

DIPLOMARBEIT

XCS in dynamischen Multiagenten-Überwachungsszenarien ohne globale Kommunikation

von

Clemens Lode

Institut für Angewandte Informatik
und Formale Beschreibungsverfahren
Universität Karlsruhe (TH)

Referent: Prof. Dr. Hartmut Schmeck
Betreuer: Dipl. Wi.-Ing. Urban Richter

Karlsruhe, 30.03.2009

Inhaltsverzeichnis

1	Einführung	1
2	Beschreibung des Szenarios	3
2.1	Definition einer Problem Instanz	4
2.2	Sichtbarkeit von Objekten	5
2.3	Kollaboration	6
2.4	Dynamik	6
2.5	Startkonfigurationen des Torus	7
2.5.1	Leeres Szenario	7
2.5.2	Szenario mit zufällig verteilten Hindernissen	8
2.5.3	Säulen Szenario	9
2.5.4	„Kreuz Szenario“	11
2.5.5	„Raum Szenario“	11
2.5.6	Schwieriges Szenario	11
2.5.7	„Irrgarten Szenario“	13
2.6	Bestimmung der Qualität eines Algorithmus	13
3	Eigenschaften der Agenten	15
3.1	Sensoren eines Agenten	16
3.1.1	Aufbau eines Sensordatenpaars	16

3.1.2	Aufbau eines Sensordatensatzes	17
3.2	Grundsätzliche Algorithmen der Agenten	19
3.2.1	Algorithmus mit zufälliger Bewegung	19
3.2.2	Einfache Heuristik	19
3.2.3	Intelligente Heuristik	21
4	Das Zielobjekt	25
4.1	Basiseigenschaften	25
4.2	Typen von Zielobjekten	26
4.2.1	Typ „Zufälliger Sprung“	26
4.2.2	Typ „Zufällige Bewegung“	27
4.2.3	Typ „Einfache Richtungsänderung“	27
4.2.4	Typ „Beibehaltung der Richtung“	28
4.2.5	Typ „Intelligentes Verhalten“	29
4.2.6	Typ „SXCS“	30
5	Ablauf der Simulation	31
5.1	Hauptschleife	31
5.2	Reihenfolge der Ausführung (<i>doOneMultiStepProblem()</i>)	31
5.3	Messung der Qualität	34
5.4	Reihenfolge der Ermittlung des <i>base reward</i>	35
5.5	Zusammenfassung	36
5.6	Implementierung eines Problemablaufs	37
6	Erste Analyse der Agenten ohne XCS	41
6.1	Statistische Merkmale	41
6.1.1	Abdeckung	42
6.2	Zielobjekt mit zufälligem Sprung	43

6.2.1	Szenario ohne Hindernisse	43
6.2.2	Säulenszenario	44
6.2.3	Zufällig verteilte Hindernisse	44
6.3	“Zufälliger Nachbar” und “Einfache Richtungsänderung”	46
6.4	“Intelligent Open” und “Intelligent Hide”	47
6.5	Always Same Direction	48
6.6	XCS	49
6.7	Zusammenfassung	49
7	XCS	51
7.1	Übersicht	56
7.2	Ablauf eines XCS	56
7.2.1	Covering	57
7.2.2	Variable <i>lastMatchSet</i>	57
7.2.3	Variable <i>actionSet</i>	57
7.3	Classifier	58
7.3.1	Der <i>condition</i> Vektor	58
7.3.2	Platzhalter im <i>condition</i> Vektor	59
7.3.3	Vergleich des <i>condition</i> Vektors mit den Sensordaten	59
7.3.4	Der <i>action</i> Wert	60
7.3.5	Der <i>fitness</i> Wert	60
7.3.6	Der <i>reward prediction</i> Wert	61
7.3.7	Der <i>reward prediction error</i> Wert	61
7.3.8	Der <i>experience</i> Wert	61
7.3.9	Der <i>numerosity</i> Wert	61
7.4	Subsummation von <i>classifier</i>	62
7.5	Genetische Operatoren	62

7.6	Bewertung der Aktionen (<i>base reward</i>)	63
8	Parameter	65
8.1	Parameter <i>max population</i> N	66
8.2	Maximalwert <i>reward</i>	69
8.3	Parameter <i>accuracy equality</i> ϵ_0	69
8.4	Parameter <i>reward prediction discount</i> γ	70
8.5	Parameter Lernrate β	71
8.6	Parameter <i>reward prediction init</i> p_i	71
8.7	Zufällige Initialisierung der <i>classifier set</i> Liste	73
8.8	Übersicht über alle Parameterwerte	74
8.9	Auswahlart der <i>classifier</i>	76
8.9.1	Auswahlart <i>tournament selection</i>	77
8.9.2	Wechsel zwischen den <i>explore</i> und <i>exploit</i> Phasen	78
9	XCS Varianten	81
9.1	Allgemeine Anpassungen und Verbesserungen	82
9.1.1	Verschiedenes, Numerosity, TODO	82
9.2	Standard XCS Multistepverfahren	83
9.3	XCS Variante für Überwachungsszenarien (SXCS)	89
9.3.1	Ereignisse	90
9.3.2	Implementierung von SXCS	91
9.3.3	Zielobjekt mit SXCS	95
10	Analyse SXCS	99
10.1	Vergleich unterschiedlicher Geschwindigkeiten des Zielobjekts	100
10.2	Zusammenfassung der bisherigen Erkenntnisse	101
10.3	Standard XCS Multistepverfahren	102

10.3.1 SXCS und Heuristiken	102
10.3.2 Vergleich Multistep / LCS	102
10.3.3 Test der verschiedenen Exploration-Modi	102
11 Kommunikation	103
11.1 Realistischer Fall mit Kommunikationsrestriktionen	103
11.2 Lösungen aus der Literatur	104
11.3 SXCS Variante mit verzögerter Reward (DSXCS)	105
11.4 Ablauf	106
11.5 Kommunikationsvarianten	111
11.5.1 Einzelne Gruppe	112
11.5.2 Gruppenbildung über Ähnlichkeit des Verhaltens der Agenten . . .	113
11.6 Bewertung Kommunikation:	116
11.6.1 Vergleich TODO	116
12 Zusammenfassung, Ergebnis und Ausblick	119
12.1 Zusammenfassung	119
12.2 Ergebnis	120
12.3 Ausblick	120
13 Verwendete Hilfsmittel und Software	123
13.1 Beschreibung des Konfigurationsprogramms	124
A Statistical significance tests	125
B Implementation	127

Abbildungsverzeichnis

2.1	„Leeres Szenario“ ohne Hindernisse	8
2.2	Szenario mit zufällig verteilten Hindernissen mit $\lambda_h = 0.05$	9
2.3	Szenario mit zufällig verteilten Hindernissen mit $\lambda_h = 0.1$	9
2.4	Szenario mit zufällig verteilten Hindernissen mit $\lambda_h = 0.2$	10
2.5	Szenario mit zufällig verteilten Hindernissen mit $\lambda_h = 0.4$	10
2.6	Startzustand des Säulen Szenarios	10
2.7	Kreuz Szenario	11
2.8	„Raum Szenario“	12
2.9	Schwieriges Szenario	12
2.10	„Irrgarten Szenario“	13
3.1	Sicht- und Überwachungsreichweite eines Agenten	18
3.2	Sich zufällig bewogender Agent	20
3.3	Agent mit einfacher Heuristik	21
3.4	Agent mit intelligenter Heuristik	24
4.1	Zielobjekt mit maximal einer Richtungsänderung	27
4.2	Bewegungsform „Beibehaltung der Richtung“: Zielobjekt das sich, wenn möglich, immer nach Norden bewegt	28

4.3	Sich intelligent verhaltendes Zielobjekt der Agenten und Hindernissen aus- weicht	29
7.1	Einfaches Beispiel zum XCS <i>multi step</i> Verfahren	53
7.2	Vereinfachte Darstellung eines <i>classifier set</i> für das Beispiel zum XCS <i>multi</i> <i>step</i> Verfahren	53
8.1	Auswirkung der Torusgröße auf die Laufzeit (leeres Szenario)	67
8.2	Auswirkung des Parameters <i>max population N</i> auf Laufzeit (leeres Szenario)	67
8.3	Verhältnis Laufzeit zu <i>max population N</i> (leeres Szenario)	68
8.4	Auswirkung des Parameters <i>max population N</i> auf Qualität (leeres Szenario)	68
8.5	Auswirkung des Parameters <i>accuracy equality</i> ϵ_0 auf die Qualität (Säulens- zenario)	70
8.6	Auswirkung des Parameters <i>learning rate</i> β auf Qualität (Säulenszenario) .	72
8.7	Auswirkung des Parameters <i>reward prediction init</i> p_i auf Qualität (Säulens- zenario)	73
9.1	Schematische Darstellung der Rewardverteilung an ActionSets	90
9.2	Schematische Darstellung der zeitlichen Rewardverteilung an und der Spei- cherung von ActionSets	92
9.3	Schematische Darstellung der Rewardverteilung an ActionSets bei einem neutralen Ereignis	93
10.1	Vergleich der Qualitäten verschiedener Algorithmen bezüglich der Geschwin- digkeit des Zielobjekts	101
11.1	Schematische Darstellung der Rewardverteilung an ActionSets bei einem neutralen Ereignis	114

11.2 Schematische Darstellung der Rewardverteilung an ActionSets bei einem neutralen Ereignis	114
11.3 Schematische Darstellung der Rewardverteilung an ActionSets bei einem neutralen Ereignis	116
11.4 Beispielhafte Darstellung der Kombination interner und externer Rewards	118
13.1 Screenshot des Konfigurationsprogramms	124

Tabellenverzeichnis

6.1	“Total Random” ohne Hindernisse	44
6.2	“Total Random” ohne Hindernisse	45
6.3	“Total Random” mit Hindernisse (12 Agenten)	45
6.4	Vergleich “Zufälliger Nachbar” und “Einfache Richtungsänderung” (12 Agenten, ohne Hindernisse)	47
6.5	Vergleich “Zufälliger Nachbar” und “Einfache Richtungsänderung” (12 Agenten, zufälliges Szenario mit $\lambda_h = 0.2$, $\lambda_p = 0.99$)	48
6.6	Vergleich “Zufälliger Nachbar” und “Einfache Richtungsänderung” (12 Agenten, Säulenszenario)	48
6.7	Vergleich “Intelligent Open” und “Intelligent Hide” (8 Agenten, ohne Hindernisse)	49
6.8	Vergleich “Intelligent Open” und “Intelligent Hide” (8 Agenten, zufälliges Szenario mit $\lambda_h = 0.2$, $\lambda_p = 0.99$)	49
6.9	Vergleich “Intelligent (Open)” und “Intelligent (Hide)” (8 Agenten, Säulenszenario)	50
8.1	Vergleichende Tests für den den Start mit und ohne zufällig gefüllten <i>classifier set</i> Listen	74
8.2	Verwendete Parameter (soweit nicht anders angegeben) und Standardparameter, TODO englisch/deutsch	75

10.1 Vergleich “Intelligent (Open)” und “Intelligent (Hide)” (8 Agenten, Säulenszenario)	99
10.2 Vergleich “Intelligent (Open)” und “Intelligent (Hide)” (8 Agenten, Säulenszenario)	100

Programmverzeichnis

3.1	Berechnung der nächsten Aktion bei der Benutzung des Algorithmus mit zufälliger Bewegung	20
3.2	Berechnung der nächsten Aktion bei der Benutzung der einfachen Heuristik	22
3.3	Berechnung der nächsten Aktion bei der Benutzung der intelligenten Heuristik	23
5.1	Zentrale Schleife für einzelne Experimente	32
5.2	Zentrale Schleife für einzelne Probleme	38
5.3	Zentrale Bearbeitung (Sensordaten und Berechnung der neuen Aktion) aller Agenten und des Zielobjekts innerhalb eines Problems	39
5.4	Zentrale Bearbeitung (Verteilung des Rewards) aller Agenten und des Zielobjekts innerhalb eines Problems	39
5.5	Zentrale Bearbeitung (Ausführung der Bewegung) aller Agenten und des Zielobjekts innerhalb eines Problems	40
7.1	Bestimmung des <i>base reward</i> Werts für Agenten	64
9.1	Korrigierte Version der <i>addNumerosity()</i> Funktion	84
9.2	Erstes Kernstück des Standard XCS Multistepverfahrens (<i>calculateReward()</i> , Bestimmung des Rewards anhand der Sensordaten), angepasst an ein dynamisches Überwachungsszenario	86

9.3	Zweites Kernstück des Multistepverfahrens (<i>collectReward()</i> - Verteilung des Rewards auf die ActionSets), angepasst an ein dynamisches Überwachungsszenario	87
9.4	Drittes Kernstück des Multistepverfahrens (<i>calculateNextMove()</i> - Auswahl der nächsten Aktion und Ermittlung des zugehörigen ActionSets), angepasst an ein dynamisches Überwachungsszenario	88
9.5	Erstes Kernstück des SXCS-Algorithmus (<i>calculateReward()</i> , Bestimmung des Rewards anhand der Sensordaten)	94
9.6	Zweites Kernstück des SXCS-Algorithmus (<i>collectReward()</i> - Verteilung des Rewards auf die ActionSets)	96
9.7	Drittes Kernstück des SXCS-Algorithmus (<i>calculateNextMove()</i> - Auswahl der nächsten Aktion und Ermittlung und Speicherung des zugehörigen ActionSets)	97
9.8	Bestimmung des <i>base rewards</i> für das Zielobjekt	98
11.1	Zweites Kernstück des verzögerten SXCS-Algorithmus (<i>collectReward()</i> - Verteilung des Rewards auf die ActionSets)	107
11.2	Auszug aus dem dritten Kernstück des verzögerten SXCS-Algorithmus (<i>calculateNextMove()</i>)	108
11.3	Viertes Kernstück des verzögerten SXCS-Algorithmus (DSXCS, Verarbeitung des Rewards, <i>processReward()</i>)	109
11.4	Verbesserte Variante des vierten Kernstück des verzögerten SXCS-Algorithmus (DXCS, Verarbeitung des Rewards, <i>processReward()</i>)	110
11.5	“Egoistische Relation“, Algorithmus zur Bestimmung des Kommunikationsfaktors basierend auf dem Verhalten des Agenten gegenüber anderen Agenten	115

Kapitel 1

Einführung

Ein aktuelles Forschungsgebiet aus dem Bereich der *learning classifier systems* (LCS) stellen die sogenannten *accuracy based* LCS (XCS) dar. In der Basis entspricht XCS einem LCS, d.h. eine Reihe von Regeln, bestehend jeweils aus einer Kondition und einer Aktion, werden mittels *reinforcement learning* schrittweise bewertet und an eine Umwelt angepasst. Die Frage nach dem Zeitpunkt der Bewertung teilt die verwendeten Algorithmen bei XCS in *single step* und *multi step* Verfahren ein. Hauptaugenmerk dieser Arbeit soll das *multi step* Verfahren sein, bei dem die Bewertung (der *reward* der Regeln erst nach einigen Schritten verfügbar ist und an zurückliegende Regeln sukzessive weitergeleitet wird um möglichst alle beteiligten Regeln an dem *reward* zu beteiligen.

Bisherige Anwendungen haben sich hauptsächlich auf statische Szenarien mit nur einem XCS oder mit mehreren Agenten mit globaler Organisation und Kommunikation beschränkt. Diese Arbeit hat sich auf die Problemstellung konzentriert, wie man XCS modifizieren sollte, damit es ein dynamisches Überwachungsszenario, mit sich bewegendem Zielobjekt und mehreren Agenten, im Vergleich zu zufälliger Bewegung möglichst gut bestehen.

Neben der Anpassung der Implementation, damit XCS für eine solche Problemstellung

anwendbar ist, wurden weitere Modifikationen durchgeführt, die in einigen Fällen zu deutlich besseren Ergebnissen als die der Standardimplementation führten.

Außerdem wurde untersucht, wie eine einfache Kommunikation ohne globale Steuereinheit stattfinden kann, um das Ergebnis weiter zu verbessern. Im Wesentlichen war dazu eine weitere Anpassung von XCS vonnöten, so dass die Implementierung auch mit (durch die Kommunikation) zeitverzögerten und externen *rewards* arbeiten konnte. Wesentliche Schlußfolgerung ist, dass sich unterschiedliche Szenarien unterschiedlich gut für Kommunikation eignen, dass Kommunikation Möglichkeiten zur Anpassung bietet um mit einer variablen, unbekannten Feldgröße besser zurecht zu kommen und, dass es Szenarien gibt, in denen Kommunikation signifikante Vorteile erbringt.

Erfolgversprechende Ansatzpunkte für weitere Forschung gibt es im Bereich der mathematischen Begründung, warum die Implementierung Vorteile erbringt, im Ausbau der Untersuchung von Kommunikation zwischen den Agenten in Verbindung mit XCS und in der Anwendung der gefundenen Ergebnisse in anderen Problemstellungen ähnlicher Natur.

Kapitel 2

Beschreibung des Szenarios

Im Wesentlichen sollen die Algorithmen, die in dieser Arbeit besprochen werden, in einem Szenario getestet werden, in dem mehrere Agenten ein sich bewegendes Zielobjekt überwachen sollen. Dies soll im folgenden als Überwachungsszenario bezeichnet werden. Die Qualität eines Algorithmus in einem solchen Überwachungsszenario wird anhand des Anteils der Zeit bewertet, in der er mit Hilfe der Agenten das Zielobjekt überwachen konnte, relativ zur Gesamtzeit (siehe Kapitel 2.6).

Verwendetes Umfeld wird ein quadratischer Torus sein, der aus quadratischen Feldern besteht. Jedes bewegliche Objekt auf einem Feld des Torus kann sich in einem Zeitschritt nur auf eines der vier Nachbarmfelder bewegen (mit Ausnahme des Zielobjekts, welches mehrere Bewegungen in einem Zeitschritt durchführen kann, Näheres dazu im Kapitel 4.1). Die Felder können entweder leer oder durch ein Objekt besetzt sein. Besetzte Felder können nicht betreten werden, eine Bewegung auf ein solches Feld schlägt ohne weitere Konsequenzen fehl.

Es gibt drei verschiedene Arten von Objekten: Unbewegliche Hindernisse, ein zu überwachendes Zielobjekt und Agenten. Sowohl das Zielobjekt als auch die Agenten bewegen sich

jeweils anhand eines bestimmten Algorithmus und bestimmter Sensordaten. Eine nähere Beschreibung der Agenten findet sich in Kapitel 3, während die Eigenschaften des Zielobjekts in Kapitel 4 beschrieben wird.

Ziel dieses Kapitels wird vor allem sein, auf Kapitel 6 vorzubereiten, in dem anhand von Tests herausgefunden werden soll, welche der hier vorgestellten Szenarien brauchbare Ergebnisse liefern kann, um zum einen das gestellte Problem an sich, als auch die jeweils erforderlichen Eigenschaften besser verstehen zu können.

Eine separate Beschäftigung mit diesen - relativ einfachen - Szenarien war notwendig, um zum einen das eigene Simulationsprogramm zu testen und zum anderen um vergleichbare Ergebnisse zu erhalten. Ein Rückgriff auf die Literatur war deshalb nicht möglich, insbesondere gibt es keine Arbeiten in Bezug auf XCS mit einer solchen Problemstellung. Zwar entspricht das Standardszenario bei XCS einem Feld, einem Agenten, Hindernissen und einem Ziel, es fehlen jedoch Arbeiten, in denen Sichtbarkeit (die Sichtweite beschränkte sich in der Literatur meist auf angrenzende Felder), Kollaboration (meist war nur ein einzelner Agenten Gegenstand der Untersuchung), Dynamik (meist gab es feste Start- und Zielpunkte) und die Messung der durchschnittlichen Qualität (meist ging es um die Anzahl der Schritte zum Ziel) gemeinsam in einem Szenario betrachtet werden.

Im folgenden sollen nun also auf diese einzelnen Punkte näher eingegangen werden und eine Abgrenzung zu Arbeiten in der Literatur aufgezeigt werden.

2.1 Definition einer Probleminstanz

Eine einzelne Probleminstanz entspricht einem Torus mit einer bestimmten Anfangsbelegung mit bestimmten Objekten und bestimmten Parametern zur Sichtbarkeit. Die An-

fangsbelegung ist über einen *random seed* Wert bestimmt. Soweit nicht anders angegeben, sollen hier Prombleninstanzen der Größe 16x16 Felder betrachtet werden, insbesondere beziehen sich die Ergebnisse der Tests auf diesen Fall.

Jedes Problem soll sich, sofern nicht anders angegeben, über 500 Zeitschritte ziehen. Ein einzelnes Experiment entspricht dem Test einer Anzahl von Probleminstanzen, die jeweils mit einer Reihe von *random seed* Werten initialisiert werden. In einem Durchlauf werden mehrere Experimente (jeweils mit unterschiedlichen Reihen an *random seed* Werten) durchgeführt. Falls nicht anders angegeben sollen die Tests jeweils über 10 Experimente mit jeweils 10 Problemen laufen.

TODO XCS, neustart etc.

2.2 Sichtbarkeit von Objekten

Der Parameter *sight range* bzw. *reward range* einer Probleminstanz bestimmt, bis zu welcher Distanz andere Objekte von einem Objekt als „gesehen“ bzw. „überwacht“ gelten, sofern die Sicht durch andere Objekte nicht versperrt ist. Der Parameter *reward range* ist relevant für die Bewertung der Qualität des Algorithmus (siehe Kapitel 2.6) und wird immer kleiner als der *sight range* Wert gewählt. Über die Sensoren kann ein Agent feststellen, ob sich Objekte in welcher der beiden Reichweiten befinden. Falls nicht anders angegeben sollen jeweils *sight range* auf 5 und *reward range* auf 2 gesetzt werden und in den Abbildungen jeweils der hellblaue Bereich den überwachten und der hell- und dunkelblaue Bereich den gesehenen Bereich darstellen.

2.3 Kollaboration

Wesentliches Hauptaugenmerk der Gestaltung der Szenarien soll Kollaboration sein, d.h. die Aufgabe soll mit Hilfe mehrerer Agenten gemeinsam gelöst werden.

TODO Literatur Definition von Kollaboration in der Literatur, Abgrenzung

Eine erfolgreiche Überwachung soll deswegen so definiert sein, dass sich ein beliebiger Agent in Überwachungsreichweite des Zielobjekts befindet. Angesichts dessen, dass diese Aufgabe auch ein einzelner Agent erfüllen kann, sofern die Geschwindigkeit des Zielobjekts kleiner oder gleich der Geschwindigkeit des Agenten ist, sollen in späteren Tests (insbesondere in Kapitel 10 beim Vergleich unterschiedlicher XCS Varianten und im Kapitel 8 beim Vergleich unterschiedlicher XCS Parameter) unterschiedliche Geschwindigkeiten getestet werden.

Bewegt sich das Zielobjekt zu schnell, werden die Agenten Schwierigkeiten haben, einen Bezug zwischen Sensordaten und eigener Aktionen zu erkennen, bewegt es sich zu langsam, wird das Problem sehr einfach, eine einzelne Regel („Bewege dich auf das Ziel zu,“) würde zur Lösung dann schon genügen.

2.4 Dynamik

Die Szenarien fallen alle unter die Kategorie „dynamisch“. Darunter soll in diesem Zusammenhang verstanden werden, dass es kein festes Ziel gibt, das erreicht werden soll oder kann, das Zielobjekt befindet sich in stetiger Bewegung, wie auch sich andere Agenten in Bewegung befinden können.

Dies ist ein wesentlicher Gesichtspunkt, dass diese Arbeit von vielen anderen unter-

scheidet, Gegenstand der Untersuchung in der Literatur sind eher statische Probleme wie z.B. 6-Multiplexer Problem und Maze1 (z.B. in [But06b]) bzw. Maze5, Maze6, Woods14 (in [BGL05])

El Fazor, Soccer

oder Probleme bei denen die Agenten globale Information besitzen

TODO

Eine nähere Diskussion zur Literatur folgt in Kapitel 7.

2.5 Startkonfigurationen des Torus

Getestet wurden eine Reihe von Szenarien (in Verbindung mit unterschiedlichen Werten für die Anzahl der Agenten, Größe des Torus und Art und Geschwindigkeit des Zielobjekts). TODO Wesentliche Rolle spielt hier die Verteilung der Hindernisse. TODO TODO

In den folgenden Abbildungen repräsentieren rote Felder jeweils Hindernisse, weiße Felder jeweils Agenten und das grüne Feld jeweils das Zielobjekt. Außerdem sind die Sicht- und Überwachungsreichweiten aus Kapitel 2.2, jeweils kreisförmig vom jeweiligen Agenten ausgehend grau den Bereich, der durch die *reward range* abgedeckt wird, und blau den Bereich, der zusätzlich noch durch die *sight range* abgedeckt wird.

2.5.1 Leeres Szenario

In Abbildung (2.1) ist ein Szenario ohne Hindernisse und mit zufälliger Verteilung der Agenten und zufälliger Position des Zielobjekts dargestellt. Im leeren Szenario soll das Verhalten der Agenten in einem Torus ohne Hindernisse untersucht werden.

TODO warum

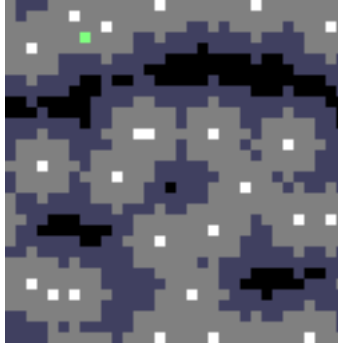


Abbildung 2.1: „Leeres Szenario“ ohne Hindernisse

2.5.2 Szenario mit zufällig verteilten Hindernissen

Zwei Parameter bestimmen das Aussehen des Szenarios mit zufällig verteilten Hindernissen, zum einen der Prozentsatz an Hindernissen an der Gesamtzahl der Felder des Torus (Hindernissanteil λ_h), zum anderen der Grad inwieweit die Hindernisse zusammenhängen (Verknüpfungsfaktor λ_p).

Bei der Erstellung des Szenarios bestimmt λ_p die Wahrscheinlichkeit für jedes einzelne angrenzende freie Feld, dass beim Verteilen der Hindernisse nach dem Setzen eines Hindernisses dort sofort ein weiteres Hindernis gesetzt wird. $\lambda_p = 0.0$ ergäbe somit eine völlig zufällig verteilte Menge an Hindernissen, während ein Wert von 1.0 eine oder mehrere stark zusammenhängende Strukturen schafft. Wird der Prozentsatz an Hindernissen λ_h auf 0.0 gesetzt, dann entspricht diesem dem oben erwähnten leeren Szenario. Ein Wert von 1.0 würde eine völlige Abdeckung des Torus bedeuten und wäre für einen Test somit unbrauchbar. Hier sollen nur geringe Werte bis 0.4 betrachtet werden, wobei später in Tests sich auf Werte bis 0.2 beschränkt wird, da bei großen Hindernissanteil die lokalen Entscheidungen einzelner Agenten zu wichtig werden, da das Zielobjekt sich oft nur in einem kleinen Bereich aufhält

TODO Zielobjektdeckung! (in 100 Schritten abgedecktes Gebiet o.ä.)

In Abbildung (2.2), Abbildung (2.3), Abbildung (2.4) und Abbildung (2.5) werden

Beispiele für zufällige Szenarien gegeben mit $\lambda_h = 0.05, 0.1, 0.2$ bzw. 0.4 und $\lambda_p = 0.01, 0.5$ bzw. 0.99 .

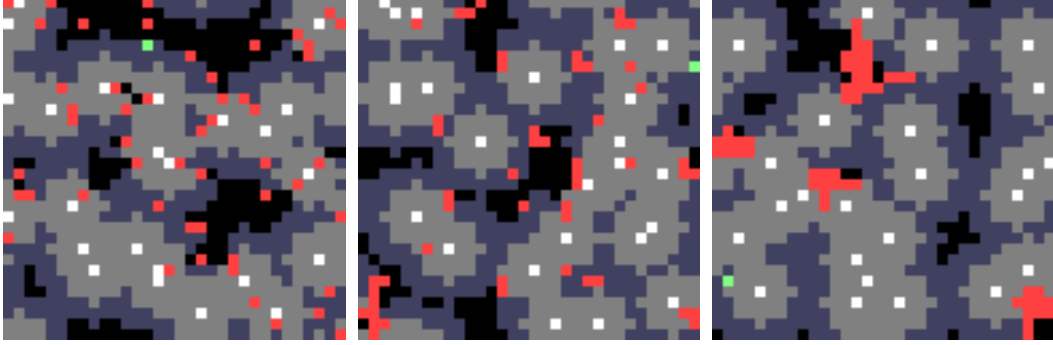


Abbildung 2.2: Szenario mit zufällig verteilten Hindernissen mit Hindernissanteil $\lambda_h = 0.05$ und Verknüpfungsfaktor $\lambda_p = 0.01, 0.5$ bzw. 0.99 .

