

# Neuronale Netze: Ein Blick in die Black Box

[Originalartikel](#)

[Backup](#)

Nina Schaaf 14. Januar 2020



Wieso sind neuronale Netze eine Black Box? &#169; Adobe: yingyaipumi / stock.adobe.com / 286542323 alt=„Wieso sind neuronale Netze eine Black Box? &#169; Adobe: yingyaipumi / stock.adobe.com / 286542323“ src=„[https://www.informatik-aktuell.de/fileadmin/\\_processed\\_/b/a/csm\\_720-AdobeStock\\_286542323\\_20e4059e6f.jpg](https://www.informatik-aktuell.de/fileadmin/_processed_/b/a/csm_720-AdobeStock_286542323_20e4059e6f.jpg)“ width=„350“ height=„196“ referrerpolicy=„no-referrer“ />

Wieso sind neuronale Netze eine Black Box? &#169; Adobe: yingyaipumi

K&#252;nstliche Intelligenz und allem voran Deep Learning ist momentan in aller Munde. Hierbei dominiert gr&#246;&#223;tenteils die Diskussion um die gesellschaftlichen Auswirkungen, welche meist zwischen Utopie und Horrorszenarien schwankt. In dieser teils sehr aufgeregten Debatte kann schnell untergehen, was eigentlich hinter den allgegenw&#228;rtigen Begriffen KI und Deep Learning steckt. Wie sind k&#252;nstliche neuronale Netze aufgebaut und wie funktionieren sie? Dar&#252;ber hinaus lohnt sich ein Blick auf eine der Herausforderungen neuronaler Netze: deren „Black-Box“-Charakter. Wieso sind neuronale Netze eine Black Box, f&#252;r welche Anwendungen ist dies eher kritisch und welche L&#246;sungsans&#228;tze existieren bereits?

Denkt man an Anwendungen mit K&#252;nstlicher Intelligenz (KI), so kommen den meisten Menschen sicherlich Sprachassistenten, wie Alexa und Google Assistant, oder das autonome Fahren in den Sinn. In der &#246;ffentlichen Debatte um KI dominieren Meldungen dieser Dom&#228;nen, seien es Fehlschl&#228;ge (ein autonomes Fahrzeug von Uber t&#246;tet eine Fu&#223;g&#228;ngerin <a href=„<https://www.informatik-aktuell.de/betrieb/kuenstliche-intelligenz/neuronale-netze-ein-blick-in-di-e-black-box.html#c30659>“>[1]</a>) oder Erfolgsmeldungen (Google Assistant bucht am Telefon einen Friseurtermin <a href=„<https://www.informatik-aktuell.de/betrieb/kuenstliche-intelligenz/neuronale-netze-ein-blick-in-di-e-black-box.html#c30659>“>[2]</a>).

An einen M&#252;nzsortierautomaten, der KI einsetzt, h&#228;tten wohl die wenigsten gedacht. Doch auch hier k&#246;nnen intelligente Systeme heute schon helfen, schwer beherrschbare Anforderungen wie Verarbeitung in Echtzeit, hohe Qualit&#228;tsanspr&#252;che, geringe Fehlertoleranzen und strikte gesetzliche Vorgaben zusammenzuf&#252;hren. So verlangt eine EU-Richtlinie von Betreibern von M&#252;nzsortierautomaten, dass die Automaten die zugef&#252;hrten M&#252;nzen zuverl&#228;ssig nach Originalen, Fremdw&#228;hrungen und Falschgeld unterscheiden k&#246;nnen. Da M&#252;nzgeld jedoch, anders als Banknoten, keine spezifischen Sicherheitsmerkmale aufweist, m&#252;ssen alternative und nichtsdestotrotz hochzuverl&#228;ssige Pr&#252;fmethoden entwickelt werden. Maschinelle Lernverfahren (ML) wie k&#252;nstliche neuronale Netze eignen sich hierf&#252;r sehr gut, da sie in der Lage sind, selbst&#228;ndig zu lernen.

Anders als bei der Entwicklung klassischer Software, die nach vom Menschen explizit festgelegten Regeln arbeitet, k&#246;nnen ML-Algorithmen die Regeln f&#252;r das L&#246;sen bestimmter Aufgaben selbst lernen und sind dadurch in der Lage, in gro&#223;en, hochdimensionalen Datens&#228;tzen Zusammenh&#228;nge und Muster zu finden, die kein Mensch je entdecken w&#252;rde. Voraussetzung hierf&#252;r ist ein ausreichend gro&#223;er Pool an Beispieldaten f&#252;r die drei zu differenzierenden F&#228;lle Original, Fremdw&#228;hrung und F&#228;lschung. Diese Daten k&#246;nnen dann daf&#252;r genutzt werden, ein neuronales Netz zu entwickeln. W&#228;hrend der sogenannten Trainingsphase lernt das

neuronale Netze selbstständig Regeln, anhand derer die drei Klassen unterschieden werden können. Ergebnis des Trainingsprozesses ist ein neuronales Netz oder Modell, das in M&#252;nzsortierautomaten neu zugeführt werden kann.

