

# Einführung in das Reinforcement Learning

Patrick Dammann

heute

## Contents

<b>1</b>	<b>Abstract</b>	<b>2</b>
<b>2</b>	<b>Markov Entscheidungs-Probleme</b>	<b>2</b>
2.1	Ein Beispiel . . . . .	2
2.2	Notation . . . . .	2
2.3	V-Values . . . . .	3
2.4	Value Iteration . . . . .	4
2.5	Policy Extraction . . . . .	5
<b>3</b>	<b>Reinforcement Learning</b>	<b>5</b>
3.1	Verbindung zu MDPs . . . . .	5
3.2	Model Based Approaches . . . . .	6
3.3	Temporal Difference Learning . . . . .	6
3.4	Q-Values/Q-Learning . . . . .	6
3.5	Exploration . . . . .	7
3.5.1	Startbedingungen . . . . .	7
3.5.2	Die $\epsilon$ -greedy Exploration . . . . .	7
3.5.3	Erkundungs-Funktionen . . . . .	8
3.6	State Features . . . . .	8
3.7	Andere Algorithmen . . . . .	9
3.7.1	SARSA . . . . .	9
3.7.2	Double Q-learning . . . . .	9
3.8	Deep Q-Learning . . . . .	9
<b>4</b>	<b>Fazit</b>	<b>9</b>
<b>5</b>	<b>Kurzgeschichte (?)</b>	<b>9</b>
<b>6</b>	<b>Quellen</b>	<b>9</b>

# 1 Abstract

## 2 Markov Entscheidungs-Probleme

Markov Entscheidungs-Probleme (kurz **MDP** für *markov decision process*) sind ein Modell für Probleme, bei denen ein Agent versucht die größtmögliche Belohnung (*reward*) zu erzielen. Er bewegt sich dazu durch eine Menge von Zuständen (*states*) indem er aus einer Menge von Aktionen (*actions*) wählt. Welchen Zustand der Agent erreicht ist nicht deterministisch, die Wahrscheinlichkeiten sind jedoch nur von der gewählten Aktion und dem aktuellen Zustand abhängig. Die Belohnung, die der Agent für einen Übergang erhält, wird durch Ausgangs- und Endzustand des Übergangs so wie die gewählte Aktion bestimmt. Der Agent startet in einem Startzustand, es kann diverse Endzustände geben, nach deren Erreichen keine Aktionen mehr ausgeführt werden können.

Ziel ist es, eine Strategie (*policy*) zu finden, welche jedem Zustand die Aktion zuordnet, durch die der höchste Gesamtgewinn (also die höchste Summe über alle erzielten Gewinne) erwartet werden kann.

### 2.1 Ein Beispiel

In diesem Beispiel werden die Zustände durch die Felder eines  $4 \times 4$ -Grids visualisiert. Der Startzustand (3,0) ist durch ein großes *S* markiert, die Zustände mit einem *X* sind nicht erreichbar.

	0	1	2	3
0				
1		X	-100	+100
2				
3	S	X	-100	-100

Die Menge der Aktionen besteht aus 4 Elementen, jedes symbolisiert einen Schritt in eine der vier Himmelsrichtungen. Beim Ausführen einer Aktion landet der Agent mit einer Chance von 0.8 ein Feld weiter in der gewählten Himmelsrichtung, so wie mit je einer Chance von 0.1 in einer der zwei orthogonalen Richtungen. Würde der Agent hierbei auf einem Feld landen, welches er nicht betreten kann (eines der mit *X* markierten, oder außerhalb des Grids), bleibt er auf seinem Feld.

Von den bunten Feldern aus führt jede Aktion mit einer Chance von 1 in den Endzustand. Der Agent erhält für diesen Übergang die auf dem Feld verzeichnete Belohnung (100 oder -100). Für jede andere Aktion wird eine Belohnung von -2 verbucht.

