



Gesten- und Sprachsteuerung für einen mobilen Roboter mittels Kinect for Windows unter Java

Studienarbeit

für die Prüfung zum

Bachelor of Science

der Angewandte Informatik

an der Dualen Hochschule Baden-Württemberg Karlsruhe

von

Simon Ebner, Volker Werling

Januar 2013

Bearbeitungszeitraum
Matrikelnummer, Kurs
Gutachter

12 Wochen
5837963 7012192, TAI10B2
Prof. Hans-Jörg Haubner

Erklärung

Wir erklären hiermit ehrenwörtlich:

1. dass wir unsere Studienarbeit mit dem Thema *Gesten- und Sprachsteuerung fuer einen mobilen Roboter mittels Kinect for Windows unter Java* ohne fremde Hilfe angefertigt haben;
2. dass wir die Übernahme wörtlicher Zitate aus der Literatur sowie die Verwendung der Gedanken anderer Autoren an den entsprechenden Stellen innerhalb der Arbeit gekennzeichnet haben;
3. dass wir unsere Studienarbeit bei keiner anderen Prüfung vorgelegt haben;
4. dass die eingereichte elektronische Fassung exakt mit der eingereichten schriftlichen Fassung übereinstimmt.

Wir sind uns bewusst, dass eine falsche Erklärung rechtliche Folgen haben wird.

Karlsruhe, Januar 2013

Simon Ebner, Volker Werling

Zusammenfassung

Inhaltsverzeichnis

1	Einleitung	7
2	Aufgabenstellung	10
3	Stand der Technik	11
3.1	Kinect	11
3.1.1	Hardware	12
3.1.2	Software	15
3.2	Java – Eclipse	16
3.2.1	OSGi	17
3.2.2	Grafische Oberfläche	17
3.3	Roboter – Ausblick	19
3.4	Mustererkennung	19
4	Kinect–Framework	21
4.1	Microsoft Kinect SDK	21
4.2	Java Frameworks	21
4.2.1	OpenNI	21
4.2.2	OpenKinect	22
4.2.3	jnect	22
4.2.4	Scoring–Modell	22
4.2.5	Analyseergebnis	22
5	Konzeption	23
5.1	Technische Vorgaben	23
5.1.1	Vorgaben durch Kinect	23
5.1.2	Vorgaben durch jnect–Framework	23

5.1.3	Service-Architektur	23
5.1.4	Abhängigkeit zu Robotertechnik	23
5.2	Fachliche Vorgaben	23
5.2.1	Vorgaben durch Mensch-Computer-Interaktion	23
5.2.2	Klassifikation der Anwendung	23
5.2.3	Analyse der Gesten- und Sprachsteuerung	23
6	Modelle zur Gesten- und Spracherkennung	24
6.1	Dynamic Time Warping	24
6.2	Künstliche neuronale Netze	25
6.3	Hidden Markov Model	25
6.3.1	Auswahl des Modells für die Arbeit	26
6.3.2	Diskrete Markov-Prozesse	26
6.3.3	Erweiterung auf Hidden Markov Model	29
6.3.4	Die drei grundlegenden Probleme eines Hidden Markov Model . .	31
6.4	Baum-Weich-Algorithmus	32
6.4.1	ToDo	32
6.5	ToDo	32
7	Gesten	33
7.1	Kreisbewegung	33
7.1.1	ToDo	33
7.2	Vorwärtsbewegung	33
7.2.1	ToDo	33
7.3	Haltesignal	33
7.3.1	ToDo	33
8	Sprachbefehle	34
8.1	ToDo	34
9	Implementierung	35
9.1	OSGI-Bundles	35
9.1.1	Gesten	35
9.1.2	Sprachbefehle	35
9.1.3	Roboterschnittstelle	35

9.1.4	Ausführbare Aktionen	35
9.2	Oberfläche	35
9.2.1	RCP–Anwendung	35
9.2.2	Grafische Darstellung	35
9.3	Ablaufbeschreibung	35
10	Ausblick - weitere Arbeiten	36
	Abbildungsverzeichnis	37
	Tabellenverzeichnis	38
	Quellcodeverzeichnis	39
	Literaturverzeichnis	40
	Glossar	42

1 Einleitung

Gestik im Sinne von kommunikativen Bewegungen stellt einen wesentlichen Teil der nonverbalen Kommunikation dar. Die Definition nach Kurtenbach und Hulteen (1990) besagt:

A gesture is a motion of the body that contains information. Waving goodbye is a gesture. Pressing a key on a keyboard is not a gesture because the motion of a finger on its way to hitting a key on a keyboard is neither observed nor significant. All that matters is which key was pressed.

Als solches bildet eine Geste abstrakte Strukturen und bildhafte Vorstellungen unmittelbar ab. Sie ist somit eine Körperbewegung, die Informationen enthält. Informationen, die herkömmliche Eingabegeräte, darunter Tastatur und Maus, nicht in dieser Form wiedergeben können. Ein Umstand, der sich besonders im Bereich der Mensch-Computer-Interaktion auswirkt.

Das Ziel der Mensch-Computer-Interaktion ist es, die Kommunikation intuitiv und unmittelbar zu gestalten. Genau an dieser Stelle setzen Gesten an, in dem sie eine bessere Schnittstelle zwischen Mensch und Maschine darstellen.

Die Spannweite von Realisierungen von Mensch-Maschine-Schnittstelle ist groß. Frühe Ansätze zeigen bereits das Potential, Gesten zur Steuerung und Nutzung von diversen Geräten zu verwenden. Das *Theremin* - ein 1919 erfundenes elektornisches Musikinstrument, das mit Händen berührungslos durch Beeinflussung eines elektromagnetischen Feldes gesteuert wird und dadurch Töne erzeugt - ist eine der ersten technischen Lösungen für eine Gestensteuerung. Eine weitere Form, der *Datenhandschuh*, 1977 von Electronic

Visualization Laboratory Labs entwickelt ¹, erregte großes Interesse und wurde unter anderem als *Power Glove* von Mattel vertrieben. Datenhandschuhe werden bis heute in den verschiedensten Bereichen eingesetzt.

