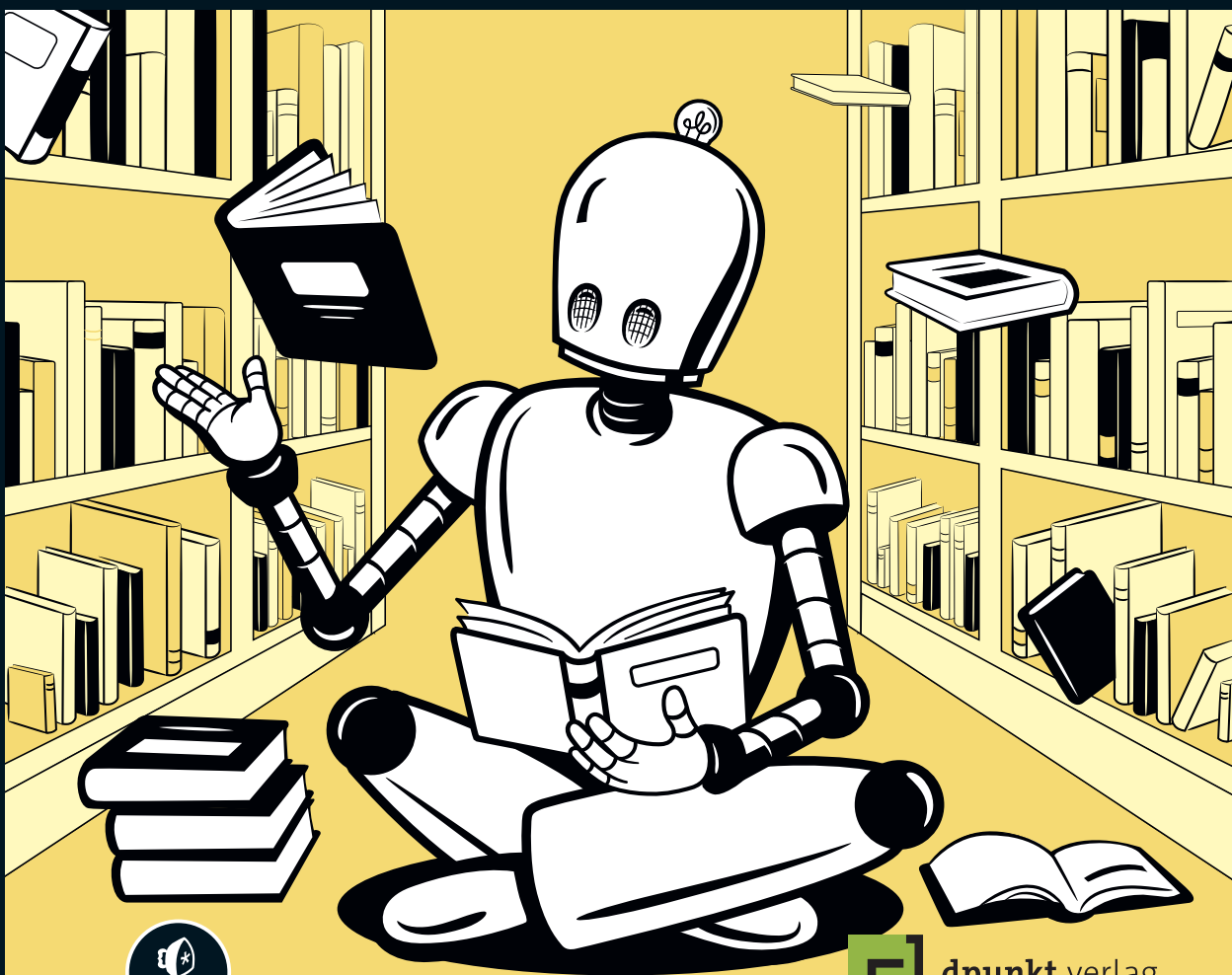


# MACHINE LEARNING & KI

ZENTRALE KONZEPTE VERSTEHEN UND ANWENDEN

**KOMPAKT**

SEBASTIAN RASCHKA



dpunkt.verlag

**Sebastian Raschka**, PhD, ist Forscher für maschinelles Lernen und KI mit einer großen Leidenschaft für Wissensvermittlung. Als Lead AI Educator bei Lightning AI brennt er dafür, KI und Deep Learning verständlich darzustellen und Menschen beizubringen, wie sie diese Technologien in großem Umfang nutzen können. Bevor er voll und ganz in Lightning AI eingestiegen ist, hatte Sebastian Raschka eine Position als Assistenzprofessor für Statistik an der University of Wisconsin-Madison inne, wo er sich auf die Erforschung von Deep Learning und maschinellem Lernen spezialisierte. Auf seiner Website (<https://sebastianraschka.com>) erfahren Sie mehr über seine Forschung. Außerdem liebt Sebastian Raschka Open-Source-Software und leistet seit über einem Jahrzehnt leidenschaftlich Beiträge dazu. Neben dem Programmieren schreibt er auch gern und ist Autor der Bestseller *Python Machine Learning* und *Machine Learning with PyTorch and Scikit-Learn* (beide bei Packt Publishing veröffentlicht).

#### *Über den Technischen Redakteur*

**Andrea Panizza** ist leitender KI-Spezialist beim Unternehmen Baker Hughes, das modernste KI/ML-Techniken einsetzt, um die technische Konstruktion zu beschleunigen, die Informationsgewinnung und -extraktion aus großen Dokumentensammlungen zu Arbeitsmappen und die unbemannte Inspektion von Anlagen mithilfe von Computer Vision zu unterstützen. Er hat einen Dokortitel in Computational Fluid Dynamics. Bevor er zu Baker Hughes kam, arbeitete er als CFD-Forscher bei CIRA (dem italienischen Zentrum für Luft- und Raumfahrtforschung).

