

Reinforcement learning

Author: Alexander Camek

Betreuer: Prof. Kramer

03.02.2004

Inhalt

Einführung in den Bereich "Learning"

Learning

Reinforcement Learning

Einführung in Reinforcement Learning

passives Lernen in bekannter Umgebung

passives Lernen in unbekannter Umgebung

aktives Lernen in unbekannter Umgebung

Exploration

Q-Values und Aktion-Wertepaare

Generalisierung in Reinforcement Learning

Beispiele

Militärische Anwendungen

Roboterkontrolle

Schach

Backgammon

Zusammenfassung

Einführung in den Bereich Learning

Was heißt eigentlich "Learning" ? Was versteht man darunter ?

Definition

Lernen bedeutet sich besser zu verhalten auf Grund von Erfahrungen

Unterteilung:

- supervised learning
- reinforcement learning

Inhalt

Einführung in den Bereich "Learning"

Learning

Reinforcement Learning

Einführung in Reinforcement Learning

passives Lernen in bekannter Umgebung

passives Lernen in unbekannter Umgebung

aktives Lernen in unbekannter Umgebung

Exploration

Q-Values und Aktion-Wertepaare

Generalisierung in Reinforcement Learning

Beispiele

Militärische Anwendungen

Roboterkontrolle

Schach

Backgammon

Zusammenfassung

Reinforcement Learning

Grundidee ist die Belohnung oder Bestrafung für das erreichte Ziel.
Egal, ob dies vom Lehrer oder vom Programm selbst kommt.

Unterschied:

- Es gibt nur Belohnung oder Bestrafung
- Es wird **NICHT** gesagt,
was schlecht war oder
wie es besser gemacht werden könnte.

Nachteil: Der Agent weiß nicht was er falsch gemacht hat!

Reinforcement Learning

Folgende Aufgabenstellungen sind dabei zu beachten:

- Die Umgebung kann zugänglich oder unzugänglich sein:
 - zugänglich: Zustände und Wahrnehmungen sind verknüpft
 - unzugänglich: innere Zustände bilden Umgebung ab
- vorgegebenes Wissen über Umgebung und Ergebnisse oder Lernen eines Modells und zusätzlichem Wissen
- Belohnung für erreichte Ziele oder für irgendeinen Zustand
- Belohnungen:
 - Teil der zu erfüllenden Aufgabe (z.B. Punkte im Tischtennis)
 - Hilfestellungen für die Aufgabe (z.B. "guter Zug")
- Art des Lernens:
 - passiv
 - aktiv

Reinforcement Learning

Der Agent lernt eine Nutzenfunktion von Zuständen.
Benutzt diese, um mittels Aktionen seinen erwarteten Nutzen zu maximieren.

Der Agent lernt eine action-value Funktion.
Diese berechnet den erwarteten Nutzen einer Aktion in einem Zustand.
(auch "Q-learning" genannt)

Vorteil:

- Kenntnisse nur über die erlaubten Bewegungen
- einfacheres Design

Nachteil:

- nicht vorausschauend, Zielrichtung nicht bekannt
- Lernfähigkeit eingeschränkt

passives Lernen in bekannter Umgebung

Geschieht durch das Generieren von Zustandsübergängen ähnlich wie beim Automaten.

Modell: M_{ij} , Wahrscheinlichkeit eines Zustandsübergangs von i nach j

Ziel:

- Erhalt des erwarteten Nutzens in einem bestimmten Zustand
- Erhalt des Nutzen einer Trainingssequenz

Definition: reward-to-go

Die Summe der Belohnung von diesem Zustand bis zum Endzustand

⇒ erwarteter Nutzen eines Zustandes = erwarteter reward-to-go dieses Zustandes

passives Lernen in bekannter Umgebung

function PASSIVE-RL-AGENT(e) **returns** an action

static: U , a table of utility estimates

N , a table of frequencies for states

M , a table of transition probabilities from state to state

$percepts$, a percept sequence (initially empty)

add e to $percepts$

increment $N[STATE[e]]$

$U \leftarrow UPDATE(U, e, percepts, M, N)$

if TERMINAL?[e] **then** $percepts \leftarrow$ the empty sequence

return the action $Observe$

passives Lernen in bekannter Umgebung

$$U \leftarrow \text{UPDATE}(U, e, \text{percepts}, M, N)$$

Das Wichtigste bei Reinforcement Learning ist der UPDATE-Algorithmus

3 Mögliche Ansätze:

- Naiveransatz mittels LMS (least mean square)
- adaptive dynamic programming (ADP)
- temporal difference (TD)

passives Lernen in bekannter Umgebung

Naiver Ansatz:

Idee:

beobachteter reward-to-go jedes Zustandes
läßt sich auf den im Moment
erwarteten reward-to-go direkt zurückschließen.

