

# Fairness und Künstliche Intelligenz: Warum Metriken nicht ausreichen

Originalartikel

Backup

<html> <header class=„article-header“><h1 class=„articleheading“>Fairness und K&#252;nstliche Intelligenz: Warum Metriken nicht ausreichen</h1><div class=„publish-info“> Isabel B&#228;r</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/18/3/2/7/1/0/7/4/shutterstock\\_1430571869.jpg-5caa03e873d791fc.jpeg](https://heise.cloudimg.io/width/700/q75.png-lossy-75.webp-lossy-75.foil1/_www-heise-de/_imgs/18/3/2/7/1/0/7/4/shutterstock_1430571869.jpg-5caa03e873d791fc.jpeg)“ srcset=„[https://heise.cloudimg.io/width/700/q75.png-lossy-75.webp-lossy-75.foil1/\\_www-heise-de/\\_imgs/18/3/2/7/1/0/7/4/shutterstock\\_1430571869.jpg-5caa03e873d791fc.jpeg\\_700w](https://heise.cloudimg.io/width/700/q75.png-lossy-75.webp-lossy-75.foil1/_www-heise-de/_imgs/18/3/2/7/1/0/7/4/shutterstock_1430571869.jpg-5caa03e873d791fc.jpeg_700w), [https://heise.cloudimg.io/width/1050/q75.png-lossy-75.webp-lossy-75.foil1/\\_www-heise-de/\\_imgs/18/3/2/7/1/0/7/4/shutterstock\\_1430571869.jpg-5caa03e873d791fc.jpeg\\_1050w](https://heise.cloudimg.io/width/1050/q75.png-lossy-75.webp-lossy-75.foil1/_www-heise-de/_imgs/18/3/2/7/1/0/7/4/shutterstock_1430571869.jpg-5caa03e873d791fc.jpeg_1050w), [https://heise.cloudimg.io/width/1500/q75.png-lossy-75.webp-lossy-75.foil1/\\_www-heise-de/\\_imgs/18/3/2/7/1/0/7/4/shutterstock\\_1430571869.jpg-5caa03e873d791fc.jpeg\\_1500w](https://heise.cloudimg.io/width/1500/q75.png-lossy-75.webp-lossy-75.foil1/_www-heise-de/_imgs/18/3/2/7/1/0/7/4/shutterstock_1430571869.jpg-5caa03e873d791fc.jpeg_1500w), [https://heise.cloudimg.io/width/2300/q75.png-lossy-75.webp-lossy-75.foil1/\\_www-heise-de/\\_imgs/18/3/2/7/1/0/7/4/shutterstock\\_1430571869.jpg-5caa03e873d791fc.jpeg\\_2300w](https://heise.cloudimg.io/width/2300/q75.png-lossy-75.webp-lossy-75.foil1/_www-heise-de/_imgs/18/3/2/7/1/0/7/4/shutterstock_1430571869.jpg-5caa03e873d791fc.jpeg_2300w)“ alt=„Artificial Intelligence, Abstract Face, Created By, Neural Network, Machine Learning“ class=„img-responsive“ referrerpolicy=„no-referrer“ /><figcaption class=„akwa-caption“>(Bild:&#160;shuttersv/Shutterstock.com)</figcaption></figure><p><strong>Klassisches Software-Testing I&#228;sst sich nicht ohne Weiteres auf KI &#252;bertragen. Model Governance und interne Audits sind n&#246;tig, um Fairness zu gew&#228;hrleisten.