Entgegen der weitläufigen Meinung sind neuronale Netze keine neue Erfindung, sondern haben eine lange Geschichte, die ihren Anfang schon in den 1940er Jahren hat. Dabei sind die frühen Lernalgorithmen, die wir heute kennen, inspiriert von der neurowissenschaftlichen Hypothese, dass die mentale Aktivität primär in Nervenzellen, den sogenannten Neuronen, stattfindet. Die Arbeitsweise künstlicher neuronaler Netze orientiert sich an den Vorgängen im menschlichen Gehirn, das aus schätzungsweise 86 Milliarden, über Dendriten verbundene Neuronen besteht (s. Abb. 1 links). Die Kommunikation zwischen Neuronen erfolgt mittels elektrischer Signale, die durch die Synapsen von einer Nervenzelle zur nächsten übertragen werden. Das Lernen im Gehirn erfolgt im Wesentlichen durch das Verstärken oder Abschwächen der Synapse

Hierdurch wird die Intensität der elektrischen Signale und somit die Stärke der Verbindung zwischen Neuronen beeinflusst. Die an einem Neuron gleichzeitig ankommenden Signale addieren sich und das Neuron „feuert“, wenn die Summe einen Grenzwert überschreitet: Ein elektrischer Impuls schießt am Axon entlang.

Abb. 1: Links: Nervenzelle (Neuron), Rechts: künstliches Neuron (Perzeptron). Quelle: Nina Schaaf

Abb. 1: Links: Nervenzelle (Neuron), Rechts: künstliches Neuron (Perzeptron). Quelle: Nina Schaaf

Die Grundeinheit eines künstlichen neuronalen Netzes ist ein einzelnes Neuron. Einfach ausgedrückt, ist dieses ein Element, das gewichtete Eingaben entgegen nimmt, verarbeitet und eine Ausgabe erzeugt. In Abb. 1 (rechts) ist ein einfaches mathematisches Modell eines Neurons, das sogenannte Perzeptron, dargestellt. Das Perzeptron ist über gewichtete Verbindungen mit einem Satz von Eingaben verbunden, wobei das Gewicht der Stärke der Synapsen entspricht. Zusätzlich erhält das künstliche Neuron eine Dummy-Eingabe mit dem Wert 1 und einem eigenen Gewicht, welches als *Bias*, sprich Achsenabschnitt, bezeichnet wird. Die Eingaben, die das Perzeptron erhält, werden auch als Eingabemerkmale bezeichnet. Bezogen auf das Eingangsbeispiel zur M&#252;nzklassifikation können einzelne Münzen etwa durch Merkmale wie ihre Dicke und ihren Durchmesser beschrieben werden. Um eine Ausgabe zu erhalten, werden die gewichteten Eingaben sowie die Dummy-Eingabe aufsummiert und darauf eine Aktivierungsfunktion *f* angewandt. Im einfachsten Fall ist die Aktivierungsfunktion eine Schwellwertfunktion, die entweder 0 oder 1 ausgibt.

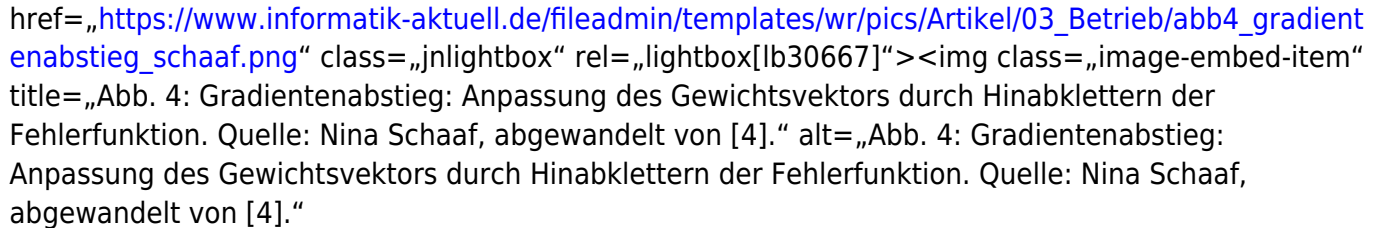
Abb. 2: Schwellwertfunktion als Aktivierungsfunktion: überschreitet die gewichtete Summe des Neurons den Schwellwert  $y \geq 0$ , so ist die Ausgabe der Aktivierungsfunktion 1,

andernfalls ist sie 0. Quelle: Nina Schaaf“ alt=„Abb. 2: Schwellwertfunktion als Aktivierungsfunktion: &#252;berschreitet die gewichtete Summe des Neurons den Schwellwert  $y \geq 0$ , so ist die Ausgabe der Aktivierungsfunktion 1, andernfalls ist sie 0. Quelle: Nina Schaaf“

src=„[https://www.informatik-aktuell.de/fileadmin/\\_processed\\_/9/9/csm\\_abb2\\_formale\\_schaaf\\_e6022c739d.png](https://www.informatik-aktuell.de/fileadmin/_processed_/9/9/csm_abb2_formale_schaaf_e6022c739d.png)“ width=„350“ height=„121“ referrerpolicy=„no-referrer“ /></a><figcaption class=„image-caption“>Abb. 2: Schwellwertfunktion als Aktivierungsfunktion: &#252;berschreitet die gewichtete