### 2.2 Notation

Zur Beschreibung eines MDPs werden einige Notationen benötigt, die hier angegeben werden. (*Zu einigen Notationen werden Beispiele zum oben genannten Beispiel angegeben.*) Zur Problemstellung selbst gehören:

- $S = \{s_1, \dots, s_n\}$  ist die Menge der **Zustände**. Die Elemente der Menge können beliebig benannt werden.
- $A = \{a_1, \dots, a_n\}$  ist die Menge der **Aktionen**. Auch diese Elemente können beliebig benannt werden.
- $T : S \times A \times S \rightarrow [0, 1]$  ist die **Übergangs-Funktion** (*transition function*).  $T(s, a, s')$  beschreibt die Wahrscheinlichkeit, mit der der Agent von Zustand  $s$  in den Zustand  $s'$  wechselt, wenn er Aktion  $a$  ausführt. Eine allgemeinere, jedoch hier nicht verwendete Schreibweise wäre  $T(s, a, s') = p(s'|s, a)$ .

$$\begin{aligned}
T(s_{0,2}, a_E, s_{0,3}) &= 0.8 && \text{(Agent wählt Osten und kommt dort an)} \\
T(s_{0,2}, a_E, s_{0,2}) &= 0.1 && \text{(Agent wählt Osten, geht nach Norden, kein valider Zug)} \\
T(s_{0,2}, a_E, s_{1,2}) &= 0.1 && \text{(Agent wählt Osten, geht nach Süden)} \\
T(s_{0,2}, a_E, s_{0,1}) &= 0.0 && \text{(Agent wählt Osten, Bewegung nach Westen unmöglich)} \\
T(s_{1,2}, \_, s_{end}) &= 1.0 && \text{(jede Aktion von bunten Feldern führt zum Endzustand)}
\end{aligned}$$

- $R : S \times A \times S \rightarrow \mathbb{R}$  ist die **Belohnungs-Funktion** (*reward function*).  $R(s, a, s')$  ordnet dem Übergang von Zustand  $s$  nach Zustand  $s'$  mit Hilfe von Aktion  $a$  eine Belohnung zu.

$$\begin{aligned}
R(s_{1,2}, \_, s_{end}) &= -100 && \text{(Belohnung von buntem Feld)} \\
R(s_{1,3}, \_, s_{end}) &= +100 && \text{(Belohnung von buntem Feld)} \\
R(s_{0,2}, a_N, s_{0,1}) &= -2 && \text{(jede Bewegung kostet 2 Belohnung)}
\end{aligned}$$

Zum Lösen der Probleme werden folgende Notationen verwendet:

- $\pi(s)$  ist eine **Policy**, die dem Zustand  $s$  die optimale Aktion zuordnet. Hierbei beschreibt  $\pi$  eine Policy im Allgemeinen,  $\pi^*$  die optimale Policy.
- $\gamma \in [0, 1]$  ist der sogenannte **Discount-Faktor**. Für den Agenten werden nach jeder Aktion alle noch erreichbaren Belohnungen mit  $\gamma$  multipliziert, damit er kürzere Wege, die zur selben Belohnung führen, bevorzugen wird. Mit  $\gamma = 1$  ist für den Agenten ein langer Weg ebenso gut wie ein kurzer, solange er die selbe Belohnung erhält. Mit  $\gamma = 0$  wird der Agent nicht mehr vorausschauend planen, da für ihn nur die nächste Belohnung Wert hat. Es gilt also, einen geeigneten Trade-Off zu finden.

## 2.3 V-Values

Um ein MDP zu lösen, lässt sich auf ein Bewertungsschema für Zustände zurückgreifen. Hierbei ergibt sich die Bewertung eines Zustandes aus der **Gesamtbelohnung, die der Agent zu erwarten hat, wenn er eine Simulation in besagtem Zustand startet, die optimale Aktion ausführt und von dort an optimal handelt**. Die Markov-Bedingung erlaubt uns in jedem Zustand anzunehmen, wir hätten die Simulation gerade erst gestartet, da es keine Faktoren gibt, die von bereits vergangenen Ereignissen abhängen.

