

3 Maschinenkunst

In diesem Kapitel stellen wir einige der Konzepte vor, die es Deep-Learning-Modellen scheinbar erlauben, Kunst zu erschaffen – ein Gedanke, den einige von uns vermutlich paradox finden. Der Philosoph Alva Noë von der University of Berkeley in Kalifornien meinte jedenfalls: »Kunst kann uns helfen, ein besseres Bild von unserer menschlichen Natur zu formen.«¹ Falls das stimmt, wie können Maschinen dann Kunst erschaffen? Oder anders formuliert: Sind die Kreationen, die von diesen Maschinen stammen, als Kunst anzusehen? Eine andere Interpretation – die uns übrigens am besten gefällt – lautet, dass diese Kreationen tatsächlich Kunst sind und dass die Programmierer Künstler sind, die ihre Deep-Learning-Modelle wie Pinsel handhaben. Wir sind nicht die einzigen, die diese Werke als wahre Kunst betrachten: Von GAN-Algorithmen (*Generative Adversarial Networks*, zu Deutsch etwa: *erzeugende gegnerische Netzwerke*) geschaffene Gemälde sind teils für mehr als 400.000 US-Dollar über den Tisch gegangen.²

In diesem Kapitel werden wir uns die hochentwickelten Konzepte hinter GANs anschauen. Sie werden Beispiele der neuartigen visuellen Werke sehen, die sie produzieren können. Wir werden eine Verbindung zwischen den latenten Räumen, die mit GANs verknüpft sind, und den Wortvektorräumen aus Kapitel 2 ziehen. Und wir werden ein Deep-Learning-Modell behandeln, das als automatisiertes Werkzeug dienen kann, um die Qualität von Fotos drastisch zu verbessern. Aber bevor es losgeht, schnappen Sie sich einen Drink ...

3.1 Eine feuchtfrohliche Nacht

Unter den Google-Büros in Montreal gibt es eine Bar namens »Les 3 Brasseurs«, zu Deutsch also »Die 3 Brauer«. Dort dachte sich Ian Goodfellow, der damals, im Jahre 2014, als PhD-Student in Yoshua Bengios renommiertem Labor (Abbildung 1–10) arbeitete, einen Algorithmus zum Herstellen realistischer aussehender

1. Noë, A. (5. Oktober 2015). What art unveils. *New York Times*.

2. Cohn, G. (25. Oktober 2018). AI art at Christie's sells for \$432.500. *New York Times*.

Bilder aus³ – eine Technik, die Yann LeCun (Abbildung 1–9) als »wichtigsten« aktuellen Durchbruch auf dem Gebiet des Deep Learning bejubelte.⁴

Goodfellows Freunde beschrieben ihm ein *generatives Modell*, an dem sie arbeiteten, das heißt, ein Computermodell, das darauf abzielt, etwas Neues zu erschaffen, sei es ein Zitat im Stil von Shakespeare, eine Melodie oder ein abstraktes Kunstwerk. In ihrem speziellen Fall versuchten die Freunde, ein Modell zu entwerfen, das fotorealistische Bilder generieren konnte, wie etwa Porträts menschlicher Gesichter. Damit dies mit dem traditionellen Machine-Learning-Ansatz einigermaßen gut funktioniert (Abbildung 1–12), müssten die Ingenieure, die das Modell entwarfen, nicht nur die entscheidenden Merkmale von Gesichtern katalogisieren und approximieren, wie Augen, Nasen und Münder, sondern auch exakt abschätzen, wie diese Merkmale relativ zueinander angeordnet werden müssten. Bislang waren die Ergebnisse wenig beeindruckend. Die generierten Gesichter waren entweder sehr unscharf oder ihnen fehlten wichtige Elemente wie die Nase oder die Ohren.

Goodfellow, dessen Kreativität möglicherweise durch das eine oder andere Bier angeregt wurde,⁵ hatte eine revolutionäre Idee: ein Deep-Learning-Modell, in dem zwei künstliche neuronale Netze (*Artificial Neural Network*, ANN) quasi im Wettstreit gegeneinander antreten. Wie in Abbildung 3–1 dargestellt wird, würde eines dieser ANN darauf programmiert werden, Fälschungen herzustellen, während das andere so programmiert würde, dass es als Detektiv agiert und die Fälschungen von den echten Bildern unterscheidet (diese würden separat angeboten werden). Diese gegnerischen Deep-Learning-Netze würden einander anstacheln: Wenn der *Generator* beim Herstellen der Fälschungen besser wird, muss der *Diskriminator* besser dabei werden, sie zu identifizieren, und so müsste der Generator noch überzeugendere Nachahmungen produzieren und so weiter. Dieser wunderbare Trainingszyklus würde schließlich zu überwältigenden neuen Bildern im Stil der echten Trainingsbilder führen, ob nun von Gesichtern oder anderen Dingen. Und das Beste an der ganzen Sache wäre, dass Goodfellows Ansatz uns der Notwendigkeit entheben würde, manuell Features in das generative Modell zu programmieren. Wie wir schon im Zusammenhang mit dem maschinellen Sehen (Kapitel 1) und der Verarbeitung natürlicher Sprache (Kapitel 2) ausgeführt haben, kümmert sich das Deep Learning automatisch um die Features der Modelle.

3. Giles, M. (21. Februar 2018). The GANfather: The man who's given machines the gift of imagination. *MIT Technology Review*.