Eine weitere Variante der Gestenerkennung ist die Bewegungskennung. Dabei wird über optische Sensoren der Körper oder einzelne Körperteile des Nutzers erkannt und somit die Steuerung durch Gesten ermöglicht. Eine frühe Lösung ist der sogenannte *Videoplace* ², entwickelt von Myron Krueger Mitte der 70'er Jahre ³. Ausgelegt als ein *Labor für künstliche Realität* war eine Person, durch ihn umgebende Projektoren und Videokameras, in der Lage seine Umgebung zu beeinflussen. Moderne Lösungen haben das Konzept der Detektion durch Kameratechnik weiterentwickelt und einer weiten Öffentlichkeit zugänglich gemacht. Darunter auch das *PLAYSTATIONEye* ⁴. Der populärste Vertreter ist *Kinect for XBOX* mit über 18 Millionen verkauften Exemplaren ⁵. Im folgenden wird mit jener Kinect gearbeitet.

In dieser Arbeit wird mit Hilfe der Kinect eine Gesten- und Spracherkennung entwickelt um eine Steuerung für einen mobilen Roboter zu implementieren. Aufgrund des Umfangs dieser Studienarbeit wird die Erarbeitung dieser Anwendung in zwei Teile aufgeteilt.

Der erste Teil der Arbeit beschäftigt sich vorrangig mit der Realisierung der Gesten- und Spracherkennung. Die eigentliche Steuerung für einen mobilen Roboter wird vorerst durch eine Schnittstelle und einer Testumgebung ersetzt. Der Fokus liegt hierbei in erster Linie auf der verwendeten Technik, in diesem Fall der Kinect zur Bewegungskennung, den verwendeten Gesten und Sprachbefehlen, sowie der Umsetzung der Software und den darin verwendeten Technologien. Der zweite Teil wird an diese Arbeit anknüpfen und im weiteren die Umsetzung der Steuerung für und den Einsatz eines mobilen Roboters beschreiben.

¹Sturman, D.J., Zeltzer, D. (January 1994). „A survey of glove-based input“. IEEE Computer Graphics and Applications 14 (1): 30–39

²Beschreibung der Einrichtung „VIDEOPLACE“. medienkunstnetz.de. Abgerufen Dezember 22, 2012

³Myron Krueger. Artificial Reality 2, Addison-Wesley Professional, 1991. ISBN 0-201-52260-8

⁴„PLAYSTATIONEye Brings Next-Generation Communication to PLAYSTATION3“. us.playstation.com. Sony Computer Entertainment America. Abgerufen Dezember 21, 2012

⁵Takahashi, Dean (Januar 9, 2012). „Xbox 360 surpasses 66M sold and Kinect passes 18M units“. venturebeat. Abgerufen Dezember 20, 2012

Die folgende Ausarbeitung spiegelt die Umsetzung des ersten Teils dieser Arbeit wieder.

2 Aufgabenstellung

3 Stand der Technik

3.1 Kinect

Kinect ist eine Gerät zur Dedektion von Bewegungen. Entwickelt wurde es von PrimeSense als Hardware zur Steuerung der Videokonsole XBOX 360 von Microsoft Corp. Nach der Ankündigung im März 2010 ¹ war die Erweiterung mit dem Erscheinungsdatum 10. November 2010 in der Ausführung *Kinect for XBOX 360* in Europa erhältlich ². Nach einem sehr erfolgreichem Verkaufsstart ³ und anhaltender Nachfrage⁴, kündigte Microsoft am 9. Januar 2012 ein weiteres Modell der Kinect, die sogenannte *Kinect for Windows* für den 1. Februar 2012 an ⁵.

Darüber hinaus veröffentlichte Microsoft bereits am 16. Juni 2011 eine erste Version ihrer *Kinect for Windows SDK*⁶. Mit diesem Software Development Kit ist es möglich auf einer Windows 7 Plattform Anwendungen zu entwickeln, die eine Kinect als Eingabegerät verwenden.

Mit der Verfügbarkeit einer Kinect-Variante, die für den Einsatz am PC ausgelegt ist und der frei zugänglichen einer breiten Öffentlichkeit als Forschungsgegenschaft zugänglich,

¹Pressemitteilung, der Veröffentlichung. Siehe Link. Microsoft.com. Abgerufen Dezember 20, 2012

²Erscheinungsdatum der *Kinect for XBOX 360*. Siehe Link. BBC UK. Abgerufen Dezember 20, 2012

³„Am schnellsten verkauftes Peripheriegerät für Spiele“. Guinnessworldrecords.com. Abgerufen Dezember 22, 2012

⁴Takahashi, Dean (Januar 9, 2012). „Xbox 360 surpasses 66M sold and Kinect passes 18M units“. venturebeat. Abgerufen Dezember 20, 2012

⁵„Ankündigung der *Kinect for Windows*“. blogs.msdn.com. Abgerufen Dezember 22, 2012

⁶„Microsoft Releases Kinect for Windows SDK Beta for Academics and Enthusiasts“. Microsoft.com. Abgerufen Dezember 21, 2012

was auch der ausschlaggebende Punkt für den Einsatz der Kinect in dieser Studienarbeit ist.

3.1.1 Hardware

Die beiden Kinect-Modelle unterscheiden sich in einigen Details⁷. Der bedeutendste Faktor ist das sogenannte *Near Mode* Feature:

Near Mode enables the depth sensor to see objects as close as 40 centimeters and also communicates more information about depth values outside the range than was previously available. There is also improved synchronization between color and depth, mapping depth to color, and a full frame API.⁷

Neben dem aktualisierten Tiefensensor unterscheiden sich die beiden Varianten auch in einer höheren Auflösung der RGB-Kamera, die für eine Gestenerkennung relevant sein kann. Aus diesem Grund wird für diese Arbeit die *Kinect for Windows* genutzt.