</strong></p><p>Der Einsatz von K&#252;nstlicher Intelligenz (KI) bringt Verantwortung mit sich. Transparenz, Erkl&#228;rbarkeit, Fairness sind dabei wesentliche Prinzipien, die ebenso gew&#228;hrleistet sein m&#252;ssen wie die hohe Leistungsf&#228;higkeit des KI-Systems. Um diese Anforderungen einzuhalten, liegt es nahe, sich an Bereichen mit einer Tradition &#252;berpr&#252;fbarer Prozesse zu orientieren. Zwar funktionieren diese Prozesse nicht fehlerlos, aber ohne sie lassen sich Sicherheitsstandards nicht verwirklichen. Am offensichtlichsten ist das in sicherheitskritischen und regulierten Branchen wie der Medizin, aber auch in der Luft- und Raumfahrt oder im Finanzwesen.</p><header class=„a-boxheader“ data-collapse-trigger=„“>Young Professionals schreiben f&#252;r Young Professionals</header><div class=„a-boxtarget a-boxcontent a-inline-textboxcontent a-inline-textboxcontent-horizontal-layout“ data-collapse-target=„“><figure class=„a-inline-textboximage-container“><img alt=„[https://heise.cloudimg.io/width/4000/q50.png-lossy-50.webp-lossy-50.foil1/\\_www-heise-de/\\_imgs/71/2/8/9/8/5/2/3/shutterstock\\_1161966886-c893e738e55765fb.jpeg](https://heise.cloudimg.io/width/4000/q50.png-lossy-50.webp-lossy-50.foil1/_www-heise-de/_imgs/71/2/8/9/8/5/2/3/shutterstock_1161966886-c893e738e55765fb.jpeg)“ srcset=„[https://heise.cloudimg.io/width/8000/q30.png-lossy-30.webp-lossy-30.foil1/\\_www-heise-de/\\_imgs/71/2/8/9/8/5/2/3/shutterstock\\_1161966886-c893e738e55765fb.jpeg\\_2x](https://heise.cloudimg.io/width/8000/q30.png-lossy-30.webp-lossy-30.foil1/_www-heise-de/_imgs/71/2/8/9/8/5/2/3/shutterstock_1161966886-c893e738e55765fb.jpeg_2x)“ class=„c1“ referrerpolicy=„no-referrer“ /></figure><div class=„a-inline-textboxcontent-container“><p class=„a-inline-textboxsynopsis“>Dieser Beitrag ist Teil einer Artikelserie, zu der heise Developer junge Entwickler:innen einl&#228;dt &#8211; um &#252;ber aktuelle Trends, Entwicklungen und pers&#246;nliche Erfahrungen zu informieren. Die Reihe „Young Professionals“ erscheint im monatlichen Rhythmus. Bist du selbst ein „Young Professional“ und willst einen (ersten) Artikel schreiben? Schicke deinen Vorschlag an die Redaktion: [developer@heise.de](mailto:developer@heise.de). Wir stehen dir beim Schreiben zur Seite.</p><ul class=„a-inline-textboxlist“><li class=„a-inline-textboxitem“><a class=„a-inline-textboxtext“ href=„<https://www.heise.de/hintergrund/Alpine.js-Das-Schweizer-Taschenmesser-fuer-dynamische-Weboberflaechen-6177628.html>“ title=„Alpine.js: Das Schweizer Taschenmesser f&#252;r dynamische