4. LeCun, Y. (28. Juli 2016). *Quora*. bit.ly/DLbreakthru

5. Jarosz, A., et al. (2012). Uncorking the muse: Alcohol intoxication facilitates creative problem solving. *Consciousness and Cognition*. 21, 487–93.

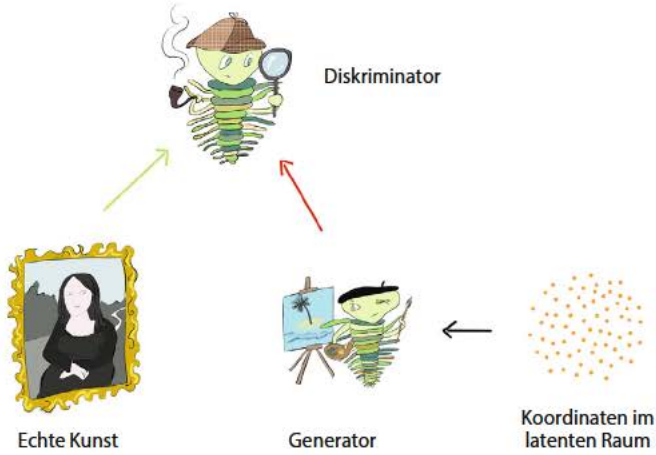


Abb. 3-1 Schematische Darstellung eines Generative Adversarial Network (GAN). Dem Diskriminator werden sowohl echte Bilder als auch Nachahmungen vorgelegt. Er hat die Aufgabe, die echten Bilder zu identifizieren. Die orange Wolke repräsentiert die Orientierungshilfe durch den latenten Raum (Abbildung 3-4), die dem Fälscher angeboten wird. Diese Lenkung kann entweder zufällig sein (wie das beim Netzwerktraining im Allgemeinen der Fall ist; siehe Kapitel 12) oder kann selektiv (während einer Erkundung nach dem Training; siehe Abbildung 3-3) erfolgen.

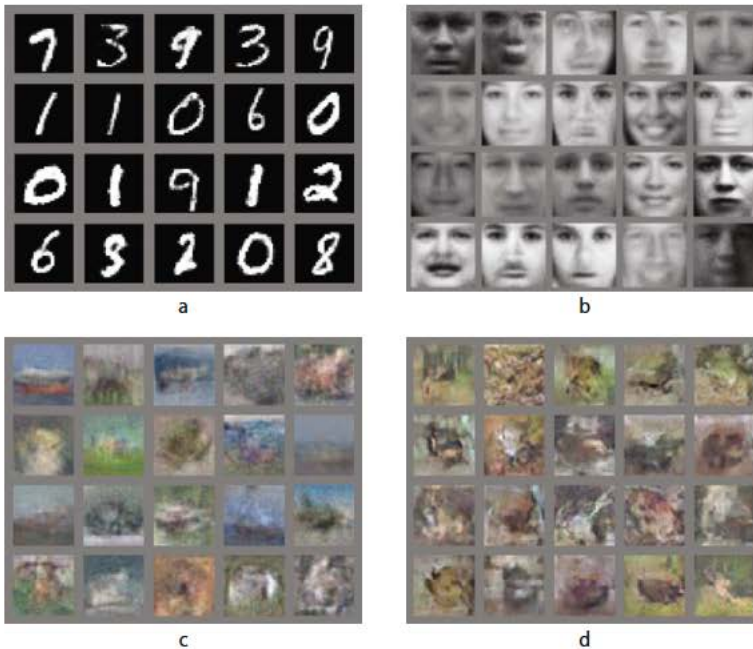


Abb. 3-2 Ergebnisse aus dem GAN-Artikel von Goodfellow und seinen Kollegen aus dem Jahre 2014

Goodfellows Freunde bezweifelten, dass sein fantasievolles Vorgehen funktionieren würde. Als er nach Hause kam und seine Freundin schlafend vorfand, arbeitete er bis in die Nacht daran, seinen Entwurf mit den zwei ANN umzusetzen. Das Ganze funktionierte beim ersten Versuch, und die erstaunliche Familie der Generative Adversarial Networks war geboren!

Im selben Jahr präsentierten Goodfellow und seine Kollegen GANs auf der angesehenen *Neural Information Processing Systems (NIPS)*⁶ *Conference* der Öffentlichkeit. Einige ihrer Ergebnisse sind in Abbildung 3–2 zu sehen. Ihr GAN erzeugte diese neuartigen Bilder, nachdem es mit (a) handgeschriebenen Ziffern⁷, (b) Fotos von menschlichen Gesichtern⁸ sowie mit (c) und (d) Fotos aus zehn unterschiedlichen Klassen (z.B. Flugzeuge, Autos, Hunde)⁹ trainiert worden war. Die Ergebnisse in (c) sind merklich weniger knackig als die in (d), weil das GAN, das die Bilder aus (d) produziert hatte, speziell für das maschinelle Sehen ausgerichtete Neuronenschichten besaß, sogenannte *Convolutional Layers* oder *Konvolutionsschichten*¹⁰, während das GAN aus (c) nur einen allgemeineren Typ verwendete.¹¹

3.2 Berechnungen auf nachgemachten menschlichen Gesichtern

Nach Goodfellows Einstieg bestimmte ein Forschungsteam unter der Leitung des amerikanischen Machine-Learning-Forschers Alec Radford die architektonischen Bedingungen für GANs, die zu einer deutlich realistischeren Bildherstellung führten. Einige der Beispiele für Porträts von nachgemachten Menschen, die von ihren *tiefen convolutional* GANs erzeugt wurden, sind in Abbildung 3–3 zu sehen. In ihrem Artikel demonstrierten Radford und seine Kollegen ziemlich geschickt die Interpolation durch den und die Berechnungen mit dem *latenten Raum*, der mit GANs¹² assoziiert wird. Untersuchen wir zunächst einmal, was latenter Raum ist, bevor wir uns mit Interpolationen und Berechnungen mit dem latenten Raum befassen.