Summe des Neurons den Schwellwert  $y \geq 0$ , so ist die Ausgabe der Aktivierungsfunktion 1, andernfalls ist sie 0. Quelle: Nina Schaaf</figcaption></figure></div><div class=„ce-bodytext“><p>&#220;berschreitet die gewichtete Summe des Neurons den Schwellwert, so ist die Ausgabe der Aktivierungsfunktion <em>1</em> &#8211; das Neuron „feuert“. Mathematisch l&#228;sst sich dieses Verhalten wie in Abb. 2 dargestellt beschreiben.</p><p>Ein Perzeptron kann eine mehrdimensionale lineare Funktion repr&#228;sentieren, wodurch sich beispielsweise lineare Klassifikationsprobleme l&#246;sen lassen. Hierbei trennt die lineare Funktion die gegebenen Daten derart in zwei Mengen auf, dass die Daten einer Menge zum Feuern des Perzeptrons f&#252;hren, w&#228;hrend die Daten der zweiten Menge das Perzeptron ruhen lassen. Im Falle der M&#252;nzklassifikation k&#246;nnte ein solcher linearer Klassifikator zwischen originalen Euro-M&#252;nzen und F&#228;lschungen unterscheiden.</p><p>Allerdings sind die meisten realen Problemstellungen nicht durch lineare Funktionen abbildbar. Beispielhaft sei hier die Unterteilung von M&#252;nzen in die <em>drei</em> Klassen Original, Fremdw&#228;hrung und F&#228;lschung genannt. Zur L&#246;sung nichtlinearer Fragestellungen k&#246;nnen deshalb mehrere Neuronen miteinander vernetzt und dadurch ein k&#252;nstliches neuronales Netz aufgebaut werden. Hierzu werden Neuronen in mehreren Schichten angeordnet und gerichtet miteinander verbunden, d. h. die Neuronen einer Schicht erhalten die Ausgaben der vorherigen Schicht als Eingabe. Die erste Schicht eines neuronalen Netzes wird als Eingabeschicht, die letzte Schicht als Ausgabeschicht bezeichnet. Die Menge der Neuronen in der Eingabeschicht wird durch die Anzahl der Eingabedaten bestimmt, die Menge der Neuronen in der Ausgabeschicht h&#228;ngt von der gew&#252;nschten Anzahl an Ergebnissen ab. Dazwischen k&#246;nnen beliebig viele verdeckte Schichten, bestehend aus beliebig vielen Neuronen, liegen.</p></div></div><div class=„ce-textpic ce-left ce-above“><div class=„ce-gallery ce-row ce-column“ data-ce-columns=„1“ data-ce-images=„1“><figure class=„image“><a href=„[https://www.informatik-aktuell.de/fileadmin/templates/wr/pics/Artikel/03\\_Betrieb/abb3\\_gewichte\\_schaaf.png](https://www.informatik-aktuell.de/fileadmin/templates/wr/pics/Artikel/03_Betrieb/abb3_gewichte_schaaf.png)“ class=„jnlightbox“ rel=„lightbox[lb30666]“><img class=„image-embed-item“ title=„Abb. 3: Lernen durch Optimierung der Gewichte eines neuronalen Netzes. Quelle: Nina Schaaf“ alt=„Abb. 3: Lernen durch Optimierung der Gewichte eines neuronalen Netzes. Quelle: Nina Schaaf“ src=„[https://www.informatik-aktuell.de/fileadmin/\\_processed\\_/3/8/csm\\_abb3\\_gewichte\\_schaaf\\_958784c04e.png](https://www.informatik-aktuell.de/fileadmin/_processed_/3/8/csm_abb3_gewichte_schaaf_958784c04e.png)“ width=„720“ height=„269“ referrerpolicy=„no-referrer“ /></a><figcaption class=„image-caption“>Abb. 3: Lernen durch Optimierung der Gewichte eines neuronalen Netzes. Quelle: Nina Schaaf</figcaption></figure></div><div class=„ce-bodytext“><p>Das „Wissen“ eines Perzeptrons, also die Regeln, nach denen es bestimmte Aufgaben l&#246;sen kann, ist in den Gewichten gespeichert. Diese sind anfangs zuf&#228;llig gew&#228;hlt und m&#252;ssen erst <em>gelernt</em> werden. Im Falle des Perzeptrons wird gelernt, indem die Gewichte des Perzeptrons derart angepasst werden, dass die gegebenen Daten m&#246;glichst fehlerfrei durch eine lineare Funktion getrennt werden. Der hierf&#252;r verwendete Perzeptron-Algorithmus passt die Gewichte nur dann an, wenn der Ausgabewert des Neurons vom Sollwert abweicht. Andernfalls bleiben die Gewichte unver&#228;ndert. Da bei mehrschichtigen neuronalen Netzen nicht mehr direkt von der Eingabe auf die Ausgabe geschlossen werden kann, ist hierf&#252;r der Perzeptron-Lernalgorithmus nicht mehr anwendbar. Der Fehler, also der Unterschied zwischen Ausgabewert und Sollwert, kann nur f&#252;r die Ausgabeschicht, nicht aber f&#252;r die verdeckten Schichten gemessen werden.</p><p>Die L&#246;sung bietet der sogenannte <em>Backpropagation</em>-Algorithmus, der es erlaubt, den Fehler von der Ausgabeschicht durch das Netzwerk zur&#252;ckzupropagieren &#8211; schematisch dargestellt in Abb. 3. Dem anfangs noch unkalibrierten Netz wird ein Trainingsdatum, bestehend aus Eingabemerkmale und

zugehöriger Sollausgabe, zugeführter. Das Netz erstellt basierend auf den Eingabemerkmale und den Gewichten eine Ausgabe  $y$ . Die Abweichung zwischen der Ausgabe des Netzes und der Sollausgabe  $y$  wird über eine Fehlerfunktion  $E$  gemessen. Der gemessene Fehler wird durch das gesamte Netz zurückgeführt, sodass darauf basierend die Gewichte des Netzes angepasst werden können. Diese Gewichtsangabe kann als Optimierungsproblem verstanden werden, bei dem es gilt,  $E$  zu minimieren.