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

Weboberfl&#228;chen“><strong>Alpine.js: Das Schweizer Taschenmesser f&#252;r dynamische Weboberfl&#228;chen [1]</strong></a></li><li class=„a-inline-textboxitem“><a class=„a-inline-textboxtext“ href=„<https://www.heise.de/hintergrund/Developer-Experience-Glueckliche-Entwickler-schreiben-besseren-Code-6150890.html>“ title=„Developer Experience: Gl&#252;ckliche Entwickler schreiben besseren Code“><strong>Developer Experience: Gl&#252;ckliche Entwickler schreiben besseren Code [2]</strong></a></li><li class=„a-inline-textboxitem“><a class=„a-inline-textboxtext“ href=„<https://www.heise.de/hintergrund/Ethik-und-Kuenstliche-Intelligenz-ein-neuer-Umgang-mit-KI-Systemen-6056059.html>“ title=„Ethik und K&#252;nstliche Intelligenz: Ein neuer Umgang mit KI-Systemen“><strong>Ethik und K&#252;nstliche Intelligenz: Ein neuer Umgang mit KI-Systemen [3]</strong></a></li><li class=„a-inline-textboxitem“><a class=„a-inline-textboxtext“ href=„<https://www.heise.de/developer/young-professionals-6065678.html>“ title=„Alle Beitr&#228;ge der Serie finden sich in der Rubrik“><strong>Alle Beitr&#228;ge der Serie finden sich in der Rubrik „Young Professionals“ [4]</strong></a></li></ul></div></div><p>&#196;hnlich wie diese Bereiche Prozesse ben&#246;tigen, um relevanten Anforderungen nachzukommen, ben&#246;tigt ein Unternehmen, das KI-Systeme einsetzt, geregelte Abl&#228;ufe, durch die es Zugriff auf Machine-Learning-Modelle (ML) kontrolliert, Richtlinien sowie gesetzliche Vorgaben umsetzt, die Interaktionen mit den Modellen und deren Ergebnissen verfolgt sowie festh&#228;lt, auf welcher Grundlage ein Modell erzeugt wurde. Insgesamt werden diese Prozesse als Model Governance bezeichnet. Model-Governance-Prozesse sind von Beginn an in jede Phase des ML-Lebenszyklus zu implementieren (Design, Development und Operations). Zur konkreten technischen Integration von Model Governance in den ML-Lebenszyklus hat die Verfasserin sich andernorts auf&#252;hrlicher ge&#228;u&#223;ert.</p><h3 class=„subheading“ id=„nav\_model0“>Model Governance: ein Muss im Regel- und Auflagenwald</h3><p>Model Governance ist nicht optional (siehe Kasten „Checkliste Model Governance“). Zum einen gibt es bereits bestehende Regularien, die Unternehmen in bestimmten Branchen erf&#252;llen m&#252;ssen. Am Beispiel des Finanzsektors l&#228;sst sich die Bedeutung von Model Governance gut illustrieren: Kreditvergabesysteme oder Zinsrisiko- und Preisbildungsmodelle f&#252;r Derivate sind risikoreich und verlangen ein hohes Ma&#223; an Kontrolle und Transparenz. Laut einer Algorithmia-Studie zu den <a href=„<https://algorithmia.com/blog/new-report-discover-the-top-10-trends-in-enterprise-machine-learning-for-2021>“ rel=„external noopener“ target=„\_blank“><strong>wichtigsten Trends im KI-Einsatz f&#252;r 2021 [5]</strong></a> ist die Mehrzahl der Unternehmen an die Erf&#252;llung rechtlicher Auflagen gebunden &#8211; 67 Prozent der Befragten m&#252;ssen mehreren Vorschriften entsprechen. Lediglich 8 Prozent gaben an, keinen gesetzlichen Vorgaben zu unterliegen.</p><p>Der Umfang der Regularien d&#252;rfte k&#252;nftig weiter zunehmen: so <a href=„<https://digital-strategy.ec.europa.eu/en/policies/european-approach-artificial-intelligence>“ rel=„external noopener“ target=„\_blank“><strong>ver&#246;ffentlichte die EU im April 2021 eine Verordnung als ersten Rechtsrahmen f&#252;r KI [6]</strong></a>, die bestehende Regularien erg&#228;nzen w&#252;rde. Der Entwurf teilt <a href=„[https://ec.europa.eu/info/strategy/priorities-2019-2024/europe-fit-digital-age/excellence-trust-artificial-intelligence\\_de](https://ec.europa.eu/info/strategy/priorities-2019-2024/europe-fit-digital-age/excellence-trust-artificial-intelligence_de)“ rel=„external noopener“ target=„\_blank“><strong>KI-Systeme in vier unterschiedliche Risikokategorien [7]</strong></a> ein („unzul&#228;ssig“, „hoch“, „begrenzt“, „minimal“). Die Risikokategorie definiert dabei Art und Umfang der Anforderungen, die an das jeweilige KI-System gestellt werden. KI-Software, die in die hohe Risikokategorie f&#228;llt, muss die strengsten Auflagen erf&#252;llen.</p><header class=„a-boxheader“ data-collapse-trigger=„>Checkliste Model Governance</header><div class=„a-boxtarget a-boxcontent a-inline-textboxcontent a-inline-textboxcontent-container“ data-collapse-target=„><p>Der Einsatz von Machine Learning bringt Verantwortung und Verpflichtungen mit sich. Um diesen Anforderungen nachzukommen, ben&#246;tigt ein Unternehmen Prozesse, durch die es</p><ul class=„boxtitel“><li>die Zugriffe auf ML-Modelle kontrolliert</li><li>Richtlinien/gesetzliche