-
6. Goodfellow, I., et al. (2014). Generative adversarial networks. *arXiv:1406.2661*.
 7. Aus LeCuns klassischem MNIST-Datensatz, den wir selbst in Teil II benutzen werden.
 8. Aus der *Toronto Face*-Datenbank der Forschungsgruppe von Hinton (Abbildung 1–15).
 9. Der CIFAR-10-Datensatz, benannt nach dem *Canadian Institute for Advanced Research*, das dessen Erschaffung unterstützt hat.
 10. Wir werden in Kapitel 10 ausführlicher auf diesen Schichttyp eingehen.
 11. Vollständig verbundene Schichten (*Dense Layers*), die in Kapitel 4 eingeführt und in Kapitel 7 ausführlicher behandelt werden.
 12. Radford et al. (2016). Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv:1511.06434v2*.

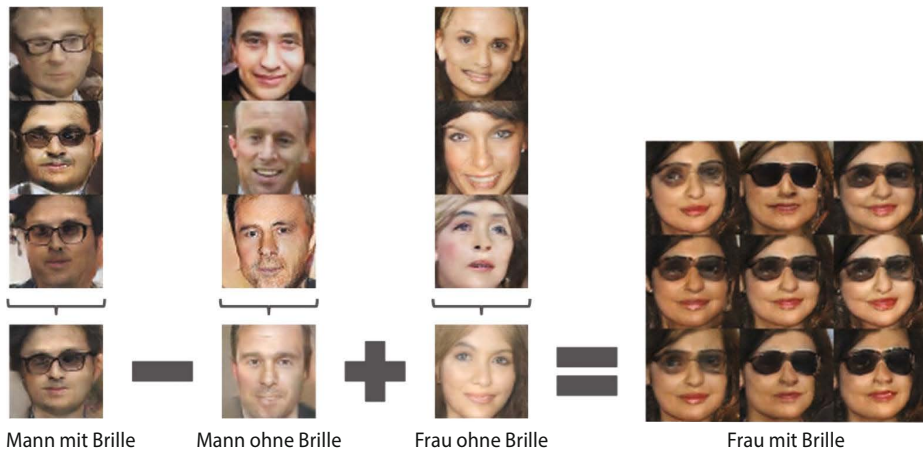


Abb. 3-3 Ein Beispiel für die Berechnungen mithilfe des latenten Raums aus Radford et al. (2016)

Die Skizze aus Abbildung 3–4 erinnert vielleicht an den Wortvektorraum aus Abbildung 2–6. Tatsächlich gibt es drei große Ähnlichkeiten zwischen latenten Räumen und Vektorräumen:

1. Während die Skizze aus Gründen der Einfachheit und Verständlichkeit nur einen dreidimensionalen Raum darstellt, sind latente Räume n -dimensional mit meist mehreren Hundert Dimensionen. Der latente Raum des GAN, das Sie in Kapitel 12 selbst herstellen werden, wird zum Beispiel $n=100$ Dimensionen besitzen.
2. Je dichter zwei Punkte im latenten Raum zueinander liegen, umso ähnlicher sind die Bilder, die diese Punkte repräsentieren.
3. Eine Bewegung durch den latenten Raum in eine bestimmte Richtung kann einer schrittweisen Änderung in einem dargestellten Konzept entsprechen. Im Fall von fotorealistischen Gesichtern könnte dies das Alter oder das Geschlecht sein.

Indem wir zwei Punkte auswählen, die weit voneinander entfernt auf einer n -dimensionalen Achse liegen, die das Alter repräsentiert, könnten wir bei einer Interpolation zwischen ihnen und einer Auswahl beliebiger Punkte von dieser interpolierten Linie einen (nachgeahmten) Mann finden, der schrittweise immer älter zu werden scheint.¹³ In unserer Skizze des latenten Raums (Abbildung 3–4) ist eine solche Achse für das »Alter« lila dargestellt. Wenn Sie die Interpolation durch einen authentischen latenten Raum eines GAN beobachten wollen, sollten

13. Ein technischer Hinweis: Wie es bei Vektorräumen der Fall ist, kann diese »Alter«-Achse (oder eine andere Richtung im latenten Raum, die ein sinnvolles Attribut repräsentiert) orthogonal zu allen anderen n Dimensionen liegen, die die Achsen des n -dimensionalen Raums bilden. Wir werden dies in Kapitel 11 näher ausführen.

