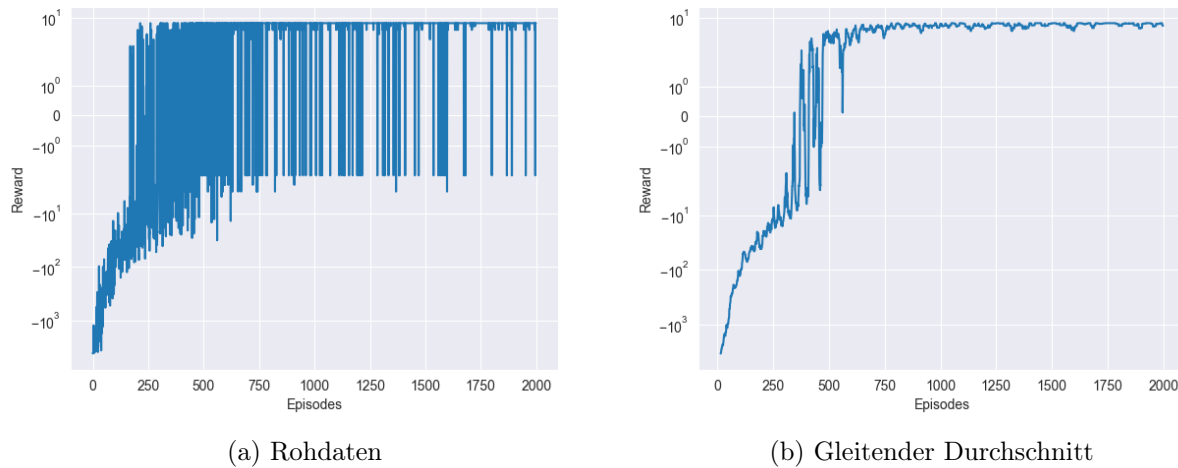


Abbildung 5: Anwendung gleitender Durchschnitt

- **Maximale Anzahl an Steps pro Episode**
Für die Festlegung dieses Hyperparameters wurde zunächst der Trainingsverlauf des initialen Wertes mit dem Trainingsverlauf von deutlich erhöhten und reduzierten Werten verglichen. So wurde der Einfluss des Parameters analysiert und eine erste Tendenz festgestellt. Darauf wurde weitere Vergleiche durchgeführt, um den optimalen Wert zu bestimmen.
- **Learning Rate**
Wie auch bei der Festlegung der maximalen Anzahl an Steps pro Episode wurde auch für die Learning Rate der initiale Wert verändert und die Auswirkungen anhand von Graphen analysiert.
- **Discount Rate**
Die Discount Rate wurde ebenfalls ermittelt, indem die Auswirkungen einer Veränderung des Hyperparameters untersucht wurden. Der initial gewählte Wert für die Discount Rate ist sehr nah am Maximum des möglichen Wertebereichs. Dementsprechend wurde ausschließlich eine Reduzierung des Parameters untersucht.
- **Exploration Rate**
Damit sich der Wert der Exploration Rate über den Verlauf des Trainings verändert, sind insgesamt vier Parameter implementiert. Die Obergrenze für den Wert, sowie der Startwert, sind initial auf Eins festgelegt. Zu Beginn hat der Agent noch kein Wissen über das Environment, eine Exploitation ist somit nicht sinnvoll. Eine Reduzierung des Wertes führt somit nur zu einem langsamen Trainingsverlauf. Die Abnahmerate und der minimale Wert sind durch Experimente, wie die vorherigen drei Hyperparameter, ermittelt worden.

3.3 Vergleichen von Algorithmen

Nach dem im vorangegangenen Kapitel die Implementation des Vergleichens der Hyperparameter beschrieben wurde, befasst sich dieser Abschnitt mit dem Vergleichen der beiden Algorithmen Q-Learning und SARSA. Um den bestmöglichen Vergleich zu gewährleisten wurden geeignete Parameter für beide Algorithmen und die jeweiligen Probleme gewählt. Diese Parameter sind den Tabellen 2

bis 4 aufgeführt. Bei der Analyse wurden alle drei Probleme, welche in 3.1 bereits erläutert wurden, berücksichtigt.