Das Diagramm zeigt eine 3D-Oberfläche, die die Fehlerfunktion darstellt. Ein Punkt auf der Oberfläche markiert den aktuellen Zustand des Gewichtsvektors. Ein Pfeil zeigt die Richtung des Gradientenabstiegs an, der zum Minimum der Fehlerfunktion führt. Die Achsen sind mit Gewichten und Fehlerwerten beschriftet.

Quelle: Nina Schaaf, abgewandelt von [4].

Dieser Vorgang lässt sich am besten vor Augen führen, wenn man sich die Fehlerfunktion als eine Oberfläche mit Höhen und Tiefen vorstellt (s. Abb. 4), die sich über den kontinuierlichen Gewichtsraum spannt, d.h. den Raum, der alle möglichen Gewichtseinstellungen abdeckt. Um nun den Punkt des minimalen Fehlers zu finden, also die Einstellung an Gewichten, für die der Fehler minimal ist, muss man sich auf den tiefsten Punkt der Oberfläche zubewegen, diese also „hinab klettern“. Hierfür wird der Gradient der Fehlerfunktion berechnet, welcher in die Richtung des größten Anstiegs des Fehlers zeigt. Soll also die Fehlerfunktion minimiert werden, so folgt man schrittweise der Richtung des negativen Gradienten, der Richtung des größten Abstiegs. Nach jedem Schritt werden die Gewichte angepasst und der Gradient wird unter Berücksichtigung der neuen Gewichte wiederholt evaluiert. Dieser Vorgang wird so lange wiederholt, bis eine Konvergenz mit dem kleinstmöglichen Fehler erreicht wird. Diese schrittweise Anpassung des Gewichtsvektors wird als Gradientenabstieg bezeichnet.

Anders als man denken könnte, verweist das Wort „Deep“ nicht etwa auf ein besonders tiefes Problem- oder Lösungsverständnis, das mit Deep Learning erreicht werden kann.

Tatsächlich bezieht sich die Tiefe auf die Struktur der verwendeten neuronalen Netze, konkret die Verwendung vieler verdeckter Schichten. Je tiefer das Netz, d. h. je mehr verdeckte Schichten es hat, umso komplexer ist die extrahierte Datenrepräsentation. Wie viele Schichten und Neuronen für eine bestimmte Problemstellung erforderlich sind, lässt sich jedoch nicht pauschal definieren. Will man beispielsweise zu wenige Neuronen, erhält man ein zu einfaches Modell, das die Daten nur unvollständig repräsentieren kann. Bei zu vielen Neuronen und Schichten hingegen „merkt“ sich das Netz alle Daten, die es während des Trainings gesehen hat, was einem Auswendiglernen gleichkommt. Dieses als **Overfitting** bezeichnete Phänomen hat den Effekt, dass das Netz schlecht generalisiert, also bei unbekannten Daten schlechte Ergebnisse liefert. Tatsächlich ist das Design der verdeckten Schichten immer noch ein äußerst aktives Forschungsgebiet; es existieren noch keine universellen theoretischen Richtlinien, denen man folgen könnte.

In der Praxis wird beispielsweise Kreuzvalidierung zur Bewertung geeigneter Architekturen eingesetzt, d. h. es wird mit verschiedenen Strukturen



