



mitp

Ian
Goodfellow

Yoshua
Bengio

Aaron
Courville

Deep Learning

Das umfassende Handbuch

Grundlagen, aktuelle Verfahren und
Algorithmen, neue Forschungsansätze



Hinweis des Verlages zum Urheberrecht und Digitalen Rechtemanagement (DRM)

Der Verlag räumt Ihnen mit dem Kauf des ebooks das Recht ein, die Inhalte im Rahmen des geltenden Urheberrechts zu nutzen. Dieses Werk, einschließlich aller seiner Teile, ist urheberrechtlich geschützt. Jede Verwertung außerhalb der engen Grenzen des Urheberrechtsgesetzes ist ohne Zustimmung des Verlages unzulässig und strafbar. Dies gilt insbesondere für Vervielfältigungen, Übersetzungen, Mikroverfilmungen und Einspeicherung und Verarbeitung in elektronischen Systemen.

Der Verlag schützt seine ebooks vor Missbrauch des Urheberrechts durch ein digitales Rechtemanagement. Bei Kauf im Webshop des Verlages werden die ebooks mit einem nicht sichtbaren digitalen Wasserzeichen individuell pro Nutzer signiert.

Bei Kauf in anderen ebook-Webshops erfolgt die Signatur durch die Shopbetreiber. Angaben zu diesem DRM finden Sie auf den Seiten der jeweiligen Anbieter.

Ian Goodfellow
Yoshua Bengio
Aaron Courville

Deep Learning

Das umfassende Handbuch

Grundlagen, aktuelle Verfahren und
Algorithmen, neue Forschungsansätze

Übersetzung aus dem Amerikanischen
von Guido Lenz



mitp

Bibliografische Information der Deutschen Nationalbibliothek

Die Deutsche Nationalbibliothek verzeichnet diese Publikation in der Deutschen Nationalbibliografie. Detaillierte bibliografische Daten sind im Internet über <http://dnb.d-nb.de> abrufbar.

ISBN 978-3-95845-701-0

1. Auflage 2018

www.mitp.de

E-Mail: mitp-verlag@sigloch.de

Telefon: +49 7953 7189 - 079

Telefax: +49 7953 7189 - 082

© 2018 mitp Verlags GmbH & Co. KG, Frechen

Dieses Werk, einschließlich aller seiner Teile, ist urheberrechtlich geschützt. Jede Verwertung außerhalb der engen Grenzen des Urheberrechtsgesetzes ist ohne Zustimmung des Verlages unzulässig und strafbar. Dies gilt insbesondere für Vervielfältigungen, Übersetzungen, Mikroverfilmungen und die Einspeicherung und Verarbeitung in elektronischen Systemen.

Die Wiedergabe von Gebrauchsnamen, Handelsnamen, Warenbezeichnungen usw. in diesem Werk berechtigt auch ohne besondere Kennzeichnung nicht zu der Annahme, dass solche Namen im Sinne der Warenzeichen- und Markenschutz-Gesetzgebung als frei zu betrachten wären und daher von jedermann benutzt werden dürften.

Authorized German translation from the English language edition, entitled Deep Learning by Ian Goodfellow, Yoshua Bengio, and Aaron Courville
ISBN 9780262035613

Original English language edition published by The MIT Press, a department of Massachusetts Institute of Technology

Copyright © 2016 Massachusetts Institute of Technology

German-language edition copyright © 2018 by mitp-Verlag. All rights reserved.

Lektorat und fachliche Prüfung: Sabine Schulz, Janina Bahlmann, Lisa Kresse

Sprachkorrektorat: Petra Heubach-Erdmann

Fachkorrektur Mathematik: Stefan Niedergriese

Fachkorrektur Deep Learning: Alexander Bresk

Coverbild: madgooch/stock.adobe.com

Satz: Dr. Joachim Schlosser, www.schlosser.info

Inhaltsverzeichnis

Website zum Buch	xi
Danksagung	xiii
Über die Fachkorrektoren zur deutschen Ausgabe	xvii
Notation	xix
1 Einleitung	1
1.1 Für wen ist dieses Buch gedacht?	10
1.2 Historische Entwicklungen im Deep Learning	12
I Angewandte Mathematik und Grundlagen für das Machine Learning	31
2 Lineare Algebra	33
2.1 Skalare, Vektoren, Matrizen und Tensoren	34
2.2 Multiplizieren von Matrizen und Vektoren	36
2.3 Einheits- und Umkehrmatrizen	38
2.4 Lineare Abhängigkeit und lineare Hülle	40
2.5 Normen	42
2.6 Spezielle Matrizen und Vektoren	43
2.7 Eigenwertzerlegung	45
2.8 Singulärwertzerlegung	48
2.9 Die Moore-Penrose-Pseudoinverse	49
2.10 Der Spuroperator	50
2.11 Die Determinante	51
2.12 Beispiel: Hauptkomponentenanalyse	51
3 Wahrscheinlichkeits- und Informationstheorie	57
3.1 Warum Wahrscheinlichkeit?	58

3.2	Zufallsvariablen	61
3.3	Wahrscheinlichkeitsverteilungen	61
3.4	Randwahrscheinlichkeit	64
3.5	Bedingte Wahrscheinlichkeit	64
3.6	Die Produktregel der bedingten Wahrscheinlichkeiten	65
3.7	Unabhängigkeit und bedingte Unabhängigkeit	65
3.8	Erwartungswert, Varianz und Kovarianz	66
3.9	Häufig genutzte Wahrscheinlichkeitsverteilungen	67
3.10	Nützliche Eigenschaften häufig verwendeter Funktionen	74
3.11	Satz von Bayes	76
3.12	Technische Einzelheiten stetiger Variablen	77
3.13	Informationstheorie	79
3.14	Strukturierte probabilistische Modelle	83
4	Numerische Berechnung	87
4.1	Überlauf und Unterlauf	87
4.2	Schlechte Konditionierung	89
4.3	Optimierung auf Gradientenbasis	90
4.4	Optimierung unter Nebenbedingungen	101
4.5	Beispiel: Lineare kleinste Quadrate	104
5	Grundlagen für das Machine Learning	107
5.1	Lernalgorithmen	108
5.2	Kapazität, Überanpassung und Unteranpassung	121
5.3	Hyperparameter und Validierungsdaten	133
5.4	Schätzer, Verzerrung und Varianz	135
5.5	Maximum-Likelihood-Schätzung	145
5.6	Bayessche Statistik	149
5.7	Algorithmen für überwachtetes Lernen	154
5.8	Algorithmen für unüberwachtetes Lernen	161
5.9	Stochastisches Gradientenabstiegsverfahren	167
5.10	Entwickeln eines Machine-Learning-Algorithmus	170
5.11	Probleme, an denen Deep Learning wächst	171
II	Tiefe Netze: Zeitgemäße Verfahren	183
6	Tiefe Feedforward-Netze	185
6.1	Beispiel: Erlernen von XOR	189
6.2	Lernen auf Gradientenbasis	195
6.3	Verdeckte Einheiten	211