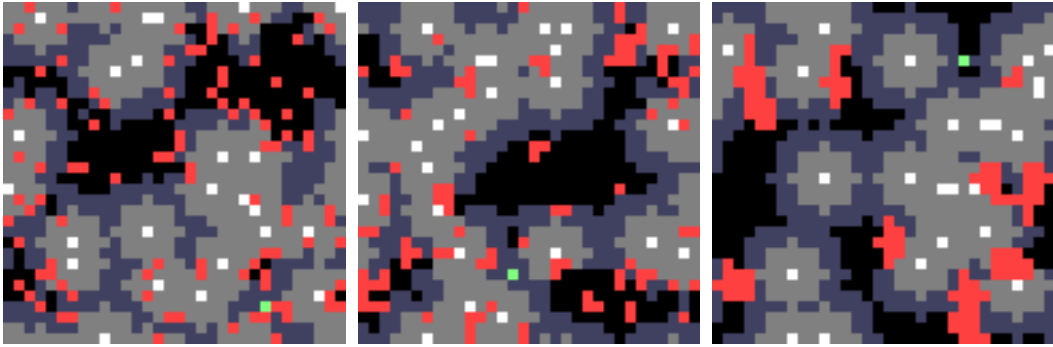


Abbildung 2.3: Szenario mit zufällig verteilten Hindernissen mit Hindernissanteil $\lambda_h = 0.1$ und Verknüpfungsfaktor $\lambda_p = 0.01, 0.5$ bzw. 0.99 .

2.5.3 Säulen Szenario

In diesem Szenario werden regelmäßig, mit jeweils 7 Feldern Zwischenraum zueinander, Hindernisse auf dem Torus verteilt. Idee ist, dass die Agenten eine kleine Orientierungshilfe besitzen sollen, aber gleichzeitig möglichst wenig Hindernisse verteilt werden. Das Zielobjekt startet an zufälliger Position, die Agenten starten mit möglichst großem Abstand zum Zielobjekt. Abbildung 2.6 zeigt ein Beispiel für den Startzustand eines solchen Szenarios, bei der das Zielobjekt sich in der Mitte und die Agenten am Rand befinden.

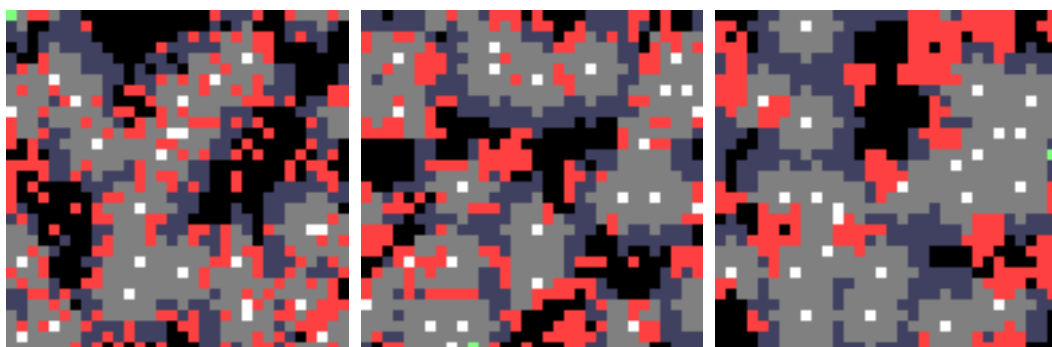


Abbildung 2.4: Szenario mit zufällig verteilten Hindernissen mit Hindernissanteil $\lambda_h = 0.2$ und Verknüpfungsfaktor $\lambda_p = 0.01, 0.5$ bzw. 0.99 .

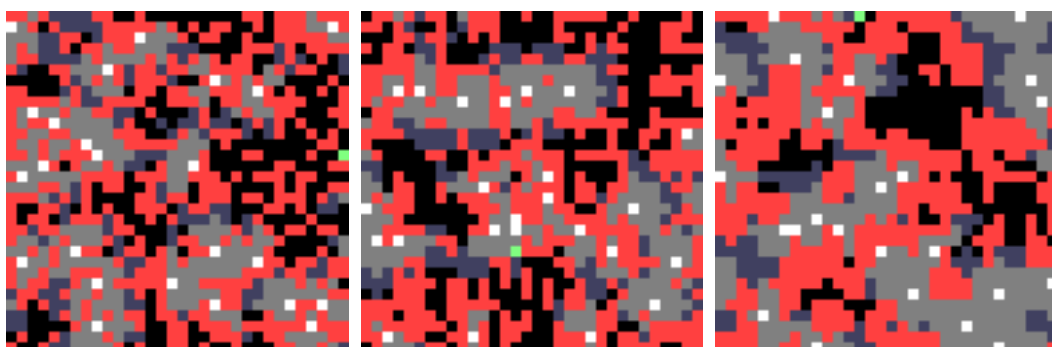


Abbildung 2.5: Szenario mit zufällig verteilten Hindernissen mit Hindernissanteil $\lambda_h = 0.4$ und Verknüpfungsfaktor $\lambda_p = 0.01, 0.5$ bzw. 0.99 .

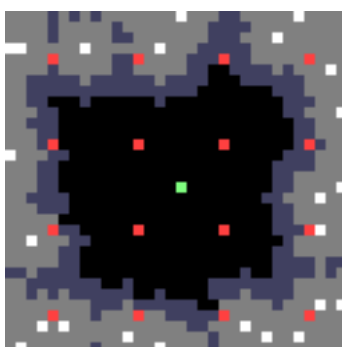


Abbildung 2.6: Startzustand des Säulen Szenarios mit regelmäßig angeordneten Hindernissen und zufälliger Verteilung von Agenten mit möglichst großem Abstand zum Zielobjekt

2.5.4 „Kreuz Szenario“

TODO raus Hier gibt eine horizontale Reihe aus Hindernissen halber Gesamtbreite welche durch eine vertikale Reihe aus Hindernissen halber Gesamthöhe in der Mitte geschnitten wird. Agenten und das Zielobjekt werden zufällig verteilt. Abbildung 2.7 zeigt ein Beispiel für ein solches Szenario.

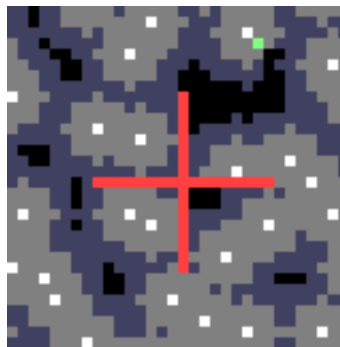


Abbildung 2.7: Kreuz Szenario mit kreuzförmiger Anordnung der Hindernisse und zufälliger Verteilung der Agenten und des Zielobjekts

2.5.5 „Raum Szenario“

Auf dem Torus wird ein Rechteck der halben Gesamthöhe und -breite des Torus erstellt, welches im Norden eine Öffnung von 4 Feldern Breite aufweist. Der Zielagent startet in der Mitte des Raums, alle Agenten starten mit maximaler Distanz zu den Hindernissen an zufälliger Position. Abbildung 2.8 zeigt ein Beispiel für eine Startkonfiguration eines solchen Szenarios.

2.5.6 Schwieriges Szenario

Hier wird der Torus an der rechten Seite vollständig durch Hindernisse blockiert, um den Torus zu halbieren. Alle Agenten starten (zufällig verteilt) am linken Rand, der Zielagent

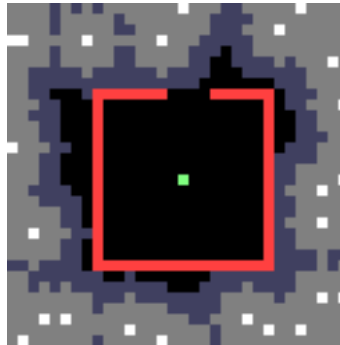


Abbildung 2.8: „Raum Szenario“ mit quaderförmiger Anordnung der Hindernisse mit Öffnung im Norden, zufälliger Verteilung der Agenten am Rand und zu Beginn im Zentrum startendem

startet auf der rechten Seite.

In regelmäßigen Abständen (7 Felder Zwischenraum) befindet sich eine vertikale Reihe von Hindernissen mit Öffnungen von 4 Feldern Breite abwechselnd im oberen Viertel und dem unteren Viertel.

Idee dieses Szenarios ist es, zu testen, inwieweit die Agenten durch die Öffnungen zum Ziel finden können. Ohne Orientierung an den Öffnungen und anderen Agenten ist es sehr schwierig, sich durch das Szenario zu bewegen. Abbildung 2.9 zeigt die Startkonfiguration des Szenarios.

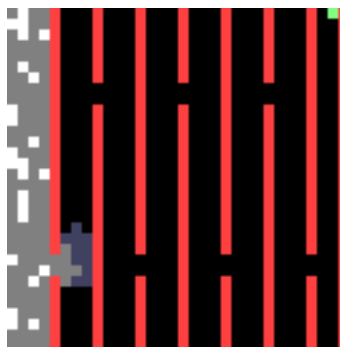


Abbildung 2.9: Schwieriges Szenario mit fester, wallartiger Verteilung von Hindernissen in regelmäßigen Abständen und mit Öffnungen, mit den Agenten mit zufälligem Startpunkt am linken Rand und mit dem Zielobjekt mit festem Startpunkt rechts oben

2.5.7 „Irrgarten Szenario“

Der Code zur Generierung der Hindernisse stammt aus [Ham04]. In den „Gängen“ des Irrgartens herrscht jeweils eine Breite von 2 Feldern. Die anfängliche Verteilung der Agenten und des Zielobjekts geschieht zufällig. Abbildung 2.10 zeigt ein Beispiel für eine Startkonfiguration eines Irrgartens.

TODO Idee oder raus

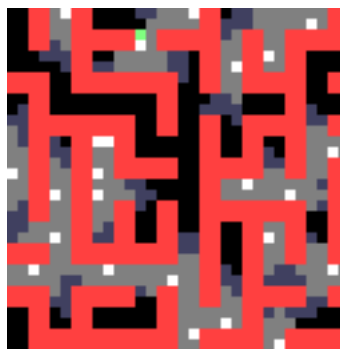


Abbildung 2.10: „Irrgarten Szenario“ mit einer Anordnung von Hindernissen in der Art, dass es nur wenige Pfade gibt.

2.6 Bestimmung der Qualität eines Algorithmus

Die Qualität eines Algorithmus zu einem Problem wird anhand des Anteils der Zeit berechnet, die er das Zielobjekt während des Problems überwachen (d.h. das Zielobjekt innerhalb einer Distanz von höchstens *reward range* halten) konnte, relativ zur Gesamtzeit.

Die Qualität eines Algorithmus zu einer Anzahl von Problemen (also einem Experiment) wird Anhand des Gesamtanteil der Zeit berechnet, die er das Zielobjekt während aller Probleme überwachen konnte, relativ zur Gesamtzeit aller Probleme.

Die Qualität eines Algorithmus entspricht dem Durchschnitt der Qualitäten des Algorithmus mehrerer Experimente.

Die Halbzeitqualität eines Algorithmus zu einem Problem entspricht dem Anteil der Zeit, die der Algorithmus das Zielobjekt während jeweils der zweiten Hälfte des Problems überwachen konnte, relativ zur halben Gesamtzeit.

Die Halbzeitqualität eines Algorithmus zu einer Anzahl von Problemen entspricht dem Anteil der Zeit, die der Algorithmus das Zielobjekt während jeweils der zweiten Hälfte des Problems überwachen konnte, relativ zur halben Gesamtzeit aller Probleme.

Die Halbzeitqualität eines Algorithmus entspricht dem Durchschnitt aller Halbzeitqualitäten des Algorithmus mehrerer Experimente.

Ein Vergleich der Qualität mit der Halbzeitqualität eines Algorithmus ermöglicht einen Einblick, wie gut sich der Algorithmus verhält, nachdem er sich auf das Problem bereits eine Zeit lang einstellen konnte.

Kapitel 3

Eigenschaften der Agenten

Ein Agent kann in jedem Schritt zwischen vier verschiedenen Aktionen wählen, die den vier Richtungen (Norden, Osten, Süden, Westen) entsprechen. Während ein Agent pro Zeiteinheit genau einen Schritt durchführen kann, kann das Zielobjekt je nach Szenario-parameter mehrere Schritte ausführen, was in Kapitel 4.1 erläutert wird.

Da wir ein Multiagentensystem auf einem diskreten Feld betrachten, werden alle Agenten werden nacheinander in der Art abgearbeitet, dass jeder Agent die aktuellen Sensordaten (siehe Kapitel 3.1) aus der Umgebung holt und auf deren Basis die nächste Aktion bestimmt.

Wurden alle Aktionen bestimmt, können die Agenten in zufälliger Reihenfolge versuchen, sie auszuführen. Ungültige Aktionen, d.h. der Versuch sich auf ein besetztes Feld zu bewegen, schlagen fehl und der Agent führt in diesem Schritt keine Aktion aus, wird aber auch nicht weiter bestraft. Eine detaillierte Beschreibung der Bewegung im Kontext anderer Agenten und Programmteile wird in Kapitel 5.2 gegeben.

Weitere Fähigkeiten eines Agenten betreffen die Kommunikation, bis Kapitel 11 soll

jedoch nur der Fall ohne Kommunikation betrachtet werden, d.h. die Agenten können untereinander keine Informationen austauschen und müssen sich alleine auf ihre Sensordaten verlassen.

3.1 Sensoren eines Agenten

Jeder Agent besitzt eine Anzahl visueller, binärer Sensoren mit begrenzter Reichweite. Jeder Sensor kann nur feststellen, ob sich in seinem Sichtbereich ein Objekt eines bestimmten Typs befindet (1) oder nicht (0). Jeder Sensor ist in eine bestimmte Richtung ausgerichtet, andere Objekte blockieren die Sicht und Sichtlinien werden durch einen einfachen Bresenham-Algorithmus bestimmt.

Zwei Sensoren, die in die selbe Richtung ausgerichtet sind und den selben Typ von Objekt erkennen, werden in diesem Zusammenhang ein Sensordatenpaar genannt (siehe Kapitel 3.1.1). Alle Sensoren, die nur gemeinsam haben, dass sie den selben Typ von Objekt erkennen, werden in einer Gruppe zusammengefasst und der Aufbau eines ganzen, aus solchen Gruppen bestehenden Sensordatensatzes soll in Kapitel 3.1.2 besprochen werden.

3.1.1 Aufbau eines Sensordatenpaares

Ein Datenpaar besteht aus zwei Sensoren, die den selben Typ von Objekt erkennen, in die selbe Richtung ausgerichtet sind und sich nur in ihrer Sichtweite unterscheiden, wodurch der Agent rudimentär die Entfernung zu anderen Objekten feststellen kann. Die Sichtweite des ersten Sensors eines Paares wird über den Parameter *sight range* bestimmt, die Sichtweite des zweiten Sensors über den Parameter *reward range* (siehe auch Kapitel 2.2). Allgemein soll *sight range* = 5.0 und *reward range* = 2.0 betragen, der überwachte Bereich ist also eine Teilmenge des sichtbaren Bereichs. In Abbildung 3.1 sind alle Sichtreichwei-

ten (heller und dunkler Bereich) und Überwachungsreichweiten (heller Bereich) für die einzelnen Richtungen dargestellt.

Anzumerken sei hier, dass wegen der gewählten Werte für beide Reichweiten ein Sensordatenpaar (01) nicht auftreten kann, da ein Objekt nicht gleichzeitig näher als 2.0 und weiter als 5.0 entfernt sein kann.

Sei $r(O_1, O_2)$ die Distanz zwischen dem Objekt, das die Sensordaten erfasst und dem nächstliegenden Objekt des Typs, den der Sensor wahrnehmen kann, dann gibt es folgende Fälle:

1. (0/0) : $r(O_1, O_2) > \textit{sight range}$ (kein passendes Objekt in Sichtweite)
2. (1/0) : $\textit{reward range} < r(O_1, O_2) \leq \textit{sight range}$ (Objekt in Sichtweite)
3. (1/1) : $r(O_1, O_2) \leq \textit{reward range}$ (Objekt in Sicht- und Überwachungsreichweite)
4. (0/1) : $\textit{reward range} \geq r(O_1, O_2) > \textit{sight range}$ (Fall kann nicht auftreten, da $\textit{reward range} < \textit{sight range}$)

3.1.2 Aufbau eines Sensordatensatzes

In einem Sensordatensatz sind jeweils 8 Sensoren zu jeweils einer Gruppe zusammengefasst, welche wiederum in 4 Richtungen mit jeweils einem Sensorenpaar aufgeteilt ist. Gleichung 3.1 stellt den allgemeinen Aufbau eines kompletten Sensordatensatzes dar, der aus den drei Gruppen der Zielobjektsensoren (z), der Agentensensoren (a) und der Hindernissensoren (h) besteht.

Seien beispielsweise im Westen und Osten sich in Überwachungsreichweite befindliche Hindernissen, im Norden außerhalb der Überwachungsreichweite aber in Sichtweite das

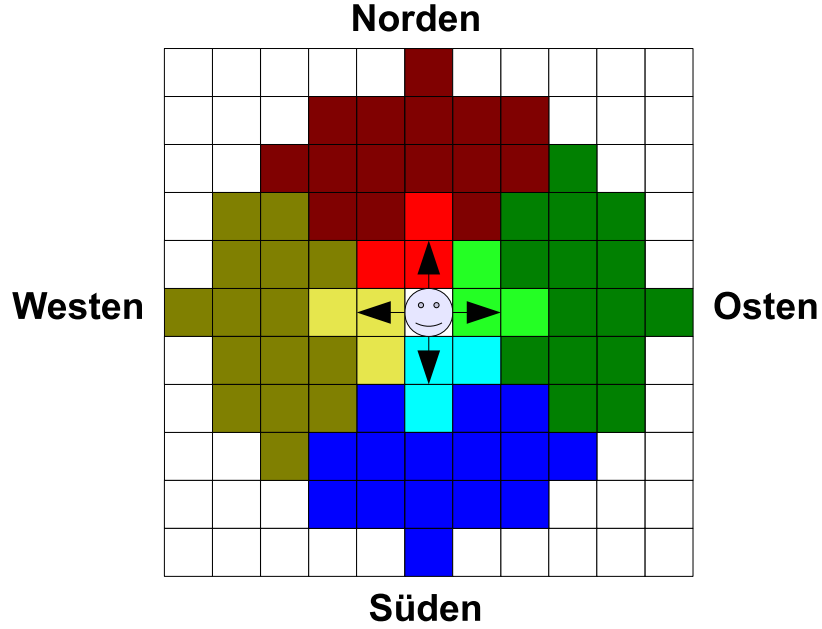


Abbildung 3.1: Sicht- (5.0, dunkler Bereich) und Überwachungsreichweite (2.0, heller Bereich) eines Agenten, jeweils für die einzelnen Richtungen

Zielobjekt und im Süden Agenten in Überwachungsreichweite des Agenten, dann ergibt sich ein Sensordatensatz $s_{Beispiel}$ wie in Gleichung 3.2 dargestellt.

$$\begin{aligned}
 \text{Sensordatensatz } s = & \underbrace{(z_{s_N} z_{r_N})(z_{s_O} z_{r_O})(z_{s_S} z_{r_S})(z_{s_W} z_{r_W})}_{\text{Erste Gruppe (Zielobjekt)}} \\
 & \underbrace{(a_{s_N} a_{r_N})(a_{s_O} a_{r_O})(a_{s_S} a_{r_S})(a_{s_W} a_{r_W})}_{\text{Zweite Gruppe (Agenten)}} \\
 & \underbrace{(h_{s_N} h_{r_N})(h_{s_O} h_{r_O})(h_{s_S} h_{r_S})(h_{s_W} h_{r_W})}_{\text{Dritte Gruppe (Hindernisse)}}
 \end{aligned} \tag{3.1}$$

$$\text{Sensordatensatz } s_{Beispiel} = (1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 1, 1) \tag{3.2}$$

3.2 Grundsätzliche Algorithmen der Agenten

Neben denjenigen Algorithmen, die auf XCS basieren und in Kapitel 7 besprochen werden, sollen hier einige, auf einfachen Heuristiken basierende, Algorithmen vorgestellt werden, um die Qualität der anderen Algorithmen besser einordnen zu können. Wesentliches Merkmal im Vergleich zu auf XCS basierenden Algorithmen ist, dass sie statische, handgeschriebene Regeln benutzen und den Erfolg oder Misserfolg ihrer Aktionen ignorieren, d.h. ihre Regeln während eines Laufs nicht anpassen.

Die in Kapitel 5.6 erwähnte und dort aufgerufene Funktion *calculateReward()* soll für die hier aufgelisteten Algorithmen also jeweils der leeren Funktion entsprechen. Im Folgenden sollen also insbesondere die Implementierungen der jeweiligen *calculateNextMove()* Funktion vorgestellt werden.

3.2.1 Algorithmus mit zufälliger Bewegung

Bei diesem Algorithmus wird in jedem Schritt eine zufällige Aktion ausgeführt. Abbildung 3.2 zeigt eine Beispielsituation, bei der der Agent jegliche Sensordaten (die 4 Agenten und das Zielobjekt, der als Stern dargestellt ist) ignoriert und eine Aktion zufällig auswählen wird.

Programm 3.1 zeigt den zugehörigen Quelltext.

3.2.2 Einfache Heuristik

Ist das Zielobjekt in Sichtweite, bewegt sich ein Agent mit dieser Heuristik auf das Zielobjekt zu, ist es nicht in Sichtweite, führt er eine zufällige Aktion aus. Abbildung 3.3 zeigt eine Beispielsituation bei der sich das Zielobjekt (Stern) im Süden befindet, der Agent mit einfacher Heuristik die anderen Agenten ignoriert und sich auf das Ziel zubewegen



Programm 3.1: Berechnung der nächsten Aktion bei der Benutzung des Algorithmus mit zufälliger Bewegung

möchte.

Programm 3.2 zeigt den zugehörigen Quelltext.

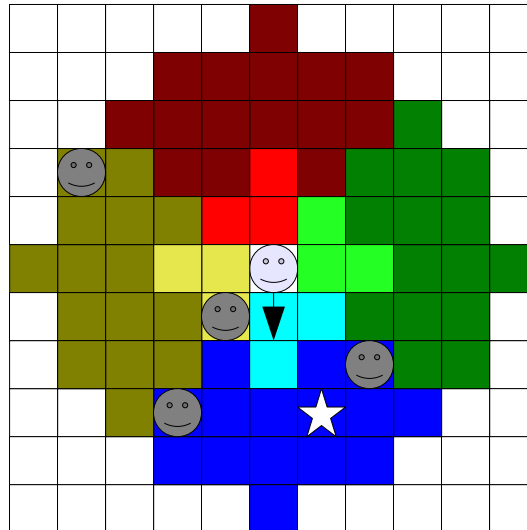


Abbildung 3.3: Agent mit einfacher Heuristik: Sofern es sichtbar ist bewegt sich der Agent auf das Zielobjekt zu.

3.2.3 Intelligente Heuristik

Ist der Zielobjekt in Sicht, verhält sich diese Heuristik wie die einfache Heuristik. Ist das Zielobjekt dagegen nicht in Sicht, wird versucht, anderen Agenten auszuweichen, um ein möglichst breit gestreutes Netz aus Agenten aufzubauen. In der Implementation heißt das, dass unter allen Richtungen, in denen kein anderer Agent gesichtet wurde, eine Richtung zufällig ausgewählt wird und falls alle Richtungen belegt (oder alle frei) sind, wird aus allen Richtungen eine zufällig ausgewählt wird. In Abbildung 3.4 ist das Zielobjekt nicht im Sichtbereich des Agenten und dieser wählt deswegen eine Richtung, in der die Sensoren keine Agenten anzeigt, in diesem Fall Norden.

Programm 3.3 zeigt den zugehörigen Quelltext.

```
1  /**
2   * Berechne nächste Aktion (einfache Heuristik)
3   */
4   private void calculateNextMove() {
5       /**
6        * Holt sich die Informationen der Gruppe der Sensoren, die auf
7        * das Zielobjekt ausgerichtet sind
8        */
9       boolean[] goal_sensor = lastState.getSensorGoal();
10      calculatedAction = -1;
11      for(int i = 0; i < Action.MAX_DIRECTIONS; i++) {
12          /**
13           * Zielagent in Sicht in dieser Richtung?
14           */
15          if(goal_sensor[2*i]) {
16              calculatedAction = i;
17              break;
18          }
19      }
20
21      /**
22       * Sonst wähle zufällige Richtung als nächste Aktion
23       */
24      if(calculatedAction == -1) {
25          calculatedAction = Misc.nextInt(Action.MAX_DIRECTIONS);
26      }
27
28  }
```

Programm 3.2: Berechnung der nächsten Aktion bei der Benutzung der einfachen Heuristik

```

1  /**
2   * Berechne nächste Aktion (intelligente Heuristik)
3   */
4  private void calculateNextMove() {
5      /**
6       * Holt sich die Informationen der Gruppe der Sensoren, die auf
7       * das Zielobjekt ausgerichtet sind
8       */
9      boolean[] goal_sensor = lastState.getSensorGoal();
10
11      calculatedAction = -1;
12      for(int i = 0; i < Action.MAX DIRECTIONS; i++) {
13          /**
14           * Zielagent in Sicht in dieser Richtung?
15           */
16          if(goal_sensor[2*i]) {
17              calculatedAction = i;
18              break;
19          }
20      }
21
22      /**
23       * Zielobjekt nicht in Sicht? Dann bewege von Agenten weg
24       */
25      if(calculatedAction == -1) {
26          calculatedAction = Misc.nextInt(Action.MAX DIRECTIONS);
27
28          boolean[] agent_sensors = lastState.getSensorAgent();
29          boolean one_free = false;
30          for(int i = 0; i < Action.MAX DIRECTIONS; i++) {
31              if(!agent_sensors[2*i]) {
32                  one_free = true;
33                  break;
34              }
35          }
36
37          if(one_free) {
38              while(agent_sensors[2*calculatedAction]) {
39                  calculatedAction = Misc.nextInt(Action.MAX DIRECTIONS);
40              }
41          }
42      }
43  }

```

Programm 3.3: Berechnung der nächsten Aktion bei der Benutzung der intelligenten Heuristik

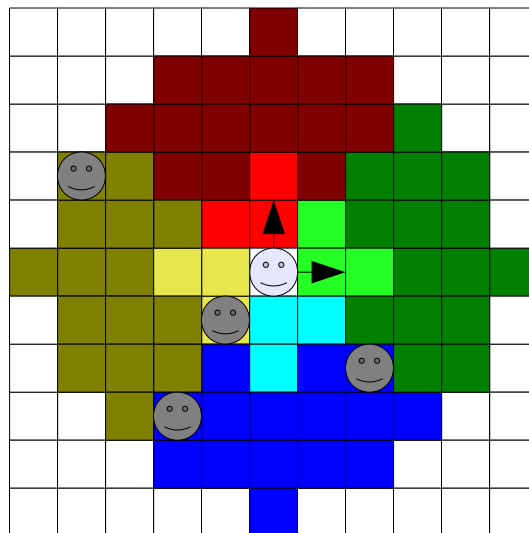


Abbildung 3.4: Agent mit intelligenter Heuristik: Falls das Zielobjekt nicht sichtbar ist bewegt sich der Agent von anderen Agenten weg.

Kapitel 4

Das Zielobjekt

Die Typen von Zielobjekten werden zum einen über ihre Geschwindigkeit und zum anderen über ihre Bewegungsart definiert. Neben der Größe des Torus und den Hindernissen trägt der Typ des Zielobjekts wesentlich zur Schwierigkeit eines Szenarios bei, da dieser die Aufenthaltswahrscheinlichkeiten des Zielobjekts unter Einbeziehung des Zustands des letzten Zeitschritts bestimmt. Die Schwierigkeit bestimmt sich über die Summe der erwarteten Aufenthaltswahrscheinlichkeiten in nicht überwachten Feldern geteilt durch die Summe der Aufenthaltswahrscheinlichkeiten in überwachten Feldern.

TODO Aufenthaltswahrscheinlichkeiten raus

4.1 Basiseigenschaften

Im wesentlichen entspricht ein Zielobjekt einem Agenten, d.h. das Zielobjekt kann sich bewegen und besitzt Sensoren. Außerdem kann sich das Zielobjekt in einem Schritt u.U. um mehr als ein Feld bewegen, was durch die durch das Szenario festgelegte Geschwindigkeit des Zielobjekts bestimmt ist. Der Wert der Geschwindigkeit kann auch gebrochene Werte annehmen, wobei in diesem Fall der gebrochene Rest dann die Wahrscheinlichkeit angibt, einen weiteren Schritt durchzuführen. Beispielsweise würde Geschwindigkeit 1.4 in 40%

der Fälle zu zwei Schritten und in 60% der Fälle zu einem einzigen Schritt führen. Die Auswertung der Bewegungsgeschwindigkeit ist relevant in Kapitel 5.2, bei der Reihenfolge der Ausführung der Aktionen der Objekte.