Da die Bewertung jedes Zustandes angibt, welche Belohnung von ihm aus zu erwarten ist, lässt sich die Bewertung wie folgt rekursiv definieren:

$$V^*(s) = \max_{a \in A} \sum_{s' \in S} T(s, a, s') [R(s, a, s') + \gamma \cdot V^*(s')]$$

In der eckigen Klammer befindet sich die unmittelbare Belohnung, wenn der Agent von Zustand  $s$  über Aktion  $a$  in Zustand  $s'$  gelandet ist, addiert zur zu erwartenden Belohnung von Zustand  $s'$  aus (welche mit  $\gamma$  multipliziert wurde, da eine Aktion vorrüber ist).

Dieser Wert wird mit der Wahrscheinlichkeit des Auftretens dieses Übergangs multipliziert und durch Aufsummieren über alle möglichen Zustände  $s'$  (man bedenke, dass  $T$  nur für Zustände, die von  $s$  über  $a$  erreicht werden können  $> 0$  ist) zum Erwartungswert für eine Aktion  $a$ .

Da  $V^*$  als erwartete Belohnung definiert ist, wenn der Agent immer optimal handelt, ist es der Erwartungswert der Aktion, von welcher die höchste Belohnung erwartet wird.

## 2.4 Value Iteration

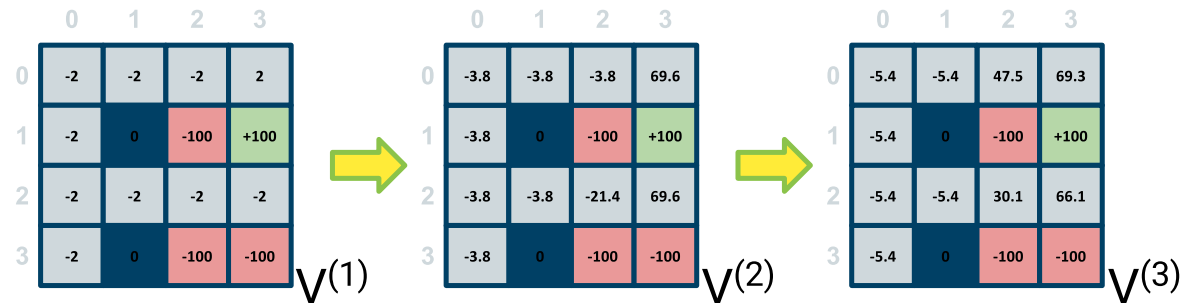
Durch die rekursive Definition von  $V^*$  lässt es sich nicht trivial berechnen. Als möglicher Lösungsansatz bietet sich daher die Value Iteration an. Hierzu wird das MDP zeitlich begrenzt und  $V^{(k)}$  definiert als *Gesamtbelohnung, die der Agent zu erwarten hat, wenn er eine Simulation in besagtem Zustand startet, die optimale Aktion ausführt und von dort an optimal handelt, **aber nur noch  $k$  Schritte übrig hat.***

$V^{(0)}$  lässt sich einfach bestimmen, da es 0 für jeden Zustand sein muss. Auch  $V^{(1)}$  folgt oft direkt aus der Problemdefinition. In unserem Beispiel wäre es  $-100$  bzw.  $+100$  für die "bunten" Zustände und  $-2$  für alle anderen, da jeder Schritt  $-2$  Belohnung bringt.

Ein beliebiges  $V^{(k)}$  lässt sich einfach berechnen. Dazu muss nur die Formel für  $V^*$  leicht angepasst werden:

$$V^{(k)}(s) = \max_{a \in A} \sum_{s' \in S} T(s, a, s') [R(s, a, s') + \gamma \cdot V^{(k-1)}(s')]$$

Da sich das hintere  $V^*$  in der Gleichung ja auf Geschehnisse bezog, die einen Zeitschritt später passieren, muss es bei der Berechnung von  $V^{(k)}$  durch ein  $V^{(k-1)}$  ersetzt werden, da der Agent im zeitlich beschränkten MDP einen Zeitschritt später selbstverständlich nur noch einen Zeitschritt weniger zur Verfügung hat.