Sie einen Blick in den Artikel von Radford und seinen Kollegen werfen und dort zum Beispiel nach den sanften Drehungen des »Photo Angle« synthetischer Schlafzimmer suchen. Zum Zeitpunkt der Entstehung dieses Buches konnte man die neuesten Entwicklungen im Bereich der GANs unter bit.ly/InterpCeleb anschauen. Dieses Video, das von Forschern des Grafikkartenherstellers Nvidia produziert wurde, bietet eine atemberaubende Interpolation durch hochwertige Porträt-»Fotografien« von B- und C-Promis.^{14, 15}

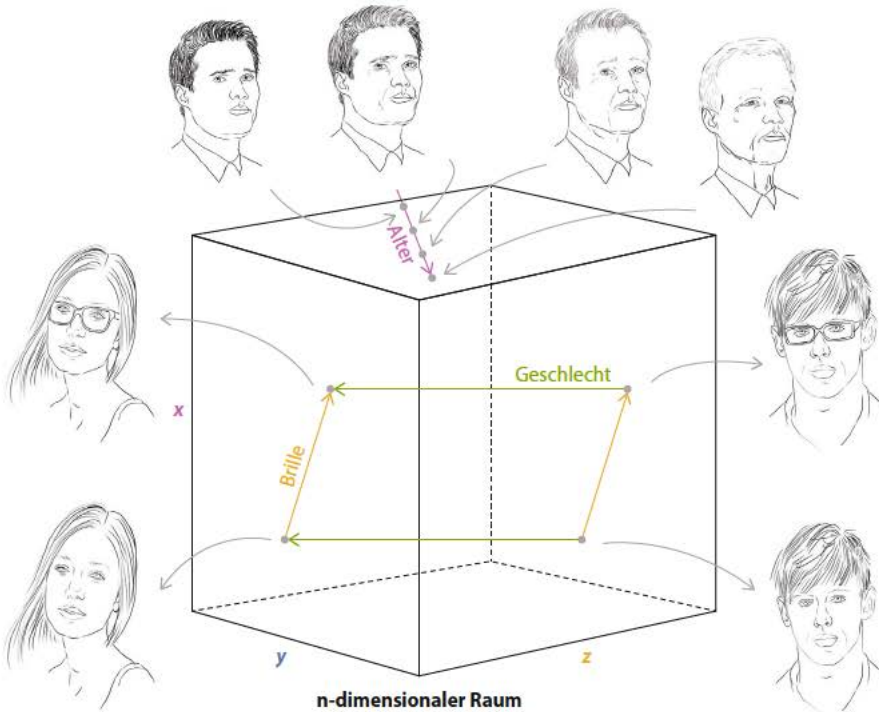


Abb. 3–4 Eine Skizze des latenten Raums von GANs (Generative Adversarial Networks). Wenn man sich an dem lila Pfeil entlangbewegt, erhält man Bilder, die die Alterung einer Person zeigen, die ansonsten annähernd gleich aussieht. Der grüne Pfeil repräsentiert das Geschlecht und der orange Pfeil das Vorhandensein einer Brille in dem Gesicht.

Wenn Sie nun mit dem Gelernten einen Schritt weiter gehen wollen, könnten Sie Berechnungen mit den Bildern anstellen, die Sie dem latenten Raum des GAN entnommen haben. Ein Punkt innerhalb des latenten Raums kann durch die Koordinaten seines Standortes dargestellt werden – der daraus resultierende Vek-

14. Karras, t. et al. (2018). Progressive growing of GANs for improved quality, stability and variation. *Proceedings of the International Conference on Learning Representations*.
15. Wenn Sie selbst einmal versuchen wollen, zwischen echten und GAN-generierten Gesichtern zu unterscheiden, besuchen Sie whichfaceisreal.com.

tor ist analog den Wortvektoren, die in Kapitel 2 beschrieben wurden. Wie bei den Wortvektoren können Sie Berechnungen mit diesen Vektoren durchführen und den latenten Raum auf semantische Weise durchlaufen. Abbildung 3–3 zeigt ein Beispiel für Berechnungen im latenten Raum aus der Arbeit von Radford und seinen Kollegen. Wir beginnen mit einem Punkt im latenten Raum ihres GANs, der einen Mann mit einer Brille repräsentiert, subtrahieren einen Punkt, der einen Mann *ohne* Brille darstellt, und addieren dann einen Punkt, der eine *Frau* ohne Brille repräsentiert. Der Punkt, den wir als Ergebnis erhalten, existiert im latenten Raum in der Nähe der Bilder, die eine Frau *mit* einer Brille darstellen. Unsere Skizze in Abbildung 3–4 verdeutlicht, wie die Beziehungen zwischen den Bedeutungen im latenten Raum gespeichert werden (auch hier wieder vergleichbar zum Wortvektorraum) und damit die Berechnungen im latenten Raum umsetzen.

3.3 Stiltransfer: Fotos in einen Monet verwandeln (und umgekehrt)

Eine besonders bezaubernde Anwendung von GANs ist der Stiltransfer. Zhu, Park und ihre Mitarbeiter aus dem *Berkeley Artificial Intelligence Research (BAIR) Lab* stellten eine neue Richtung von GAN¹⁶ vor, die erstaunliche Beispiele dafür liefert, wie Abbildung 3–5 zeigt. Alexei Efros, einer der Koautoren des Artikels, nahm im Frankreich-Urlaub Fotos auf, und die Forscher verwandelten diese Fotos dann mithilfe ihres *CycleGAN* in Bilder, die dem jeweiligen Stil des impressionistischen Malers Claude Monet, des Niederländers Vincent van Gogh, des japanischen Genres Ukiyo-E und anderer nachempfunden sind.