#### Coypright und Urheberrechte:

Die durch die dpunkt.verlag GmbH vertriebenen digitalen Inhalte sind urheberrechtlich geschützt. Der Nutzer verpflichtet sich, die Urheberrechte anzuerkennen und einzuhalten. Es werden keine Urheber-, Nutzungs- und sonstigen Schutzrechte an den Inhalten auf den Nutzer übertragen. Der Nutzer ist nur berechtigt, den abgerufenen Inhalt zu eigenen Zwecken zu nutzen. Er ist nicht berechtigt, den Inhalt im Internet, in Intranets, in Extranets oder sonst wie Dritten zur Verwertung zur Verfügung zu stellen. Eine öffentliche Wiedergabe oder sonstige Weiterveröffentlichung und eine gewerbliche Vervielfältigung der Inhalte wird ausdrücklich ausgeschlossen. Der Nutzer darf Urheberrechtsvermerke, Markenzeichen und andere Rechtsvorbehalte im abgerufenen Inhalt nicht entfernen.

**Sebastian Raschka**

# **Machine Learning und KI kompakt**

**Zentrale Konzepte verstehen und anwenden**



**dpunkt.verlag**

Sebastian Raschka

Übersetzung: Frank Langenau

Lektorat: Sandra Bollenbacher, Alissa Melitzer

Copy-Editing: Petra Heubach-Erdmann, Düsseldorf

Satz: Gerhard Alfes, mediaService, Siegen, [www.mediaservice.tv](http://www.mediaservice.tv), Birgit Bäuerlein

Herstellung: Stefanie Weidner

Cover-Illustratorin: Gina Redman

Umschlaggestaltung: Eva Hepper, Silke Braun

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:

Print 978-3-98889-031-3

PDF 978-3-98890-214-6

ePub 978-3-98890-215-3

1. Auflage 2025

Translation Copyright für die deutschsprachige Ausgabe © 2025 dpunkt.verlag GmbH

Wieblinger Weg 17

69123 Heidelberg

E-Mail: [hallo@dpunkt.de](mailto:hallo@dpunkt.de)

Copyright © 2024 by Sebastian Raschka. Title of English-language original: *Machine Learning Q and A: 30 Essential Questions and Answers on Machine Learning and AI*, ISBN 9781718503762, published by

No Starch Press Inc. 245 8th Street, San Francisco, California United States 94103.

The German-language 1st edition Copyright © 2025 by dpunkt.verlag GmbH under license by No Starch Press Inc. All rights reserved.

*Schreiben Sie uns:*

Falls Sie Anregungen, Wünsche und Kommentare haben, lassen Sie es uns wissen: [hallo@dpunkt.de](mailto:hallo@dpunkt.de).

Die vorliegende Publikation ist urheberrechtlich geschützt. Alle Rechte vorbehalten. Die Verwendung der Texte und Abbildungen, auch auszugsweise, ist ohne die schriftliche Zustimmung des Verlags urheberrechtswidrig und daher strafbar. Dies gilt insbesondere für die Vervielfältigung, Übersetzung oder die Verwendung in elektronischen Systemen.

Weiter darf der Inhalt nicht zur Entwicklung, zum Training oder zur Anreicherung von KI-Systemen, insbesondere generativen KI-Systemen, verwendet werden. Die Nutzung für Text- und Data Mining ist untersagt.

Es wird darauf hingewiesen, dass die im Buch verwendeten Soft- und Hardware-Bezeichnungen sowie Markennamen und Produktbezeichnungen der jeweiligen Firmen im Allgemeinen warenzeichen-, marken- oder patentrechtlichem Schutz unterliegen.

Alle Angaben und Programme in diesem Buch wurden mit größter Sorgfalt kontrolliert. Weder Autor noch Verlag noch Übersetzer können jedoch für Schäden haftbar gemacht werden, die in Zusammenhang mit der Verwendung dieses Buches stehen.

# Inhaltsverzeichnis

<b>Vorwort</b>	<b>xiii</b>
<b>Danksagungen</b>	<b>xv</b>
<b>Einleitung</b>	<b>xvii</b>

## Teil I Neuronale Netze und Deep Learning

<b>1</b>	<b>Einbettungen, latenter Raum und Repräsentationen</b>	<b>3</b>
1.1	Einbettungen . . . . .	3
1.2	Latenter Raum . . . . .	5
1.3	Repräsentation . . . . .	6
1.4	Übungen. . . . .	6
1.5	Referenzen . . . . .	7
<b>2</b>	<b>Selbstüberwachtes Lernen</b>	<b>9</b>
2.1	Selbstüberwachtes Lernen vs. Transferlernen . . . . .	9
2.2	Ungelabelte Daten nutzen. . . . .	11
2.3	Selbstvorhersage und kontrastives selbstüberwachtes Lernen . .	12
2.4	Übungen. . . . .	14
2.5	Referenzen . . . . .	14
<b>3</b>	<b>Few-Shot-Lernen</b>	<b>17</b>
3.1	Datensätze und Terminologie. . . . .	17
3.2	Übungen. . . . .	20
<b>4</b>	<b>Die Lotterie-Ticket-Hypothese</b>	<b>21</b>
4.1	Das Lotterie-Ticket-Trainingsverfahren . . . . .	21
4.2	Praktische Konsequenzen und Einschränkungen . . . . .	22