6.4	Architekturdesign	218
6.5	Backpropagation und andere Algorithmen zur Differentiation	225
6.6	Historische Anmerkungen	248
7	Regularisierung	253
7.1	Parameter-Norm-Strafterme	255
7.2	Norm-Strafterme als Optimierung unter Nebenbedingungen	263
7.3	Regularisierung und unterbestimmte Probleme	265
7.4	Erweitern des Datensatzes	266
7.5	Robustheit gegen Rauschen	268
7.6	Halb-überwachtes Lernen	270
7.7	Multitask Learning	271
7.8	Früher Abbruch	273
7.9	Parameter Tying und Parameter Sharing	281
7.10	Dünnbesetzte Repräsentationen	283
7.11	Bagging und andere Ensemblemethoden	285
7.12	Dropout	287
7.13	Adversarial Training	299
7.14	Tangentendistanz, Tangenten-Propagation und Mannigfaltigkeit-Tangentenklassifikator	301
8	Optimierung beim Trainieren von tiefen Modellen	305
8.1	Unterschied zwischen Lernen und reiner Optimierung . . .	306
8.2	Herausforderungen bei der Optimierung neuronaler Netze	315
8.3	Grundlegende Algorithmen	327
8.4	Verfahren zur Parameterinitialisierung	335
8.5	Algorithmen mit adaptiven Lernraten	342
8.6	Approximative Verfahren zweiter Ordnung	347
8.7	Optimierungsverfahren und Meta-Algorithmen	354
9	CNNs	369
9.1	Die Faltungsoperation	370
9.2	Motivation	374
9.3	Pooling	379
9.4	Faltung und Pooling als unendlich starke A-priori- Wahrscheinlichkeit	385
9.5	Varianten der grundlegenden Faltungsfunktion	386
9.6	Strukturierte Ausgaben	398
9.7	Datentypen	400
9.8	Effiziente Faltungsalgorithmen	402

9.9	Zufällige oder unüberwachte Merkmale	403
9.10	Die neurowissenschaftliche Basis für CNNs	405
9.11	CNNs und die Geschichte des Deep Learnings	413
10	Sequenzmodellierung: RNNs und rekursive Netze	415
10.1	Auffalten von Berechnungsgraphen	417
10.2	RNNs	420
10.3	Bidirektionale RNNs	436
10.4	Sequenz-zu-Sequenz-Architekturen	439
10.5	Tiefe RNNs	441
10.6	Rekursive neuronale Netze	443
10.7	Die Herausforderung langfristiger Abhängigkeiten	445
10.8	Echo-State-Netze	448
10.9	Leaky-Einheiten und andere Verfahren für mehrere Zeitskalen	451
10.10	Das Long Short-Term Memory und andere Gated RNNs .	453
10.11	Optimierung für langfristige Abhängigkeiten	458
10.12	Explizites Gedächtnis	462
11	Praxisorientierte Methodologie	467
11.1	Performance-Kriterien	468
11.2	Default-Baseline-Modell	471
11.3	Prüfen, ob mehr Daten gesammelt werden sollen	473
11.4	Auswählen von Hyperparametern	475
11.5	Debugging-Verfahren	485
11.6	Beispiel: Erkennen mehrstelliger Zahlen	490
12	Anwendungen	493
12.1	Deep Learning im großen Maßstab	493
12.2	Computer Vision	504
12.3	Spracherkennung	511
12.4	Verarbeitung natürlicher Sprache	514
12.5	Weitere Anwendungen	533
III	Deep-Learning-Forschung	543
13	Lineare Faktorenmodelle	547
13.1	Probabilistische PCA und Faktorenanalyse	548
13.2	Unabhängigkeitsanalyse	549
13.3	Slow Feature Analysis	553

13.4	Sparse Coding	556
13.5	Interpretation der Mannigfaltigkeit der PCA	560
14	Autoencoder	563
14.1	Untervollständige Autoencoder	564
14.2	Regularisierte Autoencoder	565
14.3	Repräsentationsleistung, Schichtgröße und Tiefe	570
14.4	Stochastische Encoder und Decoder	571
14.5	Denoising Autoencoder	573
14.6	Erlernen von Mannigfaltigkeiten mit Autoencodern	578
14.7	Contractive Autoencoder	584
14.8	Prädiktive dünnbesetzte Zerlegung	587
14.9	Anwendungen für Autoencoder	588
15	Representation Learning	591
15.1	Schichtweises unüberwachtes Pretraining mit Greedy- Algorithmen	593
15.2	Transfer Learning und Domänenadaption	602
15.3	Halb-überwachtes Separieren kausaler Faktoren	607
15.4	Verteilte Repräsentation	614
15.5	Exponentielle Verbesserungen durch Tiefe	621
15.6	Hinweise zum Aufdecken der zugrunde liegenden Ursachen	623
16	Strukturierte probabilistische Modelle für Deep Learning	627
16.1	Die Herausforderung der unstrukturierten Modellierung .	628
16.2	Verwenden von Graphen zum Beschreiben der Modellstruktur	633
16.3	Stichprobenentnahme aus graphischen Modellen	651
16.4	Vorteile der strukturierten Modellierung	653
16.5	Lernen anhand von Abhängigkeiten	654
16.6	Inferenz und approximative Inferenz	656
16.7	Der Deep-Learning-Ansatz für strukturierte probabilistische Modelle	657
17	Monte-Carlo-Verfahren	663
17.1	Stichprobenentnahme und Monte-Carlo-Verfahren	663
17.2	Importance Sampling	666
17.3	Markow-Ketten-Monte-Carlo-Verfahren	668
17.4	Gibbs-Sampling	673
17.5	Die Herausforderung, zwischen getrennten Modi zu mischen	674

18 Die Partitionsfunktion	681
18.1 Der Log-Likelihood-Gradient	682
18.2 Stochastische Maximum Likelihood und kontrastive Divergenz	684
18.3 Pseudo-Likelihood	692
18.4 Score Matching und Ratio Matching	695
18.5 Denoising Score Matching	697
18.6 Noise-Contrastive Estimation	698
18.7 Schätzen der Partitionsfunktion	701
19 Approximative Inferenz	711
19.1 Inferenz als Optimierung	713
19.2 Erwartungsmaximierung	714
19.3 MAP-Inferenz und Sparse Coding	716
19.4 Variational Inference und Variational Learning	718
19.5 Erlernte approximative Inferenz	733
20 Tiefe generative Modelle	737
20.1 Boltzmann-Maschinen	737
20.2 Restricted Boltzmann Machines	740
20.3 Deep-Belief-Netze	743
20.4 Deep Boltzmann Machines	747
20.5 Boltzmann-Maschinen für reellwertige Daten	762
20.6 Gefaltete Boltzmann-Maschinen	769
20.7 Boltzmann-Maschinen für strukturierte und sequenzielle Ausgaben	772
20.8 Weitere Boltzmann-Maschinen	773
20.9 Backpropagation durch Zufallsoperationen	775
20.10 Gerichtete generative Netze (Directed Generative Nets)	780
20.11 Ziehen von Stichproben aus Autoencodern	802
20.12 Generative stochastische Netze	806
20.13 Andere Generierungskonzepte	807
20.14 Bewerten von generativen Modellen	809
20.15 Schlussbemerkungen	812
Literaturverzeichnis	813
Abkürzungsverzeichnis	871
Index	875

Website zum Buch

Auf der Website des Verlags finden Sie Hinweise zur Übersetzung sowie Errata, sofern Fehler bekannt sind unter

www.mitp.de/700

Darüber hinaus gibt es zu diesem Buch eine englischsprachige Website. Sie enthält ergänzende Materialien auf Englisch, darunter Übungen, Vortragsfolien, Korrekturen und andere Ressourcen, die für Leser und Lehrkräfte hilfreich sind. Diese finden Sie unter

www.deeplearningbook.org

Danksagung

Dieses Buch wäre ohne die Hilfe vieler Menschen nicht entstanden.