Für MDPs, deren maximal erzielbare Belohnung nach oben beschränkt ist, wird  $V^{(k)}$  für  $k \rightarrow \infty$  konvergieren, da sich auch mit mehr Schritten kein besseres Ergebnis mehr erzielen lassen wird. Durch ein  $\gamma < 1$  kann dies  $<>$  werden.

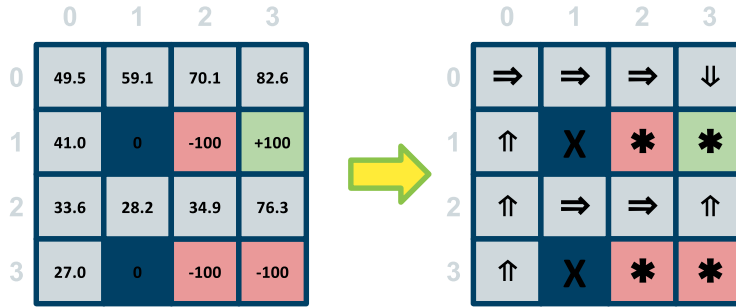
Für ein geeignet großes  $k$  gilt also  $V^{(k)} = V^* \pm \varepsilon$  mit vernachlässigbarem  $\varepsilon$  (z.B. *kleiner als Maschinengenauigkeit oder kleiner als eine vorher gewählte Toleranz*), da  $V^{(\infty)}$  genau unserer Definition für ein  $V^*$  im nicht zeitlich beschränkten MDP entspricht.

## 2.5 Policy Extraction

Mit Hilfe der optimalen V-Values  $V^*$  lässt sich jetzt eine optimale Policy für das MDP berechnen. Dazu simuliert man einen weiteren Schritt der Value Iteration. Anstatt nun aber tatsächlich das Maximum über alle  $a$  zu bestimmen, weisen wir dem aktuellen Zustand die Aktion zu, die von ihm aus die erwartete Belohnung maximieren würde.

$$\pi^*(s) = \operatorname{argmax}_{a \in A} \sum_{s' \in S} T(s, a, s') [R(s, a, s') + \gamma \cdot V^*(s')]$$

Auf diese Weise haben wir eine optimale Strategie erzeugt und somit das MDP gelöst.



## 3 Reinforcement Learning

Reinforcement Learning ist ein Machine-Learning-Konzept, in welchem ein Agent ohne anfängliche Informationen über sein Umfeld oder die Auswirkungen seines Handelns dazu trainiert werden soll, die maximale Belohnung zu erhalten.

### 3.1 Verbindung zu MDPs

Während sich zum Beispiel mit Value Iteration und Policy Extraction MDPs lösen lassen, bei denen alle Informationen vorhanden sind, bietet Reinforcement Learning eine Reihe von Algorithmen bei welchen vom Agenten keine Kenntnis von  $T(s, a, s')$  oder  $R(s, a, s')$  vorausgesetzt werden. Der Agent muss sich also selbstständig über Aktionen durch die Zustände bewegen, um die verschiedenen Übergänge und deren Belohnungen und Wahrscheinlichkeiten selbst zu erfahren.

Da der Agent erst nach jedem Übergang die erhaltene Belohnung von extern mitgeteilt bekommt, lassen sich so bei der Problemmodellierung einfache Möglichkeiten implementieren, Belohnungen zu verteilen.

### 3.2 Model Based Approaches

Die erste Möglichkeit an dieses Problem heranzugehen wäre ein **model based approach**. Das bedeutet, der Agent versucht zunächst, die fehlenden Teile des Modells (also vor allem  $T$  und  $R$ ) durch Samples zu approximieren. Sobald diese annähernd bekannt sind, lässt sich das Problem mit beliebigen MDP-Algorithmen (wie zum Beispiel Value Iteration und Policy Extraction) lösen.