Tabelle 2: Taxi Paramater

Paramater	Q-Learning	SARSA	Paramater	Q-Learning	SARSA
<i>num_episodes</i>	3000	3000	<i>num_episodes</i>	1500	1500
<i>max_steps_per_episode</i>	2000	2000	<i>max_steps_per_episode</i>	750	1000
<i>learning_rate</i>	0.4	0.1	<i>learning_rate</i>	0.4	0.05
<i>discount_rate</i>	0.6	0.99	<i>discount_rate</i>	0.1	0.99
<i>exploration_rate</i>	1	1	<i>exploration_rate</i>	1	1
<i>max_exploration_rate</i>	1	1	<i>max_exploration_rate</i>	1	1
<i>min_exploration_rate</i>	0.01	0.01	<i>min_exploration_rate</i>	0.01	0.01
<i>exploration_decay_rate</i>	0.01	0.1	<i>exploration_decay_rate</i>	0.01	0.1

Tabelle 3: Cliff Paramater

Tabelle 4: Frozen Lake Paramater

Paramater	Q-Learning	SARSA
<i>num_episodes</i>	2000	2000
<i>max_steps_per_episode</i>	500	2000
<i>learning_rate</i>	0.4	0.3
<i>discount_rate</i>	0.2	0.99
<i>exploration_rate</i>	1	1
<i>max_exploration_rate</i>	1	1
<i>min_exploration_rate</i>	0.01	0.01
<i>exploration_decay_rate</i>	0.005	0.1

Zum Vergleich der Algorithmen wurden die Rewards und Steps während des Trainings aufgenommen und anschließend in Plots anschaulich dargestellt. Dabei wurde je nach Problem eine logarithmische Skalar auf der y-Achse genutzt, um die Daten anschaulicher zu gestalten. Zusätzlich wird ein Rolling Window Ansatz gewählt, sodass immer 15 Werte zusammengefasst werden. Dies sorgt dafür das einzelne Peaks zwar weniger auffallen, dafür der Gesamttrend im Verlauf der Trainingsepisoden aber besser zu erkennen ist.

4 Ergebnisse

4.1 Optimierung der Hyperparameter

4.1.1 Anzahl an Episoden

- Q-Learning

Der in Abbildung 6a dargestellte Graph zeigt den Trainingsverlauf von Q-Learning bei dem Taxiproblem. Dort lässt sich erkennen, dass sich der Reward ab Episode 1300 kaum noch verändert. Der Agent lernt dementsprechend nicht mehr dazu. Dieser Punkt wird als optimaler Wert für die maximale Anzahl an Episode angenommen.

Um den optimalen Wert für das Cliff und Frozen Lake Problem zu ermitteln, wurden dieselben Schritte durchgeführt. Das Cliff Problem wird schneller von dem Agent erlernt, daher kann das Training schon nach etwa 1000 Episoden beendet werden. Der Trainingsverlauf des Frozen

Lake Problem hat aufgrund der besonderen Rewardstruktur einen anderen Wertebereich (vgl. Abbildung 6b). Ein Wert von etwa 750 benötigten Episoden lässt sich trotzdem ablesen.

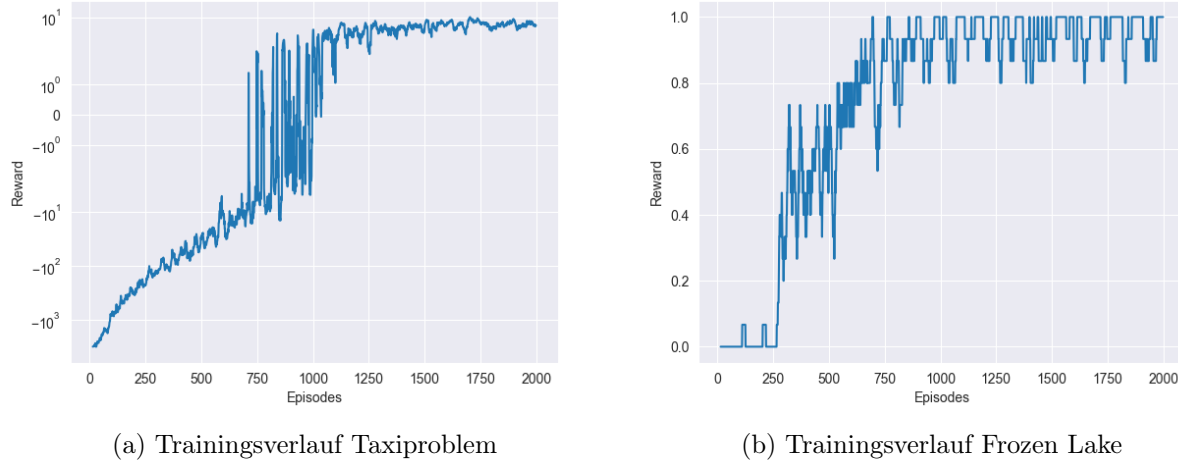


Abbildung 6: Ermittlung der optimalen Anzahl an Episoden

- SARSA

Das Verhalten von SARSA auf dem Taxiproblem ist sehr identisch zu dem von Q-Learning. Daher ist hier auch ein Wert von 1300 optimal. Bei dem Cliff Problem konvergiert der Reward ab 600 Epochen und beim FrozenLake Problem erzielt der Agent ab 2000 Epochen keinen Fortschritt mehr. In Abbildung 7 ist der Trainingsverlauf von dem Cliff Problem (Blau) und dem FrozenLake Problem (Orange) dargestellt.