---

4.3	Übungen . . . . .	23
4.4	Referenzen . . . . .	23
<b>5</b>	<b>Überanpassung mit Daten verringern</b>	<b>25</b>
5.1	Allgemeine Methoden . . . . .	25
5.2	Übungen . . . . .	28
5.3	Referenzen . . . . .	28
<b>6</b>	<b>Überanpassung durch Modellmodifikationen reduzieren</b>	<b>31</b>
6.1	Allgemeine Methoden . . . . .	31
6.2	Andere Methoden . . . . .	36
6.3	Eine Regularisierungstechnik auswählen . . . . .	37
6.4	Übungen . . . . .	37
6.5	Referenzen . . . . .	37
<b>7</b>	<b>Multi-GPU-Trainingsparadigmen</b>	<b>39</b>
7.1	Die Trainingsparadigmen . . . . .	39
7.2	Empfehlungen . . . . .	43
7.3	Übungen . . . . .	44
7.4	Referenzen . . . . .	44
<b>8</b>	<b>Der Erfolg der Transformer</b>	<b>45</b>
8.1	Der Aufmerksamkeitsmechanismus . . . . .	45
8.2	Vortraining durch selbstüberwachtes Lernen . . . . .	47
8.3	Große Anzahl von Parametern . . . . .	47
8.4	Einfache Parallelisierung . . . . .	48
8.5	Übungen . . . . .	48
8.6	Referenzen . . . . .	49
<b>9</b>	<b>Generative KI-Modelle</b>	<b>51</b>
9.1	Generative vs. diskriminative Modellierung . . . . .	51
9.2	Arten von tiefen generativen Modellen . . . . .	52
9.3	Empfehlungen . . . . .	59
9.4	Übungen . . . . .	60
9.5	Referenzen . . . . .	60
<b>10</b>	<b>Quellen der Zufälligkeit</b>	<b>61</b>
10.1	Initialisierung der Modellgewichte . . . . .	61
10.2	Sampling und Shuffling von Datensätzen . . . . .	62

10.3	Nichtdeterministische Algorithmen . . . . .	63
10.4	Verschiedene Laufzeitalgorithmen . . . . .	63
10.5	Hardware und Treiber . . . . .	64
10.6	Zufälligkeit und generative KI . . . . .	65
10.7	Übungen. . . . .	67
10.8	Referenzen . . . . .	67

## Teil II Computer Vision

<b>11</b>	<b>Die Anzahl der Parameter berechnen</b>	<b>71</b>
11.1	Wie man die Anzahl der Parameter ermittelt . . . . .	71
11.2	Praktische Anwendungen . . . . .	75
11.3	Übungen. . . . .	75
<b>12</b>	<b>Vollständig verbundene und konvolutionale Schichten</b>	<b>77</b>
12.1	Szenario: Gleiche Größen von Kernel und Eingabe. . . . .	78
12.2	Szenario: Kernel-Größe ist 1. . . . .	79
12.3	Empfehlungen . . . . .	79
12.4	Übungen. . . . .	80
<b>13</b>	<b>Große Trainingsmengen für Vision Transformer</b>	<b>81</b>
13.1	Induktive Verzerrungen in CNNs. . . . .	81
13.2	ViTs können CNNs übertreffen . . . . .	85
13.3	Induktive Verzerrungen in ViTs . . . . .	85
13.4	Empfehlungen . . . . .	87
13.5	Übungen. . . . .	87
13.6	Referenzen . . . . .	87

## Teil III Natural Language Processing

<b>14</b>	<b>Die Verteilungshypothese</b>	<b>91</b>
14.1	Word2vec, BERT und GPT . . . . .	92
14.2	Trifft die Hypothese zu? . . . . .	93
14.3	Übungen. . . . .	94
14.4	Referenzen . . . . .	94

<b>15</b>	<b>Datenvermehrung für Text</b>	<b>95</b>
15.1	Ersetzen von Synonymen . . . . .	95
15.2	Löschen von Wörtern . . . . .	96
15.3	Vertauschen von Wortpositionen . . . . .	96
15.4	Sätze mischen . . . . .	97
15.5	Rauschinjektion . . . . .	97
15.6	Rückübersetzung . . . . .	98
15.7	Synthetische Daten . . . . .	98
15.8	Empfehlungen . . . . .	99
15.9	Übungen . . . . .	99
15.10	Referenzen . . . . .	99
<b>16</b>	<b>Selbstaufmerksamkeit</b>	<b>101</b>
16.1	Aufmerksamkeit in RNNs . . . . .	101
16.2	Der Selbstaufmerksamkeitsmechanismus . . . . .	103
16.3	Übungen . . . . .	104
16.4	Referenzen . . . . .	105
<b>17</b>	<b>Encoder- und Decoder-Transformer</b>	<b>107</b>
17.1	Der ursprüngliche Transformer . . . . .	107
17.2	Encoder-Decoder-Hybride . . . . .	112
17.3	Terminologie . . . . .	112
17.4	Aktuelle Transformer-Modelle . . . . .	113
17.5	Übungen . . . . .	114
17.6	Referenzen . . . . .	114
<b>18</b>	<b>Transformer verwenden und feinabstimmen</b>	<b>117</b>
18.1	Transformer für Klassifizierungsaufgaben verwenden . . . . .	117
18.2	Kontextbezogenes Lernen, Indizierung und Prompt- Feinabstimmung . . . . .	120
18.3	Parametereffiziente Feinabstimmung . . . . .	122
18.4	Reinforcement Learning mit menschlicher Rückmeldung . . .	127
18.5	Vortrainierte Sprachmodelle anpassen . . . . .	128
18.6	Übungen . . . . .	128
18.7	Referenzen . . . . .	129

---

<b>19</b>	<b>Generative LLMs evaluieren</b>	<b>131</b>
19.1	Bewertungsmetriken für LLMs. . . . .	131
19.2	Übungen. . . . .	138
19.3	Referenzen . . . . .	138
<b>Teil IV Produktion und Deployment</b>		
<b>20</b>	<b>Zustandsloses und zustandsbehaftetes Training</b>	<b>143</b>
20.1	Zustandsloses (Re-)Training. . . . .	143
20.2	Zustandsbehaftetes Training . . . . .	144
20.3	Übungen. . . . .	145
<b>21</b>	<b>Datenzentrierte KI</b>	<b>147</b>
21.1	Datenzentrierte vs. modellzentrierte KI . . . . .	147
21.2	Empfehlungen . . . . .	149
21.3	Übungen. . . . .	150
21.4	Referenzen . . . . .	150
<b>22</b>	<b>Inferenz beschleunigen</b>	<b>151</b>
22.1	Parallelisierung. . . . .	151
22.2	Vektorisierung . . . . .	152
22.3	Schleifenkachelung. . . . .	153
22.4	Operatorfusion. . . . .	154
22.5	Quantisierung . . . . .	155
22.6	Übungen. . . . .	156
22.7	Referenzen . . . . .	156
<b>23</b>	<b>Datenverteilungsverschiebungen</b>	<b>157</b>
23.1	Kovariatenverschiebung . . . . .	157
23.2	Labelverschiebung . . . . .	158
23.3	Konzeptverschiebung . . . . .	159
23.4	Domänenverschiebung . . . . .	159
23.5	Arten von Datenverteilungsverschiebungen . . . . .	160
23.6	Übungen. . . . .	161
23.7	Referenzen . . . . .	161