Technische Daten

Eine Kinect enthält innerhalb des Gehäuses, folgende Sensoren⁸:

- Eine RGB-Kamera mit einer Auflösung von 1280x960. Dies ermöglicht eine Farbbilderfassung
- Ein Infrarot (IR) Emitter und ein IR Tiefensensor. Der Emitter emittiert infrarote Lichtstrahlen und der Tiefensensor erfasst die an den Sensor reflektierten Strahlen. Die reflektierten Strahlen werden in Tiefeninformation umgewandelt, in dem der

⁷„Informationsseite über Unterschiede der Kinect Versionen“. Microsoft.com. Abgerufen Dezember 22, 2012

⁸„Kinect for Windows Sensor Components and Specifications“. msdn.microsoft.com. Abgerufen Dezember 21, 2012

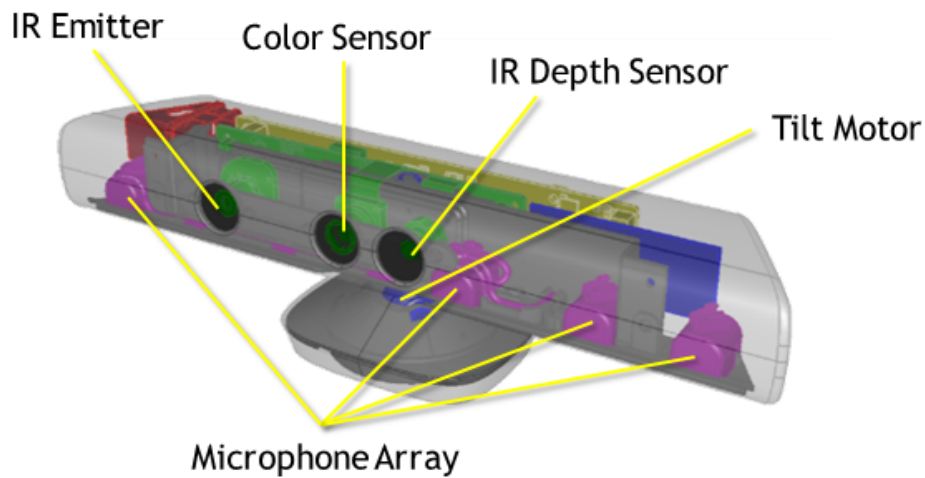


Abbildung 3.1: Schematische Ansicht der Sensoren einer *Kinect for Windows* (Quelle: Microsoft Corp.¹⁾)

Abstand zwischen Objekt und Sensor bestimmt wird. Dies ermöglicht die Erfassung von Tiefenbildern

- Ein Mikrofonarray, das vier Mikrofone zur Soundaufnahme enthält. Die Anzahl der Mikrofone ermöglicht nicht nur die Aufzeichnung von Audiodaten, sondern auch die Lokalisierung der Soundquelle und die Richtung des Audiosignals
- Ein 3-Achsen-Beschleunigungssensor, konfiguriert für einen Bereich der zweifachen Erdbeschleunigung, um die gegenwärtige Ausrichtung der Kinect zu bestimmen
- Ein Kippmotor, zur automatisierten Justierung der Sensoren

Weitere Details der technischen Spezifikation einer *Kinect for Windows* sind in Tabelle 3.1 aufgelistet:

Tabelle 3.1: Kinect for Windows – technische Spezifikation ⁸

Kinect	Spezifikation
Blickwinkel	43° vertikales, 57° horizontales Blickfeld
Vertikaler Neigebereich	$\pm 27^\circ$
Bildwiederholrate (Farb und Tiefensignal)	30 Bilder pro Sekunde (FPS)
Audioformat	16-kHz, 24-bit mono pulse code modulation (PCM)
Audioeingang	Ein Vier-Mikrofonarray mit 24-Bit Analog-Digital-Wandler (ADC)
Datensignal–Tiefensensor	640x480 16-bit, 30 Bilder pro Sekunde
Datensignal–RGB-Kamera	1280x960 16-Bit, 12 Bilder pro Sekunde 640x480 16-Bit, 30 Bilder pro Sekunde
Tiefensensorreichweite	0,4 – 4 m
<i>Skelett Tracking System</i>	Erkennung von bis zu sechs Benutzern, zwei davon trackbar/verfolgbar Verfolgung von 20 Gelenken pro aktivem Nutzer

⁸„Kinect for Windows Sensor Components and Specifications“. msdn.microsoft.com. Abgerufen Dezember 21, 2012

3.1.2 Software

Durch die Veröffentlichung eines Software Development Kits ist es möglich, die Kinect in eigene Programme einzubinden und neue Anwendungsfälle zu bearbeiten. Dazu hat Microsoft ebenfalls ein Handbuch veröffentlicht [Cor12], das Informationen und Ratschläge zum Entwurf von Anwendungen liefert.

Kinect for Windows SDK

Das *Kinect for Windows SDK* steht aktuell in der Version 1.6 ⁹ bereit. Dabei kann direkt in den Programmiersprachen C++, C# , und Visual Basic auf einer Windows 7 Plattform entwickelt werden.

Das SDK umfasst dabei folgende Funktionen ¹⁰:

- Treiber und technische Dokumentation zur Implementierung von Anwendungen, die Kinect for Windows Sensoren nutzen
- Referenz APIs und Dokumentation zur Programmierung von *managed* und *unmanaged* Code. Die APIs stellen mehrere Mediensignale bereit
- Beispiele und Best-Practice-Lösungen

Weitere Frameworks

Da Microsoft die Nutzungsmöglichkeiten seines SDK hinsichtlich verwendeter Programmiersprache und Plattform einschränkt, begannen Forscher eigene Frameworks und Treiber zur Nutzung der Kinect zu entwickeln¹¹. PrimeSense selbst veröffentlichte Treiber

⁹Downloadseite des SDK. Microsoft.com. Abgerufen Dezember 21, 2012

¹⁰„Kinect for Windows Programming Guide“. msdn.microsoft.com. Abgerufen Dezember 23, 2012

¹¹„Open Source Kinect contest has been won“. hackaday.com. November 11, 2010. Abgerufen Dezember 21, 2012

und Middleware für die Kinect ¹², als das offizielle SDK noch nicht zur Verfügung stand. Das Ziel dieser Frameworks war es zu meist, Kinect-Steuerung unter Unix-Plattformen und Programmiersprachen, wie Java, nutzbar zu machen. Ein Teil dieser Ausarbeitung ist es, ein für die Aufgabenstellung und Zielsetzung der Studienarbeit passendes Framework zu bestimmen.