experimentiert und diejenige gew&#228;hlt, welche am besten geeignet ist <a href=„<https://www.informatik-aktuell.de/betrieb/kuenstliche-intelligenz/neuronale-netze-ein-blick-in-die-black-box.html#c30659>“>[6]</a>. Eine weitere g&#228;ngige Methode ist, bereits bekannten Architekturen anderer Anwendungen zu folgen und deren Struktur zu &#252;bertragen und an einen individuellen Anwendungsfall anzupassen.</p>
<h2>Algorithmische Komplexit&#228;t &#8211; Fluch und Segen</h2>
<p>Einer der gr&#246;&#223;ten Vorteile des Deep Learning &#8211; die enorme Komplexit&#228;t neuronaler Netze und damit verbunden auch ihre F&#228;higkeit zur Approximation beliebig komplexer Funktionen &#8211; ist gleichzeitig auch einer der gr&#246;&#223;ten Kritikpunkte. Man sagt auch, tiefe neuronale Netze seien eine Black Box, da die von ihnen gelernten Zusammenh&#228;nge und Datenrepr&#228;sentationen so komplex und abstrakt sind, dass Menschen &#8211; und selbst Experten &#8211; sie nicht mehr nachvollziehen k&#246;nnen. Begr&#252;ndet ist dies darin, dass Black-Box-Modelle nicht transparent sind. Transparenz kann grunds&#228;tzlich auf drei unterschiedliche Ebenen bezogen sein: auf Modell-Ebene (Simulierbarkeit), auf Komponenten-Ebene (Unterteilbarkeit) sowie auf algorithmischer Ebene (algorithmische Transparenz) <a href=„<https://www.informatik-aktuell.de/betrieb/kuenstliche-intelligenz/neuronale-netze-ein-blick-in-die-black-box.html#c30659>“>[7]</a>. Simulierbarkeit ist dann gegeben, wenn ein Mensch alle Rechenschritte des Modells in <em>angemessener</em> Zeit auswerten kann. Diese Anforderung ist bei einem tiefen neuronalen Netz, das unter Umst&#228;nden Millionen von Gewichten hat, unm&#246;glich zu erf&#252;llen. Um dem Kriterium der Unterteilbarkeit zu entsprechen, m&#252;ssten alle Komponenten des neuronalen Netzes intuitiv verst&#228;ndlich sein, also beispielsweise der Einfluss eines jeden Parameters nachvollziehbar sein. Diese Voraussetzung ist f&#252;r neuronale Netze durch die schiere Anzahl an Parametern ebenfalls nicht gegeben. Zudem besitzen tiefe neuronale Netze in der Regel hochkomplexe Entscheidungsgrenzen, was sie folglich auch f&#252;r die Transparenz auf algorithmischer Ebene disqualifiziert.</p>
<p>Eine Frage, die immer wieder aufkommt, ist, warum wir neuronale Netze &#252;berhaupt verstehen wollen. M&#252;ssen wir wissen, wie ein Netz entscheidet oder gen&#252;gt es, den Black-Box-Charakter zu akzeptieren, solange neuronale Netze eine m&#246;glichst hohe Vorhersagegenauigkeit haben? Oftmals wird hier der Vergleich zur menschlichen Intuition gezogen. Viele Entscheidungen, die wir tagt&#228;glich treffen, erfolgen „aus dem Bauch heraus“, wir k&#246;nnen also die genauen Beweggr&#252;nde einer Entscheidung oftmals nicht solide begr&#252;nden. Sollten wir also diesen Umstand auch f&#252;r automatisiert getroffene Entscheidungen akzeptieren? Um diese Fragen zu beantworten, hilft es, einige Beispiele zu betrachten, bei denen die Erkl&#228;rbarkeit neuronaler Netze zumindest wichtig, wenn nicht unerl&#228;sslich ist.</p>
<p>Eine h&#228;ufige Anwendung in der Produktion ist die Erkennung von Ausschuss. Dies k&#246;nnen beispielsweise Werkst&#252;cke sein, die bestimmte Qualit&#228;tsanforderungen nicht erf&#252;llen oder, wie eingangs beschrieben, das Aussortieren unerw&#252;nschter Objekte wie Falschgeld. Werden f&#252;r diese Anwendungsf&#228;lle neuronale Netze eingesetzt, kann es vorkommen, dass das Netz andere Entscheidungen trifft als beispielsweise ein Mitarbeiter in der Qualit&#228;tssicherung. Entdeckt man eine solche Diskrepanz, kann es helfen, die Kriterien des neuronalen Netzes f&#252;r die spezifische Entscheidung zu hinterfragen, um eventuelle Fehler im Modell aufdecken zu k&#246;nnen.</p>
</div>
<div class=„ce-textpic ce-left ce-above“>
<div class=„ce-gallery ce-row ce-column“ data-ce-columns=„1“ data-ce-images=„1“>
<figure class=„image“>
<a href=„[https://www.informatik-aktuell.de/fileadmin/templates/wr/pics/Artikel/03\\_Betrieb/abb5\\_pneumonia\\_combined\\_schaaf.png](https://www.informatik-aktuell.de/fileadmin/templates/wr/pics/Artikel/03_Betrieb/abb5_pneumonia_combined_schaaf.png)“ class=„jnllightbox“ rel=„lightbox[lb30669]“>
<img class=„image-embed-item“ title=„Abb. 5: Patient mit Lungenentz&#252;ndung. Links: R&#246;ntgenaufnahme, Rechts: Heatmap des neuronalen Netzes. Quelle: [11]“ alt=„Abb. 5: Patient mit Lungenentz&#252;ndung. Links: R&#246;ntgenaufnahme, Rechts: Heatmap des neuronalen Netzes. Quelle: [11]“ src=„[https://www.informatik-aktuell.de/fileadmin/templates/wr/pics/Artikel/03\\_Betrieb/abb5\\_pneumonia\\_combined\\_schaaf.png](https://www.informatik-aktuell.de/fileadmin/templates/wr/pics/Artikel/03_Betrieb/abb5_pneumonia_combined_schaaf.png)“ width=„678“ height=„371“ referrerpolicy=„no-referrer“ />
</a>
<figcaption class=„image-caption“>Abb. 5: Patient mit Lungenentz&#252;ndung. Links:

Rechts: Heatmap des neuronalen Netzes. Quelle:

[11]</figcaption></figure></div><div class=„ce-bodytext“><p>Besonders relevant ist dieser Aspekt auch f&#252;r medizinische Anwendungen. Hier kann der Einsatz von KI &#196;rzt durch hochpr&#228;zise Vorhersagen bei der Diagnose von Krankheiten unterst&#252;tzen und gleichzeitig helfen, menschlichen Fehlern vorzubeugen. Eine Anwendung von tiefen neuronalen Netzen in der Medizin ist die Analyse von Thoraxr&#246;ntgenbildern, um Lungenentz&#252;ndungen zu erkennen. Mithilfe von KI k&#246;nnte eine Diagnose binnen 10 Sekunden und nicht wie bisher nach rund 20 Minuten erstellt werden <a

href=„<https://www.informatik-aktuell.de/betrieb/kuenstliche-intelligenz/neuronale-netze-ein-blick-in-di-e-black-box.html#c30659>“>[8]</a>. Eine Gruppe von Forschern hat f&#252;r ein solches Netz gepr&#252;ft, aufgrund welcher Bildbereiche das Netz ein Risiko f&#252;r

Lungenentz&#252;ndungen ausgibt <a

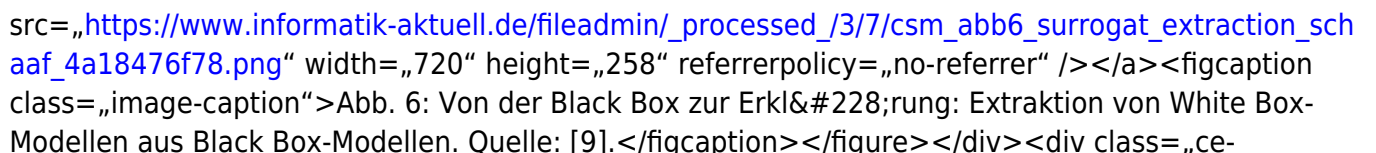
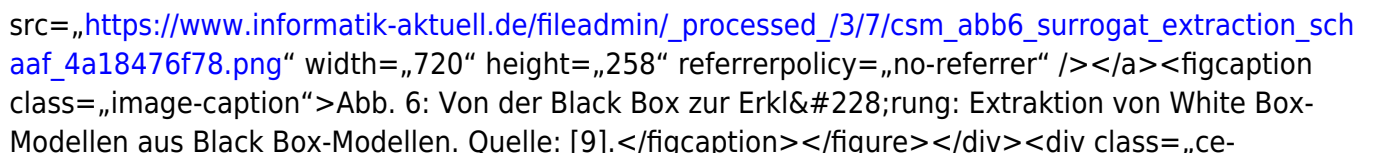
href=„<https://www.informatik-aktuell.de/betrieb/kuenstliche-intelligenz/neuronale-netze-ein-blick-in-di-e-black-box.html#c30659>“>[9]</a>. Betrachtet man das in Abb. 5 links dargestellte

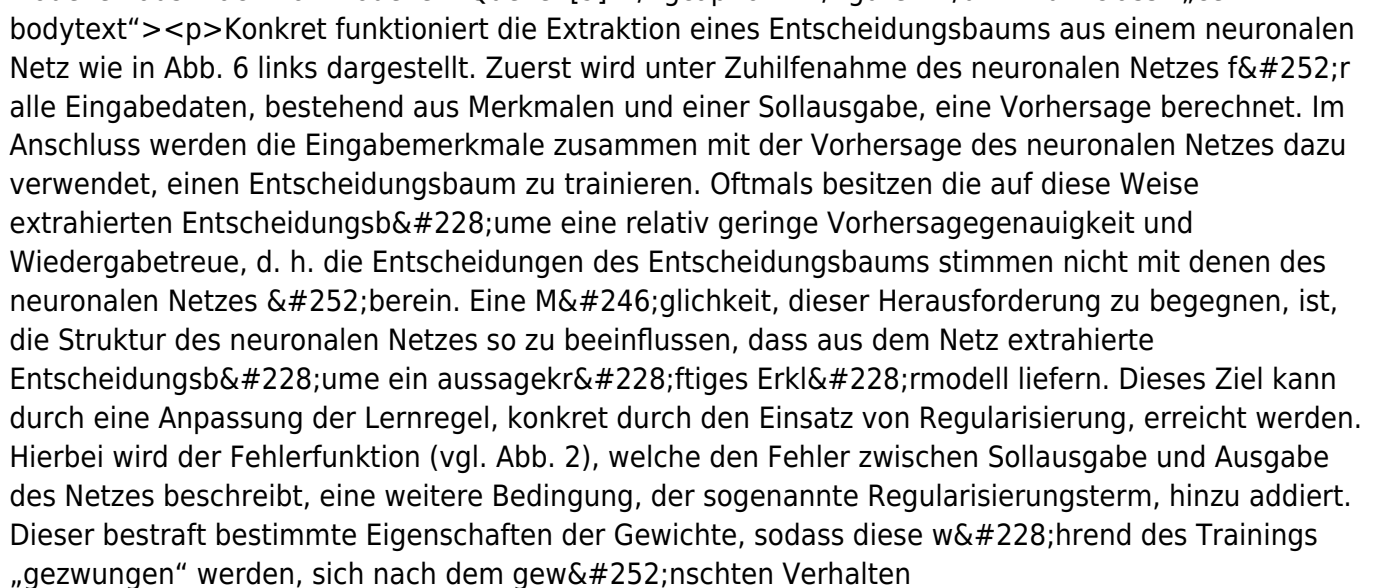
R&#246;ntgenbild, so indiziert hier der Bereich in der oberen rechten Lunge eine Lungenentz&#252;ndung. Untersucht man aber nun, welche Bildbereiche f&#252;r das Netz wirklich relevant waren (Abb. 5, rechts), so sieht man, dass dem urs&#228;chlichen Bereich hier sogar ein negativer Einfluss zugewiesen wurde &#8211; erkennbar an den negativen Werten im linken oberen Bereich des Bildes. Erstaunlicherweise werden dem Schriftzug „Portable“ und den Seitenmarkierungen in der oberen rechten Ecke, die dem verwendeten mobilen Scanner zugeordnet sind, ein positiver Einfluss zugewiesen. Es haben also nicht die tats&#228;chlich ausschlaggebenden Bildbereiche zur Risiko-Attestierung „Lungenentz&#252;ndung“ beigetragen, sondern die Indikatoren f&#252;r die Verwendung eines mobilen R&#246;ntgenger&#228;tes, welches in der Regel eingesetzt wird, wenn der Patient zu krank ist, um aus dem Bett zu kommen.</p><p>Das neuronale Netz k&#246;nnte also gelernt haben, einen Zusammenhang zwischen mobilen R&#246;ntgenaufnahmen und einem erh&#246;hten Risiko f&#252;r eine Lungenentz&#252;ndung herzustellen (Korrelation), obwohl die mobile Aufnahme nicht die eigentliche Ursache der Lungenentz&#252;ndung ist (Kausalit&#228;t). Sollen deep-learning-basierte Anwendungen tats&#228;chlich f&#252;r die computergest&#252;tzte Diagnose eingesetzt werden, so zeigt dieses Beispiel, dass eine genaue &#220;berpr&#252;fung der Entscheidungsweise oder zumindest eine intensive &#220;berpr&#252;fung des Modells unvermeidlich ist.</p><p>Auch bei