Vorgaben umsetzt</li><li>die Interaktionen mit den ML-Modellen und deren Ergebnisse verfolgt</li><li>festhält, auf welcher Grundlage ein Modell erzeugt wurde</li></ul><p><a href= „<https://algorithmia.com/blog/model-governance>“ rel= „external noopener“ target= „\_blank“><strong>Model Governance [8]</strong></a> bezeichnet diese Prozesse in ihrer Gesamtheit</p><h4>Checkliste:</h4><ul class= „boxtitel“><li>Vollständige Modelldokumentation oder Berichte. Dazu gehört auch das Reporting der Metriken durch geeignete Visualisierungstechniken und Dashboards</li><li>Versionierung aller Modelle zur Herstellung von Transparenz nach außen (Erklärbarkeit und Reproduzierbarkeit)</li><li>Vollständige Datendokumentation zur Gewährleistung hoher Datenqualität und Einhaltung des Datenschutzes</li><li>Management von ML-Metadaten</li><li>Validierung von ML-Modellen (Audits)</li><li>Laufendes Überwachen und Protokollieren von Modellmetriken</li></ul></div><p>Dazu <a href= „[https://germany.representation.ec.europa.eu/news/fur-vertrauenswürdige-künstliche-intelligenz-eu-kommission-legt-weltweit-ersten-rechtsrahmen-vor-2021-04-21\\_de](https://germany.representation.ec.europa.eu/news/fur-vertrauenswürdige-künstliche-intelligenz-eu-kommission-legt-weltweit-ersten-rechtsrahmen-vor-2021-04-21_de)“ rel= „external noopener“ target= „\_blank“><strong>zählen folgende Aspekte [9]</strong></a>: Robustheit, Sicherheit, Genauigkeit (Accuracy), Dokumentation und Protokollierung sowie angemessene Risikobewertung und Risikominderung. Weitere Anforderungen sind die hohe Qualität der Trainingsdaten, Diskriminierungsfreiheit, Nachvollziehbarkeit, Transparenz, menschliche Überwachung sowie die Erforderlichkeit einer Konformität&#228;tspr&#252;fung und der Nachweis der Konformität&#228;t mit der KI-Verordnung durch eine CE-Kennzeichnung (siehe Kasten „Plan it Legal“). Beispiele für die ML-Systeme dieser Kategorie sind private und öffentliche Dienstleistungen (wie die Bonitätsprüfung) oder Systeme, die in der Schul- oder Berufsausbildung eingesetzt werden, um über den Zugang zu Bildung und den beruflichen Werdegang einer Person zu entscheiden (beispielsweise bei der automatisierten Bewertung von Prüfungen).</p><header class= „a-boxheader“ data-collapse-trigger= „>Plan it Legal: KI-Verordnung und Konformität</header><div class= „a-boxtarget a-boxcontent a-inline-textboxcontent“ data-collapse-target= „><figure class= „a-inline-textboximage-container“><img alt= „ class= „float-center c1“ src= „[https://heise.cloudimg.io/width/5760/q50.png-lossy-50.webp-lossy-50.foil1/\\_www-heise-de\\_imgs/18/3/2/7/1/0/7/4/shutterstock\\_497068513-58db4fcaddc30e4b.jpg](https://heise.cloudimg.io/width/5760/q50.png-lossy-50.webp-lossy-50.foil1/_www-heise-de_imgs/18/3/2/7/1/0/7/4/shutterstock_497068513-58db4fcaddc30e4b.jpg)“ srcset= „[https://heise.cloudimg.io/width/11520/q30.png-lossy-30.webp-lossy-30.foil1/\\_www-heise-de\\_imgs/18/3/2/7/1/0/7/4/shutterstock\\_497068513-58db4fcaddc30e4b.jpg](https://heise.cloudimg.io/width/11520/q30.png-lossy-30.webp-lossy-30.foil1/_www-heise-de_imgs/18/3/2/7/1/0/7/4/shutterstock_497068513-58db4fcaddc30e4b.jpg) 2x“ referrerpolicy= „no-referrer“ /><figcaption class= „a-caption a-caption-textbox“>(Bild:&#160;Marian Weyo/Shutterstock.com)</figcaption></figure><div class= „a-inline-textboxcontent“><p>Die Konformität von HRKI mit der KI-Verordnung wird die Voraussetzung für die Vermarktung in der EU werden. Sie lässt sich über eine CE-Kennzeichnung nachweisen. Die EU wird zudem Standards verabschieden, bei deren Einhaltung die Konformität&#228;t mit der Verordnung anzunehmen ist.</p><p>Für die umfassenden Tests, die nach der KI-Verordnung anfallen, sollen die zuständigen Behörden Sandboxing Schemes entwickeln, also Vorgaben für sichere Testumgebungen. Die Konformität&#228;tspr&#252;fung für KI beruht auf einer ex-ante-Sicht, hat aber gleichwohl Ähnlichkeiten mit der Datenschutzfolgenabschätzung nach der DSGVO. Mehr Informationen hierzu finden sich im Blogeintrag von Dr. Benhard Freund bei planit.legal: „<a href= „<https://planit.legal/das-ki-gesetz-der-eu-entwurf-und-diskussionsstand/>“ rel= „external noopener“ target= „\_blank“><strong>Das KI-Gesetz der EU &#8211; Entwurf und Diskussionsstand [10]</strong></a>.“</p></div></div><h3 class= „subheading“ id= „nav\_konformität\_f1“>Konformität&#228;t für europäische KI-Auflagen erreichen</h3><p>Da die Verordnung nicht nur für Unternehmen in der EU ansprechende Unternehmen und Einzelpersonen gelten soll, sondern jedes Unternehmen, das KI-Dienste innerhalb der EU anbietet, hat das Gesetz einen ähnlichen Anwendungsbereich wie die DSGVO. Die Verordnung muss sowohl vom EU-Parlament gebilligt werden als auch die Gesetzgebungsverfahren