Zusätzlich dazu haben alle Arten von Bewegungen des Zielobjekts gemeinsam, dass, wenn dem Algorithmus kein freies Feld zur Verfügung steht, ein zufälliges, freies Feld in der Nähe ausgewählt und dorthin gesprungen wird. Dies kommt einem Neustart gleich und ist notwendig um eine Verfälschung des Ergebnisses zu verhindern, das daher rühren kann, dass ein oder mehrere Agenten (zusammen mit eventuellen Hindernissen) alle vier Bewegungsrichtungen des Zielobjekts blockieren.

Zu beachten ist hier, dass auch der Sprung selbst eine Verfälschung darstellt, insbesondere wenn in einem Durchlauf viele Sprünge durchgeführt werden. Falls dies passiert sollte man deshalb das Ergebnis verwerfen und z.B. andere *random seed* Werte oder einen anderen Algorithmus benutzen. Sofern nicht anders angegeben ist der Anteil solcher Sprünge jeweils unter 0.1% und wird ignoriert.

TODO weiss schon dass blockiert

4.2 Typen von Zielobjekten

4.2.1 Typ „Zufälliger Sprung“

Ein Zielobjekt dieses Typs springt zu einem zufälligen Feld auf dem Torus. Ist das Feld besetzt wird wiederholt bis ein freies Feld gefunden wurde. Mit dieser Einstellung kann die Abdeckung des Algorithmus geprüft werden, d.h. inwieweit die Agenten jeweils außerhalb der Überwachungsreichweite anderer Agenten bleiben.

Jegliche Anpassung an die Bewegung des Zielobjekts ist hier wenig hilfreich, ein Agent kann nicht einmal davon ausgehen, dass sich das Zielobjekt in der Nähe seiner Position

der letzten Zeiteinheit befindet.

4.2.2 Typ „Zufällige Bewegung“

Ein Zielobjekt dieses Typs verhält sich so wie ein Agent mit dem Algorithmus mit zufälliger Bewegung (siehe Kapitel 3.2.1). Sind alle möglichen Felder belegt, wird, wie oben beschrieben, auf ein zufälliges Feld gesprungen.

4.2.3 Typ „Einfache Richtungsänderung“

Ein Zielobjekt dieses Typs entfernt zuerst alle Richtungen, in denen sich direkt angrenzend ein Hindernis befindet. Diese Erweiterung der Sensorfähigkeiten wurde gewählt, damit das Zielobjekt nicht in Hindernissen längere Zeit steckenbleibt. Anschließend wird die Richtung entfernt, die der im letzten Schritt gewählten entgegengesetzt ist. Von den verbleibenden (bis zu) drei Richtungen wird schließlich eine zufällig ausgewählt. Sind alle drei Richtungen versperrt, wird in die entgegengesetzte Richtung zurückgegangen.

In Abbildung 4.1 sind alle Felder grau markiert, die der Zielagent innerhalb von zwei Schritten erreichen kann, nachdem er sich einmal nach Norden bewegt hat.

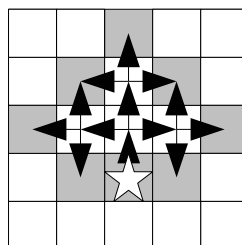


Abbildung 4.1: Zielobjekt macht pro Schritt maximal eine Richtungsänderung

4.2.4 Typ „Beibehaltung der Richtung“

Der Zielobjekt versucht, immer Richtung Norden zu gehen. Ist das Zielfeld blockiert, wählt es ein zufälliges, angrenzendes, freies Feld im Westen oder Osten. Anzumerken ist, dass dies zusätzliche Fähigkeiten darstellen, d.h. das Zielobjekt kann feststellen, ob sich direkt angrenzend ein Hindernis im Norden befindet, während normale Agenten, was die Distanz betrifft, keine Informationen darüber besitzen können.

TODO

Sind auch die Felder im Westen und Osten belegt, springt es auf ein zufälliges freies Feld in der Nähe. Schafft es der Zielobjekt innerhalb von einer bestimmten Zahl (Breite des Spielfelds) von Schritten nicht, einen weiteren Schritt nach Norden zu gehen, wird ebenfalls gesprungen, um ein „festhängen“ an einem Hindernis zu vermeiden.

In Abbildung 4.2 sind drei Situationen dargestellt, zum einen ein wiederholtes hin- und herlaufen unter den Hindernissen, der Weg links um die Hindernisse herum und der Weg rechts um die Hindernisse herum.

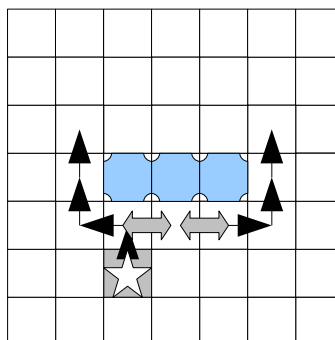


Abbildung 4.2: Bewegungsform „Beibehaltung der Richtung“: Zielobjekt bewegt sich, wenn möglich, immer nach Norden

4.2.5 Typ „Intelligentes Verhalten“

Ein Zielobjekt dieses Typs versucht bei der Auswahl der Aktion möglichst die Aktion zu wählen, bei der es außerhalb der Sichtweite der Agenten bleibt. Dazu werden alle Richtungen gestrichen, in denen ein Agent sich innerhalb der Überwachungsreichweite befindet. Außerdem werden von den verbleibenden Richtungen mit 50% diejenigen Richtungen gestrichen, in denen sich ein Agent in Sichtweite befindet. Sind alle Richtungen gestrichen worden, bewegt sich das Zielobjekt zufällig. Sind alle Richtungen blockiert, springt es wie in den anderen Varianten auch auf ein zufälliges Feld in der Nähe.

In Abbildung 4.3 wird die Richtung Süden gestrichen, da sich dort ein Agent in Überwachungsreichweite befindet. Die Richtungen Westen und Norden werden jeweils mit Wahrscheinlichkeit 50% gestrichen, da sich dort Agenten in Sichtweite befinden. Nur Richtung Osten wird als Möglichkeit sicher übrigbleiben.

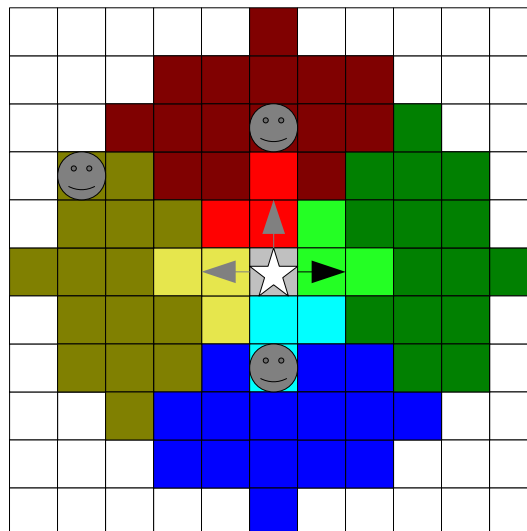


Abbildung 4.3: Zielobjekt bewegt sich mit bestimmter Wahrscheinlichkeit von Agenten und größerer Wahrscheinlichkeit von Hindernissen weg

4.2.6 Typ „SXCS“

Dieser Typ ist eine Implementierung für das Zielobjekt, das auf der SXCS Implementierung in Kapitel 9 basiert. Einziger Unterschied ist in der Art, wie das SXCS die eigenen Aktionen bewertet. Während das dort beschriebene SXCS die Nähe zum Zielobjekt belohnt, soll hier das Zielobjekt die Situationen positiv bewerten, bei denen sich keine Agenten in Überwachungsreichweite befinden. Eine genaue Beschreibung folgt im Kapitel 9, hier sollte die Idee nur der Vollständigkeit halber erwähnt werden.

TODO Pendelbewegung? TODO Fester Pfad?

Kapitel 5

Ablauf der Simulation

5.1 Hauptschleife

In der Hauptschleife (siehe Programm 5.1) wird ein Experiment mit vorgegebener Konfiguration (“*Configuration*”) durchgeführt. Dabei werden eine Anzahl von Problemen abgearbeitet, bei denen jeweils der Torus auf den Startzustand gesetzt, das eigentliche Problem berechnet und ein neuer *random seed* Wert gesetzt wird.

5.2 Reihenfolge der Ausführung (*doOneMultiStepProblem()*)

TODO vielleicht erst noch eine allgemeine Übersicht geben

Für die Berechnung eines einzelnen Problems (“*doOneMultiStepProblem()*”) stellt sich die Frage nach der Genauigkeit und der Reihenfolge der Abarbeitung, da die Simulation nicht parallel, sondern schrittweise auf einem diskreten Torus abläuft. Dies kann u.U. dazu

```

1  /**
2   * Führt eine Anzahl von Problemen aus
3   * @param experiment_nr Nummer des auszuführenden Experiments
4   */
5   public void doOneMultiStepExperiment(int experiment_nr) {
6       int currentTimestep = 0;
7
8       /**
9        * number of problems for the same population
10      */
11      for (int i = 0; i < Configuration.getNumberOfProblems(); i++) {
12
13          /**
14           * Erstellt einen neuen Torus und verteilt Agenten und das Zielobjekt neu
15          */
16          BaseAgent.grid.resetState();
17
18          /**
19           * Führe Problem aus und aktualisiere aktuellen Zeitschritt
20          */
21          currentTimestep = doOneMultiStepProblem(currentTimestep);
22
23          /**
24           * Initialisiere neuen "Random Seed" Wert
25          */
26          Misc.initSeed(Configuration.getRandomSeed() +
27                        experiment_nr * Configuration.getNumberOfProblems() + 1 + i);
28      }
29  }

```

Programm 5.1: Zentrale Schleife für einzelne Experimente

führen, dass je nach Position in der Liste abzuarbeitender Agenten die Informationen über die Umgebung unterschiedlich alt sind. Die Frage ist deshalb, in welcher Reihenfolge Sensordaten ermittelt, ausgewertet, Agenten bewegt, intern sich selbst bewertet und global die Qualität gemessen wird.

Da eine Aktion auf Basis der Sensordaten ausgewählt wird, ist die erste Restriktion, dass eine Aktion nach der Verarbeitung der Sensordaten stattfinden muss. Und da Aktionen bewertet werden sollen, also jeweils der Zustand nach der Bewegung mit dem gewünschten Zustand verglichen werden soll, ist die zweite Restriktion, dass die Bewertung einer Aktion nach dessen Ausführung stattfinden muss.

Ansonsten gibt es folgende Möglichkeiten:

1. Für alle Agenten werden erst einmal die neuen Sensordaten erfasst und sich für eine Aktion entschieden. Sind alle Agenten abgearbeitet, werden die Aktionen ausgeführt.
2. Die Agenten werden nacheinander abgearbeitet, es werden jeweils neue Sensordaten erfasst und sich sofort für eine neue Aktion entschieden.

Bei der ersten Möglichkeit haben alle Agenten die Sensordaten vom Beginn der Zeiteinheit, während bei der zweiten Möglichkeit später verarbeitete Agenten bereits die Aktionen der bereits berechneten Agenten miteinbeziehen können. Umgekehrt können dann frühere Agenten bessere Positionen früher besetzen. Da aufgrund der primitiven Sensoren nicht davon auszugehen ist, dass Agenten beginnende Bewegungen (und somit deren jeweilige Zielposition) anderer Agenten einbeziehen können, soll jeder Agent von den Sensorinformationen zu Beginn der Zeiteinheit ausgehen.

Wenn sich mehrere Agenten auf dasselbe Feld bewegen wollen, dann spielt die Reihenfolge der Ausführung der Aktionen eine Rolle. Wird die Liste der Agenten einfach linear abgearbeitet, können Agenten mit niedriger Position in der Liste die Aktion auf Basis jüngerer

Sensordaten fallen. Dies kann dazu führen, dass Aktionen von Agenten mit höherer Position in der Liste eher fehlschlagen, da das als frei angenommene Feld nun bereits besetzt ist. Da es keinen Grund gibt, Agenten mit niedrigerer Position zu bevorteilen, werden die Aktionen der Agenten in zufälliger Reihenfolge abgearbeitet.

Bezüglich der Bewegung ergibt sich hierbei eine weitere Frage, nämlich wie unterschiedliche Bewegungsgeschwindigkeiten behandelt werden sollen, da alle Agenten eine Einheitsgeschwindigkeit von maximal einem Feld pro Zeiteinheit haben, während sich das Zielobjekt je nach Szenario gleich eine ganze Anzahl von Feldern bewegen kann (siehe auch Kapitel 4.1). Die Entscheidung fiel hier auf eine zufällige Verteilung. Kann sich das Zielobjekt um n Schritte bewegen, so wird seine Bewegung in n Einzelschritte unterteilt, die nacheinander mit zufälligen Abständen (d.h. Bewegungen anderer Agenten) ausgeführt werden.

Eine weitere Frage ist, wie das Zielobjekt diese weiteren Schritte festlegen soll. Hier soll ein Sonderfall eingeführt werden, sodass das Zielobjekt in einer Zeiteinheit mehrmals (n -mal) neue Sensordaten erfassen und sich für eine neue Aktion entscheiden kann.

5.3 Messung der Qualität

Eine konkrete Antwort kann man auf diese zwei Fragen nicht geben, sie hängt davon ab, was man denn nun eigentlich erreichen möchte, also auf welche Weise die Qualität des Algorithmus bewertet wird. Der naheliegendste Messzeitpunkt ist, nachdem sich alle Agenten bewegt haben. Da wir die Agenten und das Zielobjekt in einem Durchlauf gemeinsam nacheinander bewegen, stellt sich die Frage nicht, ob wir womöglich vor der Bewegung des Zielobjekts die Qualität messen sollen. Eine Messung nach der Bewegung des Zielobjekts würde diesem erlauben, sich vor jeder Messung optimal zu positionieren, was in einer geringeren Qualität für den Algorithmus resultiert, da sich das Zielobjekt aus der Überwachungsreichweite anderer Agenten hinausbewegen kann. Letztlich ist es eine

Frage der Problemstellung, denn eine Messung nach Bewegung des Zielobjekts bedeutet letztlich, dass ein Agent einen gerade aus seiner Überwachungsreichweite heraus laufenden Zielobjekts in diesem Schritt nicht mehr überwachen kann.

Da ein wesentlicher Bestandteil die Kooperation (und somit die Abdeckung des Torus anstatt dem Verfolgen des Zielobjekts) sein soll, soll ein Bewertungskriterium sein, inwieweit der Einfluss des Zielobjekts minimiert werden soll. Auch findet, wenn man vom realistischen Fall ausgeht, die Bewegung des Zielobjekts gleichzeitig mit allen anderen Agenten statt. Die Qualität wird somit nach der Bewegung des Zielobjekts gemessen. Die Überlegung unterstreicht auch nochmal, dass es besser ist, das Zielobjekt insgesamt wie einen normalen (aber sich mehrmals bewegendem) Agenten zu behandeln.

5.4 Reihenfolge der Ermittlung des *base reward*

Keine der bisher vorgestellten Varianten machen Gebrauch von einem sogenannten *base reward*, d.h. TODO

Schließlich bleibt die Frage danach, wann geprüft werden soll, ob das Zielobjekt in Überwachungsreichweite ist, und wann sich somit ein *reward* ergeben soll. Wesentliche Punkte hierbei sind, dass der Algorithmus sich anhand der Sensordaten selbst bewertet und pro Zeitschritt die Sensordaten nur einmal erhoben werden. Letzteres folgt aus der Auslegung von XCS, der in der Standardimplementation darauf ausgelegt ist, dass der Reward jeweils genau einer Aktion zugeordnet ist. Daraus ergibt sich auch, dass der Reward von binärer Natur ist (“Zielobjekt in Überwachungsreichweite” oder “Zielobjekt nicht in Überwachungsreichweite”), weshalb Zwischenzustände für den Reward, der sich aus der mehrfachen Bewegung des Zielagenten ergeben könnte, ausgeschlossen werden soll (z.B. “War zwei von drei Schritten in der Überwachungsreichweite” $\Rightarrow \frac{2}{3}$ Reward). Insbesondere würde dies eine mehrfache Erhebung der Sensordaten erfordern.

TODO Rewarderhebung für normale Agenten irrelevant, evtl teilen und in XCS Kapitel

Für den Reward gibt es somit folgende Möglichkeiten:

1. Ermittlung der einzelnen *reward* Werte jeweils direkt nach der Ausführung einer einzelnen Aktion
2. Ermittlung aller *reward* Werte nach Ausführung aller Aktionen der Agenten und des Zielobjekts

Werden die *reward* Werte sofort ermittelt (Punkt 1), dann bezieht sich der Wert auf die veralteten Sensordaten vor der Aktion, die Aktion selbst würde bei der Ermittlung des *reward* Werts also ignoriert werden. Bei Punkt 2 müsste man bis zum neuen Zeitschritt warten, bis neue Sensordaten ermittelt wurden.

5.5 Zusammenfassung

Zusammenfassend sieht der Ablauf aller Agenten (inklusive des Zielobjekts) also wie folgt aus:

1. Bestimmen der aktuellen **Qualität**
2. Erfassung aller **Sensordaten**
3. Bestimmung der jeweiligen ***reward* Werte** für die einzelnen Objekte für den letzten Schritt
4. **Wahl der Aktion** anhand der Regeln des jeweiligen Agenten
5. **Ausführung der Aktion** (in zufälliger Reihenfolge, das Zielobjekt wiederholt Schritte 1 und 2 nach der Ausführung der Aktion)

5.6 Implementierung eines Problemablaufs

In der Schleife der Funktion zur Berechnung eines Experiments (Programm 5.1) wird die Funktion zur Berechnung des Problems (*doOneMultiStepProblem()* in Programm 5.2) aufgerufen. Dort wird, in einer weiteren Schleife über die Anzahl der maximalen Schritte, die Sicht aktualisiert (*updateSight()*), die Qualität bestimmt (*updateStatistics()*), die neuen Sensordaten und die nächste Aktion ermittelt (*calculateAgents()*, siehe Programm 5.3), der *reward* Wert ermittelt (*rewardAgents()*, siehe Programm 5.4) und schließlich die Objekte bewegt (*moveAgents()*, siehe Programm 5.5). Die konkrete Umsetzung der dort aufgerufenen Funktionen (insbesondere *calculateNextMove()* und *calculateReward()*) wird im Kapitel 9 erläutert (bzw. in Kapitel 3 was die Heuristiken betrifft, wobei *calculateReward()* dort keine Rolle spielt und eine leere Funktion aufgerufen wird).

```

1  /**
2  * Führt eine Anzahl von Schritten auf dem aktuellen Torus aus
3  * @param stepCounter Aktuelle Zeitschritt
4  * @return Der Zeitschritt nach der Ausführung
5  */
6  private int doOneMultiStepProblem(int stepCounter) {
7      /**
8       * Zeitpunkt bis zu dem das Problem ausgeführt wird
9       */
10     int steps_next_problem =
11         Configuration.getNumberOfSteps() + stepCounter;
12     for (int currentTimestep = stepCounter;
13         currentTimestep < steps_next_problem; currentTimestep++) {
14
15         /**
16          * Ermittle die Sichtbarkeit und erhebe Statistiken
17          */
18         BaseAgent.grid.updateSight();
19         BaseAgent.grid.updateStatistics(currentTimestep);
20
21         /**
22          * Ermittle neue Sensordaten und berechne Aktionen der Agenten
23          */
24         calculateAgents(currentTimestep);
25
26         /**
27          * Ermittle den Reward für alle Agenten (nach dem ersten Schritt)
28          */
29         if (currentTimestep > stepCounter) {
30             rewardAgents(currentTimestep);
31         }
32
33         /**
34          * Führe zuvor berechnete Aktionen aus
35          */
36         moveAgents();
37     }
38
39     /**
40      * Abschließende Ermittlung des Rewards
41      */
42     BaseAgent.grid.updateSight();
43     rewardAgents(steps_next_problem);
44     return steps_next_problem;
45 }

```

Programm 5.2: Zentrale Schleife für einzelne Probleme

```

1  /**
2   * Berechnet die Aktionen und führt sie in zufälliger Reihenfolge aus
3   * @param gaTimestep der aktuelle Zeitschritt
4   */
5   private void calculateAgents(final long gaTimestep) {
6
7   /**
8   * Ermittle Sensordaten und bestimme nächste Bewegung
9   */
10  for(BaseAgent a : agentList) {
11      a.acquireNewSensorData();
12      a.calculateNextMove(gaTimestep);
13  }
14  BaseAgent.goalAgent.acquireNewSensorData();
15  BaseAgent.goalAgent.calculateNextMove(gaTimestep);
16  }

```

Programm 5.3: Zentrale Bearbeitung (Sensordaten und Berechnung der neuen Aktion) aller Agenten und des Zielobjekts innerhalb eines Problems

```

1  /**
2   * Verteilt den Reward an alle Agenten
3   */
4   private void rewardAgents(final long gaTimestep) {
5       for(BaseAgent a : agentList) {
6           a.calculateReward(gaTimestep);
7       }
8       BaseAgent.goalAgent.calculateReward(gaTimestep);
9   }

```

Programm 5.4: Zentrale Bearbeitung (Verteilung des Rewards) aller Agenten und des Zielobjekts innerhalb eines Problems

```

1  /**
2  * Führt die berechnete Bewegungen der Agenten in zufälliger Reihenfolge aus
3  */
4  private void moveAgents(long gaTimestep) {
5      /**
6       * Erstelle Ausführungsliste für alle Objekte (Zielobjekt mehrfach)
7       */
8       int goal_speed = Configuration.getGoalAgentMovementSpeed();
9       ArrayList<BaseAgent> random_list =
10         new ArrayList<BaseAgent>(agentList.size() + goal_speed);
11
12       random_list.addAll(agentList);
13       for(int i = 0; i < goal_speed; i++) {
14         random_list.add(BaseAgent.goalAgent);
15       }
16
17       /**
18       * Führe die ermittelten Aktionen in zufälliger Reihenfolge aus
19       * (Zielobjekt kann mehrfach ausgeführt werden).
20       */
21       int[] array = Misc.getRandomArray(random_list.size());
22       for(int i = 0; i < array.length; i++) {
23         BaseAgent a = random_list.get(array[i]);
24         a.doNextMove();
25         if(a.isGoalAgent() && goal_speed > 1) {
26           goal_speed--;
27           a.acquireNewSensorData();
28           a.calculateNextMove(gaTimestep);
29           a.calculateReward(gaTimestep);
30         }
31       }
32     }

```

Programm 5.5: Zentrale Bearbeitung (Ausführung der Bewegung) aller Agenten und des Zielobjekts innerhalb eines Problems

Kapitel 6

Erste Analyse der Agenten ohne XCS

In diesem Kapitel sollen erste Analysen bezüglich der verwendeten Szenarien anhand des Algorithmus zufälliger Bewegung (siehe Kapitel 3.2.1), des Algorithmus mit einfacher Heuristik (siehe Kapitel 3.2.2) und des Algorithmus mit intelligenter Heuristik (siehe Kapitel 3.2.3) angefertigt werden. Die Ergebnisse aus der Analyse werden eine Grundlage für die vergleichende Betrachtung der Agenten mit XCS Algorithmen in Kapitel 10 dienen, insbesondere werden sie Anhaltspunkte dafür geben, welche Szenarien welche Eigenschaften der Algorithmen testen.

TODO: Ziel: Schwere Szenarien finden (schwierig für zufälligen, leicht für einfache heuristik)

6.1 Statistische Merkmale

Da keiner der hier vorgestellten Algorithmen lernt und somit statische Regeln besitzt, ist es nicht notwendig, die Qualitäten der Algorithmen bei verschiedener Anzahl von Zeitschritten zu betrachten und zu vergleichen, die Zahl der Zeitschritte wird somit, soweit

nicht anders angegeben, standardmäßig auf 500 festgesetzt. Außerdem sollen in den Statistiken die Werte jeweils über einen Lauf von 10 Experimenten mit jeweils 10 Problemen (siehe Kapitel 2.1) ermittelt und gemittelt werden. Die in den Tabellen jeweils angegebenen Werte sind auf zwei Stellen nach dem Komma gerundet.

Während eines Testlaufs werden eine ganze Reihe von statistischen Merkmalen erfasst. Wesentliches Merkmal zum Vergleich der Algorithmen ist der Wert der Qualität (siehe Kapitel 2.6), weitere Merkmale dienen zur Erklärung, warum z.B. ein Algorithmus bei einem Durchlauf schlechte Ergebnisse lieferte, bzw. dienen zum Testen und Finden von Fehlern oder Schwächen des Simulationsprogramms. Im Einzelnen sind hier zu nennen:

1. Anteil Sprünge des Zielagenten (siehe Kapitel 4.1), Durchläufe mit hohen Werte müssten verworfen werden
2. Anteil blockierter Bewegungen der Agenten
3. Halbzeitqualität (siehe Kapitel 2.6), größere Unterschiede zur ermittelten Qualität deuten darauf hin, dass sich der Algorithmus noch nicht stabilisiert hat und das Szenario mit höherer Schrittzahl erneut durchgeführt werden sollte
4. Abdeckung
5. Varianz der individuellen Punkte, ungefähres Maß, inwieweit einzelne Agenten an der Gesamtqualität beteiligt waren

6.1.1 Abdeckung

Die theoretisch maximal mögliche Anzahl an Felder, die die Agenten innerhalb ihrer Überwachungsreichweite zu einem Zeitpunkt haben können, entspricht der Zahl der Agenten multipliziert mit der Zahl der Felder die ein Agent in seiner Übertragungsreichweite haben kann. Ist dieser Wert größer als die Gesamtzahl aller freien Felder, wird stattdessen dieser

Wert benutzt.

Teilt man nun die Anzahl der momentan tatsächlich überwachten Felder durch die eben ermittelte maximal mögliche Anzahl an überwachten Felder, erhält man die Abdeckung, die die Agenten momentan erreichen.

6.2 Zielobjekt mit zufälligem Sprung

Im folgenden sollen alle.

In allen Szenarien mit dieser Form der Bewegung des Zielobjekts kommt es nur darauf an, dass die Agenten einen möglichst großen Bereich des Torus abdecken.

6.2.1 Szenario ohne Hindernisse

Ohne Hindernisse gibt sich ein klares Bild (siehe Tabelle 6.2), die intelligente Heuristik ist etwas besser als der des zufälligen Agenten und der einfachen Heuristik. Ein möglichst weiträumiges Verteilen auf dem Torus führt zum Erfolg, was sich auch in einem hohen Wert der Abdeckung zeigt, denn genau das wird mit dem völlig zufällig springenden Agenten getestet. Auch ist die Zahl der blockierten Bewegungen deutlich niedriger, was sich auch mit der Haltung des Abstands erklären lässt.

Die einfache Heuristik schneidet dagegen etwas schlechter als eine zufällige Bewegung ab. Zwar ist die Zahl der blockierten Bewegungen geringer, was sich dadurch erklären lässt, dass die einfache Heuristik zumindest an einem Punkt eine Sichtbarkeitsüberprüfung für die Richtung durchführt, in der sie sich bewegen möchte (nämlich wenn das Zielobjekt in Sicht ist), andererseits ist die Abdeckung etwas geringer. Dies kommt daher, dass, wenn mehrere Agenten das Zielobjekt in derselben Richtung in Sichtweite haben, mehrere Agenten sich in dieselbe Richtung bewegen. Dies beeinträchtigt die zufällige Verteilung

der Agenten auf dem Spielfeld und führt somit auch zu einer niedrigeren Abdeckung des Torus.

Bezüglich der Anzahl der Agenten ergeben sich keine Besonderheiten, mit steigender Agentenzahl steigt die Zahl der blockierten Bewegungen (aufgrund größerer Anzahl von blockierten Feldern), während die Abdeckung sinkt (aufgrund sich überlappender Überwachungsreichweiten).

Tabelle 6.1: “Total Random” ohne Hindernisse

Algorithmus	Agentenzahl	Blockierte Bewegungen	Abdeckung	Qualität
Zufällige Bewegung	8	2,82%	73,78%	32,36%
Einfache Heuristik	8	2,79%	73,22%	32,10%
Intelligente Heuristik	8	0,64%	81,26%	35,91%
Zufällige Bewegung	12	4,32%	69,55%	44,75%
Einfache Heuristik	12	4,19%	68,88%	43,86%
Intelligente Heuristik	12	1,49%	77,60%	49,49%
Zufällige Bewegung	16	5,82%	64,28%	54,55%
Einfache Heuristik	16	5,66%	63,65%	53,99%
Intelligente Heuristik	16	2,85%	71,44%	60,73%

6.2.2 Säulenszenario

Für das Säulenszenario (siehe Tabelle ??) ergeben sich erwartungsgemäß ähnliche Werte wie im Fall des Szenarios ohne Hindernisse (siehe Tabelle 6.2). Durch geringere Sicht und höhere Zahl an blockierten Bewegungen ergibt sich jeweils eine geringere Abdeckung und auch jeweils eine geringere Qualität.

