

# Wie synthetische Datensätze KI-Systeme verbessern sollen

Originalartikel

Backup

<html> <header class=„article-header“><h1 class=„articleheading“>Wie synthetische  
Datensätze KI-Systeme verbessern sollen</h1><div class=„publish-info“> Karen  
Hao</div></header><figure class=„aufmacherbild“><img  
src=„[https://heise.cloudimg.io/width/700/q75.png-lossy-75.webp-lossy-75.foil1/\\_www-heise-de\\_/imgs/18/3/1/2/1/4/1/6/ContactSheet-001-455b8139c499e9b3.jpeg](https://heise.cloudimg.io/width/700/q75.png-lossy-75.webp-lossy-75.foil1/_www-heise-de_/imgs/18/3/1/2/1/4/1/6/ContactSheet-001-455b8139c499e9b3.jpeg)“  
srcset=„[https://heise.cloudimg.io/width/700/q75.png-lossy-75.webp-lossy-75.foil1/\\_www-heise-de\\_/imgs/18/3/1/2/1/4/1/6/ContactSheet-001-455b8139c499e9b3.jpeg](https://heise.cloudimg.io/width/700/q75.png-lossy-75.webp-lossy-75.foil1/_www-heise-de_/imgs/18/3/1/2/1/4/1/6/ContactSheet-001-455b8139c499e9b3.jpeg) 700w,  
[https://heise.cloudimg.io/width/1050/q75.png-lossy-75.webp-lossy-75.foil1/\\_www-heise-de\\_/imgs/18/3/1/2/1/4/1/6/ContactSheet-001-455b8139c499e9b3.jpeg](https://heise.cloudimg.io/width/1050/q75.png-lossy-75.webp-lossy-75.foil1/_www-heise-de_/imgs/18/3/1/2/1/4/1/6/ContactSheet-001-455b8139c499e9b3.jpeg) 1050w,  
[https://heise.cloudimg.io/width/1500/q75.png-lossy-75.webp-lossy-75.foil1/\\_www-heise-de\\_/imgs/18/3/1/2/1/4/1/6/ContactSheet-001-455b8139c499e9b3.jpeg](https://heise.cloudimg.io/width/1500/q75.png-lossy-75.webp-lossy-75.foil1/_www-heise-de_/imgs/18/3/1/2/1/4/1/6/ContactSheet-001-455b8139c499e9b3.jpeg) 1500w,  
[https://heise.cloudimg.io/width/2063/q75.png-lossy-75.webp-lossy-75.foil1/\\_www-heise-de\\_/imgs/18/3/1/2/1/4/1/6/ContactSheet-001-455b8139c499e9b3.jpeg](https://heise.cloudimg.io/width/2063/q75.png-lossy-75.webp-lossy-75.foil1/_www-heise-de_/imgs/18/3/1/2/1/4/1/6/ContactSheet-001-455b8139c499e9b3.jpeg) 2063w“ alt=„ class=„img-responsive“  
referrerpolicy=„no-referrer“ /><figcaption class=„akwa-  
caption“>(Bild:&#160;Datagen)</figcaption></figure><p><strong>Deep Learning ben&#246;tigt  
gro&#223;e Informationsmengen. Aus realen Informationen abgeleitete Fake-Daten sollen  
helfen.</strong></p><p>Man kann die schwachen Stoppeln auf seiner Oberlippe sehen, die Falten  
auf seiner Stirn, die Unreinheiten seiner Haut. Er ist kein echter Mensch, aber er orientiert sich an  