Ziel:

- Ermittlung des beobachteten reward-to-go
- Ermittlung des geschätzten Nutzens

LMS erzeugt einen Nutzenschätzer mit minimalen durchschnittlichen quadratischen Fehler in Bezug auf die beobachteten Daten.

⇒ Reduzierung auf induktives Lernen

passives Lernen in bekannter Umgebung

function LMS-UPDATE($U, e, percepts, M, N$) **returns** an updated U

if TERMINAL? $[e]$ **then** $reward\text{-}to\text{-}go \leftarrow 0$

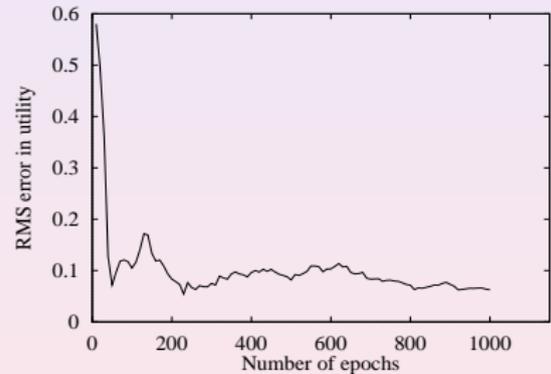
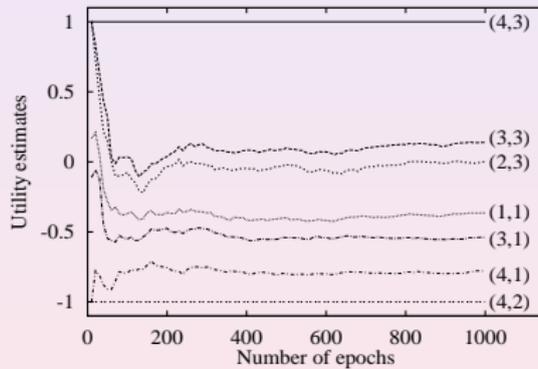
for each e_i **in** $percepts$ (starting at end) **do**

$reward\text{-}to\text{-}go \leftarrow reward\text{-}to\text{-}go + REWARD[e_i]$

$U[STATE[e_i]] \leftarrow RUNNING\text{-}AVERAGE(U[STATE[e_i]], reward\text{-}to\text{-}go, N[STATE[e_i]])$

end

passives Lernen in bekannter Umgebung



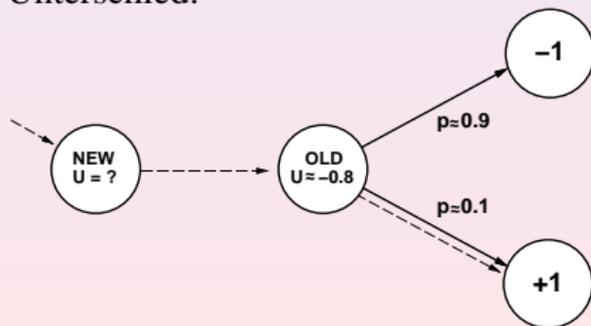
passives Lernen in bekannter Umgebung

Adaptive dynamic programming (ADP)

Idee:

- Nutzen des Wissens über die Umgebung
- Nutzen eines Modells

Unterschied:



bei LMS hätte NEW einen Nutzen von 1

passives Lernen in bekannter Umgebung

Lösung: Speicherung der Übergangswahrscheinlichkeiten in Tabelle M_{ij}

⇒ Das Problem reduziert sich zu einem wohldefinierten sequentiellen Entscheidungsproblem, sobald alle Zustände und deren Belohnungen bekannt sind

Nutzenberechnung: $U(i) = R(i) + \sum_j M_{ij}U(j)$

passives Lernen in bekannter Umgebung

Temporal difference learning (TD)

Idee:

- $U(i) = R(i) + \sum_j M_{ij}U(j)$
- ohne dies für alle möglichen Zustände zu tun

Lösung:

Nutzen der beobachteten Übergänge zur Wertanpassung der beobachteten Zustände

Beispiel:

- Übergang $i \rightarrow j \Rightarrow U(i) = -0.5$ und $U(j) = +0.5$
- Verbesserung von $U(i)$, so daß der Wert besser zum Nachfolger paßt

$$U(i) \leftarrow U(i) + \alpha(R(i) + U(j) - U(i))$$

passives Lernen in bekannter Umgebung

function TD-UPDATE($U, e, \text{percepts}, M, N$) **returns** the utility table U

if TERMINAL? $[e]$ **then**

$U[\text{STATE}[e]] \leftarrow \text{RUNNING-AVERAGE}(U[\text{STATE}[e]], \text{REWARD}[e], N[\text{STATE}[e]])$

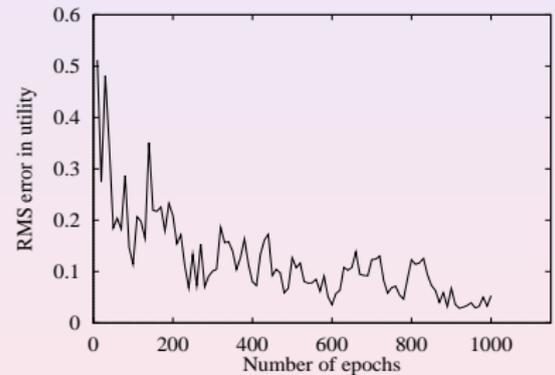
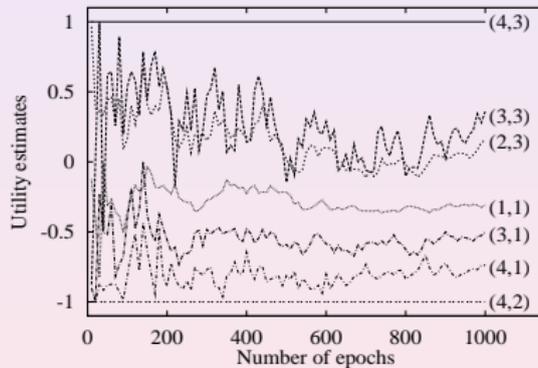
else if percepts contains more than one element **then**

$e' \leftarrow$ the penultimate element of percepts

$i, j \leftarrow \text{STATE}[e'], \text{STATE}[e]$

$U[i] \leftarrow U[i] + \alpha(N[i])(\text{REWARD}[e'] + U[j] - U[i])$

passives Lernen in bekannter Umgebung



passives Lernen in unbekannter Umgebung

Unterschied: Modell wird angepaßt

⇒ ADP-Ansatz funktioniert besser, da dieser auf dem gesamten Modellraum arbeitet

Lösung für TD-Nutzung:

- durchschnittliche Anpassung über eine große Menge von Zustandsübergängen, da Häufigkeiten jedes Nachfolgers näherungsweise proportional zu seiner Wahrscheinlichkeit ist
- Benutzung eines Umgebungsmodells zur Generierung von pseudo-Experimenten
⇒ TD → ADP

Nachteil: höherer Berechnungsaufwand

passives Lernen in unbekannter Umgebung

Anpassung des ADP-Ansatzes:

Approximation: Begrenzung der Anpassungsschritte

Heuristic: prioritized-sweeping

prioritized-sweeping:

- Hinzufügen von Wissen über Vorgänger
- Priorisierung der wahrscheinlichsten Nachfolger
- Aktualisierung des Zustands mit höchster Priorität
- Ablauf:

Step1 $U_{old} = U(i)$

Step2 $U(i) = R(i) + \sum_j M_{ij}U(j)$

Step3 Zustandspriorität O

Step4 $\Delta = |U_{old} - U(i)|$

Step5 Δ nutzen zur Prioritätsänderung bei Vorgänger

⇒ approximierter ADP \approx voller ADP

Vorteil: keine aufwendigen Iterationen am Anfang

aktives Lernen in unbekannter Umgebung

Unterschied:

- Welche Aktionen werden ausgeführt
- Was sind die Ergebnisse
- Wie wirkt das ganze auf die Belohnung

Lösung zur Wandlung passiv in aktiv:

- Modell der Übergänge wird geändert.
Zustand j wird aus i erreicht mittels Aktion a .
kurz: M_{ij}^a
- Berücksichtigung der Wahlmöglichkeiten
Ein rationaler Agent will grundsätzlich Nutzenmaximierung
Deshalb: $U(i) = R(i) + \sum_j M_{ij} U(j) \rightarrow U(i) = R(i) + \max_a \sum_j M_{ij}^a U(j)$
- Der Agent MUSZ in jedem Schritt eine Aktion ausführen

aktives Lernen in unbekannter Umgebung

function ACTIVE-ADP-AGENT(e) **returns** an action

static: U , a table of utility estimates

M , a table of transition probabilities from state to state for each action

R , a table of rewards for states

$percepts$, a percept sequence (initially empty)

$last-action$, the action just executed

add e to $percepts$

$R[STATE[e]] \leftarrow REWARD[e]$

$M \leftarrow UPDATE-ACTIVE-MODEL(M, percepts, last-action)$

$U \leftarrow VALUE-ITERATION(U, M, R)$

if TERMINAL?[e] **then**

$percepts \leftarrow$ the empty sequence

$last-action \leftarrow PERFORMANCE-ELEMENT(e)$

return $last-action$

aktives Lernen in unbekannter Umgebung

Änderungen in den bisherigen Algorithmen:

ADP:

- Änderung des M_{ij} -Modells zum M_{ij}^a -Modell
- Hinzufügen einer dritten Dimension
- Änderung auch bei der UPDATE-Funktion

TD:

- Das Modell ist dasselbe wie beim ADP
- UPDATE-Funktion muß **NICHT** geändert werden
- $TD \equiv ADP$, wenn die Anzahl der Trainingssätze gegen ∞ gehen

Exploration

last action \leftarrow PERFORMANCE-ELEMENT(e)

Frage: Was soll als Ergebnis zurückgeliefert werden?

Klar: Der Agent sollte die höchst nützlichste Aktion auswählen

Aber: Was ist mit dem Lernen?

Gedanken:

- Erhalt der Belohnung für den aktuellen Trainingsatz
- Ermöglichen des Lernens für den Erhalt zukünftiger Belohnungen

Problem: Zwiespalt zwischen sofortigem Gewinn und dem Lernen

Deshalb: Finden einse Mittelweges

Exploration

greedy-Ansatz:

Optimierung durch Auswahl von Aktionen
mit nicht starkem negativen Einfluß

wacky-Ansatz:

z.B. Boltzmann exploration

Treffen der Auswahl mittels Zufall

$ER(a)$ = erwartete Belohnung bei Ausführung von Aktion a

$$\Rightarrow P(a) = \frac{e^{ER(a)/T}}{\sum_{\hat{a} \in A} e^{ER(\hat{a})/T}}$$

Idee:

- Kombination aus beidem
- Änderung der Nutzenfunktion
- unerkundeter Zustand = hoher Nutzen

Exploration

$\Rightarrow U^+(i)$ = optimistischer Schätzer

$N(a, i)$ = Anzahl, wie oft a in i ausgeführt wurde

\Rightarrow ADP mit Wertiteration: $U^+(i) \leftarrow R(i) + \max_a f\left(\sum_j M_{ij}^a U^+(j), N(a, i)\right)$

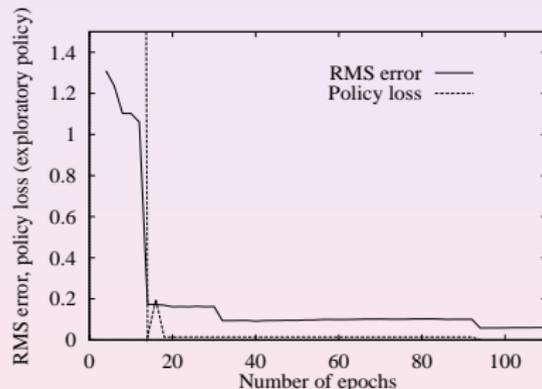
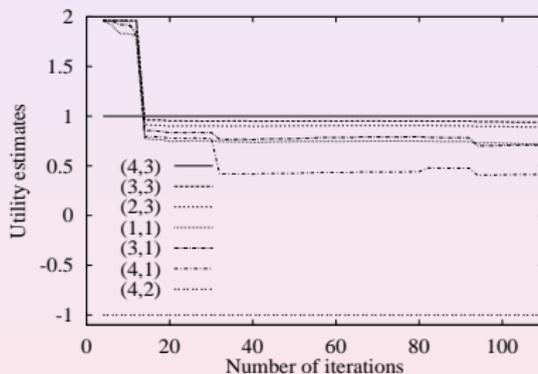
$f(u, n)$ = Explorationsfunktion,

$$\text{z.B. } f(u, n) = \begin{cases} R^+, n < N_e \\ U, \text{sonst} \end{cases}$$

R^+ = optimistischer Schätzer für bestmöglichen Gewinn

N_e = fest, Grenze wie oft Aktions-Zustandspaar ausprobiert wird.