Dieser Ansatz erfordert sehr große Mengen an Samples welche zur Auswertung gleichzeitig zur Verfügung stehen müssen. Ausserdem erfordert das Sammeln der Samples einen großen Grad an Freiheit in der Bewegung des Agenten durch die Zustände, was zum Beispiel bei Anwendungen in der Robotik oft nicht gegeben ist. Bei sehr großen (oder sogar indiskreten) Zustands-Mengen ist dieser Ansatz also nur schwer umsetzbar.

### 3.3 Temporal Difference Learning

Als weiterer Ansatz bietet sich das **temporal difference learning** an. Die Idee dahinter ist, dass mit zufällig gewählten V-Values begonnen wird und diese langsam iterativ verbessert werden.

Nach jeder Aktion (und deren Feedback) überprüft der Agent wie sich die ihm über V-Values versprochene erwartete Belohnung im Vergleich zur direkt erhaltenen Belohnung und der nun in Aussicht stehenden erwarteten Belohnung verhält und passt die V-Values des eben verlassenen Zustandes geringfügig dementsprechend an.

$$V(s) \leftarrow V(s) + \alpha [r + \gamma V(s') - V(s)]$$

Die Bewertung von Zustand  $s$  wird um einen kleinen Teil ( $\alpha < 0$ ) der Differenz zwischen versprochener ( $V(s)$ ) und erfahrener ( $r + \gamma V(s')$ ) Belohnung angepasst. Die Lernrate (*learning rate*)  $\alpha$  ist dabei ein wichtiger Hyperparameter und entscheidend für die Konvergenz des Verfahrens. Für ein sich verringerndes  $\alpha_i$  ( $\alpha$  im  $i$ -ten Zeitschritt) mit  $\sum a_i = \infty$  und  $\sum a_i^2 < \infty$  konvergiert das Verfahren in jedem Fall gegen die optimale Lösung (sofern es eine gibt). Aber auch für konstante, sehr kleine  $\alpha$  konvergiert das Verfahren in der Praxis sehr oft.

Oft lässt sich der Rechenaufwand für das Lernen auch verkleinern, in dem beim Erkunden der Umgebung zunächst viele Samples gesammelt wurden und das Lernen danach auf allen gesammelten Samples am Stück passiert. Erst danach wird die Policy aktualisiert, damit wieder neue Samples gesammelt werden können. Dies wird bis zur Konvergenz wiederholt.

### 3.4 Q-Values/Q-Learning

Das explizite Aktualisieren der Policy kann entfallen, wenn die Policy direkt implizit mitgelernt wird. Diese Idee kann durch die sogenannten Q-Values umgesetzt werden. Hierbei werden statt einem Wert pro State  $|A|$  Werte pro State gespeichert, jeweils einer pro ausführbarer Aktion.

Es ist also  $Q(s, a)$  definiert als *Gesamtbelohnung, die der Agent zu erwarten hat, wenn er eine Simulation in besagtem Zustand startet, **Aktion a ausführt** und von dort an optimal handelt*, was formal bedeutet:

$$Q^*(s, a) = \sum_{s' \in S} T(s, a, s') [R(s, a, s') + \gamma \cdot V^*(s')]$$

Hierbei lässt sich sofort erkennen, dass sowohl die V-Values als auch die gewünschte Strategie implizit mitgeliefert wird:

$$V(s) = \max_{a \in A} Q(s, a)$$

$$\pi(s) = \operatorname{argmax}_{a \in A} Q(s, a)$$

Viele beliebte Reinforcement-Learning-Algorithmen basieren auf dem direkten Lernen der Q-Values. Einer der einfachsten ist hierbei **Q-learning**, welcher nur eine Form von Temporal Difference Learning auf Q-Values darstellt. Die Update-Regel wird dabei lediglich wie folgt abgeändert:

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left[ r + \gamma \max_{a' \in A} Q(s', a') - Q(s, a) \right]$$

Selbstverständlich benötigt Q-learning aber auch mehr Samples zum konvergieren, da  $|A|$ -mal so viele Parameter gelernt werden.