ihnen &#8211; so wie Hunderttausende andere, die von Datagen hergestellt werden, einer Firma, die  
Fakes von Menschen verkauft.</p><p>Diese falschen Personen sind keine Spiel-Avatare oder  
animierte Figuren f&#252;r Filme. Sie dienen als k&#252;nstliche Daten, um damit Deep-Learning-  
Algorithmen zu f&#252;ttern. Firmen wie <a href=„<https://www.datagen.tech/technology/>“  
rel=„external noopener“ target=„\_blank“><strong>Datagen [1]</strong></a> wollen damit eine  
Alternative zum teuren und zeitaufw&#228;ndigen Sammeln von Daten aus der realen Welt anbieten.  
Das Unternehmen generiert die Informationen ma&#223; geschneidert f&#252;r den Kunden, wie und  
wann er es will &#8211; und das zu einem relativ g&#252;nstigen Preis.</p><p>Um seine  
synthetischen Menschen zu erzeugen, scannt Datagen zun&#228;chst reale Personen. Das  
Unternehmen arbeitet mit Zwischenh&#228;ndlern zusammen, die Menschen daf&#252;r bezahlen,  
in gro&#223;en Ganzk&#246;rperscannern jedes Detail von der Iris &#252;ber die  
Hautbeschaffenheit bis hin zur Kr&#252;mmung der Finger erfassen zu lassen. Aus diesen Rohdaten  
kreiert das Startup mit Hilfe einer ganzen Reihe von Algorithmen 3D-Darstellungen von K&#246;rper,  
Gesicht, Augen und H&#228;nden einer Person.</p><h3 class=„subheading“  
id=„nav\_nicht\_einfach0“>Nicht einfach „Daumen hoch“</h3><p>Das Unternehmen mit Sitz in Israel  
arbeitet nach eigenen Angaben bereits mit vier gro&#223;en US-Tech-Giganten zusammen, will aber  
nicht verraten, mit welchen. Sein wichtigster Konkurrent, <a href=„<https://synthesis.ai/>“  
rel=„external noopener“ target=„\_blank“><strong>Synthesis AI [2]</strong></a>, bietet ebenfalls  
digitale Menschen auf Abruf an. Andere Unternehmen generieren Daten f&#252;r die Finanz-,  
Versicherungs- und Gesundheitsbranche. Es gibt mittlerweile zahlreiche Firmen auf dem  
Gebiet.</p><p>Einst galten synthetische Daten im Vergleich zu realen als minderwertig. Heute  
hingegen sieht so mancher Beobachter in ihnen ein Allheilmittel. Echte Daten sind  
un&#252;bersichtlich und mit Fehlern behaftet. Neue Datenschutzbestimmungen erschweren zudem  
ihr Sammeln. Im Gegensatz dazu lassen sich aus synthetischen Daten viel leichter die  
unterschiedlichsten Datensätze erstellen. So kann man daraus zum Beispiel perfekte Gesichter  
unterschiedlichen Alters, unterschiedlicher Form und ethnischer Zugehörigkeit erzeugen. Damit