Wir danken all denen, die unser Buchkonzept unterstützt und uns dabei geholfen haben, seinen Inhalt und Aufbau zu planen: Guillaume Alain, Kyunghyun Cho, Çağlar Gülçehre, David Krueger, Hugo Larochelle, Razvan Pascanu und Thomas Rohée.

Wir danken auch den Menschen, die uns Feedback zum Text selbst gegeben haben. Einige davon haben zu vielen Kapiteln Rückmeldung gegeben: Martín Abadi, Guillaume Alain, Ion Androutsopoulos, Fred Bertsch, Olexa Bilaniuk, Ufuk Can Biçici, Matko Bošnjak, John Boersma, Greg Brockman, Alexandre de Brébisson, Pierre Luc Carrier, Sarath Chandar, Pawel Chilinski, Mark Daoust, Oleg Dashevskii, Laurent Dinh, Stephan Dreseitl, Jim Fan, Miao Fan, Meire Fortunato, Frédéric Francis, Nando de Freitas, Çağlar Gülçehre, Jurgen Van Gael, Javier Alonso García, Jonathan Hunt, Gopi Jeyaram, Chingiz Kabytayev, Lukasz Kaiser, Varun Kanade, Asifullah Khan, Akiel Khan, John King, Diederik P. Kingma, Yann LeCun, Rudolf Mathey, Matías Mattamala, Abhinav Maurya, Kevin Murphy, Oleg Mürk, Roman Novak, Augustus Q. Odena, Simon Pavlik, Karl Pichotta, Eddie Pierce, Kari Pulli, Roussel Rahman, Tapani Raiko, Anurag Ranjan, Johannes Roith, Mihaela Rosca, Halis Sak, César Salgado, Grigory Sapunov, Yoshinori Sasaki, Mike Schuster, Julian Serban, Nir Shabat, Ken Shirriff, Andre Simpelo, David Slate, Scott Stanley, David Sussillo, Ilya Sutskever, Carles Gelada Sáez, Graham Taylor, Valentin Tolmer, Massimiliano Tomassoli, An Tran, Shubhendu Trivedi, Alexey Umnov, Vincent Vanhoucke, Marco Visentini-Scarzanella, Martin Vita, David Warde-Farley, Dustin Webb, Kelvin Xu, Wei Xue, Ke Yang, Li Yao, Zygmunt Zajac und Ozan Çağlayan.

Unser Dank geht auch an jene, deren Feedback einzelnen Kapiteln gewidmet war:

- Notation: Zhang Yuanhang
- Kapitel 1, Einleitung: Yusuf Akgul, Sebastien Bratieres, Samira Ebrahimi, Charlie Gorichanaz, Brendan Loudermilk, Eric Morris, Cosmin Pârvulescu und Alfredo Solano
- Kapitel 2, Lineare Algebra: Amjad Almahairi, Nikola Banić, Kevin Bennett, Philippe Castonguay, Oscar Chang, Eric Fosler-Lussier, Andrey Khalyavin, Sergey Oreshkov, István Petrás, Dennis Prangle, Thomas Rohée, Gitanjali Gulve Sehgal, Colby Toland, Alessandro Vitale und Bob Welland
- Kapitel 3, Wahrscheinlichkeits- und Informationstheorie: John Philip Anderson, Kai Arulkumaran, Vincent Dumoulin, Rui Fa, Stephan Gouws, Artem Oboturov, Antti Rasmus, Alexey Surkov und Volker Tresp
- Kapitel 4, Numerische Berechnung: Tran Lam AnIan Fischer und Hu Yuhuang
- Kapitel 5, Grundlagen für das Machine Learning: Dzmitry Bahdanau, Justin Domingue, Nikhil Garg, Makoto Otsuka, Bob Pepin, Philip Popien, Bharat Prabhakar, Emmanuel Rayner, Peter Shepard, Kee-Bong Song, Zheng Sun und Andy Wu
- Kapitel 6, Tiefe Feedforward-Netze: Uriel Berdugo, Fabrizio Bottarel, Elizabeth Burl, Ishan Durugkar, Jeff Hlywa, Jong Wook Kim, David Krueger, Aditya Kumar Praharaaj und Sten Sootla
- Kapitel 7, Regularisierung: Morten Kolbæk, Kshitij Lauria, Inkyu Lee, Sunil Mohan, Hai Phong Phan und Joshua Salisbury
- Kapitel 8, Optimierung beim Trainieren von tiefen Modellen: Marcel Ackermann, Peter Armitage, Rowel Atienza, Andrew Brock, Tegan Maharaj, James Martens, Mostafa Nategh, Kashif Rasul, Klaus Strobl und Nicholas Turner
- Kapitel 9, CNNs: Martín Arjovsky, Eugene Brevdo, Konstantin Divilov, Eric Jensen, Mehdi Mirza, Alex Paino, Marjorie Sayer, Ryan Stout und Wentao Wu

- Kapitel 10, Sequenzmodellierung: RNNs und rekursive Netze: Gökçen Eraslan, Steven Hickson, Razvan Pascanu, Lorenzo von Ritter, Rui Rodrigues, Dmitriy Serdyuk, Dongyu Shi und Kaiyu Yang
- Kapitel 11, Praxisorientierte Methodologie: Daniel Beckstein
- Kapitel 12, Anwendungen: George Dahl, Vladimir Nekrasov und Ribana Roscher
- Kapitel 13, Lineare Faktorenmodelle: Jayanth Koushik
- Kapitel 15, Representation Learning: Kunal Ghosh
- Kapitel 16, Strukturierte probabilistische Modelle für Deep Learning: Minh Lê und Anton Varfolom
- Kapitel 18, Die Partitionsfunktion: Sam Bowman
- Kapitel 19, Approximative Inferenz: Yujia Bao
- Kapitel 20, Tiefe generative Modelle: Nicolas Chapados, Daniel Galvez, Wenming Ma, Fady Medhat, Shakir Mohamed und Grégoire Montavon
- Bibliographie: Lukas Michelbacher und Leslie N. Smith

Wir danken auch allen, die uns erlaubt haben, Abbildungen, Grafiken oder Daten aus ihren Werken zu verwenden. Wir haben diese Wiedergabe in den Bildunterschriften angemerkt.

Wir danken Lu Wang dafür, dass er pdf2htmlEX programmiert hat, mit dem wir die Webversion des Buchs erstellt haben, sowie für sein Angebot, uns bei der Verbesserung der Qualität des HTML-Codes zu unterstützen.