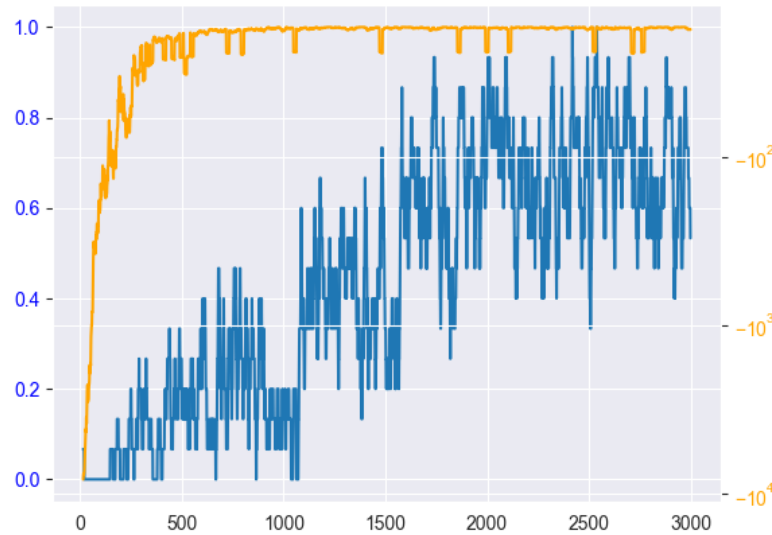
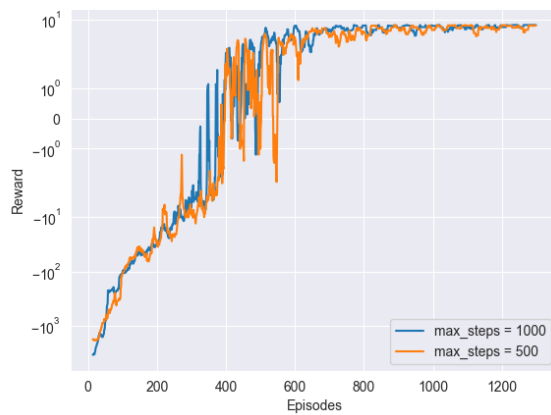


Abbildung 7: SARSA Trainingsverlauf

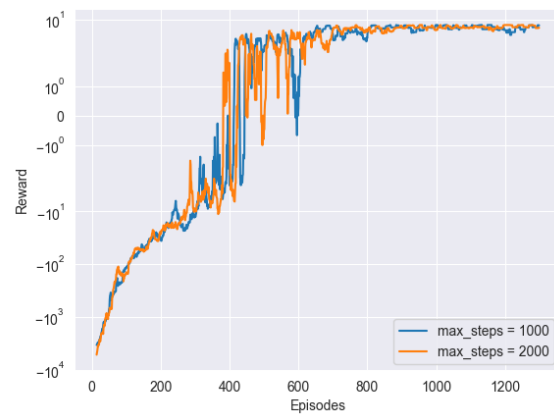
4.1.2 Maximale Anzahl an Steps pro Episode

- Q-Learning

Abbildung 8a stellt den Vergleich der Trainingsverläufe von dem Taxiproblem zwischen den anfänglichen 1000 Steps und 500 Steps dar. Diese deutliche Reduzierung der maximalen Steps hat zu Folge, dass der Agent länger ein unstabiles Verhalten aufzeigt und diese Veränderung nicht sinnvoll ist. Abbildung 8b zeigt den Effekt einer Verdopplung des Wertes auf 2000. Zu erkennen ist, dass der Agent in der Mitte des Trainings schneller lernt, diesen Vorsprung verliert er gegen Ende des Trainings jedoch wieder. Bei genauerer Betrachtung lässt sich zu dem feststellen, dass der Verlauf gegen Ende etwas weniger schwankt und somit der Agent minimal bessere Performance nach dem Training aufweist. Eine weitere Erhöhung der maximalen Anzahl an Steps führte in den Versuchen allerdings nicht zu einer Verbesserung des Ergebnisses, verursacht jedoch höhere Trainingszeiten. 2000 wurde somit als optimaler Wert des Hyperparameters für das Taxi-Problem ermittelt.



(a) Vergleich 1000 und 500



(b) Vergleich 1000 und 2000

Abbildung 8: Auswirkung der maximalen Anzahl an Episode

Für das Cliff Problem kann mit einem Wert von 750 schon das beste Ergebnis erzielt werden. Dies hängt mit der durchschnittlich niedrigeren Anzahl an Steps, die zum Lösen des Problems benötigt werden, zusammen. Bei der Untersuchung des Frozen Lake Problems konnte kein signifikanter Einfluss des Hyperparameters festgestellt werden. Um die Trainingszeit gering zu halten, wurde 500 als Wert festgelegt.