sicherheitskritischen Anwendungen, in denen neuronale Netze eingesetzt werden, ist das Verst&#228;ndnis der generellen inneren Abl&#228;ufe dieser Netze unerl&#228;sslich. Kommt etwa ein neuronales Netz bei der Programmierung eines Roboters zum Einsatz, so muss dieses System zur Abnahme einer Sicherheitspr&#252;fung unterzogen werden. Bei dieser muss die gesamte Funktionsweise des Systems ersichtlich und nachvollziehbar sein, um eine regelkonforme Arbeitsweise des Roboters garantieren zu k&#246;nnen. Besonders relevant ist dieser Aspekt, wenn Roboter mit Menschen zusammenarbeiten sollen. Des Weiteren werden neuronale Netze aktuell bereits auch f&#252;r die Anwendung von bisher regelbasiert arbeitenden Systemen wie Steuerger&#228;ten erprobt. Wie im vorherigen Fall ist auch hier die Sicherheitspr&#252;fung dieser Systeme, etwa durch den T&#220;V, eine H&#252;rde auf dem Weg, diese Verfahren in den praktischen Einsatz zu &#252;berf&#252;hren. Es m&#252;ssen also Verfahren gefunden werden, die die inneren Entscheidungswege neuronaler Netze erkl&#228;ren oder zumindest deren korrekte Funktionsweise verifizieren k&#246;nnen.</p></h2>Blick in die Black Box: Erkl&#228;rung

neuronaler Netze</h2><p>Eine M&#246;glichkeit, die generelle Arbeitsweise neuronaler Netze zu erkl&#228;ren, ist die Abbildung des Netzes mittels inh&#228;rent interpretierbarer Modelle, sogenannter White-Box-Modelle, wie Entscheidungsb&#228;umen oder Entscheidungsregeln. Entscheidungsb&#228;ume bestehen aus internen Knoten, die zu &#252;berpr&#252;fende Bedingungen definieren, und Blattknoten, die Klassen darstellen. M&#246;chte man ein Datum mithilfe eines Entscheidungsbaums klassifizieren, wird der Baum von oben nach unten traversiert, bis

man einen Blattknoten erreicht, der die Klasse codiert. Um ein Black-Box-Modell nun mithilfe eines solchen White-Box-Modells erklären zu können, wird im ersten Schritt das Black-Box-Modell, etwa ein neuronales Netz, gelernt (s. Abb. 6). Im Anschluss kann ein interpretierbares Stellvertretermodell [\[9\]](#); auch als *Surrogat* bezeichnet [\[9\]](#); aus dem neuronalen Netz extrahiert werden, welches dann zur Erzeugung von Erklärungen genutzt werden kann.

 [https://www.informatik-aktuell.de/fileadmin/templates/wr/pics/Artikel/03\\_Betrieb/abb6\\_surrogat\\_extraction\\_schaaf.png](https://www.informatik-aktuell.de/fileadmin/templates/wr/pics/Artikel/03_Betrieb/abb6_surrogat_extraction_schaaf.png) class="jnlighbox" rel="lightbox[lb30670]" data-bbox="66 260 932 328"/>  Abb. 6: Von der Black Box zur Erklärung: Extraktion von White Box-Modellen aus Black Box-Modellen. Quelle: [9]. alt="Abb. 6: Von der Black Box zur Erklärung: Extraktion von White Box-Modellen aus Black Box-Modellen. Quelle: [9]."