6.2.3 Zufällig verteilte Hindernisse

Hier ergibt sich für alle Einstellungen für λ_h und λ_p (siehe Kapitel 2.5.2) ebenfalls ein klares Bild (siehe Tabelle 6.3), die intelligente Heuristik liegt wieder vorne, gefolgt wieder von der einfachen Heuristik und der zufälligen Bewegung. Die einfache Heuristik schneidet

Tabelle 6.2: "Total Random" ohne Hindernisse

Algorithmus	Agentenzahl	Blockierte Bewegungen	Abdeckung	Qualität
Zufällige Bewegung	8	4,45%	72,11%	32,13%
Einfache Heuristik	8	4,08%	71,70%	31,99%
Intelligente Heuristik	8	2,34%	79,61%	35,29%
Zufällige Bewegung	12	5,93%	67,72%	44,44%
Einfache Heuristik	12	5,67%	67,23%	43,81%
Intelligente Heuristik	12	3,62%	75,86%	49,34%
Zufällige Bewegung	16	7,62%	62,53%	54,26%
Einfache Heuristik	16	7,23%	62,00%	53,58%
Intelligente Heuristik	16	5,18%	69,91%	60,43%

minimal besser ab als die zufällige Bewegung, die Zahl der blockierten Bewegungen scheint hier stärker ins Gewicht zu fallen. Insbesondere die intelligente Heuristik scheint Probleme mit den Hindernissen zu haben. Da Hindernisse in der Heuristik nicht beachtet werden, erzeugt die maximale Ausbreitung der Agenten einen Bewegungsdruck gegen sie.

Tabelle 6.3: "Total Random" mit Hindernisse (12 Agenten)

Algorithmus	λ_h	λ_p	Blockierte Bewegungen	Abdeckung	Qualität
Zufällige Bewegung	0.2	0.99	14.25%	57.93%	46.89%
Einfache Heuristik	0.2	0.99	11.96%	57.94%	46.99%
Intelligente Heuristik	0.2	0.99	17.76%	63.17%	51.39%
Zufällige Bewegung	0.1	0.99	9.11%	64.04%	45.75%
Einfache Heuristik	0.1	0.99	7.76%	63.82%	45.34%
Intelligente Heuristik	0.1	0.99	9.68%	70.50%	50.40%
Zufällige Bewegung	0.1	0.5	11.89%	61.82%	44.62%
Einfache Heuristik	0.1	0.5	10.27%	62.20%	44.32%
Intelligente Heuristik	0.1	0.5	12.21%	68.96%	49.16%

Agenten. Der wesentliche zweite Faktor ist hier, dass der einfache Agent, wenn er das Zielobjekt in Sicht hat, davon ausgehen kann, dass sich in dieser Richtung wahrscheinlich kein Hindernis befindet, während der zufällige Agent Hindernisse überhaupt nicht beachtet, somit öfters gegen ein Hindernis läuft und letztlich öfters stehen bleibt. Der Un-

terschied zwischen beiden Agenten ist besonders hoch in Szenarien mit größerem Anteil an Hindernissen.

Ansonsten liegt der intelligente Agent wieder eindeutig vorne, beherrscht aber besonders gut Szenarien mit hohem “Verknüpfungsfaktor” (1.0) der geringem Anteil an Hindernissen (0.1), bei denen er bis zu etwa 15% über dem Ergebnis des einfachen Agenten liegt.

Dies liegt daran, dass Szenarien mit hohem “Verknüpfungsfaktors” bedeuten, dass alle Hindernisse zusammenhängend einen großen Block bilden und somit dem Szenario ohne Hindernissen ähnlich sind, auf dem dieser Agent ja besonders gut abschneidet. In zerklüftete Szenarien hat der Algorithmus dagegen Schwierigkeiten um andere Agenten überhaupt zu Gesicht bekommen, der Vorteil der Verteilung fällt also zu einem Teil weg.

Dies bestätigt auch ein Durchlauf bei dem Behinderungen der Sicht durch Hindernisse deaktiviert sind. Hierbei erreicht der intelligente Agent im Szenario (0.4, 0.1) statt TODO evtl weg

TODO

FAZIT:

Je schneller, zufälliger

6.3 “Zufälliger Nachbar” und “Einfache Richtungsänderung”

Wesentlicher Punkt bei beiden Bewegungstypen (siehe 4.2.2, 4.2.3) ist, dass der jetzige Ort des Zielobjekts maximal zwei Felder (die Standardgeschwindigkeit des Zielobjekts in den Tests) vom Ort in der vorangegangenen Zeiteinheit entfernt ist. Somit ist ein lokales Einfangen eher von Relevanz, wenn auch das Zielobjekt grundsätzlich schneller als andere

Agenten ist.

Wesentlicher Unterschied zwischen beiden Bewegungstypen ist, dass das Zielobjekt mit Bewegungstyp “zufälliger Nachbar” bei einer Bewegungsgeschwindigkeit von 2 mit einer Wahrscheinlichkeit von $\frac{1}{4}$ auf das ursprüngliche Feld zurückkehrt, innerhalb eines Zeitschritts also stehenbleibt. Wie die Ergebnisse in Tabellen 6.5 und 6.6 zeigen (TODO vielleicht noch näher darauf eingehen), ergibt sich dadurch ein leichteres Szenario. Ein mitunter stehenbleibender Agent kann mittels Heuristiken leichter überwacht werden, während es keine signifikante Veränderung bei der zufälligen Bewegung ergibt. In weiteren Tests soll deswegen immer nur die Bewegungsform “Einfache Richtungsänderung” getestet werden.

Tabelle 6.4: Vergleich “Zufälliger Nachbar” und “Einfache Richtungsänderung” (12 Agenten, ohne Hindernisse)

Algorithmus	Blockierte Bewegungen	Abdeckung	Qualität
“Zufälliger Nachbar”			
Zufällige Bewegung	4.26%	69.41%	45.75%
Einfache Heuristik	8.24%	61.77%	80.32%
Intelligente Heuristik	5.20%	70.09%	84.20%
“Einfache Richtungsänderung”			
Zufällige Bewegung	4.23%	69.63%	48.79%
Einfache Heuristik	7.20%	62.71%	69.78%
Intelligente Heuristik	4.07%	71.24%	74.53%

6.4 “Intelligent Open” und “Intelligent Hide”

6.7 6.8 10.2

TODO: Erläuterung!

Zu beachten sei, dass im Fall von “Intelligent Hide” eine relativ große Nummer an Sprüngen des Zielobjekts (siehe Kapitel 4.1) stattgefunden hat, was die Ergebnisse etwas

Tabelle 6.5: Vergleich “Zufälliger Nachbar” und “Einfache Richtungsänderung” (12 Agenten, zufälliges Szenario mit $\lambda_h = 0.2$, $\lambda_p = 0.99$)

Algorithmus	Blockierte Bewegungen	Abdeckung	Qualität
“Zufälliger Nachbar”			
Zufällige Bewegung	14.56%	57.80%	47.32%
Einfache Heuristik	18.29%	51.22%	85.92%
Intelligente Heuristik	22.33%	57.06%	88.31%
“Einfache Richtungsänderung”			
Zufällige Bewegung	14.60%	57.78%	47.92%
Einfache Heuristik	17.11%	52.38%	78.39%
Intelligente Heuristik	21.54%	57.76%	82.31%

Tabelle 6.6: Vergleich “Zufälliger Nachbar” und “Einfache Richtungsänderung” (12 Agenten, Säulenszenario)

Algorithmus	Blockierte Bewegungen	Abdeckung	Qualität
“Zufälliger Nachbar”			
Zufällige Bewegung	5.94%	67.75%	43.37%
Einfache Heuristik	10.41%	59.85%	85.61%
Intelligente Heuristik	7.82%	68.10%	88.98%
“Einfache Richtungsänderung”			
Zufällige Bewegung	6.02%	67.60%	45.83%
Einfache Heuristik	9.54%	60.34%	76.05%
Intelligente Heuristik	6.90%	68.75%	81.28%

verzerrt, die Zahl hält sich aber noch in Grenzen (bis zu ca. 0.5% im Fall der einfachen und intelligenten Heuristik im Fall mit vielen Hindernissen).

6.5 Always Same Direction

TODO

leeres Szenario

Tabelle 6.7: Vergleich “Intelligent Open” und “Intelligent Hide” (8 Agenten, ohne Hindernisse)

Algorithmus	Abdeckung	Qualität
“Intelligent Open”		
Zufällige Bewegung	74.15%	11.32%
Einfache Heuristik	60.90%	82.86%
Intelligente Heuristik	69.62%	85.74%
“Intelligent Hide”		
Zufällige Bewegung	74.13%	12.26%
Einfache Heuristik	69.43%	55.31%
Intelligente Heuristik	74.87%	64.41%

Tabelle 6.8: Vergleich “Intelligent Open” und “Intelligent Hide” (8 Agenten, zufälliges Szenario mit $\lambda_h = 0.2$, $\lambda_p = 0.99$)

Algorithmus	Abdeckung	Qualität
“Intelligent Open”		
Zufällige Bewegung	62.54%	13.37%
Einfache Heuristik	52.23%	84.33%
Intelligente Heuristik	56.92%	85.12%
“Intelligent Hide”		
Zufällige Bewegung	62.52%	13.10%
Einfache Heuristik	50.17%	90.32%
Intelligente Heuristik	56.94%	90.45%

6.6 XCS

Wird weiter unten besprochen.

6.7 Zusammenfassung

Wie wir gesehen haben gibt es also Szenarien in denen Abdeckung kaum eine Rolle spielt und lokale Entscheidungen eine wesentliche Rolle spielen. Dies wird es erleichtern, geeignete Szenarien im Kapitel 11 “Kommunikation” zu finden.

Tabelle 6.9: Vergleich “Intelligent (Open)” und “Intelligent (Hide)” (8 Agenten, Säulenszenario)

Algorithmus	Abdeckung	Qualität
“Intelligent (Open)”		
Zufällige Bewegung	72.55%	11.58%
Einfache Heuristik	57.19%	85.58%
Intelligente Heuristik	64.26%	91.18%
“Intelligent (Hide)”		
Zufällige Bewegung	72.56%	11.78%
Einfache Heuristik	58.45%	80.98%
Intelligente Heuristik	65.65%	86.38%

Kapitel 7

XCS

Jeder Agent besitzt ein unabhängiges, sogenanntes *eXtended Classifier System* (XCS), welches einem speziellen *learning classifier system* (LCS) entspricht. Ein LCS ist ein evolutionäres Lernsystem, das aus einer Reihe von *classifier* Regeln besteht, die zusammen ein sogenanntes *classifier set* bilden (siehe Kapitel 7.1). Eine allgemeine Einführung in LCS findet sich z.B. in [But06a].

Eine wesentliche Erweiterung des LCS ist das sogenannte *accuracy-based* XCS, zuerst beschrieben in [Wil95]. Neben neuer Mechanismen zur Generierung neuer *classifier* (insbesondere im Bereich bei der Anwendung des genetischen Operators) ist im Vergleich zum LCS gibt es vor allem in der Berechnung der *fitness* Werte der *classifier* Unterschiede. Während der *fitness* Wert beim einfachen LCS lediglich auf dem *reward prediction error* Wert basierte, basiert bei XCS der *fitness* Wert auf der Genauigkeit der jeweiligen Regel. Eine ausführliche Beschreibung findet sich in [But06b].

Im einfachsten Fall, im sogenannten *single step* Verfahren erfolgt die Bewertung einzelner *classifier*, also der Bestimmung eines jeweils neuen *fitness* Werts, sofort nach Aufruf

jeder einzelnen Regel, während im sogenannten *multi step* Verfahren mehrere aufeinanderfolgende Regeln erst dann bewertet werden, sobald ein Ziel erreicht wurde.

Ein klassisches Beispiel für den Test *single step* Verfahren ist das 6-Multiplexer Problem (z.B. in [But06b]), bei dem das XCS einen Multiplexer simulieren soll, der bei der Eingabe von 2 Adressbits und 4 Datenbits das korrekte Datenbit liefert. Sind beispielsweise die 2 Adressbits auf „10“ und die 4 Datenbits auf „1101“, so soll das dritte Datenbit, also „0“ zurückgeben. Im Gegensatz zum Überwachungsszenario kann also über die Qualität eines XCS direkt bei jedem Schritt entschieden werden.

Ein klassisches Beispiel für *multi step* Verfahren ist das „Maze N “ Problem, bei dem durch ein Labyrinth mit dem kürzesten Weg von N Schritten gegangen werden muss. Am Ziel angekommen wird der zuletzt aktivierte *classifier* positiv bewertet und das Problem neugestartet. Bei den Wiederholungen erhält jede Regel einen Teil der Bewertung des folgenden *classifier*. Somit wird eine ganze Kette von *classifier* bewertet und sich der optimalen Wahrscheinlichkeitsverteilung angenähert, welche repräsentiert, welche der Regeln in welchem Maß am Lösungsweg beteiligt sind.

Als Demonstration soll das in Abbildung 7.1 dargestellte (sehr einfache) Szenario dienen. Die zum Agenten zugehörigen *classifier* sind in Abbildung 7.2 dargestellt, wobei die 4 angrenzenden Felder für jeden *classifier* jeweils die Konfiguration der Kondition darstellt und der Pfeil die Aktion (für eine genauere Beschreibung eines *classifier* siehe Kapitel *classifier:sec*). Im ersten Durchlauf werden alle *classifier* in jedem Schritt zufällig gewählt, dann erhält *classifier e*) eine positive Bewertung. Im zweiten Durchlauf erhält dann *classifier c*) einen von *classifier e*) weitergegebene positive Bewertung und *classifier e*) auf Position 3 wird mit höherer Wahrscheinlichkeit als *classifier f*) gewählt. Das geht so lange weiter bis sich für *classifier b, c, e, g* ein ausreichend großer Wert eingestellt hat und keine wesentlichen Veränderungen mehr auftreten.

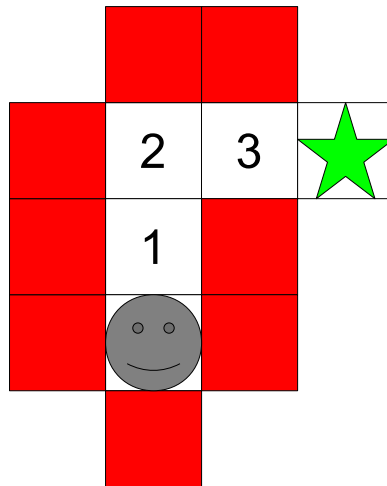


Abbildung 7.1: Einfaches Beispiel zum XCS *multi step* Verfahren

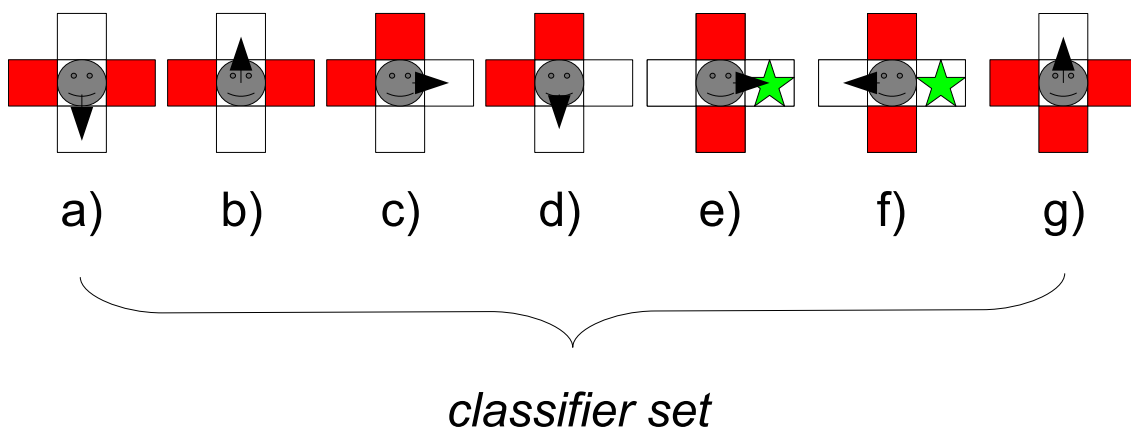


Abbildung 7.2: Vereinfachte Darstellung eines *classifier set* für das Beispiel zum XCS *multi step* Verfahren

Eine nähere Beschreibung bezüglich der Implementierung von und dem Unterschied zwischen dem *single step* und *multi step* Verfahren findet sich in [BW01].

Die in dieser Arbeit verwendete Implementierung entspricht im Wesentlichen der Standardimplementierung des *multi step* Verfahrens von [But00] (mit der algorithmischen Beschreibung des Algorithmus in [BW01]), eine

Besonderheit stellt allerdings die Problemdefinition dar, da es kein Ziel zu erreichen gibt, sondern über die Zeit hinweg ein bestimmtes Verhalten erreicht werden soll (die Überwachung des Zielobjekts). Somit gibt es auch kein Neustart des Problems und keinen festen Start- oder Zielpunkt. Zusätzlich, durch die Bewegung der anderen Agenten und des Zielobjekts, verändert sich die Umwelt in jedem Schritt, ein Lernen durch Wiederholung gemachter Bewegungsabläufe ist deswegen deutlich schwieriger.

Die meisten Implementationen und Varianten von XCS beschäftigen sich mit Szenarios, bei denen das Ziel in einer statischen Umgebung gefunden werden muss. Häufiger Gegenstand der Untersuchung in der Literatur sind insbesondere relativ einfache Probleme 6-Multiplexer Problem und Maze1 (z.B. in [But06b] [Wil95] [Wil98]), während XCS mit Problemen größerer Schrittzahl zwischen Start und Ziel Probleme hat [Bar02] [BDE⁺99]. Zwar gibt es Ansätze um auch schwierigere Probleme besser in den Griff zu bekommen (z.B. Maze5, Maze6, Woods14 in [BGL05]), indem ein Gradientenabstieg in XCS implementiert wurde. Ein konkreter Bezug zu einem dynamischen Überwachungsszenario konnte jedoch in keiner dieser Arbeiten gefunden werden.

Bezüglich Multiagentensystemen und XCS gibt es hauptsächlich Arbeiten, die auf zentraler Steuerung bzw. *OCS* [THN⁺98] basieren, also im Gegensatz zum Gegenstand dieser Arbeit auf eine übergeordnete Organisationseinheit bzw. auf globale Regeln oder

globalem Regeltausch zwischen den Agenten zurückgreifen.

Arbeiten bezüglich Multiagentensysteme in Verbindung mit LCS im Allgemeinen finden sich z.B. in [TB06], wobei es auch dort zentrale Agenten gibt, mit deren Hilfe die Zusammenarbeit koordiniert werden soll, während in dieser Arbeit alle Agenten dieselbe Rolle spielen sollen.

Vielversprechend war der Titel der Arbeit [LWB08], “Generation of Rule-based Adaptive Strategies for a Collaborative Virtual Simulation Environment”. Leider wird in der Arbeit nicht diskutiert, auf was sich der kollaborative Anteil bezog, da nicht mehrere Agenten benutzt worden sind. Auch konnte dort jeder einzelne Schritt mittels einer *reward* Funktion bewertet werden, da es globale Information gab. Dies vereinfacht ein solches Problem deutlich und macht einen Vergleich schwierig.

Eine weitere Arbeit in dieser Richtung ([HFA02]) beschreibt das „El Farol“ Bar Problem (EFBP), welches dort mit Hilfe eines Multiagenten XCS System erfolgreich gelöst wurde. Die Vergleichbarkeit ist hier auch eingeschränkt, da es sich bei dem EFBP um ein *single step* Problem handelt.

Eine der dieser Arbeit (bezüglich Multiagentensysteme) am nächsten kommende Problemstellung wurde in [ITS05] vorgestellt. Dort wurde der *base reward* (TODO) unter den (zwei) Agenten aufgeteilt, es fand also eine Kommunikation des *reward* Werts statt. Wie das Ergebnis in Verbindung mit den Ergebnissen dieser Arbeit interpretiert werden kann, wird in Kapitel 11 diskutiert.

In [KM94] wurde gezeigt, dass bei der Weitergabe des *base reward* Gruppenbildung von entscheidender Wichtigkeit ist. Nach bestimmten Kriterien werden Agenten in Grup-

pen zusammengefasst und der *base reward* anstatt an alle, jeweils nur an die jeweiligen Gruppenmitgliedern weitergegeben. Dies bestätigen auch Tests in Kapitel 11, bei der sich Agenten mit ähnelnden (was das Verhalten gegenüber anderen Agenten betrifft) *classifier set* Listen in Gruppen zusammengefasst wurden und zum Teil bessere Ergebnisse erzielt werden konnten als ohne Kommunikation.

[BD03] TODO

TODO In Kapitel 9 werden dann die Implementierungen der `calculateReward`, `calculateNextMove` beschrieben TODO Limits in Long Path Learning with XCS Gamma!

7.1 Übersicht

Ein XCS ist ein regelbasiertes evolutionäres Lernsystem, das im Wesentlichen aus folgenden Elementen besteht:

TODO ausführlicher

7.2 Ablauf eines XCS

1. Vervollständigung der *classifier* Liste (*covering*, siehe Kapitel 7.2.1)
 2. Auswahl auf die Sensordaten passender *classifier* (*matching*, siehe Kapitel 7.2.2)
 3. Bestimmung der Auswahlart der Aktion (*explore/exploit*, siehe Kapitel 8.9)
 4. Auswahl der Aktion TODO
 5. Erstellung des zur Aktion zugehörigen Liste von *classifier* (*actionSet*, siehe Kapitel 7.2.3)
- , so dass es in der Liste *classifier* deren

TODO Bei der Auswahl einer Aktion werden alle *classifier* mit *condition* Vektoren gesucht, die auf die aktuellen Sensordaten passen. Diese bilden dann das *matchSet*.

6. Im nächsten Schritt wählen wir einen *classifier* aus diesem *matchset* aus und speichern dessen Aktion.
7. Schließlich bilden wir anhand des *match set* und der gewählten Aktion das *action set*

7.2.1 Covering

TODO

7.2.2 Variable *lastMatchSet*

In der *lastMatchSet* Variable werden jeweils alle *classifier* gespeichert, die den letzten Sensordatenvektor erkannt haben. Sie entspricht dem *predictionArray* in der originalen Implementierung von XCS, dort werden nämlich außerdem Vorberechnungen zur Auswahl des nächsten *classifier* für die Bewegung durchgeführt und die Ergebnisse gespeichert.

7.2.3 Variable *actionSet*

Ein *actionSet* ist jeweils einer Zeiteinheit zugeordnet. Dort werden jeweils alle *classifier* gespeichert, die zu diesem Zeitpunkt denselben *action* Wert besitzen wie der für die Bewegung bestimmte *classifier*. In der Standardimplementierung von XCS wird jeweils nur das letzte *actionSet* gespeichert, während in SXCS eine ganze Reihe (bis zu *maxStackSize* Stück) gespeichert werden.

1. Einer Menge an Regeln, sogenannte *classifier* (siehe Kapitel 7.3), die zusammen ein *classifier set* bilden

2. Einem Mechanismus zur Auswahl einer Aktion aus dem *classifier set* (siehe Kapitel 8.9)
3. Einem Mechanismus zur Zusammenfassung aller *classifier* aus dem *classifier set* mit gleicher Aktion zu einer *action set* Liste.
4. Einem Mechanismus zur Evolution der *classifier* (mittels genetischer Operatoren, siehe Kapitel 7.5)
5. Einem Mechanismus zur Bewertung der *classifier* (mittels *reinforcement learning*, siehe Kapitel 7.6)

Während die ersten drei Punkte bei allen hier vorgestellten XCS Varianten identisch sind, gibt es wesentliche Unterschiede bei der Bewertung der *classifier*. Diese werden gesondert in Kapitel 9 im Einzelnen besprochen. Im Folgenden sollen nun die ersten drei Punkte näher betrachtet werden.

7.3 Classifier

Ein *classifier* besteht aus einer Anzahl im folgenden diskutierten Variablen die anhand der in Kapitel 8 aufgelisteten Werte initialisiert werden. Wesentliche Teile sind der *condition* Vektor (Kapitel 7.3.1) und der *action* Wert (Kapitel 7.3.4), alle restlichen Variablen dienen zur Berechnung der Wahrscheinlichkeit mit der der *classifier* ausgewählt und dessen *action* Wert ausgeführt wird.

7.3.1 Der *condition* Vektor

Der *condition* Vektor gibt die Kondition an, in welcher Situation der zugehörige *classifier* ausgewählt werden kann, d.h. welche Sensordaten von dem jeweiligen *classifier* erkannt

werden. Der Aufbau des Vektors entspricht dem Vektor der über die Sensoren erstellt wird (siehe Kapitel 3.1).

$$\underbrace{z_{sN} z_{rN} z_{sO} z_{rO} z_{sS} z_{rS} z_{sW} z_{rW}}_{\text{Erste Gruppe (Zielobjekt)}} \underbrace{a_{sN} a_{rN} a_{sO} a_{rO} a_{sS} a_{rS} a_{sW} a_{rW}}_{\text{Zweite Gruppe (Agenten)}} \underbrace{h_{sN} h_{rN} h_{sO} h_{rO} h_{sS} h_{rS} h_{sW} h_{rW}}_{\text{Dritte Gruppe (Hindernisse)}}$$

7.3.2 Platzhalter im *condition* Vektor

Neben den zu den Sensordaten korrespondierenden Werten 0 und 1 soll es noch einen dritten Zustand, den Platzhalter „#“, geben, der anzeigen soll, dass beim Vergleich zwischen Kondition und Sensordaten diese Stelle ignoriert werden soll. Eine Stelle im *condition* Vektor mit Platzhalter gilt also als äquivalent zur korrespondierenden Stelle in den Sensordaten, egal ob sie mit 0 oder 1 belegt ist. Ein Vektor, der ausschließlich aus Platzhaltern besteht, würde somit bei der Auswahl immer in Betracht gezogen werden, da er auf alle möglichen Kombinationen der Sensordaten passt. Umgekehrt können dadurch bei der Auswahl der *classifier* mehrere *classifier* auf einen gegebenen Sensordatenvektor passen. Diese bilden dann die sogenannte *match set* Liste, aus welchem dann wie in Kapitel 8.9 beschrieben der eigentliche *classifier* ausgewählt wird.

7.3.3 Vergleich des *condition* Vektors mit den Sensordaten

Beim Vergleich der Sensordaten und Daten aus dem *condition* Vektor werden immer jeweils zwei Paare verglichen. In Kapitel 3.1 wurde erwähnt, dass der Fall (0/1) in den Sensordaten nicht auftreten kann, weswegen (um die Aufgabe nicht unnötig zu erschweren) ein Datenpaar (0/1) im *condition* Vektor äquivalent zum Datenpaar (1/1) sein soll, es damit also eine gewisse Redundanz gibt. Es ergeben sich also folgende Fälle:

1. Sensorenpaar (0/0) wird erkannt von (0/0), ($\#$, 0), (0, $\#$), ($\#$, $\#$)
2. Sensorenpaar (1/0) wird erkannt von (1/0), ($\#$, 0), (1, $\#$), ($\#$, $\#$)
3. Sensorenpaar (1/1) wird erkannt von (1/1), ($\#$, 1), (1, $\#$), ($\#$, $\#$), (0/1)

Beispielsweise würden folgende Sensordaten von den folgenden *condition* Vektoren erkannt:

Sensordaten:

(Zielobjekt in Sicht im Norden, Agent im Sicht im Süden,

Hindernisse im Westen und Osten)

10 00 00 00 . 00 00 11 00 . 00 11 00 11

Beispiele für erkennende *condition* Vektoren:

10 00 00 00 . ## ## ## ## . 00 ## ## ##

. ## ## #1 00 . 00 11 ##

#0 ## ## ## . ## ## 01 ## . ## 11 ## 11

7.3.4 Der *action* Wert

Wird ein *classifier* ausgewählt, wird eine bestimmte Aktion ausgeführt, die durch den *action* Wert determiniert ist. Im Rahmen dieser Arbeit entsprechen diese Aktionsmöglichkeiten den 4 Bewegungsrichtungen, die in Kapitel 3 besprochen wurden.

7.3.5 Der *fitness* Wert

Der *fitness* Wert soll die allgemeine Genauigkeit des *classifier* repräsentieren und wird über die Zeit hinweg sukzessive an die beobachteten *reward* Werte angepasst. Der Wertebereich verläuft zwischen 0.0 und 1.0 (maximale Genauigkeit). Insbesondere eines der

frühesten Werke zu XCS [Wil95] beschäftigte sich mit diesem Aspekt der Genauigkeit.

7.3.6 Der *reward prediction* Wert

Der *reward prediction* Wert des *classifier* stellt die Höhe des *reward* Werts dar, von dem der *classifier* erwartet, dass er ihn bei der nächsten Bewertung erhalten wird.

7.3.7 Der *reward prediction error* Wert

Der *reward prediction error* Wert soll die Genauigkeit des *classifier* bzgl. des *reward prediction* Werts (die durchschnittliche Differenz zwischen *reward prediction* und *reward*) repräsentieren. U.a. auf Basis dieses Werts wird der *fitness* Wert des *classifier* angepasst.

