

# **LEARNING TO COOPERATE: REINFORCEMENT LEARNING IN THE ITERATED PRISONER'S DILEMMA**

Dokumentation

von Jonah Gräfe

Düsseldorf, 16.02.2025

Studiengang:  
Modul:  
Betreut von:

B.Sc. Data Science, AI und Intelligente Systeme  
Advances in Intelligent Systems  
Prof. Dr. Dennis Müller

# Inhaltsverzeichnis

<b>1</b>	<b>Einleitung</b>	<b>1</b>
1.1	Thema und Motivation . . . . .	1
1.2	Ziele der Arbeit . . . . .	1
<b>2</b>	<b>Hintergrund und theoretischer Rahmen</b>	<b>3</b>
2.1	Das Iterierter Gefangenendilemma . . . . .	3
2.2	Reinforcement Learning . . . . .	4
2.2.1	Q-Learning . . . . .	4
2.2.2	Deep Q-Learning . . . . .	5
<b>3</b>	<b>Methodik</b>	<b>6</b>
3.1	Aufbau des Experiments . . . . .	6
3.2	Agenten und Strategien . . . . .	6
3.3	Evaluierungsmethodik . . . . .	6
<b>4</b>	<b>Ergebnisse</b>	<b>7</b>
4.1	Daten und Beobachtungen . . . . .	7
4.2	Wichtige Erkenntnisse . . . . .	7
<b>5</b>	<b>Diskussion</b>	<b>8</b>
5.1	Interpretation der Ergebnisse . . . . .	8
5.2	Vergleich mit Erwartungen . . . . .	8
5.3	Limitationen . . . . .	8
<b>6</b>	<b>Fazit und Ausblick</b>	<b>9</b>
6.1	Zusammenfassung der Ergebnisse . . . . .	9
6.2	Ausblick . . . . .	9

# 1 Einleitung

## 1.1 Thema und Motivation

Das Iterative Gefangenendilemma (IPD) ist ein klassisches Problem der Spieltheorie, das die Herausforderungen von Kooperation und Eigennutz modelliert. Zwei Spieler müssen unabhängig voneinander entscheiden, ob sie kooperieren oder defektieren. Während Kooperation zu einem besseren kollektiven Ergebnis führt, ist aus individueller Sicht Defektion oft vorteilhafter – ein klassisches Dilemma. In der Realität tauchen ähnliche Dilemmata in vielen Bereichen auf, etwa in der Wirtschaft, der Politik und der Evolutionsbiologie. Ein zentrales Forschungsthema ist, ob und unter welchen Bedingungen Akteure langfristig kooperative Strategien entwickeln können, um dem Dilemma zu entkommen.

Mit der zunehmenden Bedeutung von künstlicher Intelligenz (KI) und Reinforcement Learning (RL) stellt sich die Frage, wie sich autonome Agenten in einem solchen Szenario verhalten. Können sie durch Lernen langfristig Kooperation aufbauen, oder werden sie egoistische Strategien bevorzugen? Ziel dieses Projekts ist es, RL-Agenten im IPD zu trainieren und zu analysieren, ob sie fähig sind, Strategien zu entwickeln, die über bloßen Eigennutz hinausgehen. Damit trägt diese Arbeit zur Diskussion über das Potenzial von RL in sozialen und spieltheoretischen Kontexten bei.

## 1.2 Ziele der Arbeit

Das Ziel dieser Arbeit ist es, zu untersuchen, wie sich RL-Agenten im IPD verhalten und ob sie in der Lage sind, kooperative Strategien zu entwickeln. Dabei stehen folgende zentrale Forschungsfragen im Fokus:

1. Können RL-Agenten lernen, langfristig zu kooperieren?
  - Entwickeln sie Strategien, die über reines Eigeninteresse hinausgehen?
  - Gibt es bestimmte Bedingungen, unter denen Kooperation wahrscheinlicher wird?
2. Welche Strategien entstehen im Laufe des Lernprozesses?
  - Verhalten sich die Agenten wie bekannte spieltheoretische Strategien (z. B. "Tit-for-Tat" oder "Always Defect")?
  - Zeigen sich neuartige, unerwartete Verhaltensmuster?
3. Wie beeinflussen verschiedene Lernparameter das Verhalten der Agenten?