Exploration



Q-Values und Aktion-Wertepaare

Definition: "action-value function"

Ordnet einem erwarteten Nutzen
eine Durchführung von einer Aktion in einem Zustand zu.

⇒ Q-values, kurz: $Q(a, i)$

Q-values sind direkt mit Nutzenwerte verbunden: $U(i) = \max_a Q(a, i)$

Wichtig:

- Genügen für Entscheidung ohne Modell
- Lernbar durch Belohnungsfeedback

Bei Korrektheit **muß** folgendes gelten: $Q(a, i) = r(i) + \sum_j M_{ij}^a \max_{\acute{a}} Q(\acute{a}, i)$

Dies kann als Updatefunktion beim ADP benutzt werden.

Q-Values und Aktion-Wertepaare

Beim TD Q-learning gilt:

$$Q(a, i) \leftarrow Q(a, i) + \alpha(R(i) + \max_{\acute{a}} Q(\acute{a}, j) - Q(a, i))$$

SARSA Methode:

$$Q(a, i) \leftarrow Q(a, i) + \alpha(R(i) + Q(\acute{a}, j) - Q(a, i))$$

Unterschied:

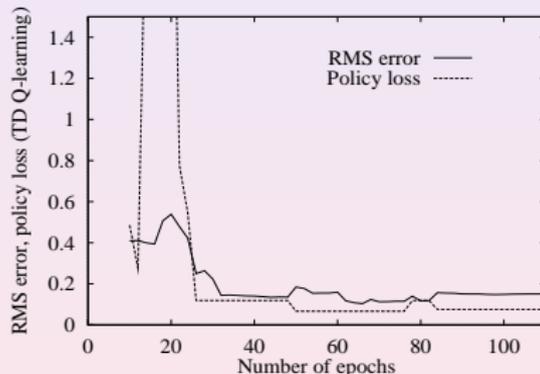
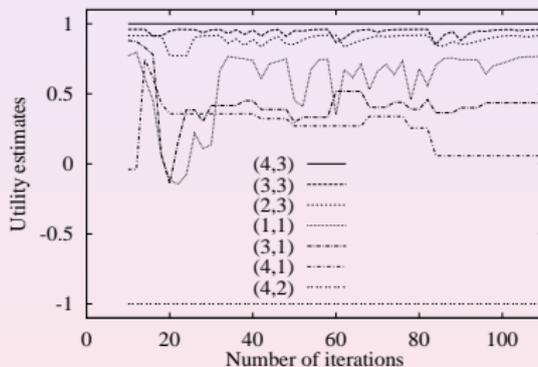
- SARSA versucht den besten Q-Value zu finden
- Q-learning lernt einen deterministischen optimalen Q-Value

Q-Values und Aktion-Wertepaare

```
function Q-LEARNING-AGENT(e) returns an action
  static: Q, a table of action values
           N, a table of state-action frequencies
           a, the last action taken
           i, the previous state visited
           r, the reward received in state i

  j ← STATE[e]
  if i is non-null then
    N[a, i] ← N[a, i] + 1
    Q[a, i] ← Q[a, i] + α(r + maxa' Q[a', j] - Q[a, i])
  if TERMINAL?[e] then
    i ← null
  else
    i ← j
    r ← REWARD[e]
  a ← arg maxa' f(Q[a', j], N[a', j])
  return a
```

Q-Values und Aktion-Wertepaare



Generalisierung in Reinforcement Learning

Stand:

- Agent(U, M, R, Q)
- interne Darstellung mittels Tabellen
- explizite Darstellung einer Ausgabe für jedes Eingabetupel

Nachteil:

bei großen Zustandsmengen steigt die Zeit für die Konvergenz und der Iteration rapide

Beispiel: Schach oder Backgammon besitzt 10^{50} bis 10^{120} Zustände

Generalisierung in Reinforcement Learning

implizite Darstellung

- Erlaubt die Berechnung der Ausgabe pro Eingabe
- z.B. Brettspiel: $U(i) = w_1 f_1(i) + w_2 f_2(i) + \dots + w_n f_n(i)$
⇒ Reduzierung von 10^{120} Werte auf n Gewichte
(bei Schach 10 Gewichte)