7.3.8 Der *experience* Wert

Der *experience* Wert des *classifier* repräsentiert die Anzahl, wie oft ein *classifier* aktualisiert wurde, also wieviel Erfahrung er sammeln konnte. Im Wesentlichen dient dieser Wert als Entscheidungshilfe, ob auf die anderen Werte des *classifier* vertraut werden kann bzw. ob der *classifier* als unerfahren gilt und somit z.B. bei Löschung und Subsumption gesondert behandelt werden muss.

7.3.9 Der *numerosity* Wert

Durch Subsumption (siehe Kapitel 7.4 und Kapitel 7.5) können *classifier* eine Rolle als *macro classifier* spielen, d.h. *classifier* die andere *classifier* in sich beinhalten. Der *nume-*

rosity Wert gibt an, wieviele andere, sogenannte *micro classifier* sich in dem jeweiligen *classifier* befinden.

7.4 Subsummation von *classifier*

Die Benutzung von den oben erwähnten Platzhaltern (Kapitel 7.3.2) erlaubt es dem XCS mehrere *classifier* zu zusammenzulegen, wodurch die Gesamtzahl der *classifier* sinkt und somit Erfahrungen, die ein XCS Agent sammelt, nicht unbedingt mehrfach gemacht werden müssen. Die dahinter stehende Annahme ist, dass es Situationen gibt, in denen der Gewinn der durch Unterscheidung zwischen zwei verschiedenen Sensordatensätzen geringer ist als die Ersparnis durch das Zusammenlegen beider *classifier*, d.h. dem Ignorieren der Unterschiede.

Besitzt ein *classifier* sowohl einen genügend großen *experience* Wert als auch einen ausreichend kleinen *reward prediction error* Wert, so kann er als sogenannter *subsumer* auftreten. Andere *classifier* (in derselben *action set* Liste, also mit gleichem *action* Wert) werden durch den *subsumer* ersetzt, sofern der von ihnen abgedeckte Sensordatenbereich eine Teilmenge des von dem *subsumer* abgedeckten Bereichs ist, der *subsumer* also an allen Stellen des *condition* Vektors entweder denselben Wert wie der zu subsummierende *classifier* oder einen Platzhalter besitzt.

7.5 Genetische Operatoren

Es werden aus der jeweiligen *action set* Liste zwei *classifier* (die Eltern) zufällig ausgewählt und zwei neue *classifier* (die Kinder) aus ihnen gebildet und in die Population

eingefügt. Dabei wird mittels *two-point crossover* ein neuer *condition* Vektor generiert und der *action* Wert auf den der Eltern gesetzt (da sie aus derselben *action set* Liste stammen, ist der Wert beider Eltern identisch). Die restlichen Werte werden standardmäßig wie in Kapitel 8 aufgelistet initialisiert. Werden Kinder in die Population eingefügt, deren *action* Wert und *condition* Vektor identisch mit existierenden *classifier* ist, werden sie stattdessen subsummiert.

Da die Sensoren und somit auch der *condition* Vektor aus drei in sich geschlossenen Gruppen bestehen, werden im Unterschied zur Standardimplementation beim *crossing over* zwei feste Stellen benutzt, die die Gruppe für das Zielobjekt, die Gruppe für Agenten und die Gruppe für feste Hindernisse voneinander trennen.

Bezeichne (z_1, a_1, h_1) bzw. (z_2, a_2, h_2) jeweils die drei Gruppen (siehe Kapitel 7.3.1) des *condition* Vektors des ersten bzw. zweiten ausgewählten Elternteils, dann können für die drei Gruppen der *condition* Vektoren (z_{1k}, a_{1k}, h_{1k}) und (z_{2k}, a_{2k}, h_{2k}) der beiden Kinder folgende Kombinationen auftreten:

$$[(z_{1k}, a_{1k}, h_{1k}), (z_{2k}, a_{2k}, h_{2k})] = [(z_1, a_1, h_1), (z_2, a_2, h_2)]$$

$$[(z_{1k}, a_{1k}, h_{1k}), (z_{2k}, a_{2k}, h_{2k})] = [(z_2, a_1, h_1), (z_1, a_2, h_2)]$$

$$[(z_{1k}, a_{1k}, h_{1k}), (z_{2k}, a_{2k}, h_{2k})] = [(z_1, a_2, h_1), (z_2, a_1, h_2)]$$

$$[(z_{1k}, a_{1k}, h_{1k}), (z_{2k}, a_{2k}, h_{2k})] = [(z_2, a_2, h_1), (z_1, a_1, h_2)]$$

7.6 Bewertung der Aktionen (*base reward*)

6 unterschiedliche Möglichkeiten, goal in reward range goal in sight range goal in reward range, kein agent in reward range goal in sight range, kein agent in sight range goal in

reward range, kein agent in sight range goal in sight range, kein agent in reward range

TODO base reward und reward unterscheiden

Programm 7.1

```
1 /**
2  * @return true Falls das Zielobjekt von diesem Agenten überwacht wird
3  *   und kein anderer Agent in dieser Richtung in
4  *   Überwachungsreichweite steht
5  */
6  public boolean checkRewardPoints() {
7      boolean[] sensor_agent = lastState.getSensorAgent();
8      boolean[] sensor_goal = lastState.getSensorGoal();
9
10     for(int i = 0; i < Action.MAX_DIRECTIONS; i++) {
11         if((sensor_goal[2*i]) && (!sensor_agent[2*i+1])) {
12             return true;
13         }
14     }
15
16     return false;
17 }
```

Programm 7.1: Bestimmung des *base reward* Werts für Agenten

Kapitel 8

Parameter

Die Einstellungen der XCS Parameter der durchgeführten Experimente entsprechen weitgehend den Vorschlägen in [BW01] (“Commonly Used Parameter Settings”). Eine Auflistung findet sich in Tabelle 8.2. Im Folgenden sollen Parameter besprochen werden, die entweder in der Empfehlung offen gelassen sind, also klar vom jeweiligen Szenario abhängen, und solche, bei denen von der Empfehlung abgewichen wurde.

Mitunter führen andere Parametereinstellungen auch zu wesentlich besseren Ergebnissen. Dies muss man aber vorsichtig bewerten, wenn die erreichte Qualität unter der des zufälligen Algorithmus liegt, da eine Auswirkung sein kann, dass der Algorithmus nicht besser lernt, sondern sich umgekehrt eher wie der zufällige Algorithmus verhält. Ein Vergleich mit der Qualität des zufälligen Algorithmus wird deswegen jeweils immer angegeben.

Anzumerken sei, dass alle Tests jeweils mit den in Tabelle 8.2 angegebenen Parameterwerten durchgeführt wurden und bei jedem Test jeweils nur der zu untersuchende Wert verändert wurde. Um synchronisierte und vergleichbare Daten zu haben, wurden die Tests deshalb in mehreren Etappen durchgeführt, die angegebenen Testergebnisse entsprechen jeweils den endgültigen Ergebnissen.

8.1 Parameter *max population N*

Der Wert von *max population N* bezeichnet die maximalen Größe der *classifier set* Liste. Bei der Wahl eines geeigneten Werts spielt insbesondere die Laufzeit eine Rolle. Für den Overhead (d.h. die Zeit, die 8 Agenten mit zufälliger Bewegung benötigen) ergab sich eine mittlere Laufzeit von 1,67s pro Experiment bei 500 Schritten (bzw. 6,50s bei 2000 Schritten), was die anfängliche Stagnierung bis $N = 32$ erklärt. Zieht man diesen von den Messwerten (siehe Abbildung 8.2) ab, erhält man im betrachteten Wertebereich einen nahezu linearen Verlauf (siehe 8.3, ab $N > 128$). Der fallende Verlauf bis 128 erklärt sich durch den Overhead des XCS Algorithmus selbst.

Größere Werte für N erlauben eine bessere Anpassung, da weniger *classifier* während eines Laufs gelöscht werden müssen und mehr Plätze zur Speicherung der Erfahrungen zur Verfügung steht. Auf der anderen Seite werden mehr Schritte benötigt um die für die jeweiligen *classifier* ausreichend Erfahrung zu sammeln (siehe Abbildung 8.4) und, wie oben demonstriert, auch eine längere Laufzeit.

Ein zu kleiner Wert erhöht dagegen die Konkurrenz zwischen den *classifier*, da einzelne *classifier* mit höherer Wahrscheinlichkeit von besseren *classifier* verdrängt und somit gelöscht werden.

Die Tests liefen auf einem T7500, 2.2 GHz in einem einzelnen Thread. Als Vergleich hierzu wurde auch der Einfluss der Kartengröße auf die Laufzeit betrachtet, wie in Abbildung 8.1 zu sehen, ist der Einfluss auf die Laufzeit im getesteten Bereich (256 - 1024) ohne Bedeutung.

In den Tests wird somit $N = 128$ gesetzt, was als ausreichender Kompromiss zwischen den erwähnten Faktoren erscheint.

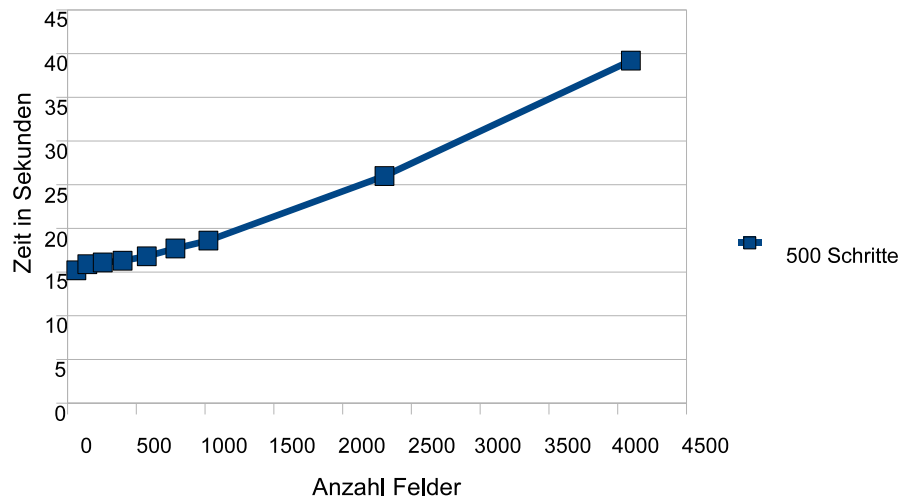


Abbildung 8.1: Darstellung der Auswirkung der Torusgröße auf die Laufzeit im leeren Szenario, zufälliger Bewegung des Zielobjekts, 8 sich zufällig bewegend Agenten

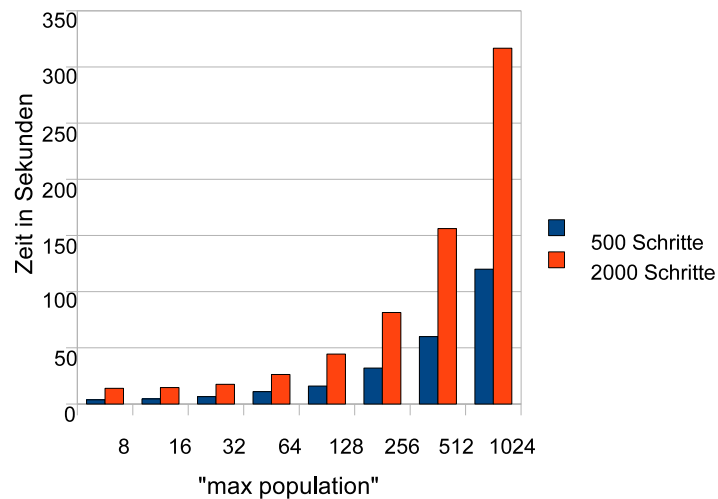


Abbildung 8.2: Darstellung der Auswirkung des Parameters *max population* N auf die Laufzeit im leeren Szenario, zufälliger Bewegung des Zielobjekts, 8 Agenten mit SXCS Algorithmus

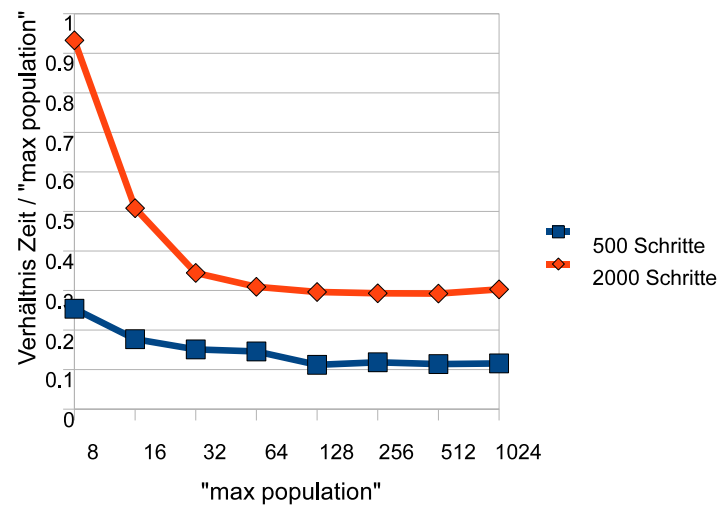


Abbildung 8.3: Darstellung der Auswirkung des Parameters *max population* N auf das Verhältnis der Laufzeit zu N im leeren Szenario, zufälliger Bewegung des Zielobjekts, 8 Agenten mit SXCS Algorithmus

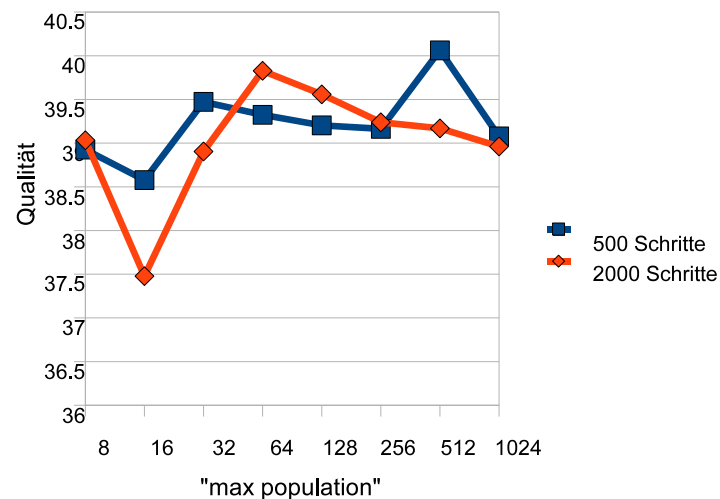


Abbildung 8.4: Darstellung der Auswirkung des Parameters *max population* N auf die Qualität im leeren Szenario, zufälliger Bewegung des Zielobjekts, 8 Agenten mit LCS Algorithmus

8.2 Maximalwert *reward*

TODO Der Wert der bei der Bewertung als *reward* vergeben wird hat lediglich ästhetische Auswirkungen und wurde auf 1.0 gesetzt. In der Standardimplementation von XCS (siehe Abbildung ??) ist der maximale *reward* äquivalent mit dem Maximalwert von ρ , da das Problem bei jedem positiven *reward* Wert neugestartet wird, also entweder der *reward* Wert aus dem letzten Schritt also immer 0 ist oder *maxPrediction* auf 0 gesetzt wurde, und $\rho = \text{reward} + \gamma \text{maxPrediction}$ gilt.

In den hier vorgestellten XCS Varianten wird dagegen der *reward* Wert absteigend, zusammen mit dem *maxPrediction* Wert, an frühere *actionSet* Listen verteilt, ρ kann also größer als 1.0 werden. In diesem Bereich ist noch Bedarf an theoretischer Forschung, in Tests haben sich Werte bis 3.0 ergeben, welche aber vom jeweiligen Szenario abhängen. Wird das Zielobjekt (z.B. wegen Hindernissen oder großen Torusdimensionen) eher selten gesehen, fällt der Wert geringer aus.

8.3 Parameter *accuracy equality* ϵ_0

TODO Der Parameter ϵ_0 gibt an, unter welchem Wert zwei *accuracy* Werte als gleich gelten sollen. Dies ist insbesondere bei der *subsummation* Funktion und der Berechnung des *accuracy* Werts von Bedeutung. In der Literatur [BW01] wird als Regel genannt, dass der Wert auf etwa 1% des Maximalwerts von ρ gesetzt werden soll, den der erwartete Reward annehmen kann. Aufgrund der Überlegungen in 8.2 wird ϵ_0 für die neuen XCS Varianten auf 0.02 gesetzt, während es für die Standardimplementation von XCS auf 0.01 gesetzt wird. Ein Testdurchlauf auf dem Säulenszenario (siehe Abbildung 8.5) ergibt aber, dass der Parameter keine besondere Auswirkung hat, weshalb der Wert auf 0.01 belassen

wird.

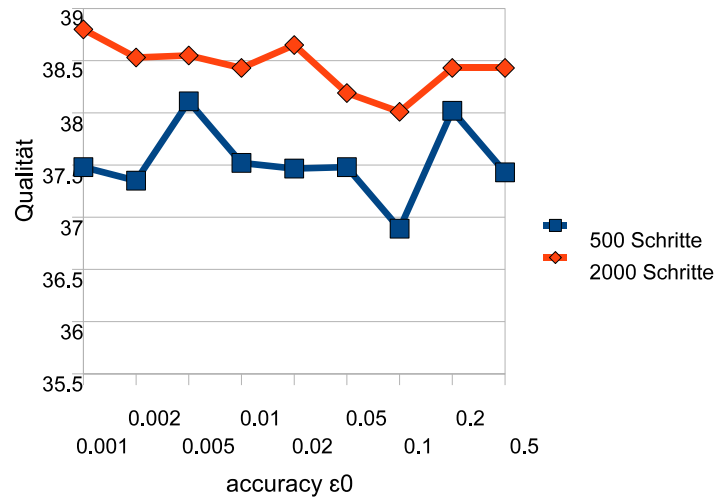


Abbildung 8.5: Auswirkung des Parameters *accuracy equality* ϵ_0 auf die Qualität im Säulenszenario, zufälliger Bewegung des Zielobjekts, 8 Agenten mit SXCS Algorithmus

8.4 Parameter *reward prediction discount* γ

dsxcs gut: alles an TODO maxpred überprüfen...

Der Einfluss von γ ist zwar vorhanden, aber sehr gering. TODO

TODO Reward prediction unnötig bei speed 2, pillar, random

TODO Abschnitt entfernen Auch für den Wert *reward prediction discount* γ hat sich ein etwas höherer Wert als sinnvoll erwiesen, als standardmäßig benutzt wird. Laut [BW01] hängt der Wert auch vom verwendeten Szenario ab. Ein höherer Wert für γ bedeutet, dass die Höhe des Werts, der über *maxPrediction* weitergegeben wird, mit zeitlichem Abstand zur ursprünglichen Bewertung mit einem *reward*, weniger schnell abfällt, wodurch eine längere Verkettung von *reward* Werten möglich ist. Umgekehrt führen zu hohe Werte für γ zu der positiven Bewertung von *classifiers* die am Erfolg gar nicht beteiligt waren, was sich negativ auf die Qualität auswirken kann.

Tabelle ?? zeigt einen Vergleich der Qualität mit dem Standardwert $\gamma = 0.71$ und dem für die in dieser Arbeit verwendeten Testszenarien gewählten Wert $\gamma = 0.95$.

TODO 0.71 lassen

auch mit Geschwindigkeit 0.1 oder so testen

TODO Tabelle prediction discount

8.5 Parameter Lernrate β

Für die Lernrate β hat sich ein etwas niedrigerer als in der Literatur angegebener Wert (0.01) als erfolgreich erwiesen. Die Lernrate bestimmt, wie stark ein ermittelter *reward* Wert den *reward prediction*, *reward prediction error*, *fitness* und *action set size* Wert pro Aktualisierung beeinflusst. TODO Auch dieser Parameter ist szenariospezifisch, über die konkrete Begründung kann nur spekuliert werden, die Schwierigkeit des Szenarios TODO

Vergleichende Tests (siehe Abbildung 8.6 mit niedrigerem bzw. höherem Wert haben zu einer etwas schlechteren Qualität geführt. TODO

8.6 Parameter *reward prediction init* p_i

In der Literatur werden Werte nahe Null bzw. 1% von ρ als Initialisierung für den *reward prediction* Wert eines *classifiers* angegeben. Im in dieser Arbeit untersuchten Fall, bei dem die Agenten nur begrenzte Sensorfähigkeiten besitzen, sich auf einem Torus frei bewegen können und keine festen Pfade suchen müssen, ist zu erwarten, dass sich die *reward prediction* Werte der einzelnen *classifier* untereinander wenig unterscheiden, während sie beispielsweise bei statischen Szenarien gegen feste, stark unterschiedliche Werte konvergieren. Beispielsweise im Einführungsbeispiel in Abbildung 7.1 würden die *reward prediction* Werte der *classifier* b), c), e) und g) eher gegen 1 und die der restlichen *classifier* gegen

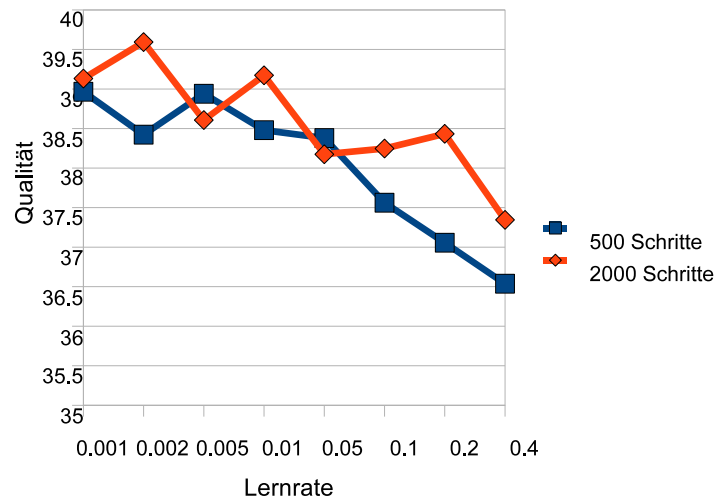


Abbildung 8.6: Auswirkung des Parameters *learning rate* β auf die Qualität im Säulenszenario, zufälliger Bewegung des Zielobjekts, 8 Agenten mit SXCS Algorithmus

0 streben.

Welchen Wert man für p_i nun als Durchschnittswert wählt, hängt vom jeweiligen Szenario ab. Beispielsweise würde ein Überwachungsszenario auf einem sehr größeren Torus mit relativ wenigen Agenten würde zu einem niedrigeren Durchschnittswert für die *reward prediction* Variable führen und umgekehrt.

Einfache Idee ist deshalb, auszunutzen, dass im jeweiligen *classifier set* bereits die Information enthalten ist, was der Durchschnittswert für das aktuelle Szenario ist, nämlich der Durchschnittswert aller *reward prediction* Werte. Dies kann man noch dadurch erweitern, dass man bei der Durchschnittsbildung nur solche *classifier* miteinbezieht, welche einen ausreichend großen *experience* Wert besitzen.

Tests haben gezeigt, dass dadurch ein deutlich schnelleres Konvergenzverhalten erreicht werden konnte

TODO Test

TODO raus unten Wählt man einen Wert, der näher am Durchschnitt der *reward*

prediction Werte der *classifier* liegt, die sich in den besten Lösungen am Ende eines Testdurchlaufs befinden, so ist zu erwarten, dass die Anzahl der benötigten Aktualisierungen des *reward prediction* Werts geringer ausfällt, das System also schneller konvergiert. Diese Überlegung wird bestätigt durch entsprechende Tests (siehe 8.7).

Wir setzen somit für SXCS den Parameter auf $p_i = 1.0$.

Zu beachten ist, dass diese Überlegung primär deswegen gilt, weil die
 TODO Standardverfahren?

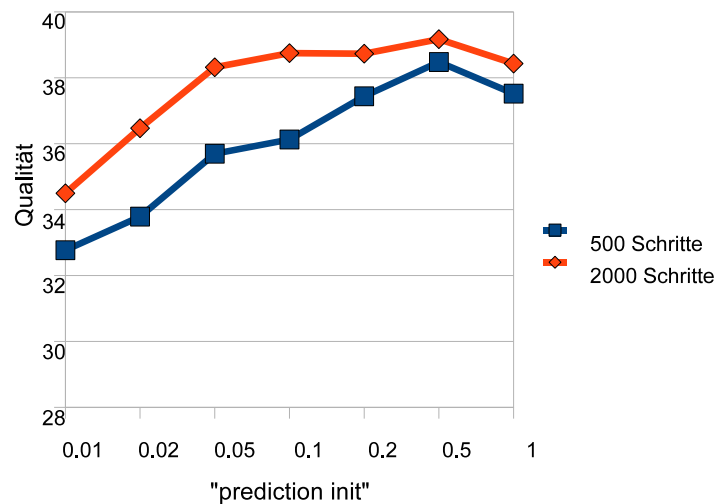


Abbildung 8.7: Darstellung Auswirkung des Parameters *reward prediction init* p_i auf die Qualität im Säulenszenario, zufälliger Bewegung des Zielobjekts, 8 Agenten mit SXCS Algorithmus

8.7 Zufällige Initialisierung der *classifier set* Liste

Normalerweise werden XCS Systeme mit leeren *classifier set* Listen initialisiert, als Option wird jedoch auch eine zufällige Initialisierung erwähnt ([But06b]), bei der zu Beginn die *classifier set* Liste mit mehreren *classifiers* mit zufälligen *action* Werten und *condition* Vektoren gefüllt wird.

Tests haben gezeigt (siehe Tabelle 8.1), dass dadurch minimal bessere Ergebnisse erzielt

werden, allerdings nur in Szenarien mit ausreichender Schrittzahl (> 100). Dies lässt sich darauf zurückführen, dass anfänglich gefüllte *classifier set* Listen die *matchSet* Listen relativ groß lassen werden, somit die Auswirkungen anfänglichen Lernens geringer ausfallen und sich die Agenten eher wie sich zufällig bewegendende Agenten verhalten. Da hier durchgeführten Tests über 500 bzw. 2000 Schritte laufen, sollen somit die *classifier set* Listen mit zufällig generierten *classifiers* gefüllt werden.

Tabelle 8.1: Vergleichende Tests für den den Start mit und ohne zufällig gefüllten *classifier set* Listen

Algorithmus	Agentenzahl	Schrittzahl	Abdeckung	Qualität
Zufälliger Agent	8	500	60.64%	61.54%
Multistep	8	500	60.64%	61.54%
LCS	8	500	60.64%	61.54%
NewLCS	8	500	60.64%	61.54%
Zufälliges Szenario	8	500	60.64%	61.54%
Säulenszenario	8	100	1.0%	1.0%
LCS Ohne Drehung	8	2000	1.0%	1.0%
LCS Mit Drehung	8	100	1.0%	1.0%
LCS Mit Drehung	8	500	1.0%	1.0%
LCS Mit Drehung	8	2000	1.0%	1.0%
Zufälliger Agent	12	500	76.03%	76.59%
Einfacher Agent	12	500	67.30%	96.86%
Intelligenter Agent	12	500	86.85%	95.08%
LCS Ohne Drehung	12	100	1.0%	1.0%
LCS Ohne Drehung	12	500	1.0%	1.0%
LCS Ohne Drehung	12	2000	1.0%	1.0%
LCS Mit Drehung	12	100	1.0%	1.0%
LCS Mit Drehung	12	500	1.0%	1.0%
LCS Mit Drehung	12	2000	1.0%	1.0%

8.8 Übersicht über alle Parameterwerte

Tabelle 8.2: Verwendete Parameter (soweit nicht anders angegeben) und Standardparameter, TODO englisch/deutsch

Parameter	Wert	Standardwert (siehe [BW01])
Max population N	128 (siehe Kapitel 8.1)	
Max value ρ	1.0 (siehe Kapitel 8.2)	[10000]
Fraction mean fitness δ	0.1	[0.1]
Deletion threshold θ_{del}	20.0	[\sim 20.0]
Subsumption threshold θ_{sub}	20.0	[20.0+]
Covering # probability $P_{\#}$	0.5	[\sim 0.33]
GAthreshold θ_{GA}	25.0	[25-50]
Mutation probability μ	0.05	[0.01-0.05]
Prediction error reduction	0.25	[0.25]
Fitness reduction	0.1	[0.1]
Reward prediction init p_i	0.5, 1.0 (siehe 8.6)	[\sim 0]
Prediction error init ϵ_i	0.0	[0.0]
Fitness init F_i	0.01	[0.01]
Condition vector	zufällig (siehe Kapitel 8.7)	[zufällig oder leer]
Numerosity	1	[1]
Experience	0	[0]
Accuracy equality ϵ_0	0.05	[1% des größten Werts]
Accuracy calculation α	0.1	[0.1]
Accuracy power ν	5.0	[5.0]
Reward prediction discount γ	0.71	[0.71]
Learning rate β	0.01 (siehe 8.5)	[0.1-0.2]
exploration probability	0.5 (siehe 7.4)	[\sim 0.5]

8.9 Auswahlart der *classifier*

In jedem Zeitschritt gilt es zu entscheiden, welche Bewegung ein Agent ausführen soll. Als Basis der Entscheidung hat ein Agent zum einen die Sensordaten und zum anderen das eigene *classifier set* zur Verfügung. Da ein Sensordatensatz von mehreren *classifier* erkannt werden kann (siehe Kapitel 7.3.2), stellt sich die Frage, welchen *classifier* (und den dazugehörigen *action* Wert, der die Bewegung bestimmt) man aus dem gebildeten *matchSet* auswählen soll. In der ursprünglichen Implementierung [But00] wurden folgende Auswahlarten benutzt:

1. *random selection* : Zufällige Auswahl eines *classifiers*
2. *roulette wheel selection* : Zufällige Auswahl eines *classifier*, mit Wahrscheinlichkeit abhängig vom Produkt seines *fitness* und *reward prediction* Werts
3. *best selection* : Auswahl des *classifiers* mit dem höchsten Produkt aus seinen *fitness* und *reward prediction* Werten

Bei einem dynamischen Überwachungsszenario ist es im Vergleich zu standardmäßigen statischen Szenarien weder nötig noch hilfreich *random selection* zu nutzen. Die Idee für diese Auswahlart in einem statischen Szenario ist, dass man möchte, dass das XCS möglichst vielen verschiedenen Situationen ausgesetzt ist. Da in einem statischen Szenario Start- und Zielposition wie auch die Hindernisse fest sind, ist es wichtig, durch *random selection* dem XCS einen gewissen Spielraum zu geben.