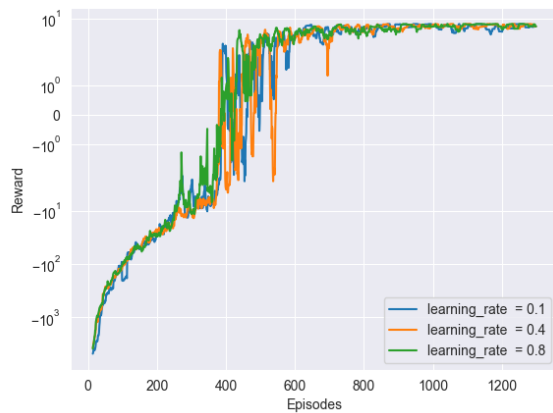
- SARSA

Die Auswirkung der maximalen Steps auf SARSA bei dem Taxiproblem sind vergleichbar mit Q-Learning. Verschiedene Versuche haben gezeigt, dass 2000 ein guter Wert für diesen Hyperparameter darstellt. Auf das Cliff Problem hat dieser Parameter kaum einen Einfluss, sodass der initiale Wert beibehalten werden kann. Bei der Lösung des Frozen Lake Problems verbessert die Erhöhung der maximalen Steps das Ergebnis. Als optimaler Wert wurde dabei erneut 2000 gewählt.

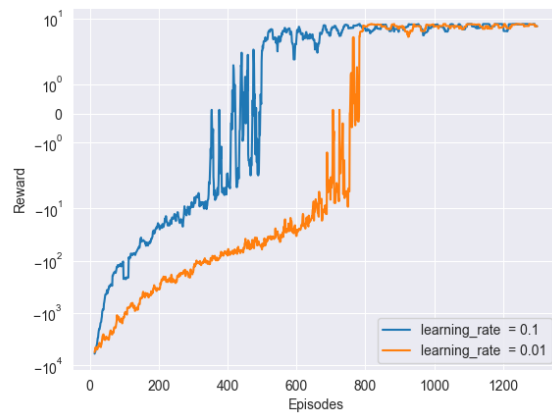
4.1.3 Learning Rate

- Q-Learning

Bei dem Vergleich verschiedener Learning Rates konnte festgestellt werden, dass Werte im Wertebereich zwischen 0.1 und 0.8 keinen Einfluss auf das Ergebnis des Trainings haben (vgl. Abbildung 9a). Es kommt zwar zu kleinen Unterschieden in der Mitte des Trainings, diese sind aber eher auf die zufällige Exploration zurückzuführen. Eine zu kleine Learning Rate, wie Abbildung 9b mit einem Wert von 0.01 veranschaulicht, führt jedoch zu einem längeren Trainingsverlauf. Für die weiteren Versuche wurde ein Wert von 0.4 festgelegt, da er sich in der Mitte des untersuchten Intervalls befindet.



(a) Einfluss Learning Rate



(b) Wahl einer zu geringen Learning Rate

Abbildung 9: Auswirkung der Learning Rate

- SARSA

Das Anheben der Learning Rate bei dem Taxiproblem führte bei SARSA sehr schnell zu einem unstabilen Verhalten. Eine Reduktion verlangsamte das Training. Somit führte keine Veränderung des Parameters zu einem bessern Ergebnis und der initiale Wert wird beibehalten. Bei dem Cliff Problem führt das Anheben erneut zu unstabilem Verhalten und die Reduzierung verbesserte das Ergebnis leicht. Somit wurde ein Wert von 0.05 gewählt. Das Frozen Lake Problem konnte mit einer leicht erhöhten Learning Rate von SARSA besser gelöst werden und es wurde ein Wert von 0.3 festgesetzt.

4.1.4 Discount Rate

- Q-Learning

Eine Reduzierung der Discount Rate auf 0.6 führte bei Q-Learning und dem Taxiproblem zu einem konstanteren Trainingsverlauf (vgl. Abbildung 10a). Die Auswirkung einer weiteren Reduzierung des Parameters auf 0.2 ist Abbildung 10b zu entnehmen und führte zur leichten Verschlechterung des Ergebnisses. Da der Wert von 0.6 die besten Ergebnisse lieferte, wird dieser Wert für die folgenden Versuche verwendet.