Vorteil:

- Kompression
- Rückschluß von besuchten Zuständen auf nicht besuchte Zustände
⇒ input generalization

Nachteil:

- Keine Funktion findbar, die der Nutzenfunktion ähnlich ist
- Größe des Suchraums = Lernzeit der Funktion

Generalisierung in Reinforcement Learning

Auswirkung auf die bisherigen Funktionen

LMS:

- Ende jeder Trainingssquenz wird ein reward-to-go mit jedem besuchten Zustand verknüpft
⇒ $\langle \text{Zustand}, \text{Belohnung} \rangle$
- Verwendung des Paares als benanntes Beispiel für einen induktiven Lernalgorithmus
⇒ Ergebniss ist eine Nutzenfunktion: $U(i)$

Generalisierung in Reinforcement Learning

TD:

- U- und/oder Q-Tabelle ersetzen durch implizite Darstellung
 - Rückgabewert der Updatefunktion ist Beispiel für Lernalgorithmus
 - gelernte Funktion wird beim nächsten Updatedurchlauf benutzt
⇒ Lernalgorithmus ist inkrementell
-
- Lernfunktion nutzt einen Gewichtsvektor \mathbf{w}
 - Gewichte werden angepaßt, so daß zeitlicher Abstand zwischen Zustand und Nachfolger verringert wird
 - neue Update-Regel: $\mathbf{w} \leftarrow \mathbf{w} + \alpha[r + U_w(i)]\nabla_{\mathbf{w}} U_w(i)$
 $U_w(i)$ = Nutzenfunktion
⇒ ähnliche Funktion kann beim Q-learning angewandt werden

Inhalt

Einführung in den Bereich "Learning"

Learning

Reinforcement Learning

Einführung in Reinforcement Learning

passives Lernen in bekannter Umgebung

passives Lernen in unbekannter Umgebung

aktives Lernen in unbekannter Umgebung

Exploration

Q-Values und Aktion-Wertepaare

Generalisierung in Reinforcement Learning

Beispiele

Militrische Anwendungen

Roboterkontrolle

Schach

Backgammon

Zusammenfassung

Militärbereich

- Simulation (z.B. Flugsimulationen)
- logistischen Planung
 - Ressourcenplanung
 - Zeitplanung
- autonomen Entscheidungen
 - Zielsuche gegen reagierende Ziele
 - Trajectoryoptimierung in wechselnden Umgebungen
 - Sensorkontrolle und dynamische Ressourcenbelegung
- automatische Entscheidungen
 - schnelle Reaktion (z.B. elektronische Kriegsführung)
 - low-level Kontrolle (z.B. Kontrolle für Roboter mit Füßen)

Roboterbereich

- Lernen des Balanzierens und der Koordination beim Laufen oder Jonglieren
- Roboterfußball (z.B. Clockwork–Orange Team)
- Kartierung von Innenräumen und der Lokalisation

Schach

Die erste Anwendung von Reinforcement Learning war das Schachprogramm von Arthur Samuel (1959;1967)

Bestandteile des Programms:

- weighted linear function zur Positionsevaluierung
- Gewichte wurden angepat
- bis zu 16 Terme genutzt

Unterschied:

- Anpassung der Gewichte mittels Unterschied zwischen aktuellem und gemerktem Wert
- Vorausschau im Suchbaum
- keine Benutzung von ermittelten Belohnungen
(Aber: Voraussetzung, da Gewichte immer positiv)

Backgammon

Neurogammon (Tesauro und Sejnowski, 1989)

- Neuronalesnetzwerk mit $Q(a,i)$
- Vorkenntnisse im Backgammonspiel
- Beispiele mit bewerteten Zügen
- Gewann die Backgammon–Olympiade 1989
- Nachteil: sehr gute Spieler konnten es besiegen

Backgammon

TD-Gammon (Tesauro, 1992)

Algorithmus:

Kombination aus $TD(\lambda)$ und nichtlinearer Funktionsapproximation

Problem:

- 30 Spielsteine und 24 Positionen
(26 mit Rahmen und dem Auenbereich)
- 20 verschiedene Wege je nach Wurf
- Reaktion auf Gegner

⇒ Spielbaum hat ber 400 Verzweigungen und zu viele Zustnde

Backgammon

Design:

- Belohnung nur am Ende des Spiels
- Neuronalesnetzwerk mit 198 Eingabeknoten und 40–80 versteckten Knoten
- Eingabe: Backgammonposition
- Ausgabe: Einschatzung der Position

Jeder Punkt im Spielfeld wird dargestellt mittels 4 Eingabeknoten pro Farbe.

kein Stein: alle 4 Knoten sind 0

ein Stein: 1.Knoten ist 1

⇒ zwei Steine: 1. und 2.Knoten ist 1

drei Steine: 1., 2. und 3.Knoten ist 1

vier Steine: 4.Knoten ist $\frac{(n-3)}{2}$ (if $n > 3$)

mit n = Anzahl der Steine auf dem Punkt.

Backgammon

- 4 weie Knoten + 4 schwarze Knoten fr 24 Positionen \Rightarrow 192 Knoten
 - + 2 Knoten fr die Anzahl der Spielsteine auf dem Spielfeldrand
(Knotenwert = $\frac{n}{2}$)
 - + 2 Knoten fr die Anzahl der entfernten Spielsteine
(Knotenwert = $\frac{n}{15}$)
 - + 2 Knoten fr die Anzeige des aktuellen Zuges
- \Rightarrow 198 Knoten

Backgammon

Berechnung im Neuronalen-Netz:

Zuerst werden die Eingaben mit ihren Gewichten multipliziert und dann bei den verdeckten Knoten aufsummiert.

Ausgabe des verdeckten Knoten j — $h(j)$ — ist eine nichtlinear Sigmoid-Funktion:

$$h(j) = \sigma\left(\sum_i w_{ij}\phi(i)\right) = \frac{1}{1+e^{-\sum_i w_{ij}\phi(i)}}$$

mit $\phi(i)$ = Wert des i .Eingabeknotens
und w_{ij} = Gewicht von Eingabe- zum verdeckten Knoten.

Der Ausgabeknoten generiert die gewichtete Summe und schickt diese durch die selbe sigmoid Funktion.

⇒ Einschtzung der Position

Backgammon

Lernen:

$$\text{Update-Regel: } \vec{\theta} = \vec{\theta} + \alpha [r + U(j) - U(i)] \vec{e}$$

$$\vec{e} = \lambda \vec{e} + \nabla_{\vec{\theta}} U(i)$$

$\vec{\theta}$ = Vektor der Gewichte (zu Beginn randomisiert)

\vec{e} = Vector der Auswahlspuren fur jede Komponente von $\vec{\theta}$

TD-Gammon lernte durch spielen mit sich selbst.

Dabei wahlte es Zuge indem jeder der 20 Wege pro Wurf und die Ergebnispositionen brucksichtigt wurden.

Danach wurde der Zug mit dem hochsten Nutzen ausgewahlt.

Backgammon

Netzgewichte wurden am Anfang zufllig gewhlt.

⇒ erste Spiele dauerten hunderte bis tausende Zge bis eine Seite durch Zufall gewann.

Nach einem dutzend Spiele stieg die Performance rapide.

TD-Gammon besiegt somit Neurogammon.

Sptere Versionen hatten eine zwei/drei-schichtige Suchprozedur.

Auswahl eines Zuges erfolgt durch berlegung, welchen Zug Gegner spielen knnte.

Zusammenfassung

Wegen seinem Potential zur Eliminierung von menschlichem Code bei Kontrollstrategien ist Reinforcement Learning eine der wichtigsten Gebiete der machine learning Forschung.

Starke Einsetzung im Bereich der autonomen Roboter.

Problem:

- Erweiterung dieser Methoden
- unzugänglichen Umgebungen

Zusammenfassung

Vielen Dank
für
Ihre Aufmerksamkeit!

Literaturhinweise

- Artificial Intelligence – A Modern Approach
Stuart Russel und Peter Norvig
- Reinforcement Learning
Sutton und Baro
- Reinforcement Learning – A Survey
Kaelbling, Littman und Moore
- Online Fitted Reinforcement Learning
Geoffrey Gordon

Literaturhinweise

- Reinforcement Learning with Perceptual Aliasing
Loonie Chrisman
- Packet Routing in Dynamically Changing Networks
Boyan und Littman
- Reinforcement Learning Methods for Military Applications
Malcom Strens
- A distributed dynamic world model for Robot soccer
Frans Groen, Jeroen Roodhart und Joost Vunderink