### 3.5 Exploration

Während der Agent die Zustände erkundet, hat er stets die Auswahl zwischen zwei generellen Lernstrategien. Er muss entscheiden, ob er bisher schlecht erkundete Wege durch den Zustandsraum besuchen möchte um mehr über sie herauszufinden, oder ob er sich an seine aktuelle Strategie hält um ihre Richtigkeit zu überprüfen und vielleicht weniger erforschte Zustände zu erforschen, die er erst durch bereits für gut befundene Entscheidungen erreicht.

Es gibt viele Methoden um einen geeigneten Trade-Off zwischen den beiden Extrema zu finden. Drei davon, welche auch beliebig untereinander kombiniert werden können, werden hier vorgestellt:

#### 3.5.1 Startbedingungen

Zu Beginn vieler Reinforcement-Learning-Algorithmen werden Startwerte zufällig initialisiert. Wenn zum Beispiel die Q-Values allesamt höher initialisiert werden als der höchste Wert, von dem man glaubt dass er von ihnen angenommen werden kann, so werden die ersten Updates sehr häufig zu Verringerungen der Werte führen. Dadurch ist der Agent in den ersten Iterationen eher geneigt Wege auszuprobieren, für die er noch keine Updates durchgeführt hat.

Selbstverständlich ist diese Technik nur in den ersten Iterationen sinnvoll und sollte mit anderen Methoden kombiniert werden.

#### 3.5.2 Die $\epsilon$ -greedy Exploration

Bei dieser Erkundungs-Technik wird der neue, namensgebende Hyperparameter  $0 \leq \epsilon \leq 1$  eingeführt. Vor jeder Aktion wirft der Agent nun eine unfaire Münze, so dass er mit einer Chance von  $(1 - \epsilon)$  nach seiner Strategie handelt, jedoch mit einer Chance von  $\epsilon$  eine völlig zufällige Aktion durchführt. Während  $\epsilon$  zur Test-Zeit für gewöhnlich auf 0 gesetzt wird, kann es zur Trainingszeit je nach Bedarf konstant sein, mit der Zeit abnehmen oder anderweitig an die Gegebenheiten angepasst werden.

### 3.5.3 Erkundungs-Funktionen

Um die Erkundung tatsächlich gezielt in unbekannte Regionen zu lenken, lässt sich der Iterationsschritt in der Berechnung der Q-Values wie folgt anpassen:

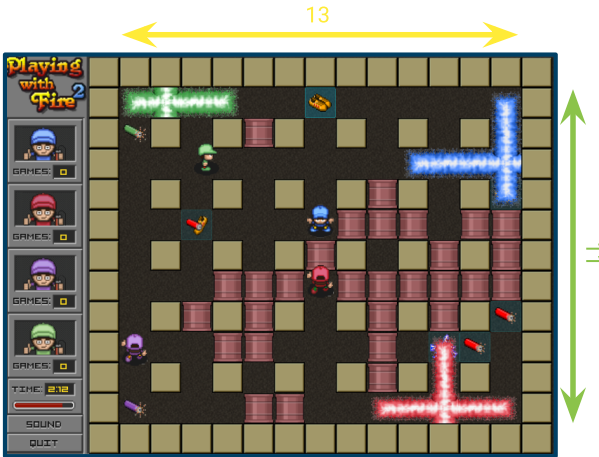
$$Q(s, a) \leftarrow Q(s, a) + \alpha \left[ r + \gamma \max_{a' \in A} (Q(s', a') + E(s', a')) - Q(s, a) \right]$$

Die Funktion  $E(s', a')$  sinkt dabei jedes Mal, wenn die Aktion  $a'$  im Zustand  $s'$  ausgeführt wird. Sie kann zum Beispiel die Form  $E(s, a) = k/n_{s,a}$  haben, wobei  $n_{s,a}$  die Anzahl der Vorkommnisse des Zustand-Aktion-Tupels  $(s, a)$  ist.