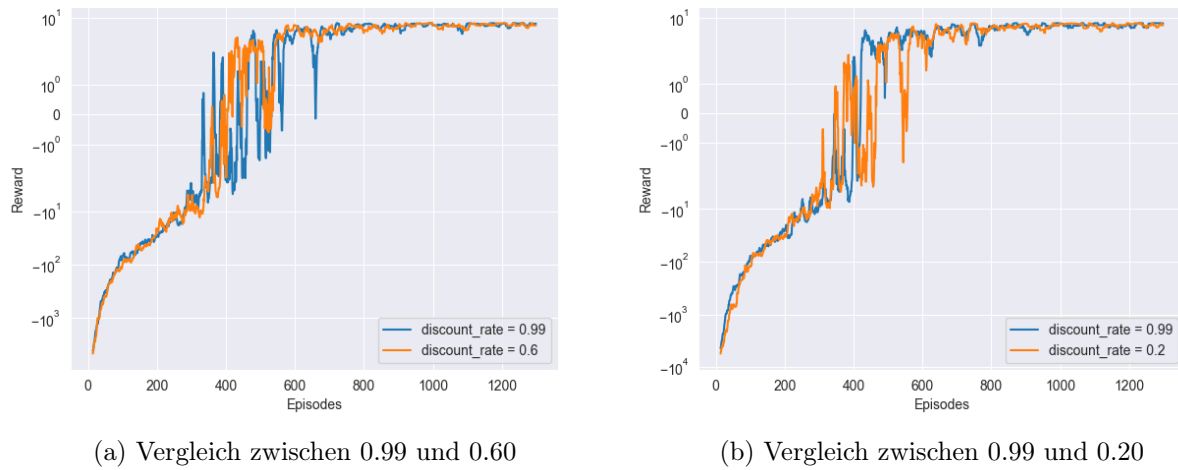


Abbildung 10: Auswirkung der Discount Rate

Bei dem Cliff Problem führt eine Absenkung des Wertes zu einer minimalen Verschlechterung der Ergebnisse, somit wird der Ausgangswert als optimal angenommen. Dieser hohe Initialwert verursachte bei dem Lake Problem häufiger einen Leistungseinbruch. Der Agent war nicht mehr in der Lage die Aufgabenstellung zu lösen. Durch eine Reduzierung des Wertes auf 0.2 konnte dieses Problem behoben werden.

- SARSA

Eine Reduzierung des Hyperparameter führte beim Taxiproblem zu deutlich schlechteren Ergebnissen. Da der initiale Wert bereits sehr nah am Maximum liegt, wurde dieser übernommen. Ähnliches Verhalten zeigte sich auch bei der Untersuchung des Cliff und Frozen Lake Problems, somit blieben auch dort die Parameter unverändert.

4.1.5 Exploration Rate

- Q-Learning

In Experimenten mit dem Taxiproblem hat sich gezeigt, dass eine Reduzierung der Exploration Decay Rate zu einem längeren Training führt, ohne dabei am Ende mehr Stabilität zu bringen (vgl. Abbildung 11a). Das Erhöhen hingegen beschleunigt das Training minimal (vgl. Abbildung 11b). Um ein schnelles Training zu erreichen und zeitlich sicherzustellen, dass der Agent alle möglichen Wege erkundet, wurde ein Wert von 0.01 für diesen Hyperparameter gewählt. Das Verhalten von Q-Learning auf dem Cliff Problem ist sehr ähnlich, daher ist dort 0.01 ein gut geeigneter Wert. Durch den hohen auf Zufall basierenden Einfluss des Environments bei dem Frozen Lake Problems ist die Wirkung der Exploration Decay Rate nur in Extremfällen feststellbar. Daher wurde keine Anpassung des Standardwertes vorgenommen.



(a) Auswirkung niedriger Exploration Decay Rate

(b) Auswirkung hoher Exploration Decay Rate

Abbildung 11: Auswirkung der Exploration Decay Rate

- SARSA

Die Reduzierung der Exploration Decay Rate führt beim Taxiproblem und Cliff Problem, zu einem sehr langsamen Training. Weitere Versuche ergaben dass 0.1 auch für SARSA ein gutes Gleichgewicht zwischen Lerngeschwindigkeit und Exploration des Environments ist. Beim Frozen Lake Problem überschattet der Zufallsfaktor des Environments jeglichen Einfluss dieses Parameters, somit bleibt der Wert unverändert.

4.2 Vergleichen von Algorithmen

Im folgenden werden die Ergebnisse beschrieben, welche bei dem Vergleichen der Algorithmen Q-Learning und SARSA entstanden sind. Dabei wurde der Vergleich an den drei Reinforcement Learning Problemen Taxi, Cliff und Frozen Lake durchgeführt. Es wurden jeweils passende Parameter für das entsprechende Problem gewählt.

4.2.1 Taxi

In der Abbildung 12 sind die Ergebnisse für den Vergleich des Q-Learning und SARSA Algorithmus für das Taxi Problem dargestellt. In der linken Abbildung 12a sind die Rewards in dem zeitlichen Verlauf über die 3000 Trainingsepisoden abgebildet. Zu erkennen ist hier, dass beide Algorithmen nach ca. 1500 Episoden Richtung Optimum konvergieren. Der Q-Learning Algorithmus trainiert jedoch ein wenig schneller, was durch die anfänglich stärkere Steigung und höheren Peaks verdeutlicht wird.

Zusätzlich ist es auffällig, dass SARSA im späteren Verlauf auch größere Peaks nach unten aufweist und generell ein wenig schlechter abschneidet als der Q-Learning Algorithmus. Somit lässt sich insgesamt für die Rewards feststellen, dass Q-Learning schneller zu einem besseren Ergebnis kommt als SARSA.