Verbreitete Lösungen, die im Kapitel 4 näher betrachtet werden, sind:

- OpenNI
- OpenKinect
- jnect

3.2 Java – Eclipse

Die Anwendung, die im Rahmen dieser Arbeit entsteht, wird auf Basis der Programmiersprache Java und der Entwicklungsumgebung Eclipse erstellt. Die Nutzung der Kinect und des Kinect SDK unterstützt nativ nicht die Entwicklung auf Basis von Java, wie in Abschnitt 3.1.2 festgestellt wurde. Da jedoch die Expertise der Autoren dieser Arbeit im Java-Umfeld liegt, wird die Anwendung dennoch in der Programmiersprache Java umgesetzt.

Als Entwicklungsumgebung kommt die Eclipse IDE zum Einsatz. Diese bietet alle Voraussetzungen zur Entwicklung von flexiblen Java-Anwendungen und bieten eine ausführliche Support-Plattform für Entwickler. Zusätzlich setzt es auf eine dynamische Supportplattform zur Modularisierung von Anwendungen¹³.

¹²Mitchell, Richard (Dezember 10, 2010). „PrimeSense releases open source drivers, middleware for Kinect“. Joystiq. Abgerufen Dezember 22, 2012

¹³„OSGi“. eclipse.org. Abgerufen Januar 1, 2013

3.2.1 OSGi

Zur Entwicklung einer komplexen Anwendung mit vielen Abhängigkeiten und Schnittstellen, wie die hier zu erarbeitende, braucht es eine entsprechend *mächtige* Software-Architektur. Dabei wird hier auf OSGi-Plattform gesetzt. OSGi¹⁴ ist ein Framework, dass Spezifikationen zur Definition eines dynamischen *Komponentenmodells* unter Java. Diese Komponenten werden *Bundles* genannt und stehen anderen Komponenten als sogenannte *Services* bereit. Die OSGi Alliance¹ selbst spezifiziert lediglich die Programmerschnittstelle. Die Umsetzung für die Eclipse IDE lautet Equinox¹⁵.

3.2.2 Grafische Oberfläche

Das Graphical User Interface (GUI) ist ein wesentlicher Bestandteil der Anwendung, da hierüber der gesamte Informationsaustausch mit dem Nutzer der Kinect stattfindet. Dabei wird auf verschiedene Techniken gesetzt, die im folgenden kurz beschrieben werden. Bei der grafischen Darstellung muss man zwischen Informationen der Anwendung und Interaktion des Nutzers mit der Software unterscheiden. Dazu werden nämlich unterschiedliche Techniken verwendet.

Rich client platform

Die Rich client platform (RCP) ist ein Werkzeug zur Entwicklung von unabhängigen Software-Komponenten, die nicht nur auf GUI-Komponenten beschränkt ist¹⁶. Dabei liegt der Fokus auf Client-Applikation, worunter auch die Anwendung dieser Arbeit fällt. Daher wird RCP im Rahmen dieser Arbeit zur Erstellung der Oberfläche der Anwendung verwendet.

¹⁴„OSGi Alliance“. osgi.org. Abgerufen Januar 1, 2013

¹⁵„equinox OSGi“. eclipse.org. Abgerufen Januar 1, 2013

¹⁶„Rich Client Platform“. wiki.eclipse.org. Abgerufen Januar 1, 2013

Die Lightweight Java Game Library

Die grafische Darstellung der Interaktion zwischen Akteur und Kinect ist ein wichtiger Bestandteil dieser Anwendung. Zur Anzeige aufwändiger grafischer Objekte wird auf die Lightweight Java Game Library (LWJGL)¹⁷ gesetzt. Damit können einfach zweidimensionale Grafiken, zum Beispiel Kreise, oder komplexe dreidimensionale Ansichten in einer Java-Anwendung angezeigt werden. Diese Java-Bibliothek stellt eine API zum Zugriff und zur Verwendung der bekannten Open Graphics Library¹⁸ bereit.

Eclipse Modeling Framework

Das Eclipse Modeling Framework (EMF)¹⁹ ist ein Framework zur Modellierung von Java-Anwendungen. Es ermöglicht aus strukturierten Datenmodellen, Java-Code zu generieren. Für die Kinect-Anwendung sind lediglich die Bereiche des Frameworks relevant, die diese *Modelle* beschreiben.

Graphical Editing Framework

Zur grafischen Darstellung und grafischen Manipulation der zuvor vorgestellten EMF-Modelle existiert das sogenannte Graphical Editing Framework (GEF)²⁰. Es stellt Methoden zur Verfügung, um Grafikeditoren und weitere Ansichten für Eclipse Anwendungen zu erstellen. Für die Studienarbeit relevante Bereiche, sind die Funktionen, die eine vereinfachte Darstellung der eben genannten EMF-Modelle ermöglichen.

¹⁷„LWJGL Lightweight Java Game Library“. lwjgl.org. Abgerufen Januar 1, 2013

¹⁸„The Industry’s Foundation for High Performance Graphics“. opengl.org. Abgerufen Januar 1, 2013

¹⁹„Eclipse Modeling Framework Project (EMF)“. eclipse.org. Abgerufen Januar 1, 2013

²⁰„GEF (Graphical Editing Framework)“. eclipse.org. Abgerufen Januar 1, 2013

3.3 Roboter – Ausblick

Ziel und Aufgabenstellung dieser Arbeit ist, wie bereits in Kapitel 2 erläutert, der Entwurf einer Steuerung für einen mobilen Roboter. Der erste Teil, der zweigeteilten Studienarbeit befasst sich dabei nicht direkt mit der eigentlichen Robotersteuerung oder der Beschreibung und Auswahl eines entsprechenden Modells. Dies ist Aufgabe und Fokus des zweiten Teils.