## Teil V Vorhersageperformance und Modellevaluierung

<b>24</b>	<b>Poisson- und ordinale Regression</b>	<b>165</b>
24.1	Übungen . . . . .	166
<b>25</b>	<b>Konfidenzintervalle</b>	<b>167</b>
25.1	Konfidenzintervalle definieren . . . . .	167
25.2	Die Methoden . . . . .	170
25.3	Empfehlungen . . . . .	175
25.4	Übungen . . . . .	175
25.5	Referenzen . . . . .	175
<b>26</b>	<b>Konfidenzintervalle vs. konforme Vorhersagen</b>	<b>177</b>
26.1	Konfidenzintervalle und Vorhersageintervalle . . . . .	177
26.2	Vorhersageintervalle und konforme Vorhersagen . . . . .	178
26.3	Vorhersagebereiche, -intervalle und -mengen. . . . .	178
26.4	Konforme Vorhersagen berechnen . . . . .	179
26.5	Beispiel für eine konforme Vorhersage . . . . .	180
26.6	Die Vorteile der konformen Vorhersagen . . . . .	181
26.7	Empfehlungen . . . . .	182
26.8	Übungen . . . . .	182
26.9	Referenzen . . . . .	183
<b>27</b>	<b>Geeignete Metriken</b>	<b>185</b>
27.1	Die Kriterien . . . . .	185
27.2	Der mittlere quadratische Fehler . . . . .	186
27.3	Der Kreuzentropieverlust . . . . .	188
27.4	Übungen . . . . .	189
<b>28</b>	<b>Das <math>k</math> in der <math>k</math>-fachen Kreuzvalidierung</b>	<b>191</b>
28.1	Kompromisse bei der Auswahl von Werten für $k$ . . . . .	192
28.2	Geeignete Werte für $k$ bestimmen . . . . .	194
28.3	Übungen . . . . .	194
28.4	Referenzen . . . . .	195
<b>29</b>	<b>Diskordanz zwischen Trainings- und Testdatensatz</b>	<b>197</b>
29.1	Übungen . . . . .	199

---

<b>30</b>	<b>Begrenzte gelabelte Daten</b>	<b>201</b>
30.1	Die Modellperformance mit begrenzten gelabelten Daten verbessern . . . . .	201
30.2	Empfehlungen . . . . .	210
30.3	Übungen. . . . .	211
30.4	Referenzen . . . . .	212
	<b>Nachwort</b>	<b>213</b>
	<b>Lösungen zu den Übungen</b>	<b>215</b>
	<b>Index</b>	<b>235</b>



---

# Vorwort

Es gibt Hunderte von Einführungsbüchern zum maschinellen Lernen, in einer Vielzahl von Stilen und Ansätzen, von theorieorientierten Perspektiven für Studenten bis hin zu geschäftsorientierten Sichtweisen für Vorstandsetagen. Diese Einführungsbücher sind unschätzbare Ressourcen für Personen, die ihre ersten Schritte in diesem Bereich machen, und sie werden dies auch in den kommenden Jahrzehnten bleiben.

Allerdings besteht der Weg zum Fachwissen nicht nur aus den Anfängen. Er führt auch über verschlungene Umwege, die steilen Anstiege und die Nuancen, die anfangs nicht offensichtlich sind. Anders ausgedrückt: Nachdem sie sich die Grundlagen angeeignet haben, stellt sich für die Lernenden die Frage: »Was kommt als Nächstes?« Genau hier, im Bereich jenseits der Grundlagen, liegt die Zweckbestimmung dieses Buches.

Sebastian führt die Leser durch ein breites Spektrum an mittleren und fortgeschritteneren Themen des angewandten maschinellen Lernens, auf die sie auf ihrem Weg zum Fachwissen wahrscheinlich stoßen werden. Man könnte sich kaum einen besseren Lehrer wünschen als Sebastian, der – ohne Übertreibung – der beste Dozent für maschinelles Lernen ist, der derzeit auf diesem Gebiet tätig ist. Auf jeder Seite vermittelt Sebastian nicht nur seine umfangreichen Fachkenntnisse, sondern teilt auch die Leidenschaft und Neugier, die wahre Kompetenz auszeichnen.

Dieses Buch richtet sich an alle Lernenden, die die erste Schwelle überschritten haben und nun tiefer einsteigen wollen. Wenn Sie dieses Buch durcharbeiten, werden Sie Ihre Fachkenntnisse um ein Vielfaches erweitern. Lassen Sie es eine Brücke zu Ihrer nächsten Phase von lohnenden Abenteuern im maschinellen Lernen sein.

Viel Glück!

*Chris Albon*

Director of Machine Learning, the Wikimedia Foundation  
San Francisco, August 2023



## Danksagungen

Ein Buch zu schreiben, ist ein enormes Unterfangen. Dieses Projekt wäre ohne die Hilfe der Open-Source- und Machine-Learning-Communities, die gemeinsam die Technologien entwickelt haben, um die es in diesem Buch geht, nicht möglich gewesen.