In der zweiten Abbildung 12b sind die Anzahl der Steps abgebildet, welche der Agent jede Episode durchläuft, bis er entweder bei der maximalen Anzahl angekommen ist oder den Passagier erfolgreich zu seinem Zielort gebracht hat. Der Verlauf der beiden Kurven von Q-Learning und SARSA ist hier sehr ähnlich. Lediglich bei ca. 400 und 1300 Episoden lassen sich leicht höhere Peaks vom SARSA

Algorithmus erkennen. Dies spricht dafür, dass SARSA in diesen Phasen etwas langsamer gelernt hat als der Q-Learning Algorithmus.

Wichtig zu erwähnen ist zusätzlich, dass beide Algorithmen nach 3000 Episoden im Test perfekt optimiert sind. Im Test nutzen die trainierten Modelle keine Exploration Rate mehr, sondern gehen lediglich den Weg mit den höchsten Q-Values. Starten beide in dem gleichen State, laufen sie fast die gleichen Wege um den Passagier abzuholen und zum Ziel zu bringen.

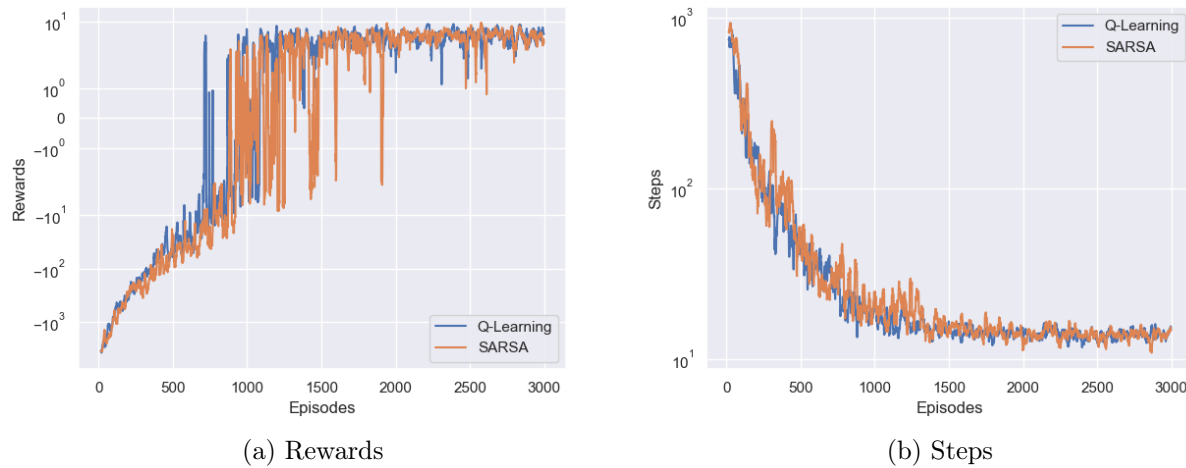


Abbildung 12: Taxi Problem über 3000 Episoden während dem Training

4.2.2 Cliff

Mit dem Cliff Problem lässt sich einer der besonderen Unterschiede zwischen Q-Learning und SARSA sehr gut zeigen. In Abbildung 13 sind wie in dem vorherigen Beispiel der Verlauf der Rewards (Abbildung 13a) und die Anzahl an Steps (Abbildung 13b) über die Episoden im Training dargestellt. Bei den Rewards ist zu erkennen dass SARSA zum einen schneller lernt, dass heißt die Kurve geht zu Beginn schneller nach oben. Zum anderen stagniert es auf einem höheren Reward zum Ende hin. Auch in der Abbildung 13b lässt sich ein Unterschied feststellen. Zu Beginn verlaufen sie sehr ähnlich und im weiteren Verlauf ist es auffällig, dass SARSA bei einer höheren Anzahl an Steps stagniert, während Q-Learning weniger Steps benötigt, aber auch mehr Schwankungen aufweist.