Wir danken Ians Ehefrau Daniela Flori Goodfellow für ihre Geduld, die sie Ian während der Abfassung des Buchs entgegengebracht hat. Und natürlich für ihre Hilfe beim Korrektorat.

Wir möchten dem Team hinter Google Brain dafür danken, dass es eine geeignete Umgebung geschaffen hat, in der wir so viel Zeit für die Arbeit an diesem Buch verbringen und gleichzeitig Feedback und Hilfe von den Kollegen annehmen konnten. Ganz besonders möchten wir Ians ehemaligem Vorgesetzten Greg Corrado und seinem derzeitigen Vorgesetzten Samy Bengio für ihre Unterstützung bei diesem Projekt danken. Und schließlich gebührt unser Dank Geoffrey Hinton, der uns bei Schwierigkeiten stets ermuntert hat, weiterzumachen.

Über die Fachkorrektoren zur deutschen Ausgabe

Fachkorrektor für den Bereich Deep Learning

Alexander Bresk ist Unternehmer, Data Engineer, Machine Teacher, Autor und Basketball-Nerd. Er arbeitet als Senior Data Engineer bei der LOVOO GmbH. Dort beschäftigt er sich mit den Themen Recommendation und Online Segmentation. Seinen Master of Science erwarb Alexander an der HTW Dresden, wo er Angewandte Informationstechnologien studierte. Seit seinem Master im Bereich des Question Answering forscht er aktiv an Verfahren zur Verarbeitung natürlicher Sprache (NLP), wobei er unter anderem Deep Learning Frameworks einsetzt. Außerdem organisiert er zusammen mit Kollegen und Freunden Meetups, Konferenzen und andere Veranstaltungen in den Bereichen Tech, Machine Learning und Startups. Mit seiner Firma Machine Rockstars berät er Unternehmen beim Einstieg in das Thema Digitalisierung sowie bei der Einführung von Machine Learning.

Alexander Bresks Website: <http://alexander.bre.sk>

Fachkorrektor für den Bereich Mathematik

Stefan Niedergriese hat Mathematik studiert und arbeitet bei einem Versicherungsunternehmen in der mathematischen Abteilung.

Notation

Dieser Abschnitt dient als Referenz. Er listet die im gesamten Buch verwendete Notation auf. Falls Ihnen einzelne der mathematischen Konzepte unbekannt sind, sollten Sie die Kapitel 2 bis 4 lesen – dort gehen wir auf die meisten davon näher ein.

Zahlen und Vektoren

a	Ein Skalar (ganzzahlig oder reell)
\mathbf{a}	Ein Vektor
\mathbf{A}	Eine Matrix
\mathbf{A}	Ein Tensor
\mathbf{I}_n	Einheitsmatrix mit n Zeilen und n Spalten
\mathbf{I}	Einheitsmatrix mit aus dem Kontext ersichtlicher Dimensionalität
$\mathbf{e}^{(i)}$	Standardbasisvektor $[0, \dots, 0, 1, 0, \dots, 0]$ mit einer 1 an der Stelle i
$\text{diag}(\mathbf{a})$	Eine quadratische Diagonalmatrix mit Diagonaleinträgen aus \mathbf{a}
a	Eine skalare Zufallsvariable
\mathbf{a}	Eine vektorwertige Zufallsvariable
\mathbf{A}	Eine matrixwertige Zufallsvariable

Mengen und Graphen

\mathbb{A}	Eine Menge
\mathbb{R}	Die Menge der reellen Zahlen

$\{0, 1\}$	Die Menge, die 0 und 1 enthält
$\{0, 1, \dots, n\}$	Die Menge aller ganzen Zahlen zwischen 0 und n
$[a, b]$	Das reellwertige Intervall der Mengen a und b
(a, b)	Das reellwertige Intervall ohne a , jedoch mit b
$\mathbb{A} \setminus \mathbb{B}$	Mengendifferenz (Komplement), also die Menge der Elemente aus \mathbb{A} , die nicht in \mathbb{B} enthalten sind
\mathcal{G}	Ein Graph
$Pa_{\mathcal{G}}(x_i)$	Die Eltern von x_i in \mathcal{G}

Indizes

a_i	Element i des Vektors \mathbf{a} , die Indexmenge beginnt mit 1
a_{-i}	Alle Elemente des Vektors \mathbf{a} bis auf das Element i
$A_{i,j}$	Element i, j der Matrix \mathbf{A}
$\mathbf{A}_{i,:}$	Zeile i der Matrix \mathbf{A}
$\mathbf{A}_{:,i}$	Spalte i der Matrix \mathbf{A}
$A_{i,j,k}$	Element (i, j, k) eines 3-D-Tensors \mathbf{A}
$\mathbf{A}_{:, :, i}$	2-D-Schnitt eines 3-D-Tensors
a_i	Element i des Zufallsvektors \mathbf{a}

Lineare Algebra

\mathbf{A}^{\top}	Transponierte der Matrix \mathbf{A}
\mathbf{A}^+	Moore-Penrose-Pseudoinverse von \mathbf{A}
$\mathbf{A} \odot \mathbf{B}$	Elementweises Produkt (Hadamard-Produkt) aus \mathbf{A} und \mathbf{B}
$\det(\mathbf{A})$	Determinante von \mathbf{A}

Analysis

$\frac{dy}{dx}$	Ableitung von y bezüglich x
$\frac{\partial y}{\partial x}$	Partielle Ableitung von y bezüglich x

$\nabla_{\mathbf{x}}y$	Gradient von y bezüglich \mathbf{x}
$\nabla_{\mathbf{X}}y$	Matrixableitungen von y bezüglich \mathbf{X}
$\nabla_{\mathbf{X}}y$	Tensor mit Ableitungen von y bezüglich \mathbf{X}
$\frac{\partial f}{\partial \mathbf{x}}$	Jacobi-Matrix $\mathbf{J} \in \mathbb{R}^{m \times n}$ von $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$
$\nabla_{\mathbf{x}}^2 f(\mathbf{x})$ oder $\mathbf{H}(f)(\mathbf{x})$	Hesse-Matrix von f im Eingabepunkt \mathbf{x}
$\int f(\mathbf{x})d\mathbf{x}$	Bestimmtes Integral über gesamten Definitionsbereich \mathbf{x}
$\int_{\mathbb{S}} f(\mathbf{x})d\mathbf{x}$	Bestimmtes Integral bezüglich \mathbf{x} über Menge \mathbb{S}

Wahrscheinlichkeits- und Informationstheorie

$a \perp b$	Die Zufallsvariablen a und b sind unabhängig
$a \perp b \mid c$	Sie sind bedingt unabhängig, wenn c zutrifft
$P(a)$	Eine Wahrscheinlichkeitsverteilung über eine diskrete Variable
$p(a)$	Eine Wahrscheinlichkeitsverteilung über eine stetige Variable oder eine Variable, deren Typ nicht angegeben wurde
$a \sim P$	Zufallsvariable a mit der Verteilung P
$\mathbb{E}_{\mathbf{x} \sim P}[f(\mathbf{x})]$ oder $\mathbb{E}f(\mathbf{x})$	Erwartung für $f(\mathbf{x})$ bezüglich $P(\mathbf{x})$
$\text{Var}(f(\mathbf{x}))$	Varianz von $f(\mathbf{x})$ unter $P(\mathbf{x})$
$\text{Cov}(f(\mathbf{x}), g(\mathbf{x}))$	Kovarianz von $f(\mathbf{x})$ und $g(\mathbf{x})$ unter $P(\mathbf{x})$
$H(\mathbf{x})$	Shannon-Entropie der Zufallsvariablen \mathbf{x}
$D_{\text{KL}}(P \parallel Q)$	Kullback-Leibler-Divergenz von P und Q
$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$	Normalverteilung über \mathbf{x} mit Mittel $\boldsymbol{\mu}$ und Kovarianz $\boldsymbol{\Sigma}$