Ich möchte mich bei den folgenden Personen für ihr ungemein hilfreiches Feedback zum Manuskript bedanken:

- Andrea Panizza, der ein hervorragender technischer Redakteur war und sehr wertvolles und aufschlussreiches Feedback geliefert hat
- Anton Reshetnikov, der ein übersichtlicheres Layout für das Flussdiagramm zum überwachten Lernen in Kapitel 30 vorgeschlagen hat
- Nikan Doosti, Juan M. Bello-Rivas und Ken Hoffman, die auf verschiedene typografische Fehler hingewiesen haben
- Abigail Schott-Rosenfield und Jill Franklin für ihre vorbildliche redaktionelle Arbeit. Ihr Geschick, die richtigen Fragen zu stellen und die Sprache zu verbessern, hat die Qualität dieses Buches erheblich gesteigert.



---

# Einleitung

Dank der rasanten Fortschritte beim Deep Learning haben sich maschinelles Lernen und künstliche Intelligenz in den letzten Jahren erheblich ausgebreitet.

Diese Entwicklung ist spannend, wenn wir davon ausgehen, dass diese Fortschritte neue Branchen schaffen, bestehende Branchen verändern und die Lebensqualität von Menschen auf der ganzen Welt verbessern werden. Andererseits kann das ständige Auftauchen neuer Techniken dazu führen, dass es schwierig und zeitaufwendig ist, mit den neuesten Entwicklungen Schritt zu halten. Dennoch ist es für Fachleute und Organisationen, die diese Technologien nutzen, unerlässlich, auf dem Laufenden zu bleiben.

Ich habe dieses Buch als Ressource für Leser und Praktiker des maschinellen Lernens geschrieben, die ihr Fachwissen auf diesem Gebiet erweitern und mehr über die Techniken erfahren möchten, die ich für nützlich und wichtig erachte, die aber in traditionellen und einführenden Lehrbüchern und Kursen oft übersehen werden. Ich hoffe, dass dieses Buch für Sie eine wertvolle Ressource ist, um neue Einblicke zu gewinnen und neue Techniken zu entdecken, die Sie in Ihrer Arbeit umsetzen können.

## An wen richtet sich dieses Buch?

Oft fühlt es sich wie eine Gratwanderung an, sich in der Welt der KI und des maschinellen Lernens zurechtzufinden, da die meisten Bücher an einem der beiden Enden angesiedelt sind: breite Einführungen für Anfänger oder tiefgründige mathematische Abhandlungen. Dieses Buch veranschaulicht und erörtert wichtige Entwicklungen in diesen Bereichen, ist dabei aber leicht verständlich und setzt keine höheren mathematischen oder programmiertechnischen Kenntnisse voraus.

Dieses Buch richtet sich an Personen, die bereits einige Erfahrung mit maschinellem Lernen haben und neue Konzepte und Techniken erlernen möchten. Es ist ideal für diejenigen, die einen Grundkurs in maschinellem Lernen oder Deep Learning absolviert oder ein entsprechendes Einführungsbuch zu diesem Thema gelesen haben. (Im gesamten Buch verwende ich *maschinelles Lernen* als Oberbegriff für maschinelles Lernen, Deep Learning und KI.)

## Was werden Sie von diesem Buch haben?

Dieses Buch ist in einem einzigartigen Frage-und-Antwort-Stil geschrieben, bei dem jedes kurze Kapitel um eine zentrale Frage zu grundlegenden Konzepten des maschinellen Lernens, des Deep Learning und der KI aufgebaut ist. Auf jede Frage folgt eine Erklärung mit mehreren Illustrationen und Abbildungen sowie Übungen, um Ihr Verständnis zu testen. Viele Kapitel enthalten auch Verweise auf weiterführende Literatur. Diese mundgerechten Informationsbrocken bieten einen unterhaltsamen Einstieg auf Ihrem Weg vom Anfänger zum Experten für maschinelles Lernen.

Das Buch deckt ein breites Themenspektrum ab. Es enthält neue Erkenntnisse über etablierte Architekturen, wie zum Beispiel Convolutional Networks (Faltungsnetze), mit denen Sie diese Technologien effektiver nutzen können. Außerdem werden erweiterte Techniken erörtert, wie etwa das Innenleben von großen Sprachmodellen (Large Language Models, LLMs) und Vision Transformers. Selbst erfahrene Forscher und Praktiker im Bereich des maschinellen Lernens werden etwas Neues finden, das sie ihrem Arsenal an Techniken hinzufügen können.

Dieses Buch macht Sie zwar mit neuen Konzepten und Ideen bekannt, ist aber weder ein Mathematik- noch Programmierbuch. Beim Lesen müssen Sie keine Beweise herleiten oder Code ausführen. Mit anderen Worten: Dieses Buch ist ein perfekter Reisebegleiter oder etwas, das Sie auf Ihrem Lieblingslesesessel mit Ihrem Morgenkaffee oder Tee lesen können.

## Wie man dieses Buch liest

Die einzelnen Kapitel dieses Buches sind in sich abgeschlossen, sodass Sie nach Belieben zwischen den Themen wechseln können. Wird ein Konzept aus einem Kapitel in einem anderen ausführlicher erklärt, habe ich Kapitelverweise eingefügt, denen Sie folgen können, um Lücken in Ihrem Verständnis zu schließen.

Allerdings sind die Kapitel in einer strategischen Reihenfolge angeordnet. Zum Beispiel bereitet das frühe Kapitel über Einbettungen die Grundlage für spätere Diskussionen über selbstüberwachtes Lernen und Few-Shot Learning. Um die Lektüre so einfach wie möglich zu gestalten und den Inhalt so umfassend wie möglich zu erfassen, empfehle ich, das Buch von Anfang bis Ende zu lesen.

Zu jedem Kapitel gibt es optionale Übungen für Leser, die ihr Verständnis testen wollen, mit einem Antwortschlüssel am Ende des Buches. Darüber hinaus finden Sie für alle Paper, auf die in einem Kapitel verwiesen wird, oder für weiterführende Literatur zum Thema des Kapitels die vollständigen Quellenangaben im Abschnitt »Referenzen« dieses Kapitels.