der einzelnen Mitgliedsstaaten passieren. Wenn das EU-Parlament die Verordnung billigt und sie die legislativen Prozesse der EU-Staaten passiert, tritt das Gesetz frühestens 2024 in Kraft. Dann müssen Hochrisikosysteme während der Entwicklung eine Konformitätsbewertung für KI-Auflagen durchlaufen, um das KI-System in einer EU-Datenbank registrieren zu lassen. Im letzten Schritt ist eine Konformitätszertifizierung notwendig, sodass KI-Systeme die notwendige CE-Kennzeichnung erhalten, damit ihre Anbieter sie in den Verkehr bringen können. Wichtig ist außerdem, dass Regulierung nicht der einzige ausschlaggebende Aspekt für Model-Governance-Prozesse ist. Denn auch Modelle, die in schweren regulierten Kontexten im Einsatz sind, müssen an Model Governance nicht vorbei [11]. Neben der Erfüllung gesetzlicher Vorgaben müssen Unternehmen wirtschaftliche Einbußen und Reputationsverluste ebenso abwenden wie juristische Schwierigkeiten. ML-Modelle, die einer Marketing-Abteilung Informationen über die Zielgruppe liefern, müssen im Betrieb an Präsentation verlieren und eine falsche Informationsgrundlage für wichtige Folgeentscheidungen bereitstellen. Somit stellen sie ein finanzielles Risiko dar. Model Governance wird also nicht nur zur Erfüllung rechtlicher Vorgaben, sondern auch zur Qualitäts sicherung von ML-Systemen und zur Minderung unternehmerischer Risiken benötigt.