Bei einem dynamischen Szenario ergibt sich dieses Problem nicht, andere Agenten und das Zielobjekt sind in stetiger Bewegung, der eigene Startpunkt ist nicht fixiert und das Problem wird bei Erreichen des Ziels nicht neugestartet. Aufgrund der Natur der Aufgabenstellung ist es in einem Überwachungsszenario außerdem wichtig, dass das XCS über eine längere Zeit hinweg eine gute Leistung liefert, also stetig gute Entscheidungen trifft,

eine zufällige Auswahl scheint also wenig hilfreich zu sein.

Außerdem wird in XCS zwischen den verschiedenen Auswahlarten hin und her geschaltet. Die Auswahlarten werden in zwei Gruppen geteilt, in die sogenannte *explore* Phase und in die *exploit* Phase. In der *exploit* Phase soll bevorzugt eine Auswahlart ausgeführt werden, die das Produkt aus den Werten *fitness* und *reward prediction* möglichst stark gewichten, beispielsweise wäre *best selection* ein Kandidat für die *exploit* Phase, während z.B. *random selection* ein Kandidat für die *explore* Phase wäre.

Dies bestätigen später auch Tests TODO Tests

8.9.1 Auswahlart *tournament selection*

Zu den oben erwähnten drei Möglichkeiten wurde in [MVBG03] eine weitere vorgestellt und in Bezug auf XCS diskutiert, die sogenannte *tournament selection*. Als Vorteile werden geringerer Selektionsdruck, höhere Effizienz, geringerer Einfluss von Störungen, wie auch Flexibilität der Anpassung über zwei Parameter, k und p , genannt. Dabei werden k *classifier* aus dem *matchSet* zufällig ausgewählt, nach ihrem Produkt aus den jeweiligen *fitness* und *reward prediction* Werten sortiert und absteigend mit Wahrscheinlichkeit p der jeweilige *classifier* ausgewählt (d.h. der erste mit p , der zweite mit $(1.0 - p) * p$, der dritte mit $(1.0 - p)^2 * p$ usw.). Mit $p = 1.0$ und $k = n$ (wobei n der Größe des *matchSets* entspricht) wäre *tournament selection* identisch mit *best selection* und mit $k = 1$ wäre es identisch mit *random selection*.

Bei der Entscheidung, welche Auswahlart jeweils für die *explore* und welche für die *exploit* Phase benutzt werden soll, ergeben sich also zwei Möglichkeiten:

Wichtig: randomize bei gleicher fitness etc.

1. *roulette wheel selection* : Zufällige Auswahl eines *classifier*, mit Wahrscheinlichkeit abhängig vom Produkt seines *fitness* und *reward prediction* Werts

2. *tournament selection* : Zufällige Wahl von k *classifiers* und daraus Wahl des jeweils besten *classifiers* mit Wahrscheinlichkeit p , Wahl des zweitbesten mit Wahrscheinlichkeit $(1.0 - p) * p$ usw. TODO

8.9.2 Wechsel zwischen den *explore* und *exploit* Phasen

In der Standardimplementierung von XCS wird zwischen jedem Problem zwischen der *explore* und der *exploit* Phase hin und hergeschaltet. Idee ist, dass man mit Hilfe der *explore* Phasen den Suchraum besser erforschen kann, dann aber zur eigentlichen Problemlösung in der *exploit* Phase möglichst direkt auf das Ziel zugeht.

Bei der Standardimplementierung für den statischen Fall ist allerdings das Erreichen eines positiven *base rewards* äquivalent mit einem Neustart des Problems. Während in der Standardimplementierung beim Neustart des Problems das gesamte Szenario (alle Agenten, Hindernisse und das Zielobjekt) auf den Startzustand zurückgesetzt wird, läuft das Überwachungsszenario weiter. Als erweiterten Ansatz soll nun deshalb eine neue Problemdefinition gelten, dass nicht das Erreichen eines positiven *base rewards* einen Phasenwechsel auslöst, sondern eine Änderung des *base rewards*, so dass mit anfänglicher *explore* Phase immer dann in die *exploit* Phase gewechselt wird, wenn das Zielobjekt in Sicht ist (bzw. umgekehrt, wenn mit der *exploit* Phase begonnen wird). Als Vergleich soll der andauernde, zufällige Wechsel zwischen der *explore* und *exploit* Phase, eine andauernde *exploit* und andauernde *explore* Phase dienen. Es sollen nun also folgende Arten des Wechsel zwischen den Phasen untersucht werden:

1. Andauernde *explore* Phase
2. Andauernde *exploit* Phase
3. Abwechselnd *explore* und *exploit* Phase (bei positivem *base reward*)

4. Abwechselnd *explore* und *exploit* Phase (bei Änderung des *base reward*, beginnend mit *explore*)
5. Abwechselnd *explore* und *exploit* Phase (bei Änderung des *base reward*, beginnend mit *exploit*)
6. In jedem Schritt zufällig entweder *explore* oder *exploit* Phase (50% Wahrscheinlichkeit jeweils)

1. Vergleich Random Explore, Roulette Wheel 2. Always Explore und Switch(exploit) schnell ausschliessen

4 verschiedene Wechsel, 2 verschiedene explore/exploit Dinger, mehrere Parametereinstellungen (p), k auf Maximum

=> 16-32 Tests

No exploration => viele ungültige Bewegungen, nicht "wegkommen" von Hindernis / stehenbleiben?

TODO SEHR WICHTIG BEI SICH WENIG BEWEGENDEN ZIELEN

Die Wahl der Auswahlart für *classifier* in Punkt (3) (in Kapitel 7.2) kann auf verschiedene Weise erfolgen. In der Standardimplementierung von XCS wird zwischen "exploit" und "explore" nach jedem Erreichen des Ziels entweder umgeschaltet oder zufällig mit einer bestimmten Wahrscheinlichkeit eine Auswahlart ermittelt. Es werden also abwechselnd ganze Probleme im "exploit" und "explore" Modus berechnet. Dies erscheint sinnvoll für die erwähnten Standardprobleme, da nach Erreichen des Ziels ein neues Problem gestartet wird und die Entscheidungen die während der Lösung eines Problems getroffen werden keine Auswirkungen auf die folgenden Probleme hat, die Probleme also nicht miteinander zusammenhängen.

Bei dem hier vorgestellten Überwachungsszenario kann nicht neugestartet werden, es gibt keine "Trockenübung", die Qualität eines Algorithmus soll deshalb davon abhängen, wie

gut sich der Algorithmus während der gesamten Berechnung, inklusive der Lernphasen, verhält. Es ist nicht möglich bei diesem Szenario zwischen *exploit* und *explore* Phasen zu differenzieren, wie dies in den Standardszenarien bei XCS der Fall ist, bei denen die Qualität nur während der *exploit* Phase gemessen wird.

Desweiteren greift auch die Idee einer reinen *explore* Phase beim Überwachungsszenario nicht, da das Szenario nicht statisch, sondern dynamisch ist. Ein zufälliges Herumlaufen kann, im Vergleich zur gewichteten Auswahl der Aktionen, dazu führen, dass der Agent mit bestimmten Situationen mit deutlich niedrigerer Wahrscheinlichkeit konfrontiert wird, da der Agent sich in Hindernissen verfängt oder das Zielobjekt ihm andauernd ausweicht. Aus diesen Gründen erscheint es sinnvoll, weitere Formen des Wechsels zwischen diesen Phasen zu untersuchen:

Möglichkeit (3.) und (4.) entspricht dem Fall in der Standardimplementierung von XCS. Dabei wird bei jedem Erreichen eines positiven Rewards zwischen “explore” und “exploit” hin und hergeschaltet, was in der Standardimplementierung dem Beginn eines neuen Problems entspricht.

TODO Umschalten bei reward, Code evtl.

TODOTESTS

TODO SWITCH EXPLORE/EXPLOIT + NEW LCS sehr gut

Kapitel 9

XCS Varianten

TODO!!

Ziel der Arbeit war es, wie man den XCS Algorithmus auf ein Überwachungsszenario anwenden kann. Notwendig dafür war es, die XCS Implementierung vollständig nachzuvollziehen, um für jeden Bestandteil entscheiden zu können, welche Rolle es bezüglich eines solchen Szenarios spielt. Für die Tests wurde nicht auf bestehende Pakete (z.B. XCSlib [Lan]) zurückgegriffen, wenn auch der Quelltext von [But00] Modell stand.

Bild mit rückwirkender Rewardvergabe

Im Vordergrund stand zum einen die grundsätzliche Frage, ob XCS in einem solchen Szenario überhaupt besser als ein Algorithmus sein kann, der sich rein zufällig verhält und wie mögliche Ansätze aussehen können, den Algorithmus zu verbessern.

Der hier entwickelte Algorithmus muss primär nicht einen Weg zum Ziel erkennen, sondern eine möglichst optimale (und auch an andere Agenten angepasste) Verhaltensstrategie finden.

In Kapitel 8 wurden mögliche Optimierungen zu den Parametern vorgestellt, in Kapitel 7.2 wurde diskutiert, in welcher Reihenfolge bei einem Multiagentensystem auf einem diskreten Torus die einzelnen Teile ausgeführt werden sollen.

Besonders die Verwaltung der Numerosity und die Verwendung des maxPrediction bereitete

Das Multistepverfahren baut darauf auf, dass die Qualität der Agenten sich sukzessive mit jeder Problemistanz verbessert, der Reward eben an immer weiter vom Ziel entfernte Aktionen TODO weitergereicht wird.

Da sich das Ziel schneller bewegt, kann eine einfache Verfolgungsstrategie nicht zum Erfolg führen. Eine einfache Implementation mit einem simplen Agenten der auf das Ziel zugeht, wenn es in Sicht ist und sich sonst wie ein sich zufällig bewegendes Agent verhält, schneidet grundsätzlich schlechter ab.

TODO!

9.1 Allgemeine Anpassungen und Verbesserungen

9.1.1 Verschiedenes, Numerosity, TODO

Durch die Benutzung von *macro classifiers* ergibt sich allerdings das programmiertechnische Problem, dass man nicht mehr direkt weiß, wieviele *micro classifiers* sich in einer Population befinden, bei jeder Benutzung des Werts der Populationsgröße müssten die *numerosity* Werte aller *classifiers* jedes Mal addiert werden. In der Standardimplementierung [But00] ist die Behandlung des *numerosity* Werts deswegen stark optimiert, jedes *classifier set* trägt eine temporäre Variable *numerositySum* mit sich, in der die aktuelle Summe gespeichert ist. Die Aktualisierung ist jedoch zum einen mangelhaft umgesetzt, zum anderen auf die Verwendung von einer einzelnen *action set* Liste optimiert, während die hier verwendete Implementierung jeweils mit bis über 100 *action set* Listen programmiert wurde, denen ein *classifier* Mitglied sein kann. Deswegen wurde die Optimierung entfernt und durch eine dezentrale Verwaltung mit einem *Observer* ersetzt, jede Änderung des *numerosity* Wertes hat also die Änderung aller *action set* Listen zur Folge, in der der

classifier Mitglied ist.

Wird also z.B. ein *micro classifier* entfernt, dann wird lediglich die Änderungsfunktion des *classifiers* aufgerufen, der dann wiederum den *numerositySum* Wert der jeweiligen Eltern anpasst. Dies macht einige Optimierungen rückgängig, erspart aber sehr viel Umstände, den *numerositySum* der Eltern immer auf den aktuellen Stand zu halten und einzelne *classifiers* zu löschen.

Positiver Nebeneffekt durch die verbesserte Struktur ist, dass man dadurch leicht auf die Menge der *action set* Listen zugreifen kann, denen ein *classifier* angehört, hierfür wurde aber im Rahmen dieser Arbeit keine Verwendung gefunden.

Ein weiteres Problem der Standardimplementierung ist, dass der *fitness* Wert eines *classifiers* als Optimierung bereits den *numerosity* Wert als Faktor enthält, während bei der Aktualisierung des *numerosity* Werts der *fitness* Wert nicht aktualisiert wurde. Das hat zur Folge, dass theoretisch *fitness* Werte von *classifiers* fast den *max population* Wert annehmen kann, wenn ein *classifier* mit *numerosity* und *fitness* Wert in der Höhe von *max population* auf einen *numerosity* Wert von 1 reduziert wird.

Dies betrifft die Funktion `public void addNumerosity(int num)` der Klasse *XClassifier* in der Datei *XClassifier.java*. Die korrigierte Fassung ist in Programm 9.1 gelistet, ein Vergleich der Qualität, mit und ohne Korrektur, ist in Abbildung ?? dargestellt.

TODO Vergleich! TODO wenig Unterschied sensoragent, evtl wieder raus

9.2 Standard XCS Multistepverfahren

Idee dieses Verfahrens ist, dass der *reward* Wert, den eine Aktion (bzw. das jeweils zugehörige *actionSet* und die dortigen *classifier*) erhält, vom erwarteten *reward* Wert der

```
1  /**
2   * Adds to the numerosity of the classifier.
3   * @param num The added numerosity (can be negative!).
4   */
5   public void addNumerosity(int num) {
6       int old_num = numerosity;
7
8       numerosity += num;
9
10      /**
11       * Korrektur der fitness
12       */
13       fitness = fitness * (double)numerosity / (double)old_num;
14
15      /**
16       * Aktualisierung der Eltern
17       */
18       for (ClassifierSet p : parents) {
19           p.changeNumerositySum(num);
20           if (numerosity == 0) {
21               p.removeClassifier(this);
22           }
23       }
24   }
```

Programm 9.1: Korrigierte Version der *addNumerosity()* Funktion

folgenden Aktion abhängt. Somit wird, rückführend vom letzten Schritt auf das Ziel, der *reward* Wert schrittweise an vorgehende Aktionen verteilt, mit der Annahme, dass dann, durch mehrfache Wiederholung des Lernprozesses, mit dem sich dadurch ergebenen Regelsatz mit höherer Wahrscheinlichkeit das Ziel gefunden wird.

Kern des Verfahrens ist die Vergabe des *base rewards*. Wird das Ziel erreicht, d.h. erhält der Algorithmus einen positiven *base reward* Wert, so wird der *reward* 1.0 an das letzte *actionSet* gegeben. Liegt kein positiver *base reward* Wert vor, so wird lediglich der für diesen Schritt erwartete *reward* Wert (nämlich der *maxPrediction* Wert) an das letzte *actionSet* gegeben.

Als Vergleich wurde das bekannte Verfahren [BW01] fast unverändert übernommen. Der wesentliche Unterschied ist, dass das Szenario bei einem positiven *base reward* nicht neugestartet wird, algorithmisch ist die Implementierung ansonsten identisch. Außerdem, wie schon in Kapitel 8.9.2 erwähnt, soll die Qualität des Algorithmus nicht nur in der *exploit* Phase gemessen werden, da ein fortlaufendes Problem und kein statisches Szenario betrachtet wird. Schließlich gibt es, neben den Parametereinstellungen im Kapitel 8, feste Schnittpunkte für das *two point crossover* beim genetischen Operator (siehe Kapitel 7.5).

TODO Programm 9.3

TODO Programm 9.4

```

1  /**
2   * Diese Funktion wird in jedem Schritt aufgerufen um den aktuellen
3   * Reward zu bestimmen, den besten Wert des ermittelten MatchSets
4   * weiterzugeben und, bei aktuell positivem Reward, das aktuelle
5   * ActionSet zu belohnen.
6   *
7   * @param gaTimestep Der aktuelle Zeitschritt
8   */
9
10 public void calculateReward(final long gaTimestep) {
11     /**
12      * checkRewardPoints liefert "wahr" wenn sich der Zielagent in
13      * Überwachungsreichweite befindet
14      */
15     boolean reward = checkRewardPoints();
16
17     if(prevActionSet != null){
18         collectReward(lastReward, lastMatchSet.getBestValue(), false);
19         prevActionSet.evolutionaryAlgorithm(classifierSet, gaTimestep);
20     }
21
22     if(reward) {
23         collectReward(reward, 0.0, true);
24         lastActionSet.evolutionaryAlgorithm(classifierSet, gaTimestep);
25         prevActionSet = null;
26         return;
27     }
28     prevActionSet = lastActionSet;
29     lastReward = reward;
30 }

```

Programm 9.2: Erstes Kernstück des Standard XCS Multistepverfahrens (*calculateReward()*, Bestimmung des Rewards anhand der Sensordaten), angepasst an ein dynamisches Überwachungsszenario

```

1  /**
2   * Diese Funktion verarbeitet den übergebenen Reward und gibt ihn an die
3   * zugehörigen ActionSets weiter.
4   *
5   * @param reward Wahr wenn der Zielagent in Sicht war.
6   * @param best_value Bester Wert des vorangegangenen ActionSets
7   * @param is_event Wahr wenn diese Funktion wegen eines Ereignisses, d.h.
8   *                einem positiven Reward, aufgerufen wurde
9   */
10
11  public void collectReward(boolean reward,
12                           double best_value, boolean is_event) {
13      double corrected_reward = reward ? 1.0 : 0.0;
14
15      /**
16       * Falls der Reward von einem Ereignis rührt, aktualisiere das
17       * aktuelle ActionSet und lösche das vorherige
18       */
19      if(is_event) {
20          if(lastActionSet != null) {
21              lastActionSet.updateReward(corrected_reward, best_value, factor);
22              prevActionSet = null;
23          }
24      }
25
26      /**
27       * Kein Ereignis, also nur das letzte ActionSet aktualisieren
28       */
29      else
30      {
31          if(prevActionSet != null) {
32              prevActionSet.updateReward(corrected_reward, best_value, factor);
33          }
34      }
35  }

```

Programm 9.3: Zweites Kernstück des Multistepverfahrens (*collectReward()* - Verteilung des Rewards auf die ActionSets), angepasst an ein dynamisches Überwachungsszenario

```

1  /**
2  * Bestimmt die zum letzten bekannten Status passenden Classifier und
3  * wählt aus dieser Menge eine Aktion. Außerdem wird das aktuelle
4  * ActionClassifierSet mithilfe der gewählten Aktion ermittelt.
5  *
6  * @param gaTimestep Der aktuelle Zeitschritt
7  */
8
9  public void calculateNextMove(long gaTimestep) {
10
11  /**
12  * Überdecke das classifierSet mit zum Status passenden Classifiern
13  * welche insgesamt alle möglichen Aktionen abdecken.
14  */
15      classifierSet.coverAllValidActions(
16          lastState, getPosition(), gaTimestep);
17
18  /**
19  * Bestimme alle zum Status passenden Classifier.
20  */
21      lastMatchSet = new AppliedClassifierSet(lastState, classifierSet);
22
23  /**
24  * Entscheide auf welche Weise die Aktion ausgewählt werden soll.
25  */
26      lastExplore = checkIfExplore(lastState.getSensorGoalAgent(),
27          lastExplore, gaTimestep);
28
29  /**
30  * Wähle Aktion und bestimme zugehöriges ActionSet
31  */
32      calculatedAction = lastMatchSet.chooseAbsoluteDirection(lastExplore);
33      lastActionSet = new ActionClassifierSet(lastState, lastMatchSet,
34          calculatedAction);
35  }

```

Programm 9.4: Drittes Kernstück des Multistepverfahrens (*calculateNextMove()* - Auswahl der nächsten Aktion und Ermittlung des zugehörigen ActionSets), angepasst an ein dynamisches Überwachungsszenario

9.3 XCS Variante für Überwachungsszenarien (SXCS)

Die Hypothese bei der Aufstellung dieser Variante des XCS-Algorithmus ist im Grunde dieselbe wie beim XCS-Multistepverfahren, nämlich dass die Kombination mehrerer Aktionen zum Ziel führt. Beim Multistepverfahren besteht die wesentliche Verbindung zwischen den *actionSet* Listen jeweils nur zwischen zwei direkt aufeinanderfolgenden *actionSets* über den *maxPrediction* Wert. In einer statischen Umgebung kann dadurch über mehrere (identische) Probleme hinweg eine optimale Einstellung (der *fitness* und der *reward prediction* Wert) für die *classifier* gefunden werden.

Bei der veränderten XCS Variante SXCS soll die Verbindung zwischen den *actionSets* zusätzlich direkt durch die zeitliche Nähe zum Ziel gegeben sein. Es wird in jedem Schritt das jeweilige *actionSet* gespeichert und aufgehoben, bis ein neues Ereignis (siehe Kapitel 9.3.1) eintritt und dann in Abhängigkeit des Alters mit einem entsprechenden *reward* Wert aktualisiert.

$r(a)$ bezeichnet den *reward* Wert für das *actionSet* mit Alter a .

Bei linearer Vergabe des *reward*:

$$r(a) = \begin{cases} \frac{a}{\text{size}(\text{ActionSet})} & , \text{ falls reward} = 1 \\ \frac{1-a}{\text{size}(\text{ActionSet})} & , \text{ falls reward} = 0 \end{cases}$$

bzw. bei quadratischer Vergabe des *reward*:

$$r(a) = \begin{cases} \frac{a^2}{\text{size}(\text{ActionSet})} & \text{ falls reward} = 1 \\ \frac{1-a^2}{\text{size}(\text{ActionSet})} & \text{ falls reward} = 0 \end{cases}$$

In Tests ergab sich für die quadratische Vergabe des *reward* ein minimal besseres Ergebnis (TODO zeigen), weitere Grafiken werden auf die lineare Vergabe des *reward*

beschränkt sein um eine verständliche Darstellung zu ermöglichen, während in den Simulationen die quadratische Vergabe des *reward* benutzt wird.

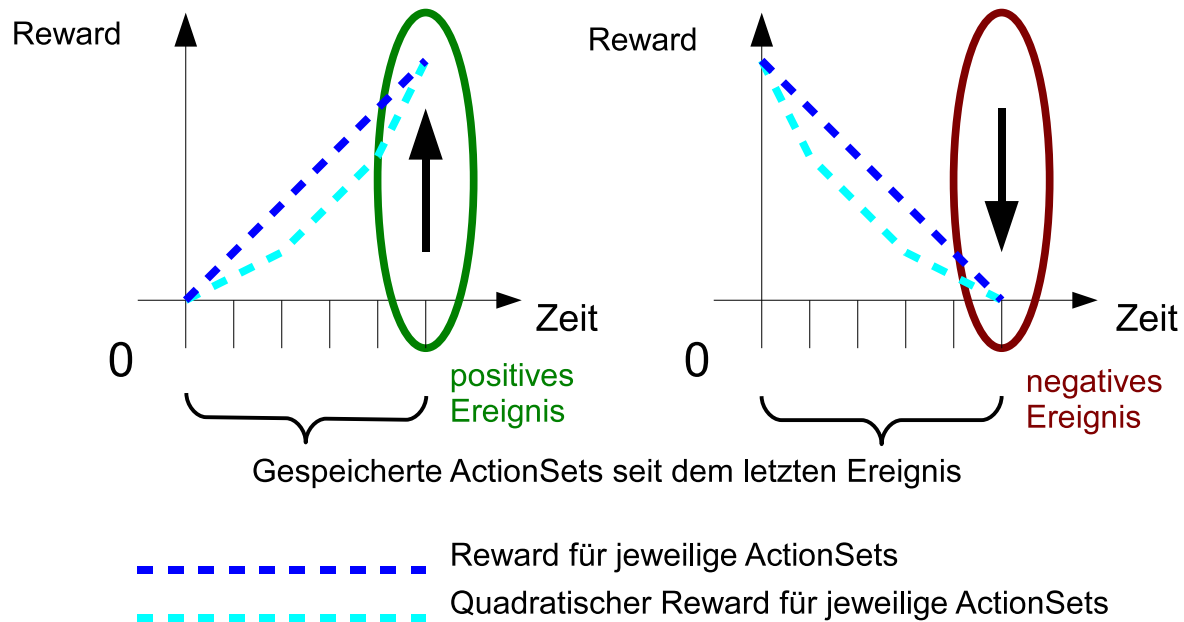


Abbildung 9.1: Schematische Darstellung der (quadratischen) Rewardverteilung an gespeicherte ActionSets bei einem positiven bzw. negativen Ereignis

9.3.1 Ereignisse

In XCS wird lediglich das jeweils letzte ActionSet aus dem vorherigen Zeitschritt gespeichert, in der neuen Implementierung werden dagegen eine ganze Anzahl (bis zu “maxStackSize”) von ActionSets gespeichert. Die Speicherung erlaubt zum einen eine Vorverarbeitung des Rewards anhand der vergangenen Zeitschritte und auf Basis einer größeren Zahl von ActionSets und zum anderen die zeitliche Relativierung eines ActionSets zu einem Ereignis. Die Classifier wird dann jeweils rückwirkend anhand des Rewards aktualisiert sobald bestimmte Bedingungen eingetreten sind.

Von einem positiven bzw. negativen Ereignis spricht man, wenn sich der Reward im Vergleich zum vorangegangenen Zeitschritt verändert hat, also wenn der Zielagent sich in

Übertragungsreichweite bzw. aus ihr heraus bewegt hat (siehe Abbildung 9.2).

Bei der Benutzung eines solchen Stacks entsteht eine Zeitverzögerung, d.h. die Classifier besitzen jeweils Information die bis zu “maxStackSize” Schritte zu alt sind. Wählen wir den Stack zu groß, nimmt die Konvergenzgeschwindigkeit und Reaktionsfähigkeit des Systems zu stark ab, wählen wir ihn zu klein, kann es sein, dass wir einen Überlauf bekommen, also “maxStackSize” Schritte lang keine Rewardänderung aufgetreten ist. Im letzteren Fall brechen wir deswegen ab, bewerten die ActionSets der ersten Hälfte des Stacks (also die $\frac{maxStackSize}{2}$ ältesten Einträge) mit dem damals vergebenem konstanten Reward (welcher dem aktuellen Reward entspricht, es ist ja keine Rewardänderung eingetreten) und nehmen sie vom Stack (siehe Abbildung 9.3). Anschließend wird normal weiter verfahren bis der Stack wieder voll ist bzw. bis eine Rewardänderung auftritt. Das Szenario mit dem maximalen Fehler wäre das, bei dem ein Schritt nach dem Abbruch eine Rewardänderung auftritt. Der Wert *maxStackSize* stellt also einen Kompromiss zwischen Zeitverzögerung bzw. Reaktionsgeschwindigkeit und Genauigkeit dar.