Funktionen

$f : \mathbb{A} \rightarrow \mathbb{B}$	Die Funktion f des Definitionsbereichs \mathbb{A} mit dem Wertebereich \mathbb{B}
$f \circ g$	Komposition der Funktionen f und g

$f(\mathbf{x}; \boldsymbol{\theta})$	Eine Funktion \mathbf{x} , parametrisiert durch $\boldsymbol{\theta}$. (Manchmal schreiben wir auch $f(\mathbf{x})$ und lassen das Argument $\boldsymbol{\theta}$ weg, um die Notation zu vereinfachen.)
$\log x$	Natürlicher Logarithmus von x
$\sigma(x)$	Logistische Sigmoidfunktion, $\frac{1}{1 + \exp(-x)}$
$\zeta(x)$	Softplus, $\log(1 + \exp(x))$
$\ \mathbf{x}\ _p$	L^p -Norm von \mathbf{x}
$\ \mathbf{x}\ $	L^2 -Norm von \mathbf{x}
x^+	Positiver Teil von x , d.h. $\max(0, x)$
$\mathbf{1}_{\text{Bedingung}}$	ist 1, wenn die Bedingung zutrifft, anderenfalls 0

Manchmal wenden wir eine Funktion f , deren Argument ein Skalar ist, auf einen Vektor, eine Matrix oder einen Tensor an: $f(\mathbf{x})$, $f(\mathbf{X})$ oder $f(\mathbf{X})$. Dies bezeichnet die elementweise Anwendung von f auf den Vektor, die Matrix oder den Tensor. Beispiel: Ist $\mathbf{C} = \sigma(\mathbf{X})$, dann gilt $C_{i,j,k} = \sigma(X_{i,j,k})$ für alle gültigen Werte von i , j und k .

Datensätze und Verteilungen

p_{data}	Die datengenerierende Verteilung
\hat{p}_{data}	Die empirische Verteilung, die anhand der Trainingsdatensmenge definiert wurde
\mathbb{X}	Eine Sammlung mit Trainingsbeispielen
$\mathbf{x}^{(i)}$	Das i -te Beispiel (Eingabe) eines Datensatzes
$\mathbf{y}^{(i)}$ oder $\mathbf{y}^{(i)}$	Der mit $\mathbf{x}^{(i)}$ für das überwachte Lernen verknüpfte Zielwert
\mathbf{X}	Die Matrix $m \times n$ mit dem Eingabebeispiel $\mathbf{x}^{(i)}$ in Zeile $\mathbf{X}_{i,:}$

1

Einleitung

Schon lange träumen Erfinder von einer Maschine, die selbstständig denken kann. Dieser Wunsch lässt sich bis zu den alten Griechen zurückverfolgen. Die Sagengestalten Pygmalion, Dädalus und Hephaistos sind allesamt berühmte Erfinder. Und Galatea, Talos sowie Pandora lassen sich sämtlich als künstliche Lebensformen betrachten (*Ovid und Martin*, 2004; *Sparkes*, 1996; *Tandy*, 1997).

Als man – über 100 Jahre vor dem Bau des ersten programmierbaren Rechners – über die Möglichkeit solcher Computer nachdachte, fragten sich die Menschen, ob diese Maschinen wohl eines Tages intelligent werden würden (*Lovelace*, 1842). Heute ist die **Künstliche Intelligenz** (KI) eine erfolgreiche Disziplin mit einer Vielzahl praktischer Anwendungen und aktiver Forschungsbereiche. Wir verlassen uns auf »intelligente« Software, um Routineaufgaben zu erledigen, Sprache oder Bilder zu erfassen, medizinische Diagnosen zu stellen und uns bei grundlegender wissenschaftlicher Forschung unterstützen zu lassen.

In der Anfangszeit der Künstlichen Intelligenz hat die Forschung schnell Probleme in Angriff genommen und gelöst, die für Menschen schwierig zu lösen, aber relativ unkompliziert für Computer sind – dabei geht es um Probleme, die durch eine Reihe formaler mathematischer Regeln beschrieben werden können. Die wahre Herausforderung an eine Künstliche Intelligenz besteht jedoch darin, Aufgaben zu erledigen, mit denen Menschen keinerlei Probleme haben, obwohl sie schwer in Worte zu fassen sind – solche Dinge also, die wir intuitiv erledigen und die wir nebenbei tun. Gute Beispiele dafür sind das Erkennen der gesprochenen Sprache oder von Gesichtern in einer Fotografie.

In diesem Buch geht es um die Lösung solch intuitiver Aufgabenstellungen. Die Lösung besteht darin, Computern zu ermöglichen, aus Erfahrungen zu lernen und die Welt so erfassen, als ob sie aus hierarchischen Konzepten besteht, wobei jedes Konzept wiederum definiert ist durch die Beziehung zu einfacheren Konzepten. Indem Computer Wissen aus Erfahrung sammeln, ist es nicht notwendig, mit menschlicher Arbeitskraft formal das Wissen genau zu beschreiben, das der Computer zum Lösen solcher Aufgaben braucht. Die Hierarchie der Konzepte ermöglicht dem Computer das Erlernen komplexer Konzepte, indem er diese aus einfacheren zusammensetzt. Wenn wir diese Modelle grafisch darstellen, blicken wir von oben auf viele Schichten, die sich weit in die Tiefe erstrecken. Deshalb nennen wir diese Methode der Künstlichen Intelligenz **Deep Learning**.

Viele der frühen Erfolge der Künstlichen Intelligenz fanden in relativ sterilen und formalen Umgebungen statt. Das hatte den Vorteil, dass die Computer nur wenig Wissen über die echte Welt benötigten. Ein Beispiel ist IBMs Deep Blue. Dieses Schachprogramm besiegte im Jahr 1997 den Weltmeister Garri Kasparow (*Hsu, 2002*). Beim Schach handelt es sich um ein recht einfaches Universum mit nur 64 Feldern und 32 Figuren, deren Bewegungsmöglichkeiten zudem durch klare Regeln eingeschränkt sind. Die Entwicklung einer erfolgreichen Schachstrategie ist eine echte Meisterleistung. Das liegt aber nicht daran, dass es so kompliziert wäre, einem Computer die Regeln und zulässigen Züge beizubringen. Tatsächlich lassen sich die Regeln in eine sehr kurze Liste formaler Anweisungen fassen, die durch den Programmierer im Vorfeld eingegeben werden.