Für die Zwecke der Arbeiten an diesem Teil und der dabei zu erstellenden Anwendung, wird lediglich ein Grundgerüst und eine Schnittstelle erarbeitet, um auf die Kinect-Steuerung zuzugreifen. Siehe hierzu Abschnitt 9.1.3.

3.4 Mustererkennung

Die *Kinect for Windows* und die dazu gehörende *Kinect for Windows SDK* besitzen bereits Algorithmen zur Detektion von Bewegungen einzelner Körperteile und der Erkennung von Sprache. Dabei muss festgehalten werden, dass zwischen der Detektion der Bewegung eines Gelenkes und der Erkennung einer Geste eine große Diskrepanz besteht. Die Bewegung an sich macht noch keine Geste aus. Zur *Gestenerkennung* sind Verfahren und Algorithmen notwendig, die über den Funktionsumfang der Kinect hinausgehen. Hierzu sind Verfahren der Mustererkennung notwendig, die auf mathematischen Modellen beruhen. Dabei haben sich zwei Richtungen herausgebildet. Eine Variante nutzt hierbei sogenannte künstliche neuronale Netze. Beschrieben und umgesetzt wurde dieser Ansatz unter anderem von Chun Zhu and Weihua Sheng [ZS09].

Eine weitere Variante setzt das sogenannte Hidden Markov Model (HMM) ein, das von Leonard E. Baum [BP66] postuliert wurde. Dabei handelt es sich um ein stochastisches Modell, indem ein System durch eine Markov-Kette (Siehe Abschnitt 6.3.2) modelliert wird. Eine der ersten Arbeiten zum Einsatz in der Gestenerkennung schrieb Lawrence R. Rabiner [Rab89]. Dieser Ansatz findet häufige Verwendung in Bereichen der Spracherkennung, Gestenerkennung und in der Mensch-Maschine-Kommunikation. Aus diesem

Grund wurde entschieden, ein HMM zu erstellen. Das stochastische Modell dahinter wird in Kapitel 6 und die Umsetzung in der Implementierung in Kapitel 9 beschrieben.

4 Kinect–Framework

4.1 Microsoft Kinect SDK

asdf

4.2 Java Frameworks

adsf

4.2.1 OpenNI

Beschreibung

adsf

Technische Daten

SWOT–Analyse

4.2.2 OpenKinect

Beschreibung

Technische Daten

SWOT–Analyse

4.2.3 Kinect

Beschreibung

Technische Daten

SWOT–Analyse

4.2.4 Scoring–Modell

4.2.5 Analyseergebnis

adsasfd

5 Konzeption

5.1 Technische Vorgaben

5.1.1 Vorgaben durch Kinect

5.1.2 Vorgaben durch jnect-Framework

5.1.3 Service-Architektur

5.1.4 Abhängigkeit zu Robotertechnik

5.2 Fachliche Vorgaben

5.2.1 Vorgaben durch Mensch-Computer-Interaktion

5.2.2 Klassifikation der Anwendung

5.2.3 Analyse der Gesten- und Sprachsteuerung

6 Modelle zur Gesten- und Spracherkennung

Es existieren diverse Verfahren zur Gesten- und Spracherkennung. Die bekanntesten Vertreter sind künstliche neuronale Netze, Dynamic time warping (DTW) und HMM. Im Folgenden werden dennoch kurz die diversen Eigenschaften und Algorithmen hinter den weiteren Modell erörtert, um eine fundierte Entscheidung über die Auswahl eines dieser Modelle für diese Arbeit und das zu erarbeitende Programm zu treffen.

6.1 Dynamic Time Warping

DTW ist ein Algorithmus zur Messung von Ähnlichkeiten zwischen einer Eingabesequenz und einem gegebenen Muster, die sich im zeitlichen Ablauf oder der Geschwindigkeit unterscheiden. Dabei würden zum Beispiel die Gemeinsamkeiten im Laufmuster erkannt werden, selbst wenn eine Person in einem Video langsam läuft und diese Person in einer weiteren Aufnahme schneller laufen würde, oder gar eine Beschleunigung oder Verlangsamung während des Zeitraumes der Beobachtung stattfindet.

Formalisiert versteht man unter DTW eine vorlagenbasierte Erkennungstechnik auf Grundlage von dynamischer Programmierung [LK99, S. 963]. Trotz der Erfolge bei Aufgaben mit kleinem Vokabular, braucht DTW eine große Anzahl an Vorlagen für ein grösseres Umfeld an Variation. Weiterhin kann es keine undefinierte Muster verarbeiten.

Takahashi et al. [NO96] haben einen sogenannten *spotting algorithm* vorgeschlagen Referenz [FB93] [KM12]

6.2 Künstliche neuronale Netze

6.3 Hidden Markov Model

Das HMM ist ein stochastisches Modell, in dem ein System durch eine Markov-Kette mit unbeobachteten Zuständen modelliert wird. Die Theorie dazu wurde von Baum [BP66], 1966 veröffentlicht.

Der versteckte (englisch *hidden*) Prozess ist eine Markov-Kette und besteht aus Zuständen und Übergangswahrscheinlichkeiten. Das Modell ist wie folgt definiert:

Definition

Ein HMM $\lambda = (S, A, O, B, \pi)$ ist gegeben durch

- $S = S_1, \dots, S_n$ – Menge aller Zustände,
- $A = a_{ij}$ – Übergangsmatrix zwischen den Zuständen, wobei a_{ij} die Wahrscheinlichkeit angibt, dass Zustand S_i in Zustand S_j gewechselt wird,
- $O = O_1, \dots, O_m$ – Menge der möglichen Beobachtungen (*Emissionen*),
- $B = b_{ij}$ – Beobachtungsmatrix, wobei b_{ij} die Wahrscheinlichkeit angibt, im Zustand S_i die Beobachtung $O_j \in O$
- π – Anfangswahrscheinlichkeitsverteilung mit $\pi(i)$ Wahrscheinlichkeit, dass S_i der Startzustand ist.