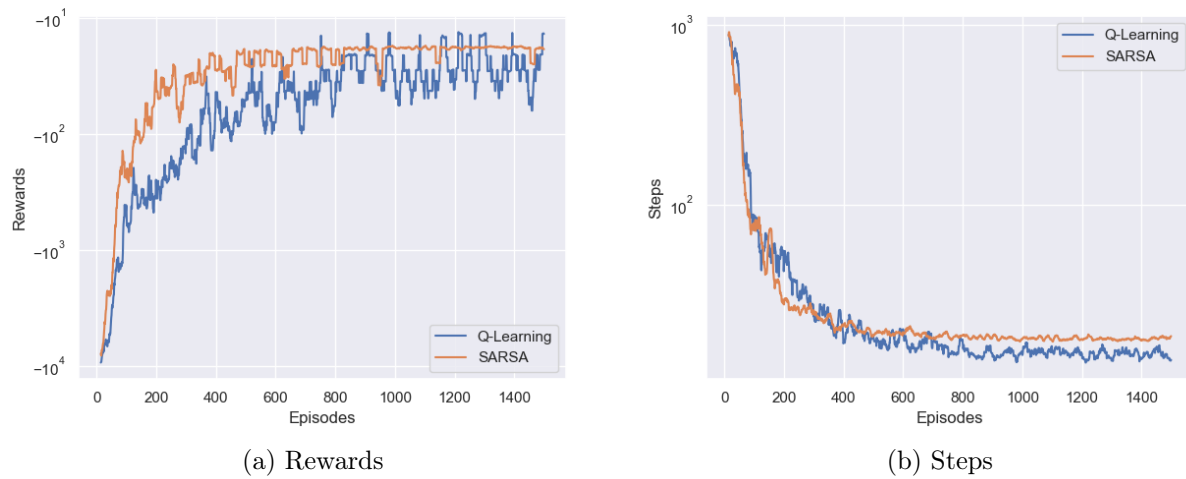


Abbildung 13: Cliff Problem über 1500 Episoden während dem Training

Dies liegt daran, dass SARSA ein on-policy Algorithmus ist. Mit anderen Worten, es berücksichtigt bei der Berechnung des Q-Values, dass es eine Wahrscheinlichkeit gibt, zu der der Agent im nächsten Schritt explorieren statt exploiten wird. Gerade bei dem Cliff Problem ist dies besonders spannend, da eine Exploration, welche in der Klippe endet, sehr fatal sein kann. In der Abbildung 14 sind die Wege dargestellt, welche die Algorithmen wählen. SARSA entscheidet sich aufgrund der on-policy Strategie für den “safe path” während Q-Learning für den “optimal path” geht. Beide Wege haben durchaus ihre Vor- und Nachteile und können je nach Anwendungsfall den subjektiv besseren Weg darstellen. Aus den eben genannten Gründen nennt man Q-Learning auch einen *greedy*-Algorithmus, da er sich für den möglichst kürzesten Weg entscheidet, während SARSA etwas vorsichtiger, also weniger *greedy*, agiert.

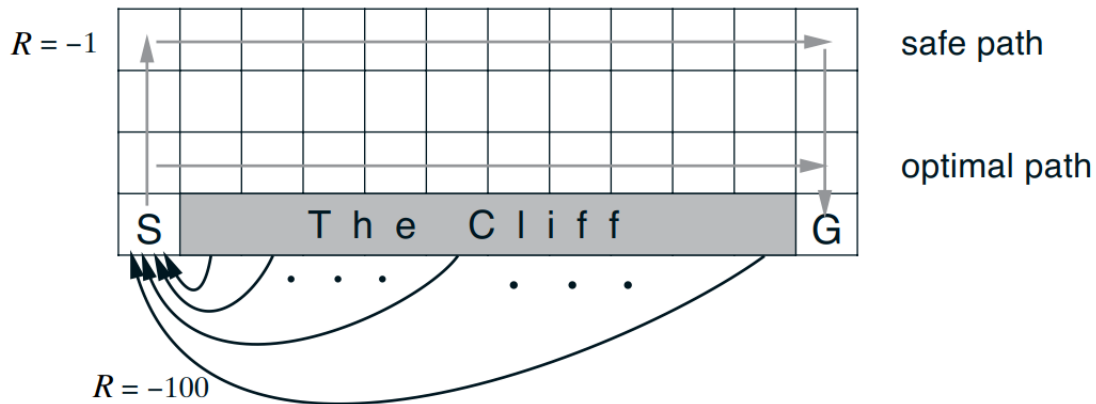


Abbildung 14: Cliff Q-Learning vs SARSA Pfad

4.2.3 Frozen Lake

Das Frozen Lake Problem hat die Besonderheit, dass es nicht zu 100% optimiert werden kann. Dies ist auch in den Abbildungen der Rewards (Abb. 15a) und der Steps (Abb. 15b) zu erkennen. Beide Gra-

phen verlaufen zwar jeweils ähnlich und stagnieren auf einem Niveau unter dem Optimum. Bezüglich der Vergleichen der Algorithmen lässt sich hier kein sichtbarer Unterschied der Performance aufzeigen. Bei den Rewards brauchen beide Algorithmen ca. 500 Episoden, um das Zielfeld regelmäßiger zu finden. Nach 2000 Episoden schaffen beide es regelmäßig zum Ziel. Anhand des Graphen 15b ist zu erkennen, dass beide Algorithmen zu Beginn schnell in eines der Löcher fallen und somit die Episoden frühzeitig enden. Über die Zeit schaffen sie es aber länger auf dem Frozen Lake zu laufen. Optimalerweise würde die Anzahl der Steps gegen Ende wieder etwas sinken, dies ist aber in den ersten 2000 Episoden nicht der Fall.