## 1 Einleitung

- Welche Rolle spielen Faktoren wie Belohnungsfunktionen, Discount-Faktor oder Explorationsstrategien?
- Wie wirken sich unterschiedliche Trainingsumgebungen auf die Ergebnisse aus?

Welche Rolle spielen Faktoren wie Belohnungsfunktionen, Discount-Faktor oder Explorationsstrategien? Um diese Fragen zu beantworten, werden verschiedene RL-Algorithmen auf das IPD angewendet, die Trainingsverläufe analysiert und die resultierenden Strategien miteinander verglichen. Durch die gewonnenen Erkenntnisse soll ein tieferes Verständnis dafür geschaffen werden, wie lernende Agenten spieltheoretische Herausforderungen bewältigen und ob sie sich aus dem klassischen Dilemma befreien können.

## 2 Hintergrund und theoretischer Rahmen

### 2.1 Das Iterierter Gefangenendilemma

Das Gefangenendilemma ist eines der bekanntesten Probleme der Spieltheorie und beschreibt eine Situation, in der zwei Spieler unabhängig voneinander entscheiden müssen, ob sie kooperieren oder defektieren. Die Auszahlung für jeden Spieler hängt sowohl von der eigenen Entscheidung als auch von der des Gegenübers ab. Die klassische Auszahlungsmatrix sieht dabei folgendermaßen aus:

	Spieler B kooperiert	Spieler B defektiert
Spieler A kooperiert	(R, R) Belohnung für Kooperation	(S, T) A verliert, B gewinnt
Spieler A defektiert	(T, S) A gewinnt, B verliert	(P, P) Bestrafung für gegenseitige Defektion

Dabei gilt üblicherweise:

- $T$  (Temptation)  $>$   $R$  (Reward)  $>$   $P$  (Punishment)  $>$   $S$  (Sucker's payoff)
- $2R > T + S$ , sodass sich gegenseitige Kooperation langfristig mehr lohnen würde als wechselseitige Ausnutzung.

Im einmaligen Gefangenendilemma ist die dominante Strategie, zu defektieren, da dies in jedem individuellen Fall die höhere Auszahlung sichert – unabhängig von der Entscheidung des Gegenspielers. Dies führt jedoch zu einem sozial suboptimalen Ergebnis.

Im iterierten Gefangenendilemma (IDG) wird das Spiel jedoch mehrfach hintereinander gespielt, sodass frühere Entscheidungen zukünftige Interaktionen beeinflussen können. Dadurch eröffnen sich neue Möglichkeiten für kooperative Strategien, bei denen Agenten versuchen, durch wechselseitige Zusammenarbeit langfristig höhere Erträge zu erzielen. Bekannte Strategien aus der Spieltheorie für das IDG sind beispielsweise: "Tit-for-Tat" (Spiele das, was dein Gegner in der vorherigen Runde getan hat). "Always Defect" (Immer defektieren, um kurzfristig die höchste Auszahlung zu sichern). "Grim Trigger" (Kooperiere, aber falls der Gegner einmal defektiert, defektiere für immer).

Die zentrale Forschungsfrage im IPD lautet daher: Ist es möglich, langfristige Kooperation zu etablieren, oder führt Eigennutz immer zu gegenseitiger Defektion? Diese Fragestellung ist besonders relevant für Reinforcement Learning-Agenten, da sie ihre Strategie durch wiederholte Interaktion und Belohnungsmechanismen erlernen.

## 2.2 Reinforcement Learning

Reinforcement Learning (RL) ist ein Teilbereich des maschinellen Lernens, bei dem ein Agent durch Interaktion mit einer Umgebung eine optimale Strategie erlernt. Der Lernprozess basiert auf einem Belohnungssystem: Der Agent führt Aktionen aus, erhält daraufhin Belohnungen oder Bestrafungen und passt sein Verhalten entsprechend an. RL-Probleme werden typischerweise als Markov-Entscheidungsprozesse (MDP) modelliert, bestehend aus:

- Zustand (State,  $S$ ): Die aktuelle Situation der Umgebung.
- Aktion (Action,  $A$ ): Eine Entscheidung, die der Agent treffen kann.
- Belohnung (Reward,  $R$ ): Eine Rückmeldung, die die Qualitäten der gewählten Aktion bewertet.
- Übergangsmodell ( $P(s'|s, a)$ ): Wahrscheinlichkeiten, dass der Zustand  $s'$  nach einer Aktion  $a$  im Zustand  $s$  entsteht.
- Policy ( $\pi(s)$ ): Die Strategie des Agenten zur Auswahl von Aktionen.