Unter bit.ly/cycleGAN können Sie Beispiele für den umgekehrten Fall entdecken (Monet-Gemälde, die in fotorealistische Bilder umgewandelt wurden) sowie:

- n Sommerszenen, die in Winterszenen verwandelt wurden, und umgekehrt
- n Körbe mit Äpfeln, aus denen Körbe mit Orangen wurden, und umgekehrt
- n Flache Fotos niedriger Qualität, die auf einmal aussehen, als wären sie mit High-End-(Spiegelreflex-)Kameras aufgenommen worden
- n ein Video eines Pferdes, das sich in ein Zebra verwandelt
- n ein Video einer am Tage aufgenommenen Fahrt, die in eine Nachtfahrt konvertiert wird

16. Sie werden auch »CycleGANs« genannt, weil sie über viele Zyklen des Netzwerktrainings die Konsistenz des Bildes erhalten. Zhu, J.-Y., et al. (2017). Unpaired image-to-image translation using cycle-consistent adversarial networks. *arXiv:1703.10593*.

10 Maschinelles Sehen

Willkommen zu Teil III, liebe Leser! In Teil I gaben wir Ihnen einen Überblick über spezielle Anwendungen des Deep Learning. Mit der grundlegenden Theorie, die wir anschließend in Teil II behandelt haben, sind Sie nun bestens ausgestattet, um sich einer Reihe von Anwendungen aus verschiedenen Bereichen zu widmen. Dies wird vor allem in Form von Codebeispielen geschehen. In diesem Kapitel untersuchen Sie z. B. Convolutional Neural Networks und setzen diese für Aufgaben aus dem Gebiet des maschinellen Sehens (Machine Vision) ein. Außerdem behandeln wir im weiteren Verlauf von Teil III praktische Beispiele für:

- n Recurrent Neural Networks zur Verarbeitung natürlicher Sprache (in Kapitel 11)
- n Generative Adversarial Networks für eine visuelle Kreativität (in Kapitel 12)
- n Deep Reinforcement Learning für eine sequenzielle Entscheidungsfindung in komplexen, veränderlichen Umgebungen (in Kapitel 13)

10.1 Convolutional Neural Networks

Ein *convolutional*, zu Deutsch »faltendes«, neuronales Netz – auch *Convolutional Neural Network*, ConvNet oder CNN genannt – ist ein künstliches neuronales Netz, das ein oder mehrere Konvolutionsschichten oder *Convolutional Layers* aufweist. Dieser Schichttyp erlaubt es einem Deep-Learning-Modell, effizient räumliche Muster zu verarbeiten. Wie Sie in diesem Kapitel aus erster Hand erfahren werden, eignen sich solche Netzwerke daher besonders gut für Anwendungen aus dem Bereich des maschinellen Sehens oder der Computer Vision.

10.1.1 Die zweidimensionale Struktur der visuellen Bilddarstellung

In unseren bisherigen Codebeispielen mit den handgeschriebenen MNIST-Ziffern wandelten wir die Bilddaten in ein eindimensionales Array aus Zahlen um, die wir dann in die vollständig verbundene verborgene Schicht eingaben. Genauer gesagt, begannen wir mit 28×28 Pixel großen Graustufenbildern und verwandelten diese in eindimensionale Arrays mit 784 Elementen.¹ Obwohl dieser Schritt im Kontext eines vollständig verbundenen Netzwerks notwendig war (wir mussten die 784 Pixelwerte egalisieren, damit sie jeweils in ein Neuron der ersten verborgenen Schicht eingegeben werden konnten), bringt die Reduzierung eines zweidimensionalen Bildes in eine Dimension einen beträchtlichen Strukturverlust mit sich. Wenn Sie mit einem Stift eine Ziffer auf Papier zeichnen, stellen Sie sich das Ganze nicht als eine kontinuierliche lineare Abfolge von Pixeln vor, die von oben links nach unten rechts verlaufen. Würden wir z.B. hier eine MNIST-Ziffer als 784 Pixel langen Strom aus Graustufen ausdrucken, dann wetten wir, dass Sie nicht in der Lage wären, die Ziffer zu erkennen. Menschen nehmen visuelle Informationen nämlich in einer zweidimensionalen Form wahr,² und unsere Fähigkeit, zu erkennen, was wir anschauen, ist von Natur aus mit den räumlichen Beziehungen zwischen den Formen und Farben verbunden, die wir wahrnehmen.

10.1.2 Berechnungskomplexität

Neben dem Verlust der zweidimensionalen Struktur beim Reduzieren eines Bildes müssen wir noch eine zweite Sache bedenken, wenn wir Bilder in ein vollständig verbundenes Netzwerk leiten: die Berechnungskomplexität. Die MNIST-Bilder sind sehr klein – 28×28 Pixel mit nur einem *Kanal*. (Es gibt nur einen Farb»kanal«, weil MNIST-Ziffern monochromatisch sind; um vollfarbige Bilder darzustellen, sind wenigstens drei Kanäle – üblicherweise Rot, Grün und Blau – erforderlich.) Wenn man MNIST-Bildinformationen in eine vollständig verbundene Schicht übergibt, entspricht dies 785 Parametern pro Neuron: 784 Gewichte für jedes der Pixel plus der Bias des Neurons. Würden wir jedoch ein mäßig großes Bild verarbeiten – sagen wir ein 200×200 Pixel großes vollfarbiges RGB-Bild³ – nimmt die Anzahl der Parameter ganz drastisch zu. In diesem Fall hätten wir drei Farbkanäle, jeweils mit 40.000 Pixeln, was insgesamt 120.001 Parametern pro Neuron in einer vollständig verbundenen Schicht entspricht.⁴ Bei einer bescheidenen

-
1. Erinnern Sie sich daran, dass die Pixelwerte durch 255 dividiert wurden, um alles in den Bereich $[0 : 1]$ zu skalieren.
 2. Na ja ... dreidimensional, aber wir wollen die Tiefe bei dieser Diskussion einfach einmal ignorieren.
 3. Die erforderlichen *Rot*-, *Grün*- und *Blau*-Kanäle für ein vollfarbiges Bild.
 4. $200 \text{ Pixel} \times 200 \text{ Pixel} \times 3 \text{ Farbkanäle} + 1 \text{ Bias} = 120.001 \text{ Parameter}$.