Ist sich dann ein Programm zur Gesichtserkennung entwickeln, das für alle Bevölkerungsgruppen funktioniert. Aber synthetische Daten haben auch Nachteile. Spiegeln sie die Realität nicht richtig wider, könnte das zu schlechteren Ergebnissen führen als weniger genaue Daten aus der realen Welt; oder zumindest zu den Problemen führen, die diese haben. „Ich mag hier nicht einfach ein „Daumen hoch“ setzen und sagen: Oh, das wird so viele Probleme lösen“, sagt Cathy O’Neil, eine Datenwissenschaftlerin und Gründerin der auf die überprüfung von Algorithmen spezialisierte Firma ORCAA. „Denn bei der Methode werden auch viele Dinge nicht beachtet.“

**Realistisch, nicht real**

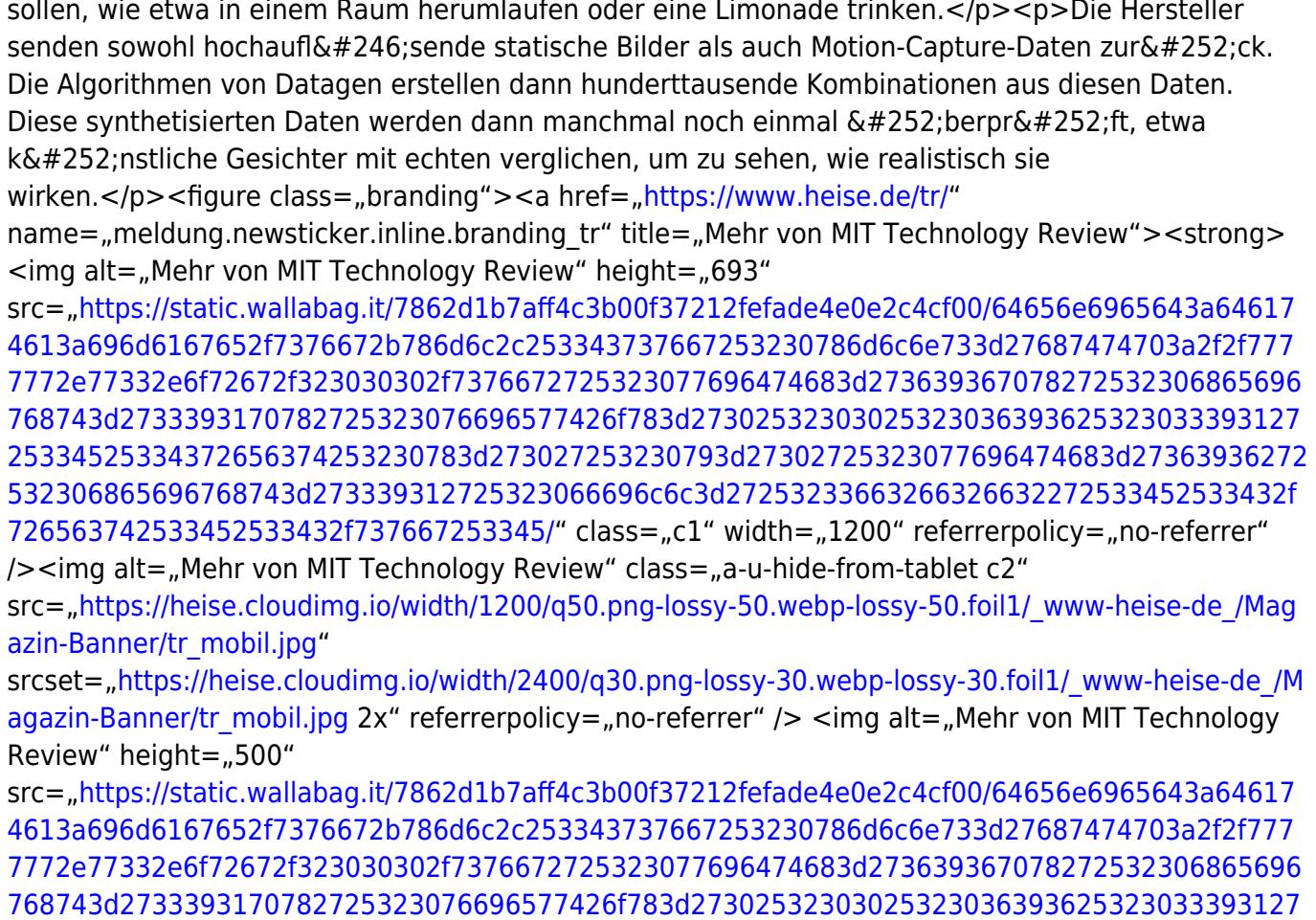
Deep Learning brauchte schon immer viele Daten. Aber in den letzten Jahren hat sich gezeigt, dass deren Qualität wichtiger ist als ihre Menge. Selbst kleine Mengen richtiger, sauber zugeordneter Daten verbessern die Leistung eines KI-Systems mehr als die zehnfache Menge schlecht aufbereiteter Daten, selbst wenn man leistungsfähigere Algorithmen verwendet.

Das sollten Unternehmen bei der Entwicklung ihrer KI-Modelle berücksichtigen, sagt Ofir Chakon, CEO und Mitbegründer von Datagen. Heute sammeln die Firmen zumeist erst einmal so viele Daten wie möglich und optimieren dann ihre Algorithmen. Stattdessen sollten sie das Gegenteil tun: Den selben Algorithmus verwenden, aber die Qualität ihrer Daten verbessern. Doch reale Daten für solches iteratives Experimentieren zu sammeln, ist zu kostspielig und zeitintensiv. An dieser Stelle kommt Datagen ins Spiel. Mit einem Generator für synthetische Daten können Teams Dutzende von neuen Datensätzen pro Tag erstellen und sie testen, um herauszufinden, welche die Realität am besten abbilden.