Die Ironie liegt gerade darin, dass die abstrakten und formalen Aufgaben, die von uns Menschen höchste Geistesleistungen erfordern, für den Computer am unteren Ende der Schwierigkeitsskala stehen. Schon lange können Computer auch die besten Großmeister im Schach schlagen. Doch erst seit kurzer Zeit vermögen sie es, auch bei Objekt- und Spracherkennung mit den Menschen gleichzuziehen. Wir alle greifen im Alltag auf einen gewaltigen Erfahrungsschatz zurück, mit dem wir uns in der Welt zurechtfinden. Ein Großteil unseres Wissens ist subjektiv und intuitiv, lässt sich also nur schwer formal ausdrücken. Damit Computer auf intelligente Weise agieren können, benötigen sie genau dieses Wissen und diese Erfahrungen. Eine der größten Herausforderungen auf dem Gebiet der Künstlichen Intelligenz besteht darin, dem Computer dieses informale und nicht eindeutige Wissen zu vermitteln.

In diversen Projekten der Künstlichen Intelligenz wurde versucht, unser Weltwissen fest in formalen Sprachen zu codieren. Mithilfe logischer Inferenzregeln kann ein Computer auf dieser Basis Schlüsse ziehen. Das ist die sogenannte **wissensbasierte** Herangehensweise an die Künstliche

Intelligenz. Allerdings hat keines dieser Projekte den großen Durchbruch gebracht. Ein berühmtes Beispiel ist *Cyc* (*Lenat und Guha, 1989*). *Cyc* ist eine Inferenzmaschine samt Datenbank mit Aussagen in der Sprache *CycL*. Diese Aussagen werden von menschlichen Supervisoren recht umständlich eingegeben. Für uns Menschen ist es ein großes Problem, formale Regeln aufzustellen, die komplex genug sind, um die Welt korrekt zu beschreiben. Das zeigte sich, als *Cyc* einen Bericht über Fred, der sich morgens rasierte, gründlich missverstand (*Linde, 1992*). Die Inferenzmaschine kam zu dem Schluss, dass der Bericht inkonsistent sei. Warum? Nun, *Cyc* wusste, dass Menschen keine elektrischen Bauteile enthalten. Aber da Fred einen elektrischen Rasierapparat in der Hand hielt, ging das System davon aus, dass »FredBeimRasieren« zum Teil aus Elektronik bestand. Daher wollte die Inferenzmaschine wissen, ob Fred beim Rasieren noch ein Mensch ist.

Die Schwierigkeiten, denen sich Systeme gegenüber sehen, die sich auf stark codiertes Wissen stützen, legen nahe, dass KI-Systeme in der Lage sein müssen, ihr eigenes Wissen zu erwerben, indem sie aus Rohdaten Muster extrahieren. Diese Fähigkeit wird als **Machine Learning** bezeichnet. Das Aufkommen des Machine Learnings versetzte Computer in die Lage, Probleme anzugehen, für deren Lösung Wissen über die Realität erforderlich ist, und dafür scheinbar subjektive Entscheidungen zu treffen. Ein einfacher Machine-Learning-Algorithmus namens **logistische Regression** kann entscheiden, ob ein Kaiserschnitt angeraten ist (*Mor-Yosef et al., 1990*). Ein einfacher Machine-Learning-Algorithmus namens **Naive Bayes** kann erwünschte E-Mails und Spam voneinander trennen.

Die Leistung dieser einfachen Machine-Learning-Algorithmen hängt stark von der **Repräsentation** (also der Darstellung oder Aufbereitung) der Ausgangsdaten ab. Ein Beispiel: Wenn logistische Regression verwendet wird, um Empfehlungen zu einem Kaiserschnitt anzugeben, untersucht nicht das KI-System selbst die Patientin. Stattdessen gibt ein Arzt relevante Daten in das System ein, zum Beispiel ob eine Gebärmutternarbe vorhanden ist. Jede in der Repräsentation der Patientin enthaltene Angabe wird als **Merkmal** bezeichnet. Die logistische Regression »lernt«, wie jedes dieser Merkmale der Patientin mit unterschiedlichen Resultaten korreliert. Sie kann jedoch keinerlei Einfluss darauf nehmen, wie die Merkmale definiert werden. Würde der logistischen Regression eine MRT-Aufnahme der Patientin anstelle des formalisierten Arztberichts zur Verfügung gestellt, könnte das System damit keine nützliche Vorhersage machen. Die einzelnen Pixel des MRTs haben kaum eine Korrelation mit möglichen Komplikationen, die während der Entbindung auftreten können.

Diese Abhängigkeit von der Repräsentation ist ein allgemeines Phänomen, das sich durch die gesamte Informatik und sogar unseren Alltag zieht. In der Informatik erfolgt ein Vorgang wie das Durchsuchen einer Datensammlung exponentiell schneller, wenn die Datensammlung auf intelligente Weise strukturiert und indiziert ist. Menschen können problemlos mit arabischen Ziffern rechnen, doch mit römischen Zahlenzeichen fällt es den meisten sehr viel schwerer. Kein Wunder also, dass die Auswahl der Repräsentation oder Darstellung einen gewaltigen Effekt auf die Leistungsfähigkeit von Machine-Learning-Algorithmen hat. Abbildung 1.1 zeigt ein einfaches Beispiel.

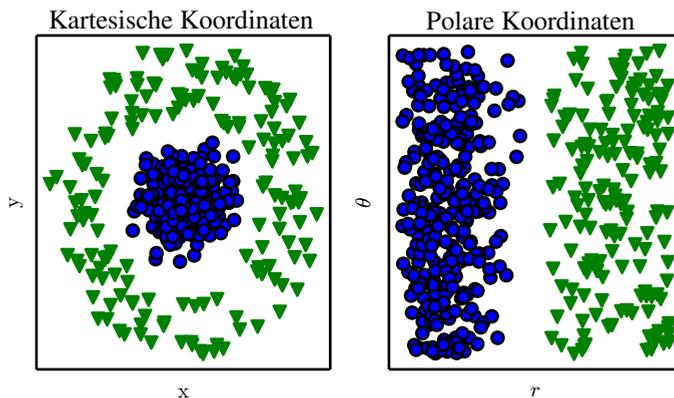


Abbildung 1.1: Beispiel für unterschiedliche Darstellungen: Zwei Datenkategorien sollen durch eine Gerade in einem Punktdiagramm voneinander getrennt werden. Im linken Diagramm werden die Daten mithilfe kartesischer Koordinaten dargestellt. Das macht die Lösung der Aufgabe unmöglich. Im rechten Diagramm sind die Daten mithilfe polarer Koordinaten dargestellt und eine einfache Senkrechte reicht aus, um die Aufgabe zu lösen. (Abbildung gemeinsam mit David Warde-Farley erstellt.)

Viele Aufgaben der Künstlichen Intelligenz können gelöst werden, indem man die richtige Zusammenstellung von Merkmalen konzipiert, die für eine Aufgabe eine Rolle spielen, und diese Merkmale einem einfachen Machine-Learning-Algorithmus zur Verfügung stellt. Um beispielsweise einen Sprecher anhand des Klangs seiner Stimme zu identifizieren, kann sich das Volumen seines Stimmapparats als nützliches Merkmal erweisen. Dieses Merkmal liefert einen guten Hinweis darauf, ob es sich beim Sprecher um einen Mann, eine Frau oder ein Kind handelt.