6.3.1 Auswahl des Modells für die Arbeit

Rabiner et al. [Rab89, S. 257f], sowie Lee und Kim [LK99, S. 961] sprechen von den guten Erfolgen des HMM als stochastisches Werkzeug zur Modellierung von Gesten. Daher wird für diese Arbeit eine Implementierung mittels HMM erarbeitet.

6.3.2 Diskrete Markov-Prozesse

Eine Markov-Kette ist ein Stochastischer Prozess. Dabei unterscheidet man zwischen Markov-Ketten in diskreter und in stetiger Zeit. Weiter spricht man von Markov-Ketten, sofern ein diskreter Zustandsraum vorhanden ist, von *Markov-Prozessen* im Allgemeinen wenn ein stetiger Zustandsraum vorliegt.

Bei einem HMM werden ausschließlich diskrete Markov-Ketten verwendet, die wie folgt definiert sind.

Definition

Gegeben sei ein System, das zu jeder gegebenen Zeit t beschrieben werden kann, als sich in einer Menge von N verschiedenen Zuständen S_1, S_2, \dots, S_N befindend. Abbildung 6.1 zeigt eine solche Markov-Kette mit fünf Zuständen S_1 bis S_5 . In regelmäßigen diskreten Zeitabständen nimmt das System eine Zustandsänderung vor, gemäß einer Wahrscheinlichkeitsmenge, die mit dem Zustand verknüpft ist. Die Zeitinstanzen, die mit den Zustandsänderungen verknüpft sind, werden durch $t = 1, 2, \dots$ gekennzeichnet und der aktuelle Zustand zum Zeitpunkt t als q_t . Für den Fall, einer diskreten Markov-Kette erster Ordnung ist diese Charakterisierung ausreichend, da

$$P[q_t = S_j | q_{t-1} = S_i, q_{t-2} = S_k, \dots] = P[q_t = S_j | q_{t-1} = S_i]. \quad (6.1)$$

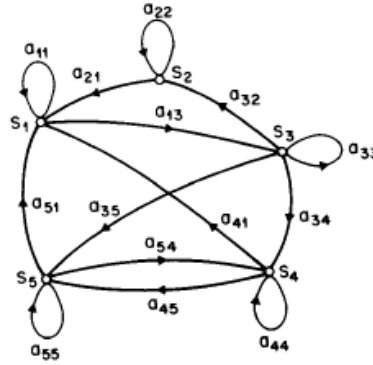


Abbildung 6.1: Markov-Kette mit fünf Zuständen (S_1 bis S_5) mit Verbindungen (Quelle: [Rab89])

Weiterhin werden ausschließlich die Prozesse betrachtet, in denen die rechte Seite der Gleichung 6.1 zeitunabhängig ist, was somit zur Menge der Übergangswahrscheinlichkeiten a_{ij} der Form

$$a_{ij} = P[q_t = S_j | q_{t-1} = S_i], 1 \leq i, j \leq N \quad (6.2)$$

führt, wobei die Zustandsübergangskoeffizienten folgende Eigenschaften erfüllen:

$$a_{ij} \geq 0 \quad (6.3a)$$

$$\sum_{j=1}^N a_{ij} = 1 \quad (6.3b)$$

Ein solcher Stochastischer Prozess stellt aber noch kein HMM dar.

Um ein besseres Verständnis für diese Prozesse und deren Eigenschaften zu erlangen, wird im folgenden anhand eines Beispiels aus einer Anleitung von Lawrence Rabiner [Rab89] näher auf das Thema eingegangen.

Man betrachte ein einfaches Markov Modell, dass mittels drei Stati das Wetter beschreibt. Diese drei Zustände sind folgende:

Status 1: Regnerisch

Status 2: Bewölkt

Status 3: Sonnig

Es wird gefordert, dass das Wetter am Tag t durch einen der drei oberen Zustände beschrieben wird. Weiter wird die Matrix A der Wahrscheinlichkeiten der Zustandsübergänge wie folgt definiert:

$$\mathbf{A} = a_{ij} = \begin{pmatrix} 0,4 & 0,3 & 0,3 \\ 0,2 & 0,6 & 0,2 \\ 0,1 & 0,1 & 0,8 \end{pmatrix}$$

Das Wetter am ersten Tag ($t = 1$) sei als sonnig (Status 3) gegeben. Stellt man nun die Frage, nach der Wahrscheinlichkeit für eine Folge von Wetterzuständen, oder anders formuliert, für einen Wetterverlauf der nächsten sieben Tage, wie „sonnig, sonnig, regnerisch, regnerisch, sonnig, bewölkt, sonnig...“, so kann man hierzu formal eine Sequenz aus mehreren Zuständen betrachten. Hierzu sei die Beobachtungssequenz O definiert als $O = S_3, S_3, S_1, S_1, S_3, S_2, S_3$, was mit $t = 1, 2, \dots, 8$ übereinstimmt. Die Wahrscheinlichkeit kann ausgedrückt werden als,

$$\begin{aligned} P(O|\text{Modell}) &= P[S_3, S_3, S_3, S_1, S_1, S_3, S_2, S_3|\text{Modell}] \\ &= P[S_3] \cdot P[S_3|S_3] \cdot P[S_3|S_3] \cdot P[S_1|S_3] \cdot P[S_1|S_1] \cdot P[S_3|S_1] \cdot P[S_2|S_3] \cdot P[S_3|S_2] \\ &= \pi_3 \cdot a_{33} \cdot a_{33} \cdot a_{33} \cdot a_{31} \cdot a_{11} \cdot a_{13} \cdot a_{32} \cdot a_{23} \\ &= 1 \cdot (0,8)(0,8)(0,1)(0,4)(0,3)(0,1)(0,2) \\ &= 1,536 \times 10^{-4}, \end{aligned}$$

wobei

$$\pi_i = P[q_1 = S_i], \quad 1 \leq i \leq N \quad (6.4)$$

die Wahrscheinlichkeit für den initialen Status darstellt.