**Hunderttausende Kombinationen**

Um die Realitätssicherheit der Daten zu gewährleisten, gibt Datagen seinen Lieferanten detaillierte Anweisungen, wie viele Personen in jeder Alters- und Gewichtsklasse sowie ethnischer Zugehörigkeit zu scannen sind. Hinzu kommt eine Liste von Aktionen, die sie aufzuzeigen sollen, wie etwa in einem Raum herumlaufen oder eine Limonade trinken.

Die Hersteller senden sowohl hochauflösende statische Bilder als auch Motion-Capture-Daten zurück. Die Algorithmen von Datagen erstellen dann hunderttausende Kombinationen aus diesen Daten. Diese synthetisierten Daten werden dann manchmal noch einmal überprüft, etwa künstliche Gesichter mit echten verglichen, um zu sehen, wie realistisch sie wirken.



25334525334372656374253230783d273027253230793d27302725323077696474683d27363936272  
532306865696768743d273339312725323066696c6c3d2725323366326632272533452533432f  
726563742533452533432f737667253345" class="c3" width="1830" referrerpolicy="no-referrer"/>  
/>  
[3]</strong></a></figure><p>Datagen generiert beispielsweise Gesichtsausdrücke zur Überwachung der Aufmerksamkeit von Fahrerinnen und Fahrgästen in smarten Autos. Oder Körperbewegungen, um Kunden in kassenlosen Geschäften zu verfolgen, sowie Iris und Handbewegungen, um die Augen- und Hand-Tracking-Funktionen von VR-Headsets zu verbessern. Laut dem Unternehmen dienen seine Daten bereits zur Entwicklung von Bilderkennungssystemen (Computer Vision), die mehrere Millionen Nutzerinnen und Nutzer verwenden.</p><h3 class="subheading" id="nav\_von3">Von Fahrzeuginspektion bis zur Medizin</h3><p>Es werden nicht nur synthetische Menschen in Massenproduktion geschaffen. Das Startup <a href="https://www.click-ins.com/" rel="external noopener" target="\_blank"><strong>Click-Ins</strong></a> verwendet zum Beispiel synthetische Daten für automatische Fahrzeuginspektionen. Mithilfe von Design-Software werden alle Automarken und -modelle, die die KI erkennen muss, neu erstellt und dann mit verschiedenen Farben, Schichten und Verformungen unter verschiedenen Lichtverhältnissen und vor verschiedenen Hintergründen gerendert. Auf diese Weise kann das Unternehmen sein KI-Modell aktualisieren, wenn Autohersteller neue Fahrzeugvarianten auf den Markt bringen. Außerdem wird so kein Datenschutz verletzt in Ländern wie Deutschland, in denen Nummernschilder als private Informationen gelten und daher nicht in Fotos zum Training der KI enthalten sein dürfen.</p><p><a href="https://mostly.ai/industries/" rel="external noopener" target="\_blank"><strong>Mostly.ai</strong></a> arbeitet wiederum mit Finanz-, Telekommunikations- und Versicherungsunternehmen zusammen, um Tabellen mit synthetischen Kundendaten bereitzustellen. Auf diese Weise können die Unternehmen den Aufbau ihrer Kundendatenbank mit externen Dienstleistern auf rechtskonforme Weise teilen. Denn selbst wenn Daten anonymisiert werden, schreibt dass zuweilen nicht ausreichend die Privatsphäre der Menschen. Mit den synthetischen Daten lassen sich Datensätze mit den gleichen statistischen Eigenschaften wie die der echten Daten eines Unternehmens generieren. So können auch Daten simuliert werden, die das Unternehmen noch gar nicht hat, zum Beispiel von hypothetischen zukünftigen Kundengruppen oder Szenarien betrügerischer Aktivitäten.</p><p>Proponenten synthetischer Daten sagen, dass diese auch bei der Bewertung der Fähigkeiten von KI helfen können. Ein Beispiel zeigten Suchi Saria, Professorin für maschinelles Lernen und Gesundheitswesen an der Johns Hopkins University, und ihre Mitautoren, in einer kürzlich auf einer KI-Konferenz veröffentlichten Arbeit: Mit Hilfe von Techniken zur Datengenerierung ließen sich verschiedene Patientengruppen aus einem einzigen Datensatz extrapoliieren. Das könnte nützlich sein, wenn ein Unternehmen etwa nur Daten von der eher jugendlichen Bevölkerung von New York City vorliegen hat, aber verstehen möchte, wie seine KI bei einer älteren Bevölkerung mit einer höheren Diabetesprävalenz funktioniert. Um solche medizinische KI-Systeme zu testen, findet Saria nun ihr eigenes Unternehmen namens Bayesian Health.</p><h3 class="subheading" id="nav\_datenschutz4">Datenschutz nicht automatisch gewährleistet</h3><p>Doch gibt es um synthetische Daten einen ungültigen Hype? In Sachen Datenschutz bedeutet die Tatsache, dass die Daten 'synthetisch' sind und nicht direkt den realen Benutzerdaten entsprechen, nicht, dass sie keine sensiblen Informationen über reale Personen enthalten", sagt Aaron Roth, Professor für Computer- und Informationswissenschaften an der University of Pennsylvania. Es habe sich gezeigt, dass einige Datengenerierungstechniken Bilder oder Texte aus ihren (echten) Trainingsdaten einfach nur