Für viele Aufgaben besteht jedoch die Schwierigkeit darin zu wissen, welche Merkmale wichtig sind und extrahiert werden müssen. Beispiel: Ein Programm soll Fotografien prüfen und entscheiden, ob darauf ein Pkw zu sehen ist. Wie wir wissen, haben Pkws Räder. Ein Merkmal könnte also das

Vorhandensein von Rädern sein. Aber wie lässt sich ein Rad in Form von Pixelwerten beschreiben? Schwierig, nicht wahr? Ein Rad weist eine einfache geometrische Form auf, aber im Bild gibt es viele andere Elemente, die eine Erkennung erschweren: Schatten, Sonnenblendungen auf Metallteilen des Rads, Kotflügel oder andere Objekte, die Teile des Rads verdecken usw.

Eine Lösung besteht darin, Machine Learning zu verwenden, um nicht nur aus dem Mapping zwischen Repräsentation und Ausgabe, sondern auch aus der Repräsentation (oder Darstellung) selbst Erkenntnisse zu gewinnen. Dieser Ansatz wird als **Representation Learning** bezeichnet. Erlernete Repräsentationen führen häufig zu einer sehr viel besseren Leistung gegenüber individuell angepassten Repräsentationen. Und sie befähigen KI-Systeme dazu, sich bei minimaler menschlicher Intervention schnell an neue Aufgaben anzupassen. Ein Representation-Learning-Algorithmus kann eine gutes Set von Merkmalen für eine einfache Aufgabe in wenigen Minuten erkennen. Bei komplexen Aufgaben kann er einige Stunden oder Monate benötigen. Die händische Konzipierung von Merkmalen für eine komplexe Aufgabe dagegen erfordert jede Menge menschliche Zeit und Mühe – ein Forschungsteam könnte Jahrzehnte damit beschäftigt sein.

Ein Musterbeispiel für einen Representation-Learning-Algorithmus ist der **Autoencoder**. Ein Autoencoder ist eine Kombination aus einer **Encoder**-Funktion, die Eingabedaten in eine andere Repräsentation umwandelt, und einer **Decoder**-Funktion, die diese neue Repräsentation wieder zurück in das Ausgangsformat umwandelt. Autoencoder werden darauf trainiert, so viele Daten wie möglich zu erhalten, wenn eine Eingabe die Encoder- und Decoder-Funktion durchläuft. Gleichzeitig werden sie aber auch darauf trainiert, der neuen Repräsentation einige nützliche Eigenschaften zu verleihen. Je nach Art des Autoencoders liegt das Augenmerk auf unterschiedlichen Eigenschaften.

Beim Entwickeln von Merkmalen oder Algorithmen zum Erlernen von Merkmalen ist es üblicherweise unser Ziel, die **Faktoren der Variation** zu separieren, die die beobachteten Daten erklären. In diesem Kontext bezieht sich das Wort »Faktor« rein auf die unterschiedlichen Einflussquellen. Der Begriff hat nichts mit einer Multiplikation zu tun. Diese Faktoren sind häufig keine direkt beobachteten Größen, sondern entweder unbeobachtete Objekte oder unbeobachtete Kräfte in der physischen Welt, die ihrerseits die einer Beobachtung zugänglichen Größen beeinflussen. Es kann sich auch um Konstrukte des menschlichen Denkens handeln, die nützliche verallgemeinerte Erklärungen oder vermutete Ursachen für die beobachteten Daten liefern. Sie können sich diese Faktoren als Konzepte oder Abstraktionen vorstellen, die dabei helfen, der großen Variabilität der Daten einen Sinn zu

geben. Beim Analysieren einer Sprachaufzeichnung gehören zu den Faktoren der Variation zum Beispiel das Alter und das Geschlecht des Sprechers, der Dialekt und die verwendeten Wörter. Beim Analysieren einer Abbildung eines Fahrzeugs wiederum gehören dessen Position und Farbe, der Betrachtungswinkel und die Helligkeit des Sonnenlichts zu den Faktoren der Variation.

Bei vielen KI-Anwendungen in der Praxis liegt eine große Schwierigkeit darin, dass viele der Faktoren der Variation sich auf jedes einzelne Datenelement auswirken, das wir beobachten können. Die einzelnen Pixel des Fotos eines roten Autos können bei Nacht nahezu schwarz erscheinen. Die Form der Fahrzeugumrisse ist abhängig vom Blickwinkel. Die meisten Anwendungen erfordern es von uns, dass wir die Faktoren der Variation *separieren* und diejenigen verwerfen, die nicht von Belang sind.

Das Extrahieren solch übergeordneter, abstrakter Merkmale aus den Rohdaten kann sich als extrem schwierig erweisen. Viele der Faktoren der Variation wie ein Dialekt können nur bestimmt werden, wenn ein tiefes, nahezu menschliches Verständnis der Daten vorhanden ist. Wenn es genauso schwierig ist, eine Repräsentation zu erzielen wie das ursprüngliche Problem zu lösen, scheint uns das Representation Learning auf den ersten Blick nicht hilfreich zu sein.

Deep Learning löst dieses zentrale Problem beim Representation Learning, indem Repräsentationen eingesetzt werden, die in Form von anderen simpleren Repräsentationen dargestellt werden. Deep Learning befähigt den Computer, komplexe Konzepte aus einfacheren Konzepten zu konstruieren. Abbildung 1.2 zeigt, wie ein Deep-Learning-System das Konzept einer Personenfotografie darstellen kann, indem es einfachere Konzepte kombiniert wie hier Ecken und Umrisse, die wiederum über Kantenverläufe definiert werden.

Ein Musterbeispiel für ein Deep-Learning-Modell ist ein tiefes Feedforward-Netz, auch **mehrschichtiges Perzeptron** (engl. *multilayer perceptron*, MLP) genannt. Ein mehrschichtiges Perzeptron ist eine mathematische Funktion, die eine Menge von Eingabewerten entsprechenden Ausgabewerten zuordnet. Die Funktion wird aus vielen einfacheren Funktionen zusammengesetzt. Dabei können Sie sich vorstellen, dass jede Anwendung einer anderen mathematischen Funktion eine neue Repräsentation der Eingaben liefert.