Auch die Frage nach der Wahrscheinlichkeit wie lange das Modell in einem gegebenen

und bekannten Status genau d Tage verbleibt, kann mittels HMM beantwortert werden. Dabei kann die Wahrscheinlichkeit durch die Beobachtungssequenz

$$O = S_{i_1}, S_{i_2}, S_{i_3}, \dots, S_{i_d}, S_{j_{d+1}} \neq S_i$$

ermittelt werden:

$$P(O|Modell, q_1 = S_i) = (a_{ii})^{d-1}(1 - a_{ii}) = p_i(d) \quad (6.5)$$

Die Größe $p_i(d)$ ist die (diskrete) Wahrscheinlichkeitsdichtefunktion der Dauer d im Status i . Diese ist charakteristisch für die Dauer eines Status in einer Markov-Kette. Basierend auf $p_i(d)$ kann man die erwartete Zahl der Betrachtungen (Dauer) in einem Status berechnen, unter der Voraussetzung, dass in diesem Status gestartet wird, mit

$$\bar{d}_i = \sum_{d=1}^{\infty} d p_i(d) \quad (6.6a)$$

$$= \sum_{d=1}^{\infty} d (a_{ii})^{d-1} (1 - a_{ii}) = \frac{1}{1 - a_{ii}}. \quad (6.6b)$$

Daraus ergeben sich die Anzahl aufeinanderfolgender Sonnentage, gemäß dem Modell $\frac{1}{0,2} = 5$, die Anzahl aufeinanderfolgender bewölkter Tage entsprechend mit 2,5 und die Anzahl aufeinanderfolgender Regentage mit 1,67.

6.3.3 Erweiterung auf Hidden Markov Model

Um Markov Prozesse auf Probleme anwenden zu können, in denen sich die Beobachtungen nicht auf Zustand beziehen, muss das bisher dargestellte Konzept erweitert werden. Dazu wird der Fall hinzugefügt, dass eine Beobachtung durch eine Wahrscheinlichkeitsfunktion des Zustandes ausgedrückt wird. Nach Rabiner [Rab89] besteht das resultierende Modell aus einem doppelt eingebetteten stochastischen Prozess mit einem darunterliegenden, nicht beobachtbaren (versteckten, daher *hidden*), stochastischen Prozess und

wird als HMM bezeichnet.

Ein HMM ist wie folgt charakterisiert:

1. Sei N die Anzahl aller Zustände im Modell, dann sind die verschiedenen Zustände definiert als $S = S_1, S_2, \dots, S_N$ und der Zustand zum Zeitpunkt t ist definiert als q_t .
2. Sei M die Anzahl aller verschiedenen Beobachtungssymbole pro Zustand, dann sind die individuellen Symbole definiert als $V = v_1, v_2, \dots, v_M$.
3. Sei $A = a_{ij}$ die Übergangsmatrix, wobei

$$a_{ij} = P[q_{t+1} = S_j \mid q_t = S_i], \quad 1 \leq i, j \leq N. \quad (6.7)$$

Für den speziellen Fall, dass jeder beliebige Zustand jeden beliebigen anderen Zustand in einem Schritt erreichen kann, ist $a_{ij} > 0$ für alle i, j . Für alle weiteren Typen eines HMM ist $a_{ij} = 0$ für mindestens ein (i, j) Paar.

4. Sei $B = b_j(k)$ die Beobachtungsmatrix im Zustand j , wobei

$$b_j(k) = P[v_k \text{ zum Zeitpunkt } t \mid q_t = S_j], \quad \begin{aligned} 1 \leq j \leq N \\ 1 \leq k \leq M. \end{aligned} \quad (6.8)$$

5. Sei $\pi = \pi_i$ die Anfangswahrscheinlichkeitsverteilung, wobei

$$\pi_i = P[q_1 = S_i], \quad 1 \leq i \leq N. \quad (6.9)$$

Sind entsprechende Werte für N, M, A, B und π gesetzt, so kann das HMM als Generator verwendet werden, um eine Beobachtungssequenz

$$O = O_1 O_2 \dots O_T \quad (6.10)$$

wie folgt zu erhalten:

1. Wähle einen Startzustand $q_1 = S_i$ entsprechend der Anfangswahrscheinlichkeitsverteilung π .
2. Setze $t = 1$.
3. Wähle $O_t = v_k$ entsprechend der Beobachtungsmatrix für den Zustand S_i , zum Beispiel $b_i(k)$.
4. Gehe in einen neuen Zustand $q_{t+1} = S_j$ entsprechend der Beobachtungsmatrix für den Zustand S_i über, zum Beispiel a_{ij} .
5. Setze $t = t + 1$; kehre zu Schritt 3 zurück, falls $t < T$; beende andernfalls die Prozedur.

Nach Rabiner [Rab89, S. 5] kann diese Prozedur sowohl dazu verwendet werden, Beobachtungen zu generieren und als Modell, wie eine gegebene Beobachtungssequenz von einem entsprechendem HMM erstellt wurde.

Eine vollständige Spezifikation eines HMM benötigt die Beschreibung zweier Modellparameter (N und M), die Spezifikation von Beobachtungssymbolen und die Beschreibung der drei Wahrscheinlichkeitsgrößen A, B und π . Meist wird hierfür die kompakte Notation

$$\lambda = (A, B, \pi) \tag{6.11}$$

verwendet, um die Vollständigkeit aller Parameter des Modells anzuzeigen.

6.3.4 Die drei grundlegenden Probleme eines Hidden Markov Model ¹

Verwendet man die oben Form eines HMM, die in Abschnitt refsubsec:HMM beschrieben ist, so ergeben sich drei grundlegende Probleme, die vor einem Einsatz von HMM in einer Anwendung gelöst werden müssen. Diese lauten wie folgt:

¹Die Ausführungen basieren auf den Vorträgen von Jack Ferguson von IDA an den Bell Laboratories