abkupfern.</p><p>Das mag für eine Firma wie DataGen in Ordnung sein, deren synthetische Daten nicht dazu gedacht sind, die Identität der Personen zu verbergen, denn die haben dem Scan zugestimmt. Aber es wäre eine schlechte Nachricht für Unternehmen, die in der Methode eine Möglichkeit sehen, sensible Finanz- oder Patientendaten zu schützen.</p><p>Bisherige Forschung legt nahe, dass insbesondere die Kombination von zwei Techniken für synthetische Daten &#8211; die sogenannte <a href=„<https://www.heise.de/hintergrund/Kollektiver-Datenschutz-Was-dahinter-steckt-und-warum-er-nuetig-ist-6054822.html>“><strong>Differential Privacy [6]</strong></a> und <a href=„<https://www.heise.de/hintergrund/Neuronales-Netz-erzeugt-atraktive-Gesichter-fuer-jeden-Geschmack-6018436.html>“><strong>Generative Adversarial Networks [7]</strong></a> (GANs) &#8211; für guten Schutz sorgen kann, sagt Bernease Herman, Datenwissenschaftler am University of Washington eScience Institute. Skeptiker befürchten jedoch, dass dieser Aspekt im Marketing-Jargon der Anbieter synthetischer Daten verloren geht. Denn die sprechen nicht immer offen darüber, welche Techniken sie verwenden.</p><h3 class=„subheading“ id=„nav\_voreingenommenheit5“>Voreingenommenheit nicht ausgeschlossen</h3><p>Bislang deutet wenig darauf hin, dass synthetische Daten zuverlässig vor Voreingenommenheit schützen. Ist ein Datensatz verzerrt, so kommt man durch ein „Hochrechnen“ nicht zu wirklich repräsentativen Daten. Die Rohdaten von DataGen enthalten zum Beispiel proportional weniger ethnische Minderheiten, was bedeutet, dass weniger reale Datenpunkte verwendet werden, um Fake-Menschen aus diesen Gruppen zu erzeugen. „Wenn Ihre Gesichter mit dunklerer Hautfarbe keine besonders guten Annahmen an reale Gesichter sind, dann löschen Sie das Problem nicht wirklich“, sagt O’Neil.</p><p>Zum anderen führen perfekt ausbalancierte Datensätze nicht automatisch zu perfekt fairen KI-Systemen, sagt Christo Wilson, außerordentlicher Professor für Informatik an der Northeastern University. Wenn ein Kreditkartenanbieter versucht, einen KI-Algorithmus zur Bewertung potenzieller Kreditnehmer zu entwickeln, würde er nicht alle marginalen Diskriminierungen beseitigen, indem er einfach Weiße genauso wie Schwarze in seinen Daten repräsentiert. Diskriminierung kann sich immer noch durch Unterschiede zwischen Bewerbern aus verschiedenen Gruppen einschleichen.</p><p>Um die Sache weiter zu verkomplizieren, zeigen erste Forschungsergebnisse, dass es in manchen Fällen gar nicht möglich ist, mit synthetischen Daten sowohl die Privatsphäre zu schützen als auch eine faire KI zu entwickeln. In einer kürzlich auf einer KI-Konferenz veröffentlichten Arbeit versuchten Forscher der Universität Toronto und des Vector-Instituts <a href=„<https://dl.acm.org/doi/10.1145/3442188.3445879>“ rel=„external noopener“ target=„\_blank“><strong>dies mit Röntgenaufnahmen der Brust [8]</strong></a> zu erreichen. Sie fanden heraus, dass sie nicht in der Lage waren, ein akkurate medizinisches KI-System zu erstellen, als sie versuchten, einen synthetischen Datensatz durch Kombination von Differential Privacy und GANs zu erstellen.</p><p>All dies heißt nicht, dass synthetische Daten nicht verwendet werden sollten. In der Tat kann das durchaus notwendig werden. Da die Aufsichtsbehörden KI-Systeme auf ihre Rechtskonformität hin überprüfen müssen, können sie marginalerweise nur so bedarfsgerechte, gezielte Testdaten generieren, so O’Neil. Aber dadurch ist es noch wichtiger, die Grenzen der Methode zu hinterfragen. „Synthetische Daten werden mit der Zeit wahrscheinlich besser werden“, sagt sie, „aber nicht von alleine.“</p><p>()</p><hr /><p><strong>URL dieses Artikels:</strong><br /><small><code><https://www.heise.de/-6071301></code></small></p><p><strong>Links in diesem Artikel:</strong><br /><small><code><strong>[1]</strong>&#160;<https://www.datagen.tech/technology/></code></small><br /><small><code><strong>[2]</strong>&#160;<https://synthesis.ai/></code></small><br /><small><code><strong>[3]</strong>&#160;<https://www.heise.de/tr/></code></small><br /><small><code><strong>[4]</strong>&#160;<https://www.click-ins.com/></code></small><br /><small><code><strong>[5]</strong>&#160;<https://mostly.ai/industries/></code></small><br />