Das Konzept des Erlernens der richtigen Repräsentation für die Daten stellt eine Perspektive auf das Deep Learning dar. Eine weitere Perspektive rückt die Tiefe in den Mittelpunkt: Sie ermöglicht dem Computer das Erlernen eines mehrstufigen Computerprogramms. Jede Schicht der

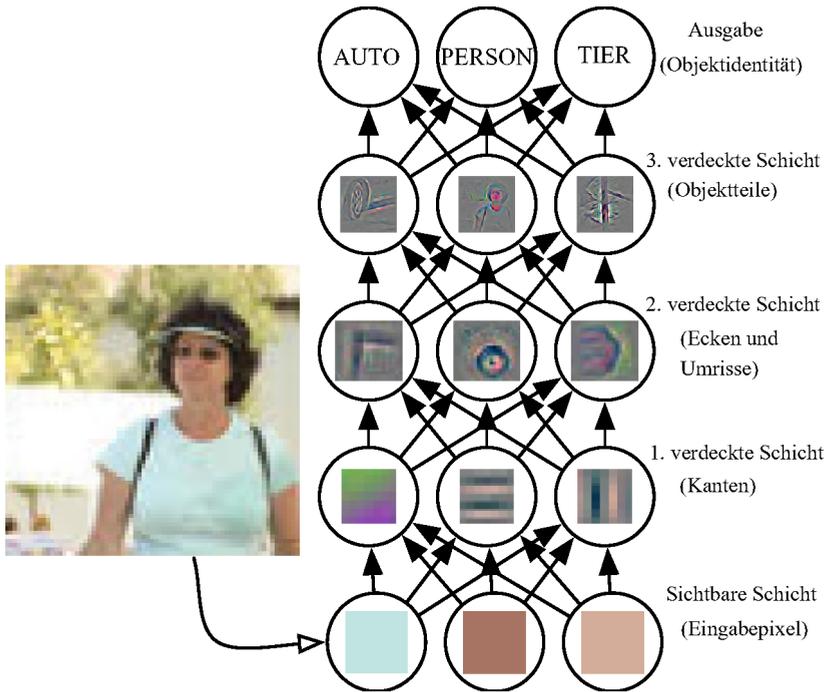


Abbildung 1.2: Darstellung eines Deep-Learning-Modells. Für den Computer ist es nicht einfach, die Bedeutung roher Sensoreingaben zu verstehen. Ein Beispiel ist dieses Foto, das als Ansammlung von Pixelwerten dargestellt wird. Die Funktion zur Zuordnung einer Pixelmenge zu einer Objektidentität ist sehr kompliziert. Das Erlernen oder Bewerten dieser Zuordnung erscheint unmöglich, wenn man das Problem direkt angehen will. Mit Deep Learning wird die Schwierigkeit jedoch vermindert, indem die erwünschte komplizierte Zuordnung in eine Reihe verschachtelter einfacherer Zuordnungen aufgeteilt wird. Jede davon beschreibt eine andere Schicht des Modells. Die Eingabe wird als **sichtbare Schicht** dargestellt. Die Bezeichnung wurde gewählt, weil sie die Variablen enthält, die wir beobachten können. Darauf werden in einer Reihe von **verdeckten Schichten** immer abstraktere Bildmerkmale extrahiert. Diese Schichten werden als »verdeckt« bezeichnet, weil ihre Werte nicht in den Daten enthalten sind, sondern das Modell selbst entscheiden muss, welche Konzepte zum Erklären der Beziehungen in den beobachteten Daten hilfreich sind. Die gezeigten Abbildungen visualisieren die Art von Merkmalen, die durch jede verdeckte Einheit dargestellt werden. Anhand der Pixel können in der ersten Schicht durch Vergleich der Helligkeitswerte benachbarter Pixel die Kanten bestimmt werden. Auf Basis der in der ersten verdeckten Schicht ermittelten Kanten erfolgt in der zweiten verdeckten Schicht die Suche nach Ecken und verlängerten Konturen, die als Aneinanderreihung von Kanten erkannt werden. Nachdem die zweite verdeckte Schicht das Bild anhand von Ecken und Konturen beschrieben hat, sucht die dritte verdeckte Schicht nach kompletten Teilen bestimmter Objekte. Dazu wird auf bestimmte Sammlungen von Konturen und Ecken geprüft. Im Ergebnis steht eine Bildbeschreibung in Form von Objektteilen zur Verfügung, die dann zum Erkennen der Objekte im Foto verwendet werden kann. (Abbildungen mit freundlicher Genehmigung von *Zeiler und Fergus (2014)*)

Repräsentation wird dabei als Speicherzustand des Computers nach der parallelen Ausführung einer weiteren Anweisungsgruppe betrachtet. Netze mit einer größeren Tiefe können mehr Anweisungen nacheinander ausführen. Sequenzielle Anweisungen bieten mehr Leistungsfähigkeit, denn spätere Anweisungen können auf die Ergebnisse früherer Anweisungen zurückgreifen. Bei dieser Perspektive auf das Deep Learning codieren nicht alle Informationen in den Aktivierungen einer Schicht notwendigerweise Faktoren der Variation, die die Eingabe erklären. Die Repräsentation speichert außerdem Informationen über den Zustand, mit deren Hilfe ein Programm ausgeführt werden kann, das die Eingabe auswertet. Diese Zustandsinformationen ähneln einem Zähler oder Pointer in einem klassischen Computerprogramm. Sie haben nichts mit dem eigentlichen Inhalt der Eingabe zu tun, sondern helfen dem Modell, die Verarbeitung zu strukturieren.

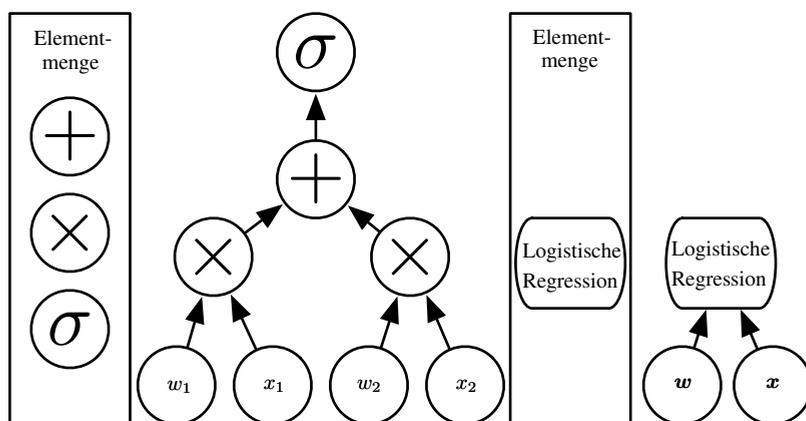


Abbildung 1.3: Abbildung des Berechnungsverlaufs (Graphen) für eine Eingabe, bei dem jeder Knoten auf dem Weg zur Ausgabe eine Berechnung durchführt. Die Tiefe ist der längste Pfad von der Eingabe zur Ausgabe. Sie richtet sich allerdings danach, was als möglicher Berechnungsschritt definiert ist. Die hier dargestellte Berechnung ist die Ausgabe eines logistischen Regressionsmodells $\sigma(\mathbf{w}^T \mathbf{x})$ mit σ als logistischer Sigmoidfunktion. Bei Verwendung von Addition, Multiplikation und logistischen Sigmoidfunktionen als Elemente unserer Programmiersprache hat dieses Modell eine Tiefe von 3. Gilt die logistische Regression als eigenes Element, hat das Modell die Tiefe 1.

Es gibt zwei wesentliche Möglichkeiten zum Messen der Modelltiefe. Bei der ersten wird die Anzahl sequenzieller Anweisungen bestimmt, die zum Bewerten der Architektur ausgeführt werden müssen. Das entspricht quasi dem längsten Pfad in einem Ablaufplan, der die Berechnung der einzelnen Ausgaben des Modells anhand der Eingaben beschreibt. So wie zwei äquivalente Computerprogramme abhängig von der Programmiersprache