Das Ziel ist es, eine Optimale Policy  $\pi^*$  zu lernen, die die kumulierte zukünftige Belohnung maximiert. Dafür gibt es verschiedene Methoden, darunter Q-Learning und Deep Q-Learning.

### 2.2.1 Q-Learning

Q-Learning ist ein wertbasierter RL-Algorithmus, der darauf abzielt, die Q-Werte für jede Zustands-Aktions-Kombination zu erlernen. Der Q-Wert  $Q(s, a)$  repräsentiert die erwartete zukünftige Belohnung, wenn der Agent in Zustand  $s$  Aktion  $a$  wählt und danach der optimalen Strategie folgt. Die Aktualisierung der Q-Werte erfolgt iterativ mit der Bellman-Gleichung:

$$Q(s, a) \leftarrow Q(s, a) + \alpha(r + \gamma \max_{a'} Q(s', a') - Q(s, a)) \quad (2.1)$$

Hierbei sind:

- $\alpha$  die Lernrate, die bestimmt, wie stark neue Informationen alte Werte überschreiben.
- $\gamma$  der Diskontfaktor, der die Gewichtung zukünftiger Belohnungen bestimmt.
- $r$  die unmittelbare Belohnung nach der Aktion  $a$ .

Da Q-Learning tabellarisch arbeitet, ist es nur für kleine Zustandsräume geeignet, da die Q-Tabelle bei vielen möglichen Zuständen und Aktionen schnell zu groß wird. In komplexeren Umgebungen ist daher eine neuronale Netzarchitektur notwendig - hier kommt Deep Q-Learning (DQN) ins Spiel.

### 2.2.2 Deep Q-Learning

Deep Q-Learning (DQN) erweitert Q-Learning durch den Einsatz eines neuronalen Netzwerks zur Approximation der Q-Werte, anstatt eine explizite Tabelle zu speichern. Dies erlaubt das Lernen in hochdimensionalen Zustandsräumen, die tabellarische Methoden überfordern würden. Die Hauptbestandteile von DQN sind:

1. Neuronales Netz als Q-Funktion
  - Das Netz nimmt den Zustand  $s$  als Eingabe und gibt geschätzte Q-Werte für alle möglichen Aktionen  $a$  aus.
  - Die Gewichte des Netzwerks werden durch Gradientenabstieg und einen Mean-Squared-Error (MSE)-Loss aktualisiert.
2. Erfahrungsspeicher (Experience Replay)
  - Anstatt direkt mit den neuesten Erfahrungen zu trainieren, werden vergangene Erfahrungen  $(s, a, r, s')$  in einem Speicher abgelegt.
3. Zielnetzwerk (Target Network)
  - Zusätzlich zum Hauptnetzwerk existiert eine separate Kopie, die in regelmäßigen Abständen aktualisiert wird.
  - Dies verhindert zu starke Schwankungen in den Q-Werten und stabilisiert das Training.

Die Aktualisierungsregel für DQN basiert auf dem Mean-Squared-Error zwischen dem vorhergesagten und dem Ziel-Q-Wert:

$$L = \left( r + \gamma \max_{a'} Q_{\text{target}}(s', a') - Q_{\text{current}}(s, a) \right)^2 \quad (2.2)$$

DQN ist besonders mächtig für komplexe Umgebungen, in denen eine tabellarische Q-Funktion nicht mehr praktikabel ist.

## 3 Methodik

### 3.1 Aufbau des Experiments

Um zu untersuchen, wie Reinforcement-Learning-Agenten im Iterierten Gefangenendilemma (IPD) agieren, wurde eine experimentelle [Python-Umgebung](#) implementiert. Diese Umgebung ermöglicht es, Spiele zu simulieren, RL-Agenten zu trainieren, und diese dann zu evaluieren.

### 3.2 Agenten und Strategien

### 3.3 Evaluierungsmethodik



## 4 Ergebnisse

### 4.1 Daten und Beobachtungen

### 4.2 Wichtige Erkenntnisse

## 5 Diskussion

### 5.1 Interpretation der Ergebnisse

### 5.2 Vergleich mit Erwartungen

### 5.3 Limitationen

## 6 Fazit und Ausblick

### 6.1 Zusammenfassung der Ergebnisse

### 6.2 Ausblick