/><small><code><strong>[6]</strong>&#160;[<https://www.heise.de/hintergrund/Kollektiver-Datenschutz-Was-dahinter-steckt-und-warum-er-noetig-ist-6054822.html>](https://www.heise.de/hintergrund/Kollektiver-Datenschutz-Was-dahinter-steckt-und-warum-er-noetig-ist-6054822.html)</code></small><br/><small><code><strong>[7]</strong>&#160;[<https://www.heise.de/hintergrund/Neuronales-Netz-erzeugt-attractive-Gesichter-fuer-jeden-Geschmack-6018436.html>](https://www.heise.de/hintergrund/Neuronales-Netz-erzeugt-attractive-Gesichter-fuer-jeden-Geschmack-6018436.html)</code></small><br/><small><code><strong>[8]</strong>&#160;[<https://dl.acm.org/doi/10.1145/3442188.3445879>](https://dl.acm.org/doi/10.1145/3442188.3445879)</code></small><br/><small><code><strong>[9]</strong>&#160;[<mailto:bsc@heise.de>](mailto:bsc@heise.de)</code></small><br/></p><p class=„printversioncopyright“><em>Copyright &#169; 2021 Heise Medien</em></p></html>

From:

<https://schnipsl.qgelm.de/> - Qgelm

Permanent link:

<https://schnipsl.qgelm.de/doku.php?id=wallabag:wb2wie-synthetische-datenstze-ki-systeme-verbessern-sollen>

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