Da  $E(s, a)$  gegen 0 konvergiert, und fälschlicherweise gut klassifizierte Aktionen oft gewählt und im Wert verringert werden, sollte beim Konvergieren der Q-Learning-Methode der Wert der Erkundungsfunktion keinen merklichen Einfluss mehr auf die gelernte Strategie haben.

### 3.6 State Features

Gerade in realistischen Anwendungen wird die Anzahl verschiedener Zustände schnell unüberschaubar groß. Eine ungefähre Vorstellung lässt sich über diese grobe, obere Abschätzung der möglichen Zustände eines Flash-Spiels vermitteln:



Man stelle sich ein Szenario vor, in dem der Reinforcement-Learning-Agent lernen soll, ein einfaches Spiel im Stil des bekannten Videospiele "Bomberman" zu gewinnen.

Das Spielfeld in diesem Beispiel besteht aus  $11 \cdot 13$  Feldern, welche jeweils 15 verschiedene Zustände annehmen können. Diese Zustände sind: leer, zerstörbarer Stein, unzerstörbarer Stein, Feuer von je einem der 4 Spieler, Dynamit von je einem der 4 Spieler oder eines von 4 sammelbaren Power-Ups. Allein dadurch bekommt bereits  $15^{11 \cdot 13} \approx 1.5 \cdot 10^{168}$  verschiedene Zustände, bereits mehr als das Quadrat aller Atome im Universum.

Wenn man jetzt noch betrachtet, dass in jedem dieser Zustände jeder der vier Spieler auf jedem der  $44 \cdot 44$  Pixel jedes Feldes stehen könnte und möglicherweise die Zeit, die auf dem Timer auf die Sekunde genau abgelesen werden kann, die Spielweise des Agenten beeinflussen sollte und somit neue



Zustände generiert, kommt man sogar auf eine Kardinalität von  $15^{11 \cdot 13} \cdot (13 \cdot 11 \cdot 44^2)^4 \cdot (3 \cdot 60) \approx 1.6 \cdot 10^{192}$ .

Natürlich ist dies nur eine obere Abschätzung, da einige der Zustände sicherlich unerreichbar sind. Die ungefähre Größenordnung zeigt jedoch gut auf, wie schnell die Zustandsmenge unermesslich groß werden kann.

Nicht nur, dass der Agent für jeden dieser Zustände eine eigene Aktion für die Strategie bzw. einen Q-Value lernen muss, er kann auch nicht über verschiedene Zustände hinweg abstrahieren und somit von Gelerntem über Zustände auf ähnliche Zustände schließen.

Aus diesem Grund lohnt es sich, statt eines Lookup-Tables für die Q-Values eine Q-Funktion zu approximieren. Eine Möglichkeit hierfür ist, nur bestimmte Features aus den Zuständen zu extrahieren und mit diesen über eine gewichtete Summe die Q-Values anzunähern. Dies könnte beispielsweise diese Form haben:

$$Q(s, a) = w_1 f_1(s, a) + w_2 f_2(s, a) + \dots + w_n f_n(s, a)$$

An Stelle der Q-Values selbst werden nun die Gewichte zu ihrer Berechnung gelernt. Ein Update nachdem in Zustand  $s$  Aktion  $a$  ausgeführt wurde könnte dabei so aussehen:

$$w_i \leftarrow w_i + \alpha \left[ r + \gamma \max_{a' \in A} Q(s', a') - Q(s, a) \right] \cdot f_i(s, a)$$

Natürlich hängt hierbei der Erfolg der Methode sehr stark von den gewählten Features ab.

## 3.7 Andere Algorithmen

### 3.7.1 SARSA

### 3.7.2 Double Q-learning

## 3.8 Deep Q-Learning

## 4 Fazit

## 5 Kurzgeschichte (?)

## 6 Quellen