Das Buch ist in fünf Hauptteile gegliedert, die sich mit den wichtigsten Themen des maschinellen Lernens und der KI in der heutigen Zeit befassen.

**Teil I: Neuronale Netze und Deep Learning** behandelt Fragen zu tiefen neuronalen Netzen und Deep Learning, die nicht spezifisch für einen bestimmten Teilbereich sind. Zum Beispiel erörtern wir Alternativen zum überwachten Lernen und Techniken, die Überanpassung – ein häufiges Problem bei Modellen für maschinelles Lernen, die auf praktische Probleme mit begrenzter Datenmenge angewendet werden – vermeiden oder reduzieren sollen.

- **Kapitel 1: Einbettungen, Latenter Raum und Repräsentationen** beschäftigt sich mit den Unterschieden und Ähnlichkeiten zwischen Einbettungsvektoren, latenten Vektoren und Repräsentationen. Es wird verdeutlicht, wie diese Konzepte dazu beitragen, Informationen im Kontext des maschinellen Lernens zu codieren.
- **Kapitel 2: Selbstüberwachtes Lernen** konzentriert sich auf selbstüberwachtes Lernen, eine Methode, die es neuronalen Netzen ermöglicht, große, nicht gelabelte Datensätze in überwachter Art und Weise zu nutzen.
- **Kapitel 3: Few-Shot Learning** stellt Few-Shot Learning vor, eine spezialisierte Technik des überwachten Lernens, die auf kleine Trainingsdatensätze zugeschnitten ist.
- **Kapitel 4: Die Lotterieticket-Hypothese** untersucht die Idee, dass zufällig initialisierte neuronale Netze kleinere, effiziente Teilnetze enthalten.
- **Kapitel 5: Überanpassung mit Daten reduzieren** setzt sich mit dem Problem der Überanpassung im maschinellen Lernen auseinander und erörtert Strategien, bei denen Datenvermehrung und nicht gelabelte Daten im Mittelpunkt stehen, um Überanpassung zu verhindern.
- **Kapitel 6: Überanpassung mit Modellmodifikationen verringern** erweitert die Betrachtungen zur Überanpassung und konzentriert sich dabei auf modellbezogene Lösungen wie Regularisierung, die Entscheidung für einfachere Modelle und Ensemble-Techniken.
- **Kapitel 7: Multi-GPU-Trainingsparadigmen** erläutert verschiedene Trainingsparadigmen für Multi-GPU-Setups inklusive Daten- und Modellparallelität, um das Modelltraining zu beschleunigen.
- **Kapitel 8: Der Erfolg der Transformer** untersucht die beliebte Transformer-Architektur, wobei es insbesondere um Features wie Attention-Mechanismen, einfache Parallelisierung und hohe Parameteranzahlen geht.
- **Kapitel 9: Generative KI-Modelle** bietet einen umfassenden Überblick über tiefe generative Modelle, die dafür gedacht sind, verschiedene Medienformen zu erzeugen, darunter Bilder, Text und Audio. Erörtert werden die Stärken und Schwächen der einzelnen Modelltypen.
- **Kapitel 10: Quellen des Zufalls** behandelt die verschiedenen Quellen des Zufalls beim Training von tiefen neuronalen Netzen, die zu inkonsistenten und nicht reproduzierbaren Ergebnissen sowohl beim Training als auch bei der Inferenz führen können. Während Zufälligkeiten unbeabsichtigt auftreten können, ist es auch möglich, sie vom Konzept her absichtlich einzubringen.

**Teil II: Computer Vision** konzentriert sich auf Themen, die hauptsächlich mit Deep Learning zu tun haben, aber spezifisch für Computer Vision sind. Viele davon betreffen CNNs und Vision Transformer.

- **Kapitel 11: Die Anzahl der Parameter berechnen** erläutert das Verfahren, mit dem sich die Parameter in einem CNN bestimmen lassen. Dies ist nützlich, um den Speicherbedarf eines Modells abzuschätzen.
- **Kapitel 12: Vollständig verbundene Schichten und Convolutional Layer** veranschaulicht die Umstände, unter denen Faltungsschichten nahtlos vollständig verbundene Schichten ersetzen können. Dies ist nützlich, um die Hardware zu optimieren oder Implementierungen zu vereinfachen.
- **Kapitel 13: Große Trainingsdatensätze für Vision Transformer** untersucht die Gründe, warum Vision Transformer im Vergleich zu herkömmlichen CNNs umfangreichere Trainingsdatensätze benötigen.

**Teil III: Verarbeitung natürlicher Sprache (Natural Language Processing, NLP)** behandelt Themen rund um die Verarbeitung von Text, von denen viele mit Transformer-Architekturen und Selbstaufmerksamkeit zu tun haben.

- **Kapitel 14: Die Verteilungshypothese** befasst sich mit der Verteilungshypothese, einer linguistischen Theorie, die besagt, dass Wörter, die in den gleichen Kontexten vorkommen, tendenziell ähnliche Bedeutungen haben. Diese Eigenschaft ist nützlich für das Training von Modellen des maschinellen Lernens.
- **Kapitel 15: Datenvermehrung für Text** beleuchtet die Bedeutung der Vermehrung von Text, einer Technik, mit der sich die Datensatzgröße künstlich erhöhen lässt. Dies kann hilfreich sein, um die Modellperformance zu verbessern.
- **Kapitel 16: Selbstaufmerksamkeit** stellt mit Selbstaufmerksamkeit einen Mechanismus vor, der es jedem Segment der Eingabe in ein neuronales Netz erlaubt, auf andere Teile zu verweisen. Selbstaufmerksamkeit ist ein entscheidender Mechanismus in modernen Sprachmodellen.
- **Kapitel 17: Encoder- und Decoder-Transformer** beschreibt die Nuancen von Encoder- und Decoder-Transformer-Architekturen und erläutert, welche Art von Architektur für die jeweilige Sprachverarbeitungsaufgabe am nützlichsten ist.
- **Kapitel 18: Vortrainierte Transformer verwenden und feinabstimmen** erläutert verschiedene Methoden zur Feinabstimmung vortrainierter LLMs und erörtert ihre Stärken und Schwächen.
- **Kapitel 19: Generative LLMs evaluieren** listet bekannte Evaluierungsmetriken für Sprachmodelle wie Perplexity, BLEU, ROUGE und BERTScore auf.