[https://www.informatik-aktuell.de/fileadmin/\\_processed\\_/3/7/csm\\_abb6\\_surrogat\\_extraction\\_schaaf\\_4a18476f78.png](https://www.informatik-aktuell.de/fileadmin/_processed_/3/7/csm_abb6_surrogat_extraction_schaaf_4a18476f78.png) width="720" height="258" referrerpolicy="no-referrer" data-bbox="66 328 932 589"/>  Abb. 6: Von der Black Box zur Erklärung: Extraktion von White Box-Modellen aus Black Box-Modellen. Quelle: [9].

Konkret funktioniert die Extraktion eines Entscheidungsbaums aus einem neuronalen Netz wie in Abb. 6 links dargestellt. Zuerst wird unter Zuhilfenahme des neuronalen Netzes für alle Eingabedaten, bestehend aus Merkmalen und einer Sollausgabe, eine Vorhersage berechnet. Im Anschluss werden die Eingabemerkmale zusammen mit der Vorhersage des neuronalen Netzes dazu verwendet, einen Entscheidungsbaum zu trainieren. Oftmals besitzen die auf diese Weise extrahierten Entscheidungsbaume eine relativ geringe Vorhersagegenauigkeit und Wiedergabetreue, d. h. die Entscheidungen des Entscheidungsbaums stimmen nicht mit denen des neuronalen Netzes überein. Eine Möglichkeit, dieser Herausforderung zu begegnen, ist, die Struktur des neuronalen Netzes so zu beeinflussen, dass aus dem Netz extrahierte Entscheidungsbaume ein aussagekräftiges Erklärungsmodell liefern. Dieses Ziel kann durch eine Anpassung der Lernregel, konkret durch den Einsatz von Regularisierung, erreicht werden. Hierbei wird der Fehlerfunktion (vgl. Abb. 2), welche den Fehler zwischen Sollausgabe und Ausgabe des Netzes beschreibt, eine weitere Bedingung, der sogenannte Regularisierungsterm, hinzu addiert. Dieser bestraft bestimmte Eigenschaften der Gewichte, sodass diese während des Trainings „gezwungen“ werden, sich nach dem gewünschten Verhalten

auszurichten.

Oblicherweise wird Regularisierung dazu verwendet, das zuvor beschriebene Overfitting zu vermeiden. Es existieren jedoch bereits Ansätze, die Regularisierung einsetzen, um die Struktur von neuronalen Netzen so zu beeinflussen, dass extrahierte Surrogate bestimmten Kriterien (z. B. Größe, Zuverlässigkeit etc.) entsprechen. Der Tree Regularizer <https://www.informatik-aktuell.de/betrieb/kuenstliche-intelligenz/neuronale-netze-ein-blick-in-die-black-box.html#c30659> [10] schätzt die Größe des Entscheidungsbaums, während der L1-orthogonal Regularisierer in [\[11\]](https://www.informatik-aktuell.de/betrieb/kuenstliche-intelligenz/neuronale-netze-ein-blick-in-die-black-box.html#c30659) neuronale Netze mit wenigen, aber orthogonalen Gewichten favorisiert. Beides führt zu neuronalen Netzen mit Entscheidungsgrenzen, die von Entscheidungsbaumen gut angenähert werden können. Die Baume sind gleichzeitig von geringer Größe, wodurch sie für den Menschen leicht verständlich sind.

## Blick in die Glaskugel: wohin geht die Reise?

Schon heute existiert eine Vielzahl an Methoden und Techniken, die uns Menschen erlauben, einen ersten Blick in die Black Box zu werfen. Doch sind wir noch lange nicht am Ziel angekommen. Stand heute gibt es bereits eine Vielzahl an Unternehmen, die die enormen Potenziale der künstlichen Intelligenz nutzt, sei es in hochrelevanten Anwendungen, wie der computergestützten Diagnose im Gesundheitswesen oder dem autonomen Fahren, aber auch zur Anfangserkennung des Falschgeldgebrauchs. Auf der anderen Seite stehen viele Firmen hier

aufgrund bisher nicht erkannter Potenziale, fehlenden Know-hows oder gesetzlicher Hürden erst am Anfang. Doch mit steigender Verbreitung der KI-Anwendungen wird auch der Bedarf von Erklärbarkeit der eingesetzten Algorithmen steigen. Hier gilt es, bereits existierende Verfahren in die Anwendung zu integrieren und neue Methoden zu entwickeln. Auch die Frage, welche Art der Erklärung für welche Zielgruppe am besten geeignet ist, wird hierbei noch eine große Rolle spielen. Schließlich hat ein Endanwender ein anderes Interesse an einem KI-Verfahren als der Experte, der es programmiert hat und selbst besser verstehen will.

From:  
<https://schnipsel.qgelm.de/> - Qgelm

Permanent link:  
[https://schnipsel.qgelm.de/doku.php?id=wallabag:wb2neuronale-netze\\_ein-blick-in-die-black-box](https://schnipsel.qgelm.de/doku.php?id=wallabag:wb2neuronale-netze_ein-blick-in-die-black-box)

Last update: **2025/06/27 11:17**