6.4 Baum-Weich-Algorithmus

6.4.1 ToDo

6.5 ToDo

7 Gesten

7.1 Kreisbewegung

7.1.1 ToDo

7.2 Vorwärtsbewegung

7.2.1 ToDo

7.3 Haltesignal

7.3.1 ToDo

8 Sprachbefehle

8.1 ToDo

9 Implementierung

9.1 OSGI–Bundles

9.1.1 Gesten

9.1.2 Sprachbefehle

9.1.3 Roboterschnittstelle

9.1.4 Ausführbare Aktionen

9.2 Oberfläche

9.2.1 RCP–Anwendung

9.2.2 Grafische Darstellung

9.3 Ablaufbeschreibung

10 Ausblick - weitere Arbeiten

Abbildungsverzeichnis

3.1	Schematische Ansicht der Sensoren einer Kinect for Windows	13
6.1	Markov-Kette mit fünf Zuständen	27

Tabellenverzeichnis

3.1 Technische Spezifikation der Kinect for Windows	14
---	----

Quellcodeverzeichnis

Literaturverzeichnis

- [BP66] E. Leonard Baum and Ted Petrie. Statistical inference for probabilistic functions of finite state markov chains. *The Annals of Mathematical Statistics*, Vol. 37, No. 6:1554–1563, 1966.
- [Cor12] Microsoft Corp., editor. *HUMAN INTERFACE GUIDELINES*, volume v1.5.0. Microsoft Corp., 2012.
- [FB93] Janet Finlay and Russell Beale. Neural networks and pattern recognition in human-computer interaction. *SIGCHI Bull.*, 25(2):25–35, April 1993.
- [KM12] Matthew Kober, Jens; Glisson and Michael Mistry. Playing catch and juggling with a humanoid robot. 2012.
- [LK99] Hyeon-Kyu Lee and Jin H. Kim. An hmm-based threshold model approach for gesture recognition. *IEEE Trans. Pattern Anal. Mach. Intell.*, 21(10):961–973, 1999.
- [NO96] T. Nishimura and R. Oka. Spotting recognition of human gestures from time-varying images. In *Automatic Face and Gesture Recognition, 1996., Proceedings of the Second International Conference on*, pages 318 –322, oct 1996.
- [Rab89] Lawrence R. Rabiner. A tutorial on hidden markov models and selected applications in speech recognition. *Proceedings of the IEEE*, Vol. 77, No. 2:257–286, 1989.

- [ZS09] Chun Zhu and Weihua Sheng. Online hand gesture recognition using neural network based segmentation. *Intelligent Robots and Systems, 2009. IROS 2009. IEEE/RSJ International Conference on*, 1:2415–2420, 2009.

Glossar

Bewegungsdetektion

Unter Bewegungsdetektion versteht man Methoden des Maschinellen Sehens, die auf die Erkennung von Fremdbewegung im Erfassungsbereich eines optischen Detektors (also dem Blickfeld der Maschine) abzielen. 8

Datenhandschuh

Der Datenhandschuh ist ein Eingabegerät in Form eines Handschuhs. Durch Bewegungen der Hand und Finger erfolgt eine Orientierung im virtuellen Raum. 7, 8

Eclipse

Eclipse (von englisch eclipse ‚Sonnenfinsternis‘, ‚Finsternis‘, ‚Verdunkelung‘) ist ein quelloffenes Programmierwerkzeug zur Entwicklung von Software verschiedenster Art. Ursprünglich wurde Eclipse als integrierte Entwicklungsumgebung (IDE) für die Programmiersprache Java genutzt, aber mittlerweile wird es wegen seiner Erweiterbarkeit auch für viele andere Entwicklungsaufgaben eingesetzt. Für Eclipse gibt es eine Vielzahl sowohl quelloffener als auch kommerzieller Erweiterungen. 16

Framework

Ein Framework (englisch für Rahmenstruktur) ist ein Programmiergerüst, das in der Softwaretechnik, insbesondere im Rahmen der objektorientierten Softwareentwicklung sowie bei komponentenbasierten Entwicklungsansätzen, verwendet wird. 17

Kinect

Kinect ist eine Plattform zur Steuerung der Videospielkonsole Xbox 360, die seit Anfang November 2010 verkauft wird. 8

Mensch-Computer-Interaktion

Die Mensch-Computer-Interaktion (englisch *Human-Computer Interaction, HCI*) als Teilgebiet der Informatik beschäftigt sich mit der benutzergerechten Gestaltung von interaktiven Systemen und ihren Mensch-Maschine-Schnittstellen. 7

Mensch-Maschine-Schnittstelle

Die Benutzerschnittstelle (nach Gesellschaft für Informatik, Fachbereich Mensch-Computer-Interaktion auch Benutzungsschnittstelle) ist die Stelle oder Handlung, mit der ein Mensch mit einer Maschine in Kontakt tritt. 7

Mikrofonarray

Beamforming ist ein Verfahren zur Positionsbestimmung von Quellen in Wellenfeldern (z. B. Schallfeldern). Entsprechende Vorrichtungen werden auch akustische Kamera, Mikrofonarray oder akustische Antenne genannt. 13, 14

Open Graphics Library

OpenGL (Open Graphics Library) ist eine Spezifikation für eine plattform- und programmiersprachenunabhängige Programmierschnittstelle zur Entwicklung von

2D- und 3D-Computergrafik. Der OpenGL-Standard beschreibt etwa 250 Befehle, die die Darstellung komplexer 3D-Szenen in Echtzeit erlauben. 18

Programmierschnittstelle

Eine Programmierschnittstelle (englisch application programming interface, API; deutsch Schnittstelle zur Anwendungsprogrammierung) ist ein Programmteil, der von einem Softwaresystem anderen Programmen zur Anbindung an das System zur Verfügung gestellt wird, plural=API. 17

Skelett Tracking System

Unter Motion Capture, oder auch Skelett Tracking System, wörtlich Bewegungs-Erfassung, versteht man ein Tracking-Verfahren, das es ermöglicht, jede Art von Bewegungen so zu erfassen und in ein von Computern lesbares Format umzuwandeln, dass diese die Bewegungen analysieren, aufzeichnen, weiterverarbeiten und zur Steuerung von Anwendungen verwenden können. 14

Software Development Kit

Ein Software Development Kit (SDK) ist eine Sammlung von Werkzeugen und Anwendungen, um eine Software zu erstellen, meist inklusive Dokumentation. Mit diesem ist es Softwareentwicklern möglich, eigene darauf basierende Anwendungen zu erstellen. 11, 15