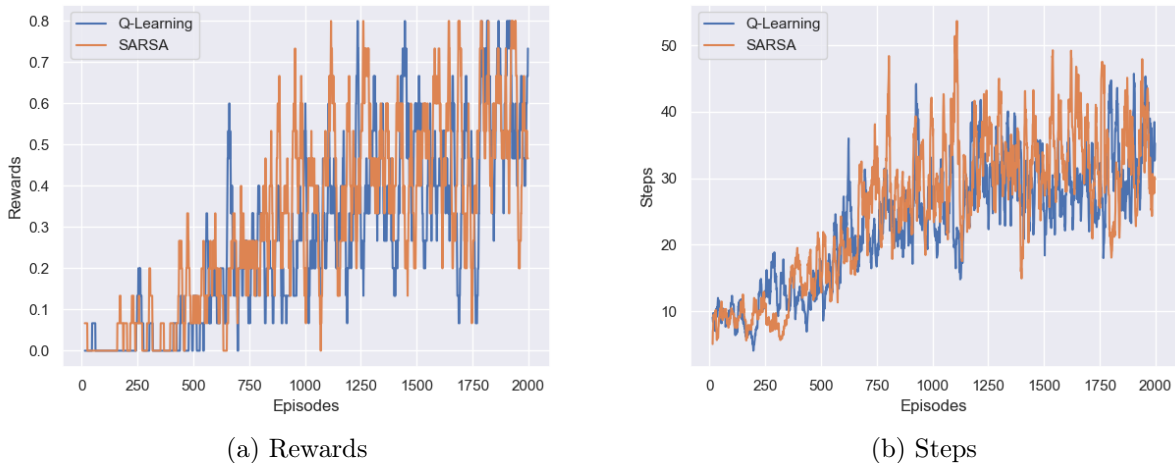


Abbildung 15: Frozen Lake Problem über 2000 Episoden während dem Training

Dies kann daran liegen, dass das Environment durch den Zufallsparameter sehr unvorhersehbar agiert. Dies ist auch daran zu erkennen, dass trotz des gleitenden Durschnitts, die Graphen sehr stark schwanken.

5 Zusammenfassung und Ausblick

Im Rahmen der Arbeit konnte sich ausgiebig mit den Grundlagen des Reinforcement Learnings auseinandergesetzt werden. Es wurde anhand von drei einfacheren Szenarien zwei klassische Algorithmen für das Temporal Difference Learning getestet. Für die beiden Algorithmen gibt es eine Vielzahl von Parametern, welche je nach Anwendungsfall passend gewählt werden müssen. Im Rahmen der Arbeit konnten die Einflüsse dieser Parameter aufgezeigt werden. Für die Temporal Difference Learning Algorithmen Q-Learning und SARSA konnte ein ausführlicher Vergleich anhand der drei Probleme durchgeführt werden. Hierbei wurde festgestellt, dass die Algorithmen je nach Problem zu unterschiedlichen Ergebnissen führen. So hat das Q-Learning beim Taxi Problem etwas besser abgeschnitten, während SARSA den wohl sichereren Weg bei dem Cliff Problem gewählt hat.

In weiteren Arbeiten ist es von Interesse, sich mit komplexeren Problemen und dem Thema Deep Reinforcement Learning zu beschäftigen. Dies ermöglicht das Anwenden von Reinforcement Learning auf weitaus komplexere Herausforderungen, welche nur schwer mit den in dieser Arbeit genutzten Algorithmen lösbar sind.

A Selbstständigkeitserklärung Joshua Henjes

Ich erkläre, dass

- o alle sinngemäßen Übernahmen aus Arbeiten Dritter mit der Quellenangabe gekennzeichnet sind,
- o alle wörtlichen Übernahmen von Textpassagen aus Arbeiten Dritter durch Anführungszeichen und ausführliche Angabe der Belegstelle als Zitat gekennzeichnet sind,
- o die vorliegende Arbeit selbständig unter Verwendung der im experimentellen Teil genannten Methoden angefertigt wurde und
- o Primärdaten von Experimenten der Arbeit unverändert und in geeigneter Form beigefügt sind.

Ort und Datum

Unterschrift

B Selbstständigkeitserklärung Bjarne Seen

Ich erkläre, dass

- o alle sinngemäßen Übernahmen aus Arbeiten Dritter mit der Quellenangabe gekennzeichnet sind,
- o alle wörtlichen Übernahmen von Textpassagen aus Arbeiten Dritter durch Anführungszeichen und ausführliche Angabe der Belegstelle als Zitat gekennzeichnet sind,
- o die vorliegende Arbeit selbständig unter Verwendung der im experimentellen Teil genannten Methoden angefertigt wurde und
- o Primärdaten von Experimenten der Arbeit unverändert und in geeigneter Form beigelegt sind.

Ort und Datum

Unterschrift