**Teil IV: Produktion und Deployment** behandelt Fragen zu praktischen Szenarios wie zum Beispiel die Steigerung der Inferenzgeschwindigkeiten und verschiedene Arten von Verteilungsverschiebungen.

- **Kapitel 20: Zustandsloses und zustandsbehaftetes Training** unterscheidet zwischen zustandslosen und zustandsbehafteten Trainingsmethodiken, die beim Deploying von Modellen verwendet werden.
- **Kapitel 21: Datenzentrierte KI** untersucht datenzentrische KI, bei der die Verfeinerung von Datensätzen zur Verbesserung der Modellperformance im Vordergrund steht. Dieser Ansatz steht im Gegensatz zum herkömmlichen modellzentrierten Ansatz, bei dem es in erster Linie um die Verbesserung von Modellarchitekturen oder -methoden geht.
- **Kapitel 22: Die Inferenz beschleunigen** stellt Techniken vor, um die Geschwindigkeit der Modellinferenz zu erhöhen, ohne die Modellarchitektur anzupassen oder die Genauigkeit zu beeinträchtigen.
- **Kapitel 23: Verschiebungen in der Datenverteilung:** Nach dem Deployment können KI-Modelle mit Diskrepanzen zwischen Trainingsdaten und realen Datenverteilungen konfrontiert werden, die als Verschiebungen in der Datenverteilung bekannt sind. Diese Verschiebungen können die Modellperformance verschlechtern. In diesem Kapitel werden gängige Verschiebungen wie Kovariantenverschiebung, Konzeptdrift, Labeldrift und Domänenverschiebung kategorisiert und näher beleuchtet.

**Teil V: Vorhersageperformance und Modellevaluierung** vertieft verschiedene Aspekte, um die Vorhersageleistung zu optimieren, indem man zum Beispiel die Verlustfunktion ändert, k-fache Kreuzvalidierung einrichtet und mit begrenzt gelabelten Daten arbeitet.

- **Kapitel 24: Poisson- und ordinale Regression** verdeutlicht die Unterschiede zwischen Poisson- und ordinaler Regression. Poisson-Regression kommt für abzählbare Daten infrage, die einer Poisson-Verteilung folgen, beispielsweise die Anzahl der Erkältungen, die sich Personen in einem Flugzeug zugezogen haben. Im Gegensatz dazu eignet sich die ordinale Regression für geordnete kategoriale Daten, bei denen keine äquidistanten Kategorien angenommen werden, wie zum Beispiel beim Schweregrad von Krankheiten.
- **Kapitel 25: Konfidenzintervalle** beschreibt Methoden, mit denen sich Konfidenzintervalle für Klassifikatoren des maschinellen Lernens erstellen lassen. Es wird erläutert, welchen Zweck Konfidenzintervalle haben und wie sie Parameter unbekannter Populationen schätzen. Außerdem lernen Sie Techniken wie Intervalle der Normal-Approximation, Bootstrapping und Retraining mit verschiedenen zufälligen Startwerten kennen.
- **Kapitel 26: Konfidenzintervalle vs. konforme Vorhersagen** erörtert die Unterscheidung zwischen Konfidenzintervallen und konformen Vorhersagen und beschreibt Letztere als Tool, um Vorhersageintervalle zu erstellen, die tatsächliche Ergebnisse mit einer bestimmten Wahrscheinlichkeit abdecken.

- **Kapitel 27: Geeignete Metriken** konzentriert sich auf die wesentlichen Eigenschaften einer geeigneten Metrik in Mathematik und Informatik. Es wird untersucht, ob die beim maschinellen Lernen häufig verwendeten Verlustfunktionen, wie der mittlere quadratische Fehler und der Kreuzentropieverlust, diese Eigenschaften erfüllen.
- **Kapitel 28: Das  $k$  in  $k$ -facher Kreuzvalidierung** untersucht die Rolle des  $k$  in der  $k$ -fachen Kreuzvalidierung sowie die Vor- und Nachteile bei der Wahl eines großen  $k$ .
- **Kapitel 29: Diskordanz zwischen Trainings- und Testdatensatz** befasst sich mit dem Szenario, in dem ein Modell mit dem Testdatensatz besser performt als mit dem Trainingsdatensatz. Bietet Strategien an, um Diskrepanzen zwischen Trainings- und Testdatensätzen zu entdecken und zu behandeln. Dabei lernen Sie das Konzept der adversarialen Validierung kennen.
- **Kapitel 30: Begrenzte gelabelte Daten** stellt verschiedene Techniken vor, um die Modellleistung in Situationen zu verbessern, in denen die Daten begrenzt sind. Behandelt werden Datenbeschriftung, Bootstrapping und Paradigmen wie Transferlernen, aktives Lernen und multimodales Lernen.

## Online-Ressourcen

Auf GitHub habe ich optionales Ergänzungsmaterial mit Codebeispielen für bestimmte Kapitel bereitgestellt, damit Sie Ihre Lernerfahrung verbessern können (siehe <https://github.com/rasbt/MachineLearning-QandAI-book>). Diese Materialien sind als praktische Erweiterungen und zur Vertiefung der im Buch behandelten Themen gedacht. Sie können sie parallel zu den einzelnen Kapiteln verwenden oder erst nach dem Lesen des Buches erkunden, um Ihr Wissen zu festigen und zu erweitern.

Nun genug der Vorrede, gehen wir in medias res!