Anzahl an Neuronen in der Schicht – sagen wir, 64 – wären das fast 8 Millionen Parameter allein für die erste verborgene Schicht unseres Netzwerks.⁵ Außerdem hat das Bild nur 200×200 Pixel – das sind kaum 0,4 MP⁶, während die meisten modernen Smartphones Kamerasensoren mit 12 MP oder mehr haben. Im Allgemeinen müssen Machine-Vision-Aufgaben nicht auf hochaufgelösten Bildern ausgeführt werden, um erfolgreich zu sein, aber der Punkt ist doch klar: Bilder können sehr viele Datenpunkte enthalten, und wenn man diese naiverweise in einer vollständig verbundenen Struktur verarbeiten möchte, explodieren die Anforderungen an die Rechenkapazitäten des neuronalen Netzes geradezu.

10.1.3 Konvolutionsschichten

Konvolutionsschichten (Convolutional Layers) bestehen aus Gruppen von *Kernen*, die auch als *Filter* oder Filterkernel bezeichnet werden. Jeder dieser Kernel ist ein kleines Fenster (ein sogenanntes *Patch*), das das Bild von oben links nach unten rechts scannt (technisch ausgedrückt: Der Filter führt eine *Konvolution* (Faltung) durch). Abbildung 10–1 verdeutlicht diese *Konvolutionsoperation*.

Kernel bestehen aus Gewichten, die – wie in vollständig verbundenen Schichten – durch Backpropagation gelernt werden. Ihre Größe ist unterschiedlich, typischerweise sind sie 3×3 groß, und so werden wir das auch in den Beispielen in diesem Kapitel nutzen.⁷ Für die monochromatischen MNIST-Ziffern würde dieses 3×3 -Pixel-Fenster aus $3 \times 3 \times 1$ Gewichten bestehen – d.h., aus neun Gewichten für insgesamt 10 Parameter. (Wie ein künstliches Neuron in einer vollständig verbundenen Schicht hat auch jeder Konvolutionsfilter einen Bias-Term *b*.) Zum Vergleich: Würden wir mit vollfarbigen RGB-Bildern arbeiten, dann hätte ein Kernel, der die gleiche Anzahl an Pixeln abdeckt, dreimal so viele Gewichte – $3 \times 3 \times 3$ für insgesamt 27 Gewichte und 28 Parameter.

5. $64 \text{ Neuronen} \times 120.001 \text{ Parameter pro Neuron} = 7.680.064 \text{ Parameter}$.

6. Megapixel.

7. Eine andere typische Größe ist 5×5 , wobei Kernel, die größer sind als diese, nur selten verwendet werden.

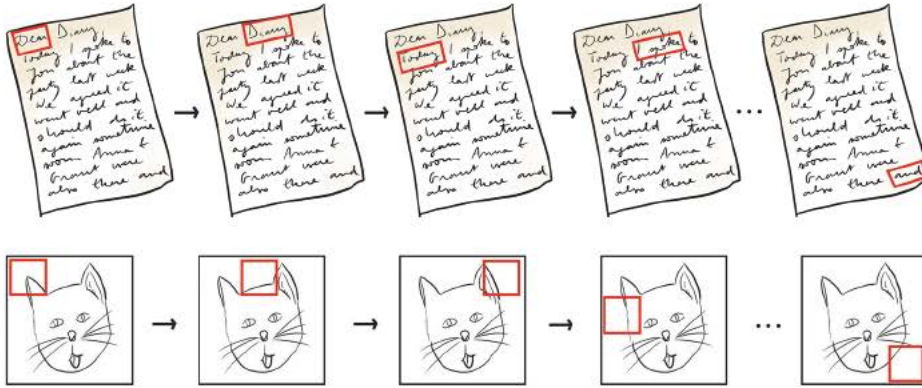


Abb. 10-1 Wenn wir die Seite eines Buches lesen, das in Englisch (oder Deutsch) geschrieben ist, beginnen wir in der obersten linken Ecke und lesen nach rechts. Immer wenn wir das Ende einer Textzeile erreichen, machen wir in der nächsten Zeile weiter. Auf diese Weise erreichen wir irgendwann die untere rechte Ecke und haben damit alle Wörter auf der Seite gelesen. Analog beginnt der Kernel in einer Konvolutionsschicht in einem kleinen Fenster aus Pixeln in der oberen linken Ecke eines Bildes. Der Kernel scannt von der oberen Reihe nach unten, und zwar von links nach rechts, bis er schließlich die untere rechte Ecke erreicht hat. Damit hat er dann alle Pixel des Bildes gescannt.