Ein Ereignis tritt auf, wenn:

1. Positive Rewardänderung (Zielagent war im letzten Zeitschritt nicht in Überwachungsreichweite) \Rightarrow positives Ereignis (mit reward = 1)
2. Negative Rewardänderung (Zielagent war im letzten Zeitschritt in Überwachungsreichweite) \Rightarrow negatives Ereignis (mit reward = 0)
3. Überlauf des Stacks (keine Rewardänderung in den letzten “maxStackSize” Schritten), Zielagent ist in Überwachungsreichweite \Rightarrow neutrales Ereignis (mit reward = 1)
4. Überlauf des Stacks (keine Rewardänderung in den letzten “maxStackSize” Schritten), Zielagent ist nicht in Überwachungsreichweite \Rightarrow neutrales Ereignis (mit reward = 0)

9.3.2 Implementierung von SXCS

TODO Erläuterung

TODO Programm 9.5

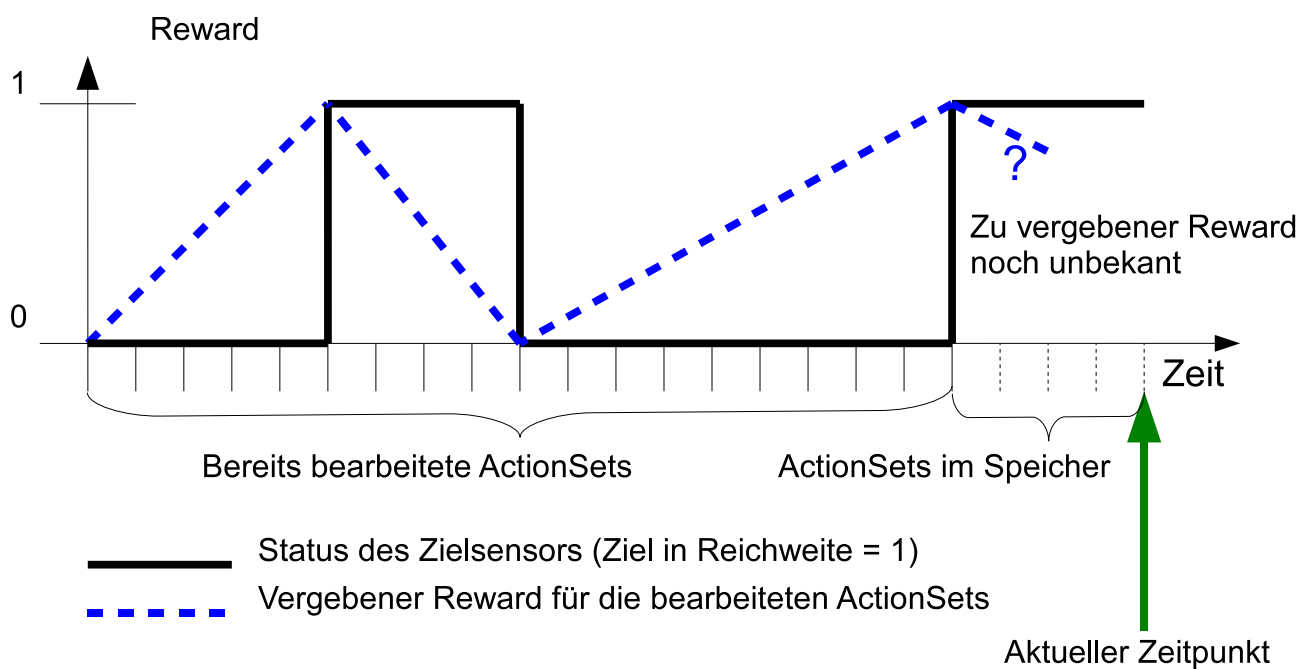


Abbildung 9.2: Schematische Darstellung der zeitlichen Rewardverteilung an ActionSets nach mehreren positiven und negativen Ereignissen und der Speicherung der letzten ActionSets

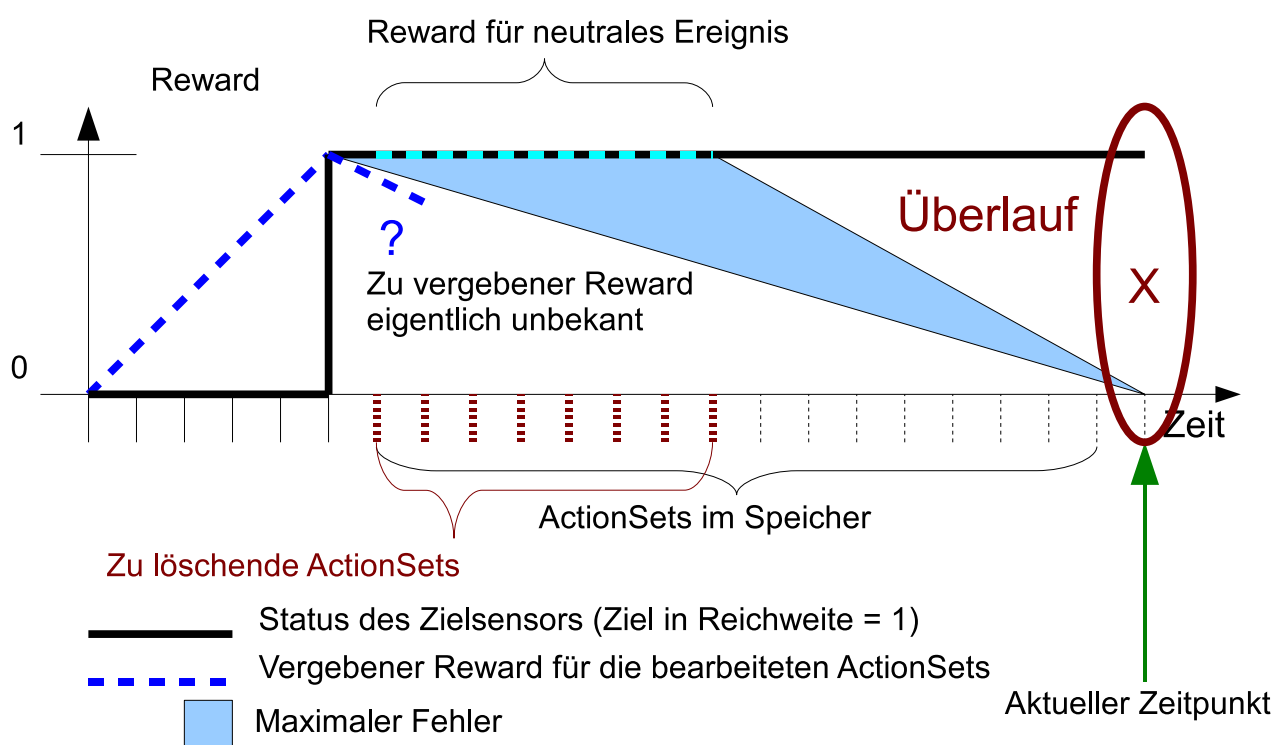


Abbildung 9.3: Schematische Darstellung der Rewardverteilung an ActionSets bei einem neutralen Ereignis

```

1  /**
2   * Diese Funktion wird in jedem Schritt aufgerufen um den aktuellen
3   * Reward zu bestimmen und positive, negative und neutrale Ereignisse
4   * den besten Wert des ermittelten MatchSets weiterzugeben und, bei
5   * aktuell positivem Reward, das aktuelle ActionSet zu belohnen.
6   *
7   * @param gaTimestep Der aktuelle Zeitschritt
8   */
9
10 public void calculateReward(final long gaTimestep) {
11     /**
12      * checkRewardPoints liefert "wahr" wenn sich der Zielagent in
13      * Überwachungsreichweite befindet
14      */
15     boolean reward = checkRewardPoints();
16
17     if (reward != lastReward) {
18         int start_index = historicActionSet.size() - 1;
19         collectReward(start_index, actionSetSize, reward, 1.0, true);
20         actionSetSize = 0;
21     }
22     else
23
24     if (actionSetSize >= Configuration.getMaxStackSize())
25     {
26         int start_index = Configuration.getMaxStackSize() / 2;
27         int length = actionSetSize - start_index;
28         collectReward(start_index, length, reward, 1.0, false);
29         actionSetSize = start_index;
30     }
31
32     lastReward = reward;
33 }

```

Programm 9.5: Erstes Kernstück des SXCS-Algorithmus (*calculateReward()*, Bestimmung des Rewards anhand der Sensordaten)

TODO Programm 9.6

TODO Programm 9.7

9.3.3 Zielobjekt mit SXCS

Wie bereits in Kapitel 4.2.6 erwähnt, soll hier eine Implementierung von SXCS für das Zielobjekt diskutiert werden. Bis auf die Funktion *checkRewardPoints()* (siehe Kapitel 7.6) ist die Implementierung für das Zielobjekt identisch. Die abgeänderte Version ist in Programm 9.8 aufgelistet.

```

1  /**
2   * Diese Funktion verarbeitet den übergebenen Reward und gibt ihn an die
3   * zugehörigen ActionSets weiter.
4   *
5   * @param reward Wahr wenn der Zielagent in Sicht war.
6   * @param best_value Bester Wert des vorangegangenen ActionSets
7   * @param is_event Wahr wenn diese Funktion wegen eines Ereignisses, d.h.
8   *                einem positiven Reward, aufgerufen wurde
9   */
10
11 public void collectReward(
12     boolean reward, double best_value, boolean is_event) {
13     double corrected_reward = reward ? 1.0 : 0.0;
14     /**
15      * Wenn es kein Event ist, dann gebe den Reward weiter wie beim
16      * Multistepverfahren
17      */
18     double max_prediction = is_event ? 0.0 :
19         historicActionSet.get(start_index+1).getMatchSet().getBestValue();
20
21     /**
22      * Aktualisiere eine ganze Anzahl von ActionSets
23      */
24     for(int i = 0; i < action_set_size; i++) {
25
26         /**
27          * Benutze aufsteigenden bzw. absteigenden Reward bei einem positiven
28          * bzw. negativen Ereignis
29          */
30         if(is_event) {
31             corrected_reward = reward ?
32                 calculateReward(i, action_set_size) :
33                 calculateReward(action_set_size - i, action_set_size);
34         }
35         /**
36          * Aktualisiere das ActionSet mit dem bestimmten Reward und
37          * gebe bei allen anderen ActionSets den Reward weiter wie
38          * beim Multistepverfahren
39          */
40         ActionClassifierSet action_classifier_set =
41             historicActionSet.get(start_index - i);
42         action_classifier_set.updateReward(
43             corrected_reward, max_prediction, factor);
44
45         max_prediction =
46             action_classifier_set.getMatchSet().getBestValue();
47     }
48 }

```

Programm 9.6: Zweites Kernstück des SXCS-Algorithmus (*collectReward()* - Verteilung des Rewards auf die ActionSets)


```

1  /**
2   * Bestimmt die zum letzten bekannten Status passenden Classifier und
3   * wählt aus dieser Menge eine Aktion. Außerdem wird das aktuelle
4   * ActionClassifierSet mithilfe der gewählten Aktion ermittelt.
5   * Im Vergleich zur originalen Multistepversion wird am Schluß noch
6   * das ermittelte ActionSet gespeichert.
7   *
8   * @param gaTimestep Der aktuelle Zeitschritt
9   */
10
11  public void calculateNextMove(long gaTimestep) {
12
13  /**
14   * Überdecke das classifierSet mit zum Status passenden Classifiern
15   * welche insgesamt alle möglichen Aktionen abdecken.
16   */
17      classifierSet.coverAllValidActions(
18          lastState, getPosition(), gaTimestep);
19
20  /**
21   * Bestimme alle zum Status passenden Classifier.
22   */
23      lastMatchSet = new AppliedClassifierSet(lastState, classifierSet);
24
25  /**
26   * Entscheide auf welche Weise die Aktion ausgewählt werden soll,
27   * wähle Aktion und bestimme zugehöriges ActionSet
28   */
29      lastExplore = checkIfExplore(lastState.getSensorGoalAgent(),
30                                  lastExplore, gaTimestep);
31
32      calculatedAction = lastMatchSet.chooseAbsoluteDirection(lastExplore);
33      lastActionSet = new ActionClassifierSet(lastState, lastMatchSet,
34                                              calculatedAction);
35
36  /**
37   * Speichere das ActionSet und passe den Stack bei einem Überlauf an
38   */
39      actionSetSize++;
40      historicActionSet.addLast(lastActionSet);
41      if (historicActionSet.size() > Configuration.getMaxStackSize()) {
42          historicActionSet.removeFirst();
43      }
44  }

```

Programm 9.7: Drittes Kernstück des SXCS-Algorithmus (*calculateNextMove()* - Auswahl der nächsten Aktion und Ermittlung und Speicherung des zugehörigen ActionSets)

```
1  /**
2   * @return true Falls das Zielobjekt von keinem Agenten überwacht wird
3   */
4  @Override
5  public boolean checkRewardPoints() {
6      boolean[] sensor_agent = lastState.getSensorAgent();
7
8      for(int i = 0; i < Action.MAX_DIRECTIONS; i++) {
9          if(sensor_agent[2*i+1]) {
10             return false;
11         }
12     }
13
14     return true;
15 }
```

Programm 9.8: Bestimmung des *base rewards* für das Zielobjekt

Kapitel 10

Analyse SXCS

Tabelle 10.1: Vergleich “Intelligent (Open)” und “Intelligent (Hide)” (8 Agenten, Säulenszenario)

Algorithmus	Abdeckung	Qualität
“Intelligent (Open)”		
Zufällige Bewegung	72.55%	11.58%
XCS	71.35%	13.98%
SXCS	72.10%	13.50%
“Intelligent (Hide)”		
Zufällige Bewegung	72.56%	11.78%
XCS	71.33%	14.27%
SXCS	72.05%	13.90%

TODO

TODO auch sich langsam bewegende analysieren! Und auch stehenbleibende : z.B. im Raumszenario.

Geschwindigkeit 2 problematisch, Geschwindigkeit 1 ok?

TODO classifier ausgeben

Tabelle 10.2: Vergleich “Intelligent (Open)” und “Intelligent (Hide)” (8 Agenten, Säulenszenario)

Algorithmus	Abdeckung	Qualität
“Intelligent (Open)”		
Zufällige Bewegung	72.55%	11.58%
XCS	71.35%	13.98%
SXCS	72.10%	13.50%
“Intelligent (Hide)”		
Zufällige Bewegung	72.56%	11.78%
XCS	71.33%	14.27%
SXCS	72.05%	13.90%

10.1 Vergleich unterschiedlicher Geschwindigkeiten des Zielobjekts

In Abbildung 10.1 ist ein Vergleich der unterschiedlichen Geschwindigkeiten des Zielobjekts dargestellt. XCS (mit 500 Schritten) macht bei keiner Geschwindigkeit Lernfortschritte, die Qualität pendelt zwischen 31.69% und 33.40%, also in etwa identisch mit der zufälligen Bewegung. Die SXCS Implementierung scheint dagegen die geringere Geschwindigkeit ausgenutzt zu haben und ist dadurch in der Lage das Zielobjekt besser zu verfolgen. Mit 500 Schritten ist die Qualität abnehmend von 39.64% (Geschwindigkeit 1.0) bis 35.96% (Geschwindigkeit 2.0), im Fall mit 2000 Schritten erhöht sich dieser Bereich leicht auf 40.15% bis 37.71%.

Auch bei den Heuristiken zeichnet sich ein klares Bild ab, bei niedrigen Geschwindigkeiten ist die Ausbreitung der Agenten auf dem Feld (intelligente Heuristik) weniger wichtig als die konstante Verfolgung des Zielobjekts, während bei höheren Geschwindigkeiten die Verteilung auf dem Feld wichtiger wird.

Bester Agent nach 20000 Schritten (Zielgeschwindigkeit 2.0, SXCS, 2000 Schritte)

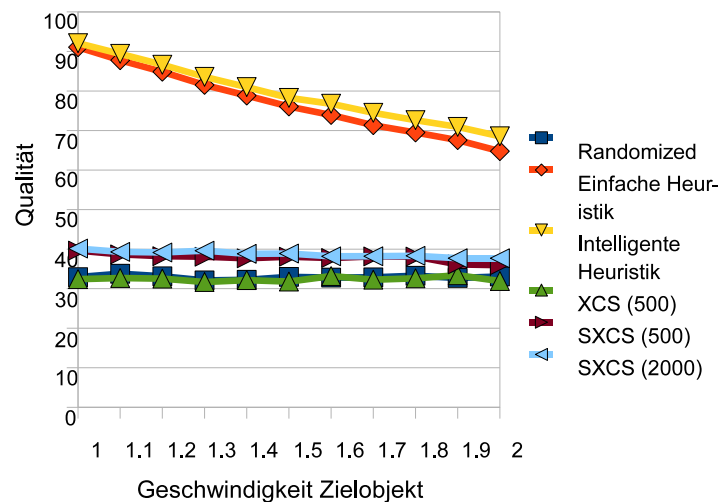


Abbildung 10.1: Vergleich der Qualitäten verschiedener Algorithmen bezüglich der Geschwindigkeit des Zielobjekts

#0#####.###0#0##.#0#0###0-S : [Fi: 0.38] [Ex: 00450.0] [Pr: 0.74] [PE: 0.38]

....

TODO

10.2 Zusammenfassung der bisherigen Erkenntnisse

Algorithmen mit Ergebnissen die unter dem des zufälligen Algorithmus liegt, sind unbrauchbar und nicht vergleichbar. “Verbesserungen”, die die Qualität des Algorithmus näher an das Ergebnis des zufälligen Algorithmus bringen, sind in Wirklichkeit Veränderungen, die den Algorithmus eher zufällige Entscheidungen treffen lassen, und keine tatsächlichen Lernerfolge.

SXCS sehr gut bei NO DIRECTION CHANGE und speed 1!

nicht geschafft: Pillar, one direction change, speed 2, XCS ...besser... weil zufälliger

10.3 Standard XCS Multistepverfahren

10.3.1 SXCS und Heuristiken

erst multistep... mit random vergleichen

In allen Tests erreichten die Heuristiken deutlich bessere Ergebnisse. Diesen Nachteil hat sich LCS in diesen Szenarien durch deutlich überlegene Flexibilität erkauft. Ein Großteil der eingehenden Informationen ist für die Auswertung nicht relevant und lokale Information ist zu ungenau. Bei einer komplexeren Implementierung mit Distanzen

Insbesondere der Vergleich mit dem intelligenten Agenten, der anderen Agenten ausweicht, zeigt, dass die LCS Agenten unmöglich ein solches globales Ziel erreichen können, es ist also kein emergentes Verhalten zu beobachten. Dies ist dadurch zu begründen, dass bei der Berechnung des Rewards keine Information außer der eigenen, lokalen Information der Abstand zu anderen Agenten nicht Teil der Berechnung des Rewards ist, noch gibt keine eingebaute Heuristik. Man könnte zwar

TODO statistical value:Error in predictions!

10.3.2 Vergleich Multistep / LCS

Szenarien, Parameter.

10.3.3 Test der verschiedenen Exploration-Modi

Prediction Error sehr hoch, da dynamisches

Kapitel 11

Kommunikation

Einführung, Kommunikationsbeschränkungen (nur Reward weitergeben)

Vergleich Agentenzahl (1, 2, 3, 4, 5, 6, 7, 8)

reward all equally besser als reward none Unterscheidung interner und externer reward

11.1 Realistischer Fall mit Kommunikationsrestriktionen

Bisher wurde der Fall betrachtet, dass Kommunikation mit beliebiger Reichweite stattfinden kann. Dies ist natürlich kein realistisches Szenario. Geht man jedoch davon aus, dass die Kommunikationsreichweite zumindest ausreichend groß ist um nahe Agenten zu kontaktieren, so kann man argumentieren, dass man dadurch ein Kommunikationsnetzwerk aufbauen kann, in dem jeder Agent jeden anderen Agenten - mit einer gewissen Zeitverzögerung - erreichen kann. Bei ausreichender Agentenzahl relativ zur freien Fläche fallen dadurch nur vereinzelte Agenten aus dem Netz, was der Effektivität der Agentengruppe erwartungsgemäß nur geringfügig schadet (TODO zeigen?) Stehen die Agenten nicht indirekt andauernd miteinander in Kontakt (mit anderen Agenten als Proxy), son-

dern muss die Information zum Teil durch aktive Bewegungen der Agenten transportiert werden, tritt eine Zeitverzögerung auf. Auch kann die benötigte Bandbreite die verfügbare übersteigen, was ebenfalls Zeit benötigt. Im realistischen Fall ist also davon auszugehen, dass jede Kommunikation erst mit einer gewissen Verzögerung ausgeführt wird, weshalb für Kommunikation nur der zuvor besprochene verzögerte LCS Algorithmus in Frage kommt.

pg. 286 Zentralisierung der Daten

TODO bei Faktorberechnung Ranking

11.2 Lösungen aus der Literatur

Da wir ein Multiagentensystem betrachten, stellt sich natürlich die Frage nach der Kommunikation. In der Literatur gibt es Multiagentensysteme die auf Learning Classifier Systemen aufbauen, wie z.B. TODO Literatur. Alle Ansätze in der Literatur erlauben jedoch globale Kommunikation, z.T. Gibt es globale Classifier auf die alle Agenten zurückgreifen können, z.T. gibt es globale Steuerung.

Verteilung des rewards an alle - soccer

TODO Einordnen In [KM94] gezeigt, Gruppenbildung (rationality, grade 2 confusion) soccer!

[THN⁺98] OCS, centralized control system

In dieser Arbeit betrachte ich das Szenario ohne globale Steuerung oder globale Classifier, also mit der Restriktion einer begrenzten, lokalen Kommunikation. Geht man davon aus, dass über die Zeit hinweg jeder Agent indirekt mit jedem anderen Agenten in Kontakt treten kann, Nachrichten also mit Zeitverzögerung weitergeleitet werden können, ist eine Form der globalen, wenn auch zeitverzögerten, Kommunikation möglich. TODO Eine spezielle Implementierung für diesen Fall werde ich weiter unten besprechen TODO

11.3 SXCS Variante mit verzögerter Reward (DSXCS)

Eine hilfreiche Voraussetzung für Kommunikation ist, wenn die dadurch möglicherweise entstehende Verzögerung vom jeweiligen Algorithmus unterstützt wird. Während weiter oben

Realistischer Fall

Drei Werte weitergeben... Egoismus Faktor, Reward und Timestamp

Der wesentliche Unterschied zur ersten XCS Variante SXCS ist, dass jeglicher ermittelter *reward* Wert und der jeweils zugehörige Faktor lediglich erst einmal zusammen mit den jeweiligen *actionSets* in einer Liste (*historicActionSet* TODO Bezeichnung) gespeichert werden und in jedem Schritt immer nur die *classifiers* des *actionSets* des ältesten Eintrags in der *historicActionSet* Liste aktualisiert wird. Somit haben wir also eine zeitlich beliebig verzögerbare Aktualisierungsfunktion, welche uns erlaubt, mehrere gleichzeitig stattgefundenen (aber erst verzögert eintreffende, wegen z.B. Kommunikationsschwierigkeiten) Ereignisse zusammen auszuwerten. Dies ist eine wesentliche Voraussetzung für Kommunikation zwischen den Agenten. TODO

Wann immer ein *base reward* Wert an einen Agenten verteilt wird, kann es sinnvoll sein, diesen *base reward* an andere Agenten weiterzugeben. Dies wurde z.B. in einem ähnlichen Szenario in [ITS05] festgestellt, bei dem zwei auf XCS basierende Agenten gegen bis zu zwei anderen (zufälligen) Agenten eine vereinfachte Form des Fußballs spielen. Das in dieser Arbeit besprochene Szenario ist wesentlich komplexer, was d

Die Funktion *calculateReward()* ist identisch mit der in Kapitel ?? besprochenen Funktion bei der SXCS Variante ohne verzögerten *reward*.

In der Funktion *processReward()* werden die gespeicherten *reward* und *factor* ausgewertet. In der Implementation in Programm ?? werden einfach alle nacheinander auf das

action set angewendet, während in der verbesserten Version in Programm ?? nur der *reward* Wert aus dem Paar mit dem größten Produkt aus den *reward* und *factor* Werten für die Aktualisierung benutzt wird. In beiden Implementationen werden außerdem Einträge mit sowohl einem *reward* als auch *factor* Wert von 1.0 ignoriert, sie wurden bereits in Programm 11.1 ausgewertet.

TODO Programm 11.2

TODO Programm 11.3

TODO Programm 11.4

11.4 Ablauf

TODO wann weitergabe des rewards

Jeder Reward, der aus einem normalen Ereignis generiert wird, wird unter Umständen an alle anderen Agenten weitergegeben. Wie ein solches sogenanntes “externes Ereignis” von diesen Agenten aufgefasst wird, hängt von der jeweiligen Kommunikationsvariante ab, die in (11.5) besprochen werden.

Durch eine gemeinsame Schnittstelle erhält jeder Agent den Reward zusammen mit dem Kommunikationsfaktor. Dabei ergibt sich das Problem, dass sich Rewards überschneiden können, da jeder Reward sich rückwirkend auf die vergangenen ActionClassifierSets auswirken kann. Auch können mehrere externe Rewards eintreffen als auch ein eigener lokaler Reward aufgetreten sein. Würden die Rewards nach ihrer Eingangsreihenfolge abgearbeitet werden, kann es passieren, dass das selbe ActionClassifierSet sowohl mit einem hohen als auch einem niedrigen Reward aktualisiert wird. Da das globale Ziel ist, den Zielagenten durch *irgendeinen* Agenten zu überwachen, ist es in jedem einzelnen Zeitschritt nur relevant, dass ein *einzelner* Agent einen hohen Reward produziert bzw. weitergibt um die eigene Aktion als zielführend zu bewerten.

Befindet sich das Ziel beispielsweise gerade in Überwachungsreichweite mehrerer Agen-

```

1  /**
2   * Diese Funktion verarbeitet den übergebenen Reward und gibt ihn an die
3   * zugehörigen ActionSets weiter. Wesentlicher Unterschied zum LCS ohne
4   * Verzögerung ist, dass maxPrediction erst bei der endgültigen
5   * Verarbeitung des historicActionSets ermittelt wird.
6   *
7   * @param reward Wahr wenn der Zielagent in Sicht war.
8   * @param best_value Bester Wert des vorangegangenen actionSets
9   * @param is_event Wahr wenn diese Funktion wegen eines Ereignisses, d.h.
10  *      einem positiven Reward, aufgerufen wurde
11  */
12
13  public void collectReward(
14      boolean reward, double best_value, boolean is_event) {
15      double corrected_reward = reward ? 1.0 : 0.0;
16
17      /**
18       * Aktualisiere eine ganze Anzahl von Einträgen im historicActionSet
19       */
20      for(int i = 0; i < action_set_size; i++) {
21
22          /**
23           * Benutze aufsteigenden bzw. absteigenden Reward bei einem positiven
24           * bzw. negativen Ereignis
25           */
26          if(is_event) {
27              corrected_reward = reward ?
28                  calculateReward(i, action_set_size) :
29                  calculateReward(action_set_size - i, action_set_size);
30          } else {
31              if(corrected_reward == 1.0 && factor == 1.0) {
32                  historicActionSet.get(start_index - i).
33                      rewardPrematurely(
34                          historicActionSet.get(start_index - i + 1).getBestValue());
35              }
36          }
37
38          /**
39           * Füge den ermittelten Reward zum historicActionSet
40           */
41          historicActionSet.get(start_index - i).
42              addReward(corrected_reward, factor);
43
44      }
45  }

```

Programm 11.1: Zweites Kernstück des verzögerten SXCS-Algorithmus (*collectReward()* - Verteilung des Rewards auf die ActionSets)

```

1  /**
2  *
3  * Der erste Teil der Funktion ist identisch mit dem calculateNextMove
4  * der SXCS Variante ohne Kommunikation. Der Zusatz ist, dass beim
5  * Überlauf die im HistoricActionSet gespeicherte Rewards verarbeitet
6  * werden
7  */
8
9  public void calculateNextMove(long gaTimestep) {
10
11    // ...
12
13    /**
14     * HistoryActionSet voll? Dann verarbeite den dort gespeicherten Reward
15     */
16     if (historicActionSet.size() > Configuration.getMaxStackSize()) {
17         HistoryActionClassifierSet first = historicActionSet.pop();
18         last.processReward(historicActionSet.getFirst().getBestValue());
19     }
20 }

```

Programm 11.2: Auszug aus dem dritten Kernstück des verzögerten SXCS-Algorithmus (*calculateNextMove()*)

ten und verliert ein anderer Agent das Ziel aus der Sicht, sollte der Agent (und alle anderen Agenten), der das Ziel in Sicht hat, deswegen nicht bestraft werden, da das globale Ziel ja weiterhin erfüllt wurde.

TODO überlegen ob das noch Sinn macht, inwieweit das erklärt werden musws

Gebe keinen Reward an andere Agenten weiter. Es ist nicht relevant, ob ein Agent das Ziel aus den Augen verliert oder nicht, es ist nur relevant, ob der Zielagent weiterhin von anderen Agenten beobachtet wird. Ein Sonderfall ist, wenn im vorherigen Schritt der Zielagent nicht in Sichtweite eines anderen Agenten stand, also in diesem Schritt auf einmal mehrere Agenten den Zielagenten sehen können. In diesem Fall gibt nur der erste Agent den Reward weiter und setzt ein Flag.