# Neuronale Netze und Deep Learning



---

# 1 Einbettungen, latenter Raum und Repräsentationen

**Beim Deep Learning sind Begriffe wie *Einbettungsvektoren*, *Repräsentationen* und *latenter Raum* gebräuchlich. Was haben diese Konzepte gemeinsam und wie unterscheiden sie sich?**

---

Auch wenn diese drei Begriffe oft synonym verwendet werden, können wir zwischen ihnen feine Unterscheidungen treffen:

- Einbettungsvektoren sind Repräsentationen von Eingabedaten, bei denen ähnliche Elemente nahe beieinanderliegen.
- Latente Vektoren sind Zwischenrepräsentationen von Eingabedaten.
- Repräsentationen sind codierte Versionen der ursprünglichen Eingabedaten.

Die folgenden Abschnitte untersuchen die Beziehung zwischen Einbettungen, latenten Vektoren und Repräsentationen. Außerdem erfahren Sie, wie sie jeweils im Kontext des maschinellen Lernens Informationen codieren.

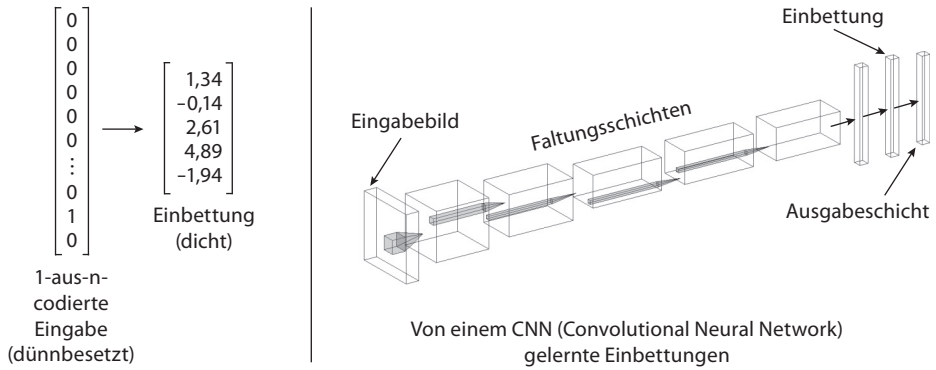
## 1.1 Einbettungen

Einbettungsvektoren – kurz *Einbettungen* – codieren relativ hochdimensionale Daten in relativ niedrigdimensionale Vektoren.

Mithilfe von Einbettungsmethoden können wir einen kontinuierlichen dichten (nicht-dünnbesetzten) Vektor aus einer (dünnbesetzten) 1-aus-n-Codierung (englisch One-hot encoding) erzeugen. Die *1-aus-n-Codierung* ist eine Methode, um kategoriale Daten als binäre Vektoren darzustellen, wobei jede Kategorie auf einen Vektor abgebildet wird, der an der Position, die dem Index der Kategorie entspricht, eine 1, und an allen anderen Positionen eine 0 enthält. Damit ist sichergestellt, dass die kategorialen Werte so dargestellt werden, dass bestimmte Algorithmen für maschinelles Lernen sie verarbeiten können. Wenn wir zum Beispiel eine kategoriale Variable Farbe mit den drei Kategorien Rot, Grün und Blau haben, stellt die 1-aus-n-Codierung Rot als  $[1, 0, 0]$ , Grün als  $[0, 1, 0]$  und Blau als  $[0, 0, 1]$  dar. Diese 1-aus-n-codierten kategorialen Variablen lassen sich dann in kon-

tinuierliche Einbettungsvektoren abbilden, indem die gelernte Gewichtsmatrix einer Einbettungsschicht oder eines Moduls verwendet wird.

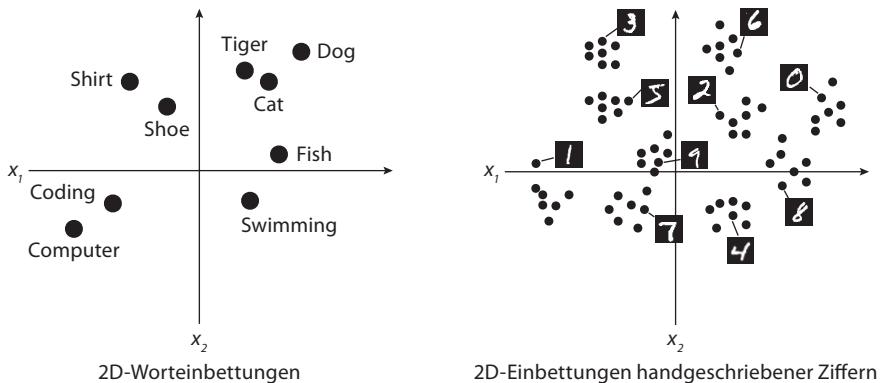
Einbettungsmethoden eignen sich auch für dichte Daten wie Bilder. Zum Beispiel können die letzten Schichten eines Convolutional Neural Networks (CNN) Einbettungsvektoren liefern, wie Abbildung 1-1 veranschaulicht.



**Abb. 1-1** Eine Eingabeeinbettung (links) und eine Einbettung von einem neuronalen Netz (rechts)

Um technisch korrekt zu sein, könnten alle Ausgaben der Zwischenschicht eines neuronalen Netzes Einbettungsvektoren liefern. Je nach Trainingsziel kann auch die Ausgabeschicht nützliche Einbettungsvektoren erzeugen. Der Einfachheit halber assoziiert das Convolutional Neural Network in Abbildung 1-1 die vorletzte Schicht mit Einbettungen.

Es ist möglich, dass Einbettungen eine höhere oder niedrigere Anzahl von Dimensionen haben als die ursprüngliche Eingabe. Mithilfe von Einbettungsmethoden für extreme Ausdrücke lassen sich beispielsweise Daten in zweidimensionale dichte und kontinuierliche Darstellungen für Visualisierungszwecke und Clustering-Analysen codieren, wie Abbildung 1-2 zeigt.



**Abb. 1-2** Abbildung von Wörtern (links) und Bildern (rechts) auf einen zweidimensionalen Merkmalsraum