Wie in Abbildung 10-1 dargestellt ist, besetzt der Kernel während seiner Faltungsoperation diskrete Positionen entlang dem Bild. Bleiben wir für unsere Erklärung einmal bei der Kernel-Größe 3×3 . Während der Forwardpropagation wird eine multidimensionale Variante der »wichtigsten Gleichung in diesem Buch« $w \cdot x + b$ (vorgestellt in Abbildung 6-7) an jeder Position berechnet, die der Kernel während der Konvolution entlang dem Bild einnimmt. Wenn wir das 3×3 -Fenster aus Pixeln und den 3×3 -Kernel in Abbildung 10-2 als Eingaben x bzw. Gewichte w annehmen, dann können wir die Berechnung der gewichteten Summe $w \cdot x$ demonstrieren, bei der Produkte elementweise berechnet werden. Die Berechnung basiert dabei auf der Ausrichtung der vertikalen und horizontalen Stellen. Es hilft, wenn man sich vorstellt, dass der Kernel über die Pixelwerte gelegt ist. Die Berechnung ist hier:

$$\begin{aligned}
 w \cdot x &= 0,01 \cdot 0,53 + 0,09 \cdot 0,34 + 0,22 \cdot 0,06 \\
 &+ -1,36 \cdot 0,37 + 0,34 \cdot 0,82 + -1,59 \cdot 0,01 \\
 &+ 0,13 \cdot 0,62 + 0,69 \cdot 0,91 + 1,02 \cdot 0,34 \\
 &= -0,3917
 \end{aligned}$$

Gleichung 10-1

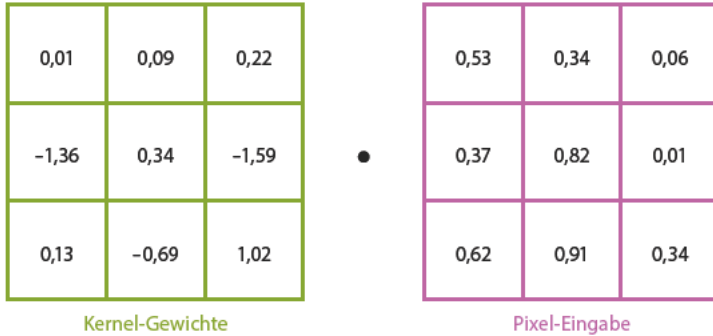


Abb. 10-2 Ein 3×3 -Kernel und ein 3×3 -Pixelfenster

Jetzt addieren wir entsprechend Gleichung 7-1 noch einen Bias b (sagen wir, 0,19), um z zu erhalten:

$$\begin{aligned}
 z &= w \cdot x + b \\
 &= -0,39 + b \\
 &= -0,39 + 0,20 \\
 &= -0,19
 \end{aligned}$$

Gleichung 10-2

Mit z können wir schließlich einen Aktivierungswert a berechnen, indem wir z an die Aktivierungsfunktion unserer Wahl übergeben, etwa an die Tanh-Funktion oder an die ReLU-Funktion.

Sie sehen, dass sich die grundlegenden Rechenoperationen im Vergleich zu den künstlichen Neuronen aus den Kapiteln 6 und 7 nicht geändert haben. Faltdende Kernel haben **Gewichte**, **Eingaben** und einen **Bias**. Daraus wird mithilfe unserer wichtigsten Gleichung eine gewichtete Summe erzeugt und das daraus resultierende z wird an irgendeine nichtlineare Funktion übergeben, um eine **Aktivierung** zu erzeugen. Geändert hat sich, dass es nicht für *jede* Eingabe ein Gewicht gibt, sondern einen diskreten Kernel mit 3×3 Gewichten. Diese Gewichte ändern sich nicht, während der Kernel die Konvolution durchführt, sondern sie werden über *alle* Eingaben verwendet. Auf diese Weise kann eine Konvolutionsschicht um Größenordnungen weniger Gewichte haben als eine vollständig verbundene Schicht. Ein weiterer wichtiger Punkt ist, dass die Ausgaben dieses Kernels (alle Aktivierungen) genau wie die Eingaben in einem zweidimensionalen Array angeordnet sind. Wir werden uns das gleich genauer anschauen, aber zuvor ...

10.1.4 Mehrere Filter

Typischerweise haben wir mehrere Filter in einer Konvolutionsschicht. Jeder Filter erlaubt es dem Netzwerk, eine Repräsentation der Daten in einer bestimmten Schicht auf einzigartige Weise zu lernen. Falls z. B. analog zu den einfachen Zellen von Hubel und Wiesel im biologischen visuellen System (Abbildung 1–5) die erste verborgene Schicht in unserem Netzwerk eine Konvolutionsschicht ist, dann könnte sie einen Kernel enthalten, der optimal auf vertikale Linien reagiert. Das bedeutet: Immer wenn er über eine vertikale Linie im Eingabebild gleitet, erzeugt er einen großen Aktivierungswert (a). Weitere Kernel in dieser Schicht können lernen, andere einfache räumliche Eigenschaften darzustellen, wie etwa horizontale Linien und Farbübergänge. (Schauen Sie sich z. B. das Bild unten links in Abbildung 1–17 an.) Auf diese Weise kamen die Kernel zu ihrer Bezeichnung *Filter*: Sie überfliegen das Bild und filtern den Ort bestimmter Features heraus. Dabei erzeugen sie hohe Aktivierungen, wenn sie einem Muster, einer Form und/oder einer Farbe begegnen, auf deren Erkennung sie spezialisiert sind. Man könnte sagen, dass sie als Marker fungieren, die ein zweidimensionales Array aus Aktivierungen erzeugen, die andeuten, *wo* das besondere Feature dieses Filters im Originalbild existiert. Aus diesem Grund wird die Ausgabe eines Kernels als *Aktivierungs-Map* bezeichnet.