Ziel verschwindet aus Sicht War der Zielagent von keinem anderen Agenten in Sicht, dann hat sich der Zielagent hiermit aus der Sichtweite aller Agenten bewegt. Somit haben alle Agenten versagt und der negative Reward wird weitergegeben.

```
1  /**
2   * Zentrale Routine des HistoryActionSets zur Verarbeitung aller
3   * eingegangenen Rewards bis zu diesem Punkt.
4   */
5
6   public void processReward(double max_prediction) {
7
8   /**
9   * Finde das größte reward / factor Paar TODO Verbessern
10  */
11   for(RewardHelper r : reward) {
12   /**
13   * Dieser Eintrag wurde schon in collectReward() verwertet
14   */
15   if(r.reward == 1.0 && r.factor == 1.0) {
16   continue;
17   }
18   /**
19   * Aktualisiere den Eintrag mit den entsprechenden Werten und dem
20   * übergebenen maxPrediction Wert
21   */
22   actionClassifierSet.updateReward(r.reward, max_prediction, r.factor);
23   }
24   }
```

Programm 11.3: Viertes Kernstück des verzögerten SXCS-Algorithmus (DSXCS, Verarbeitung des Rewards, *processReward()*)

```

1  /**
2   * Zentrale Routine des HistoryActionSets zur Verarbeitung aller
3   * eingegangenen Rewards bis zu diesem Punkt.
4   */
5
6   public void processReward(double max_prediction) {
7
8       double max_value = 0.0;
9       double max_reward = 0.0;
10
11      /**
12       * Finde das größte reward / factor Paar TODO Verbessern
13       */
14      for(RewardHelper r : reward) {
15          /**
16           * Dieser Eintrag wurde schon in collectReward() verwertet
17           */
18          if(r.reward == 1.0 && r.factor == 1.0) {
19              return;
20          }
21
22          if(r.reward * r.factor > max_value) {
23              max_value = r.reward * r.factor;
24              max_reward = r.reward;
25          }
26      }
27      /**
28       * Aktualisiere den Eintrag mit dem ermittelten Wert und dem
29       * übergebenen maxPrediction Wert
30       */
31      actionClassifierSet.updateReward(max_reward, max_prediction, 1.0);
32  }

```

Programm 11.4: Verbesserte Variante des vierten Kernstück des verzögerten SXCS-Algorithmus (DXCS, Verarbeitung des Rewards, *processReward()*)

Selbiges wenn das Ziel in Sicht kommt und von keinem anderen Agenten in Sicht ist. Die Agenten waren offensichtlich erfolgreich und können belohnt werden.

TODOTODOTODOTODO Ist kein Event aufgetreten und leeren wir die Hälfte des Stacks ist es nicht sinnvoll, einen 0-Reward weiterzugeben, da zwangsläufig immer mehrere Agenten eine längere Zeit den Zielagenten nicht sehen, selbst wenn sie sich optimal verteilen / bewegen. TODO

Dies zeigt auch der Test: TODO

Ist kein Event aufgetreten und haben wir einen 1-Reward vorliegen, dann stellt sich die Frage, ob bereits andere Agenten diesen Reward weitergereicht haben. Befinden sich andere Agenten in Reichweite soll nur ein Agent den Reward weiterreichen. TODO Test

11.5 Kommunikationsvarianten

Allen hier vorgestellten Kommunikationsvarianten ist gemeinsam, dass sie einen Kommunikationsfaktor berechnen, nach denen sie den externen Reward, den ihnen ein anderer Agent übermittelt hat, bewerten. Der Kommunikationsfaktor gewichtet alle Verwendungen des Parameters β (welcher die Lernrate bestimmt). Ein Faktor von 1.0 hieße, dass der externe Reward wie ein normaler Reward behandelt wird, ein Faktor von 0.0 hieße, dass externe Rewards deaktiviert sein sollen. Die Idee ist, dass unterschiedliche Agenten unterschiedlich stark am Erfolg des anderen Agenten beteiligt sind, da ohne Kommunikation jeder Agent versuchen wird, selbst den Zielagenten möglichst in die eigene Überwachungsreichweite zu bekommen, anstatt mit anderen Agenten zu kooperieren, also das Gebiet des Grids möglichst großräumig abzudecken.

Gruppenbildung

11.5.1 Einzelne Gruppe

Mit dieser Variante wird der Kommunikationsfaktor fest auf 1.0 gesetzt und es werden alle Rewards in gleicher Weise weitergegeben. Dadurch wird zwischen den Agenten nicht diskriminiert, was letztlich bedeutet, dass zwar zum einen diejenigen Agenten korrekt mit einem externen Reward belohnt werden, die sich zielführend verhalten, aber zum anderen eben auch diejenigen, die es nicht tun. Deren Classifier werden somit zu einem gewissen Grad zufällig bewertet, denn es fehlt die Verbindung zwischen Classifier und Reward.

Letztlich ist eine Zusammenlegung der Rewards im Grunde mit einer Zusammenlegung aller Sensoren zu vergleichen, Tatsächlich nur ein einzelner Agent?

In Tests (TODO) haben sich dennoch in bestimmten Fällen mit “Reward all equally” deutlich bessere Ergebnisse gezeigt als im Fall ohne Kommunikation. Dies ist wahrscheinlich darauf zurückzuführen, dass in diesen Fällen die Kartengröße und Geschwindigkeit des Zielagenten relativ zur Sichtweite und Lerngeschwindigkeit zu groß war, die Agenten also annahmen, dass ihr Verhalten schlecht ist, weil sie den Zielagenten relativ selten in Sicht bekamen. Eine Weitergabe des Rewards an alle Agenten kann hier also zu einer Verbesserung führen, dabei ist der Punkt aber nicht, dass Informationen ausgetauscht werden, sondern, dass obiges Verhältnis zugunsten der Sichtweite gedreht wird. Für die Auswahl geeigneter Tests sollten die Szenario-Parameter also möglichst so gewählt werden, dass “Reward all equally” keinen signifikanten Vorteil gegenüber “No external reward” bringt. Blickt man auf diesen Sachverhalt aus einer etwas anderen Perspektive ist es auch einleuchtend. Es scheint offensichtlich, dass es relevant ist, ob das Spielfeld z.B. 100x100 oder nur 10x10 Felder groß ist, wenn es darum geht, das Verhalten über die Zeit hinweg zu bewerten. In den Algorithmus für die Kommunikation bzw. für die Rewardvergabe müsste man deshalb einen weiteren (festen) Faktor einbauen, der zu Beginn in Abhängigkeit von Größe des zu überwachenden Feldes berechnet wird. Dies soll aber nicht Teil der Arbeit

werden. TODO

TODO Idee: Verteilt man den Reward an alle Agenten mit gleichem Faktor heisst das letztlich, dass jeder Agent in jedem Zeitschritt den selben Rewardwert erhält. Dann bildet das System der Agenten im Grunde als gemeinsames System von Agenten mit gemeinsamen Sensoren und gemeinsame, ClassifierSet TODO

11.5.2 Gruppenbildung über Ähnlichkeit des Verhaltens der Agenten

Eine weitere Variante berechnet erst einmal für jeden Agenten einen “Egoismus-Faktor”, indem grob die Wahrscheinlichkeit ermittelt wird, dass ein Agent, wenn sich ein anderer Agent in Sicht befindet, sich in diese Richtung bewegt. “Egoismus”-Faktor, weil ein großer Faktor bedeutet, dass der Agent eher einen kleinen Abstand zu anderen Agenten bevorzugt, also wahrscheinlich eher auf eigene Faust versucht, den Zielagenten in Sicht zu bekommen anstatt ein möglichst großes Gebiet abzudecken.

Die Hypothese ist, dass Agenten mit ähnlichem Egoismus-Faktor auch einen ähnlichen Classifiersatz besitzen und der Reward nicht an alle Agenten gleichmäßig weitergegeben wird, sondern bevorzugt an ähnliche Agenten.

Damit gäbe es einen Druck in Richtung eines bestimmten Egoismus-Faktors. TODO

Der Vorteil gegenüber den anderen Verfahren liegt darin, dass der Kommunikationsaufwand hier nur minimal ist, neben dem *reward* muss lediglich der Egoismus Faktor übertragen und pro Zeitschritt nur einmal berechnet werden.

Ein Problem dieser Variante kann sein, dass der Ansatz das Problem selbst schon löst, indem er kooperatives Verhalten belohnt, unabhängig davon, ob Kooperation für das Problem sinnvoll ist.

Die Variante müsste also zum einen in

schlecht abschneiden TODO

TODO rewardrange = kommrage, Vereinfachung

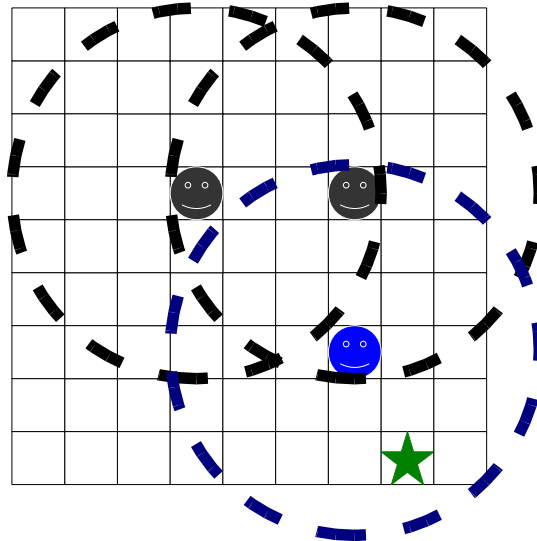


Abbildung 11.1: Schematische Darstellung der Rewardverteilung an ActionSets bei einem neutralen Ereignis

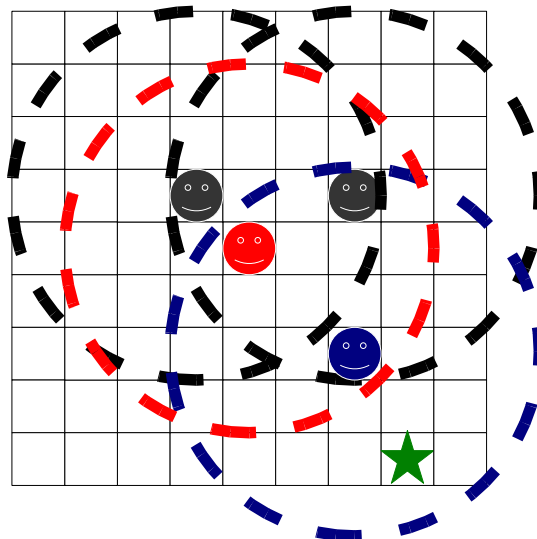


Abbildung 11.2: Schematische Darstellung der Rewardverteilung an ActionSets bei einem neutralen Ereignis

TODO Programm 11.5

```

1  /**
2   * Relation of this classifier set (the active agent classifier set,
3   * e.g. the set that received a reward) to another classifier set
4   * @param other The other set we want to compare with
5   * @return degree of relationship (0.0 – 1.0)
6   */
7  public double checkEgoisticDegreeOfRelationship(
8      final MainClassifierSet other) {
9      double ego_factor =
10         getEgoisticFactor() - other.getEgoisticFactor();
11      if(ego_factor == 0.0) {
12          return 0.0;
13      }
14      return 1.0 - ego_factor * ego_factor;
15  }
16
17  public double getEgoisticFactor() throws Exception {
18      double factor = 0.0;
19      double pred_sum = 0.0;
20      for(Classifier c : getClassifiers()) {
21          if(!c.isPossibleSubsumer()) {
22              continue;
23          }
24          factor += c.getEgoFactor();
25          pred_sum += c.getFitness() * c.getPrediction();
26      }
27      if(pred_sum > 0.0) {
28          factor /= pred_sum;
29      } else {
30          factor = 0.0;
31      }
32      return factor;
33  }

```

Programm 11.5: “Egoistische Relation“, Algorithmus zur Bestimmung des Kommunikationsfaktors basierend auf dem Verhalten des Agenten gegenüber anderen Agenten

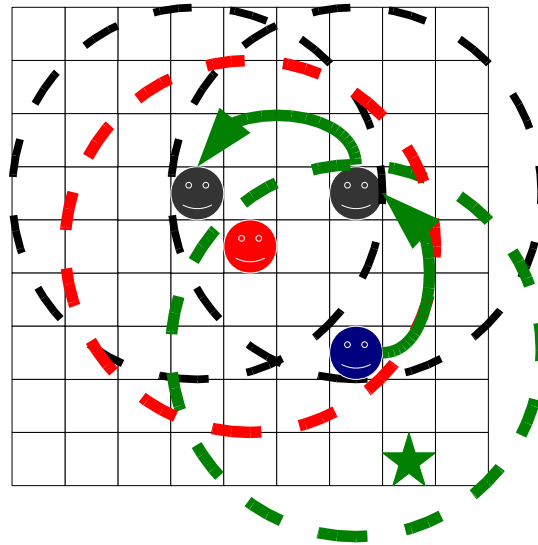


Abbildung 11.3: Schematische Darstellung der Rewardverteilung an ActionSets bei einem neutralen Ereignis

11.6 Bewertung Kommunikation:

Die Vorteile, die man durch Kommunikation erzielen kann, hängt stark von dem Szenario ab. Beispielsweise in dem Fall, bei dem zufällige Agenten bereits fast 100% Abdeckung erreichen, also so viele Agenten auf dem Feld sind, dass der Gewinn durch Absprache minimal ist. Auch ist, weil wir nur mit Binärsensoren arbeiten, die Sensorik gestört, wenn sich sehr viele Agenten auf dem Feld befinden, weil die Sensoren sehr oft gesetzt sind und somit wenig Aussagekraft haben. Erweiterungen wie zusätzliche Sensoren die die Abstände bestimmen würde hier wahrscheinlich klarere Ergebnisse liefern.

Umgekehrt ist der Einfluss bei sehr wenigen Agenten gering. TODO Vergleich unterschiedliche Agentenanzahl, unterschiedliche Kommunikationsmittel Vergleich mit LCS?

11.6.1 Vergleich TODO

Old LCS Agent New LCS Agent

Multistep LCS Agent Dieser Algorithmus stellt eine Implementation des Standard XCS

Algorithmus dar. Unterschied zur Standardimplementation ist, dass die Problemistanz bei Erreichen des temporären Ziels (d.h. den Zielagenten in Sicht zu bekommen) nicht tatsächlich neugestartet wird. Events, wie bei den neuen LCS Implementationen gibt es nicht, ist das Ziel in Sicht wird Reward 1.0 weitergegeben.

Sight range/Kommunikationsrange

LCS Agenten schneiden auch ohne Kommunikation (bei ausreichender Anzahl von Schritten) immer besser ab als zufällige Agenten.

TODOGrafiken

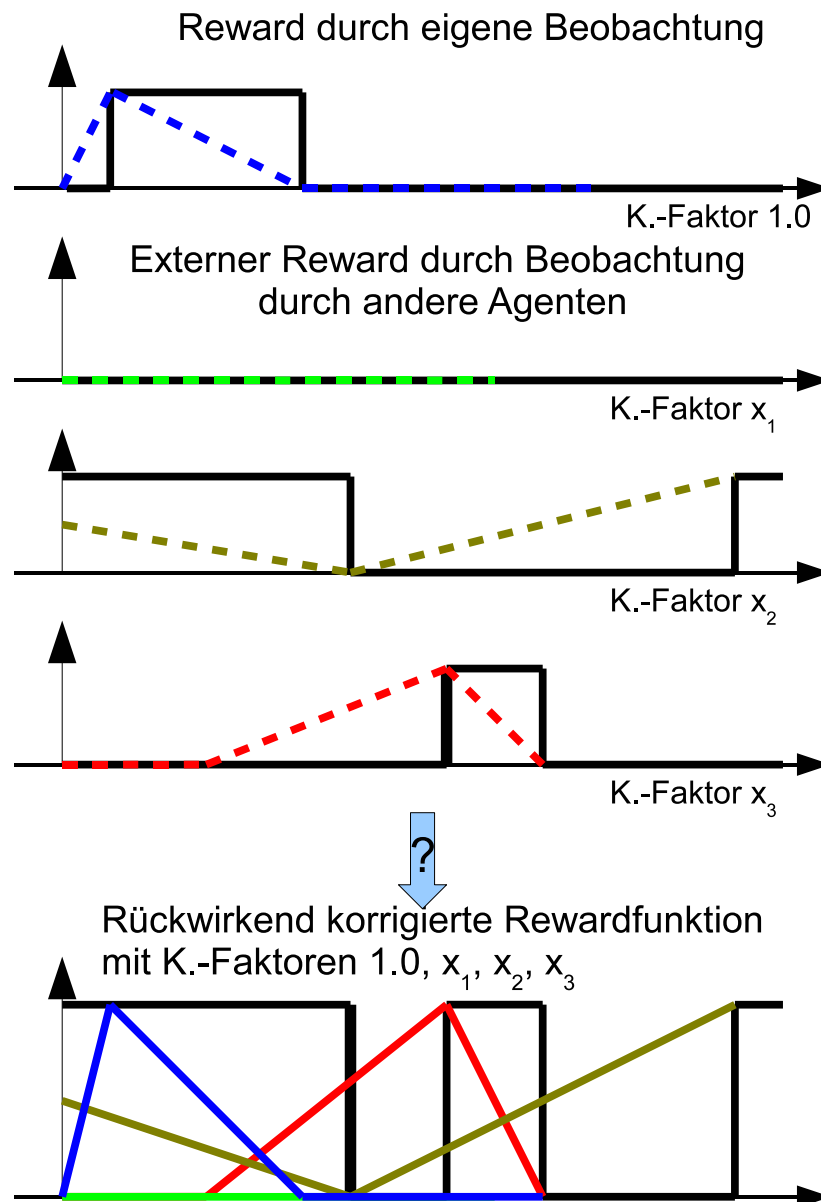


Abbildung 11.4: Beispielhafte Darstellung der Kombination interner und externer Rewards

Kapitel 12

Zusammenfassung, Ergebnis und Ausblick

12.1 Zusammenfassung

Zu Beginn wurde auf die Szenariodefinition und die Fähigkeiten der Agenten eingegangen. Anhand von Beispielen heuristischer Agenten wurden einige Grundeigenschaften der präsentierten Szenarien als Vorbereitung für die Analyse der Learning Classifier Systeme bestimmt. Nach der Einführung in LCS, der Beschreibung des Standardverfahren XCS und der angepassten Implementierung für Überwachungsszenarios konnten dann umfangreiche Tests ausgeführt werden.

von der Möglichkeit zur Kommunikation eine angepasste Implementierung für verzögerten Reward definiert auf Basis dessen dann mehrere Varianten für die Weitergabe des Rewards vorgestellt, analysiert und verglichen wurden.

12.2 Ergebnis

Das wesentliche Ergebnis ist, dass die Implementierung des XCS auf Überwachungsszenarios ausgeweitet werden kann ohne wesentliche Veränderungen am Algorithmus vorzunehmen. Während sich die Qualität der resultierenden Agenten im Allgemeinen über dem zufälligen Agenten befindet, ist die Effizienz der Implementierung, im Vergleich zu einfachen Heuristiken, sehr gering. Mit der verwendeten Implementierung hat XCS Probleme, eine optimale Regelmenge zu finden bzw. zu halten. Eine Regel wie z.B. „laufe auf das Ziel zu, wenn es in Sicht ist“, ist als Heuristik sehr erfolgreich, bei dauerhafter Überwachung ohne Kommunikation läuft es aber eher auf ein Verfolgungsszenario hinaus. Aufgrund andauerndem Lernens TODO

Die alleinige Anpassung des XCS Multistepverfahrens, dass ein neues Problem gestartet wird, wann immer sich das Ziel in Überwachungsreichweite befand führte nicht zum Erfolg, die Ergebnisse waren nicht besser als ein sich zufällig bewogender Agent.

Erst durch Verknüpfung des Rewards mit dem zeitlichen Abstand zu einer Änderung des Zustands führte zu deutlich besseren Ergebnissen.

TODO Desweiteren wurde untersucht, inwiefern sich der Austausch an minimaler Information unter den Agenten, ohne zentrale Steuerung oder globalem Regeltausch, auf die Qualität auswirkt. Zwar gab es vereinzelt positive Effekte, diese waren jedoch auf andere Faktoren zurückzuführen.

12.3 Ausblick

Ein

Weitere Untersuchungen sind nötig um zu bestimmen, inwiefern Kommunikation, beispielsweise mit einer größeren Zahl an besseren Sensoren, zu einem besseren Ergebnis

führen kann. TODO

Vom theoretischen Standpunkt ist noch zu klären, warum genau der zeitliche Abstand zum Erfolg geführt hat und wo die Grenzen hierfür liegen.

Erschwerung, mehr Kollaboration TODO aus verschiedenen Richtungen betrachten? Mehrere Agenten notwendig?

Probiert, aber verworfen:

Während der Arbeit wurden auch einige Ansätze probiert aber mangels Erfolgsaussichten wieder verworfen. Ursprünglich wurde das Szenario auf Basis von Rotation konzipiert. Die Annahme war, dass ein Agent, der für einen Satz an Sensordaten eine optimales *classifier set* gefunden hat, dieses *classifier set* auch für Sensordaten eines um 90, 180 und 270 Grad gedrehten Szenarios (mit entsprechend 90, 180 und 270 Grad gedrehter Aktion des jeweiligen *classifier*) optimal sei. Aufgrund der deutlichen Komplexitätssteigerung des Programms, der niedrigeren Laufzeit und mangels konkreter Qualitätssteigerungen gegenüber dem Ansatz ohne Rotation wurde diese Idee jedoch fallengelassen. Möglicherweise könnte man durch Hinzunahme eines weiteren Bits im *condition* Vektor, das bestimmt, ob dieser *classifier* gleichzeitig auch die drei rotierten Szenarien erkennen kann, die Leistung des Systems verbessern, dies bedarf aber weiterer Untersuchung und geht am eigentlichen Thema dieser Arbeit vorbei.

Abnehmende Exploration LITERATUR Intelligent Exploration Method to Adapt Exploration Rate in XCS, Based on Adaptive Fuzzy Genetic Algorithm An Adaptive Approach for the Exploration-Exploitation Dilemma for Learning Agents

Im Bereich der Kommunikation wurde neben der „egoistischen Relation“ (siehe Kapitel 11.5.2) auch weitere Verfahren ausprobiert, mit welchen versucht wurde, gleichartige Gruppen zu finden. Hier wurden ganze *classifier set* Listen unterschiedlicher Agenten miteinander auf Ähnlichkeit geprüft um daraus einen Faktor zu berechnen, der (wie bei der

„egoistischen Relation“) Einfluss auf die Weitergabe des *reward* Werts haben sollte. Der dadurch deutlich erhöhte Kommunikations- und Berechnungsaufwand lag jedoch in keinem Verhältnis zu eventuell beobachteten Qualitätsverbesserungen, im Gegenteil wurden eher Qualitätsverschlechterungen bemerkt. Die Ergebnisse mit dem Test der „egoistischen Relation“ zeigen jedoch, dass hier zumindest etwas Potential stecken könnte und für bestimmte Szenarien die zwei Grundideen, dass sich die Agenten zum einen an die Größe des Szenarios anpassen und zum anderen der *reward* möglichst nur an sich ähnlich verhaltende Agenten weitergegeben wird, nicht ganz falsch sein können. Genauere, insbesondere theoretische, Untersuchungen sind hier aber nötig.

Kapitel 13

Verwendete Hilfsmittel und Software

Zu Beginn stellte sich die Frage, welche Software zu benutzen ist, da es sich um ein recht komplexe Problemstellung handelt. Begonnen habe ich mit der YCS Implementierung [Bul03]. Sie ist in der Literatur wenig vertreten, die Implementierung bot aber einen guten Einstieg in das Thema, da sie sich auf das Wesentliche eines LCS beschränkte und keine Optimierungen enthielt.

Der nächste Schritt war zu entscheiden, auf welchem System die Agenten simuliert werden sollen. Implementierungen wie

Unter einer Reihe von vorhandenen Implementierungen entschied ich mich für eine eigene Implementation.

Wesentlicher Grund war die Unerfahrenheit mit den Lösungen (und der damit verbundenen Einarbeitungszeit) wie auch Überlegungen bzgl. der Geschwindigkeit, dem Speicherverbrauch und der Kompatibilität. TODO

Das Programm und die zugehörige Oberfläche zum Erstellen von Test-Jobs wurden in Netbeans 6.5 programmiert.

Grafiken wurden mittels GnuPlot erstellt.

Grafiken der Grid-Konfiguration wurden im Programm mittels GifEncode TODO erste
* @version 0.90 beta (15-Jul-2000) * @author J. M. G. Elliott (tep@jmge.net)

Wesentlicher Bestandteil der Konfigurationsoberfläche war auch eine Automatisierung der Erstellung von Konfigurationsdateien, Batchdateien (für ein Einzelsystem und für JoSchKA) zum Testen einer ganzen Reihe von Szenarien und auch GnuPlot Skripts.

Speicherverbrauch

Speicherung der Agentenpositionen und des Grids verbrauchen fast keinen Speicher
TODO Wesentlicher Faktor waren die LCS Systeme mit ihren ClassifierSets TODO

OpenOffice

L^{Ed} Latex

13.1 Beschreibung des Konfigurationsprogramms

Abbildung 13.1: Screenshot des Konfigurationsprogramms

Anhang A

Statistical significance tests

This is the first appendix

Anhang B

Implementation

The second appendix...

Literaturverzeichnis

- [Bar02] A. Barry. The stability of long action chains in xcs, 2002.
- [BD03] Alwyn Barry and Claverton Down. Limits in long path learning with xcs. In *Proc. GECCO 2003, Genetic and Evolutionary Computation Conference, 2003*, pages 1832–1843. Springer-Verlag, 2003.
- [BDE⁺99] W. Banzhaf, J. Daida, A. E. Eiben, M. H. Garzon, V. Honavar, M. Jakiela, and R. E. Smith. Extending the representation of classifier conditions, part i: From binary to messy coding. In *Proceedings of the Genetic and Evolutionary Computation Conference*, pages 337–344. Morgan Kaufmann, 1999.
- [BGL05] M. V. Butz, D. E. Goldberg, and P. L. Lanzi. Gradient descent methods in learning classifier systems: improving xcs performance in multistep problems. *IEEE Transactions on Evolutionary Computation*, 9(5):452–473, Oct. 2005.
- [Bul03] Larry Bull. A simple accuracy-based learning classifier system. Technical report, Learning Classifier Systems Group Technical Report UWELCSG03-005, 2003.
- [But00] Martin V. Butz. Xcs classifier system in java, 2000.
- [But06a] Martin V. Butz. *Simple Learning Classifier Systems*, chapter 4, pages 31–50. Springer, 2006.

- [But06b] Martin V. Butz. *The XCS Classifier System*, chapter 4, pages 51–64. Springer, 2006.
- [BW01] Martin V. Butz and Stewart W. Wilson. An algorithmic description of XCS. *Lecture Notes in Computer Science*, 1996:253–272, 2001.
- [Ham04] Carol Hamer, 2004.
- [HFA02] Luis Miramontes Hercog, Terence C. Fogarty, and London Se Aa. Social simulation using a multi-agent model based on classifier systems: The emergence of vacillating behaviour in the „el farol“ bar problem. In *Proceedings of the International Workshop in Learning Classifier Systems 2001*. Springer-Verlag, 2002.
- [ITS05] Hiroyasu Inoue, Keiki Takadama, and Katsunori Shimohara. Exploring xcs in multiagent environments. In *GECCO '05: Proceedings of the 2005 workshops on Genetic and evolutionary computation*, pages 109–111, New York, NY, USA, 2005. ACM.
- [KM94] S. Kobayashi K. Miyazaki, M. Yamamura. On the rationality of profit sharing in reinforcement learning. In *Proceedings of the 3rd International Conference on Fuzzy Logic, Neural Nets and Soft Computing*, pages 285–288, 1994.
- [Lan] P. L. Lanzi. The xcs library.
- [LWB08] A. Lujan, R. Werner, and A. Boukerche. Generation of rule-based adaptive strategies for a collaborative virtual simulation environment. In *Proc. IEEE International Workshop on Haptic Audio visual Environments and Games HAVE 2008*, pages 59–64, 18–19 Oct. 2008.

- [MVBG03] K. Sastry M. V. Butz and D. E. Goldberg. Tournament selection: Stable fitness pressure in xcs. In *Lecture Notes in Computer Science*, pages 1857–1869, 2003.
- [TB06] J.-M. Nigro T. Benouhiba. An evidential cooperative multi-agent system. *Expert Systems with Applications*, 30(2):255–264, 2006.
- [THN⁺98] K. Takadama, K. Hajiri, T. Nomura, M. Okada, S. Nakasuka, and K. Shimohara. Learning model for adaptive behaviors as an organized group of swarm robots. *Artificial Life and Robotics*, 2(3):123–128, 1998.
- [Wil95] Stewart W. Wilson. Classifier fitness based on accuracy. *Evolutionary Computation*, 2(3):149–175, 1995.
- [Wil98] Stewart W. Wilson. Generalization in the xcs classifier system. In *Genetic Programming 1998: Proceedings of the Third Annual Conference*, pages 665–674. Morgan Kaufmann, 1998.

Erklärung

Ich versichere hiermit wahrheitsgemäß , die Arbeit bis auf die dem Aufgabensteller bereits bekannte Hilfe selbständig angefertigt, alle benutzten Hilfsmittel vollständig und genau angegeben und alles kenntlich gemacht zu haben, was aus Arbeiten anderer unverändert oder mit Abänderungen entnommen wurde.

Karlsruhe, 30. März 2009,

Clemens Lode