Analog zu den hierarchischen Repräsentationen des biologischen visuellen Systems (Abbildung 1–6) erhalten nachfolgende Konvolutionsschichten diese Aktivierungs-Maps als Eingaben. Je tiefer das Netzwerk wird, umso komplexer werden die Kombinationen dieser simplen Features, auf die die Filter in diesen Schichten reagieren, d. h., umso abstrakter werden die räumlichen Muster, die repräsentiert werden. So baut sich eine Hierarchie aus einfachen Linien und Farben nach und nach zu komplexen Texturen und Formen auf (siehe die Einzelbilder unten in Abbildung 1–17). Auf diese Weise haben spätere Schichten innerhalb des Netzwerks die Fähigkeit, ganze Objekte zu erkennen oder sogar ein Bild einer Dänischen Dogge von dem eines Yorkshire Terriers zu unterscheiden.

Die Anzahl der Filter in der Schicht ist genauso wie die Anzahl der Neuronen in einer vollständig verbundenen Schicht ein Hyperparameter, den wir selbst konfigurieren können. Wie bei den anderen Hyperparametern, die wir bereits besprochen haben, gibt es auch hier einen Goldlöckchen-haften Mittelweg für die Anzahl der Filter. Mithilfe dieser Faustregeln können Sie die nötigen Filter für ein bestimmtes Problem festlegen:

- n Eine größere Anzahl an Kernen erlaubt die Identifizierung komplexerer Features. Sie sollten daher die Komplexität der Daten und des zu lösenden Problems in Betracht ziehen. Natürlich geht eine größere Anzahl an Kernen zulasten der Berechnungseffizienz.
- n Wenn ein Netzwerk mehrere Konvolutionsschichten besitzt, kann die optimale Anzahl der Kernel von Schicht zu Schicht stark schwanken. Denken Sie daran, dass frühere Schichten einfache Features erkennen, während spätere Schichten komplexe Neukombinationen dieser einfachen Features identifizieren. Berücksichtigen Sie dies beim Entwurf des Netzwerks. Wie wir bei den Codebeispielen für CNNs weiter hinten in diesem Kapitel sehen werden, besteht ein gern genommener Ansatz darin, in späteren Konvolutionsschichten im Verhältnis mehr Kernel zu verwenden als in früheren Schichten.
- n Streben Sie wie immer danach, die Berechnungskomplexität zu minimieren: Versuchen Sie, die kleinstmögliche Anzahl an Kernen zu benutzen, die niedrige Kosten für Ihre Validierungsdaten verursacht. Falls eine Verdoppelung der Anzahl an Kernen (etwa von 32 auf 64 auf 128) in einer Schicht die Validierungskosten Ihres Modells merklich senkt, sollten Sie möglicherweise den höheren Wert nehmen. Erhöht eine Halbierung der Anzahl (etwa von 32 auf 16 auf 8) die Validierungskosten des Modells nicht, dann sollten Sie den niedrigeren Wert nutzen.

10.1.5 Ein Beispiel für Konvolutionsschichten

Konvolutionsschichten (Faltungsschichten) sind eine nichttriviale Abweichung von den einfacheren vollständig verbundenen Schichten aus Teil II. Damit Sie verstehen, wie die Pixelwerte und Gewichte gemeinsam Feature-Maps erzeugen, zeigen wir in den Abbildungen 10–3 bis 10–5 ein ausführliches Beispiel und erläutern auch die Berechnungen, die damit einhergehen. Stellen Sie sich zu Beginn vor, wir falten über ein einzelnes RGB-Bild von 3×3 Pixeln Größe. In Python werden diese Daten in einem $[3,3,3]$ -Array gespeichert, wie oben in Abbildung 10–3 zu sehen ist.⁸

8. Wir geben zu, dass das RGB-Beispiel des Baumes deutlich mehr als nur neun Pixel hat, aber wir hatten Probleme, ein passendes Farbbild zu finden, das 3×3 Pixel groß ist.

Schlüsselkonzepte

Hier sind die Schlüsselkonzepte aus diesem Buch aufgelistet. Das letzte Konzept – das wir in diesem Kapitel behandelt haben – ist **lila** hervorgehoben.

- n Parameter:
 - Gewicht w
 - Bias b
- n Aktivierung a
- n künstliche Neuronen:
 - Sigmoid
 - Tanh
 - ReLU
 - linear
- n Eingabeschicht
- n verborgene Schicht
- n Ausgabeschicht
- n Schichttypen:
 - vollständig verbunden (»dense«)
 - Softmax
 - Konvolutionsschicht (convolutional)
 - De-Convolutional
 - Max-Pooling
 - Upsampling
 - Flatten-
 - Einbettungsschicht
 - RNN
 - (bidirektional-)LSTM
 - Verkettungsschicht
- n Kosten-(Verlust-)funktionen:
 - quadratisch (mittlere quadratische Abweichung)
 - Kreuzentropie
- n Forwardpropagation
- n Backpropagation
- n instabile (besonders verschwindende) Gradienten
- n Glorot-Gewichtsinitialisierung
- n Batch-Normalisierung
- n Dropout

- n Optimierer:
 - stochastischer Gradientenabstieg
 - Adam
- n Optimierer-Hyperparameter:
 - Lernrate η
 - Batch-Größe
- n word2vec
- n GAN-Komponenten:
 - Diskriminator-Netzwerk
 - Generator-Netzwerk
 - Adversarial-Netzwerk
- n Deep-Q-Learning