**Model Governance als Herausforderung**

Die sich abzeichnenden EU-Vorgaben, bestehende Regelungen und Unternehmensrisiken machen es notwendig, Model-Governance-Prozesse von Beginn an zu implementieren. Die Bedeutung von Model Governance ergibt sich für viele Unternehmen allerdings oft erst dann, wenn ML-Modelle in die Produktion gehen und in Einklang mit gesetzlichen Regelungen stehen sollen. Dazu kommt, dass der abstrakte Charakter rechtlicher Vorgaben Unternehmen vor die Herausforderung der praktischen Umsetzung stellt: So geben [Algorithmia](https://algorithmia.com/blog/model-governance) [12] nach der bereits zitierten 56 Prozent der Befragten die Implementierung von Model Governance als eine der größten Herausforderungen an, um ML-Anwendungen langfristig erfolgreich in Produktion zu bringen. Dazu passen auch die Zahlen der „State of AI in 2021“-Studie mit Blick auf die Risiken Künstlicher Intelligenz: 50 Prozent der befragten Unternehmen geben die Einhaltung gesetzlicher Vorschriften als Risikofaktor an, andere hoben Mangel bei Erklärbarkeit (44 Prozent der Befragten), Reputation (37 Prozent), Gerechtigkeit und Fairness (30 Prozent) als relevante Risikofaktoren hervor.

**Audits als standardisierte Prozesse im Model-Governance-Framework**

Ein wichtiger Bestandteil von Model Governance sind Audits [13] als Werkzeuge, um zu prüfen, ob KI-Systeme den Unternehmensrichtlinien, Branchenstandards oder Vorschriften entsprechen. Dabei gibt es interne und externe Audits. Die im Artikel „[Ethik und Künstliche Intelligenz: ein neuer Umgang mit KI-Systemen \[14\]](https://www.heise.de/hintergrund/Ethik-und-Künstliche-Intelligenz-ein-neuer-Umgang-mit-KI-Systemen-6056059.html)“ auf *Heise* von der Verfasserin besprochene Studie „Gender Shades“ ist ein Beispiel für einen externen Auditprozess: Sie prüfte Gesichtserkennungssysteme gegenüber Anbieter hinsichtlich ihrer Genauigkeit bezüglich des Geschlechtes und der Ethnie und konnte dabei eine abweichende Präsentation des Modells je nach Ethnie und Geschlecht feststellen.

Dieser Blick von außen ist aber limitiert, da externe Prüfprozesse nur Zugang zu Modellergebnissen, aber nicht zu den zugrundeliegenden Trainingsdaten oder Modellversionen besitzen. Das sind wertvolle Quellen, die Unternehmen in einem internen Auditprozess einbeziehen müssen. Diese Prozesse sollen eine kritische Reflexion

&#252;ber die potenziellen Auswirkungen eines Systems erm&#246;glichen. Zun&#228;chst sind jedoch an dieser Stelle Grundlagen &#252;ber KI-Systeme zu kl&#228;ren. </p><h3 class=„subheading“ id=„nav\_besonderheiten4“>Besonderheiten von KI-Systemen</h3><p>Um KI-Software pr&#252;fen zu k&#246;nnen, ist es wichtig zu <a href=„<https://christophm.github.io/interpretable-ml-book/what-is-machine-learning.html>“ rel=„external noopener“ target=„\_blank“><strong>verstehen, wie Machine Learning funktioniert [15]</strong></a>: Maschinelles Lernen besteht aus einer Reihe von Methoden, die Computer verwenden, um Vorhersagen oder Verhaltensweisen auf der Grundlage von Daten zu treffen und zu verbessern. Um diese Vorhersagemodele aufzubauen, m&#252;ssen ML-Modelle eine Funktion finden, die zu einer bestimmten Eingabe eine Ausgabe (Label) erzeugt. Daf&#252;r ben&#246;tigt das Modell Trainingsdaten, die zu den Eingabedaten die jeweils passende Ausgabe enthalten. Dieses Lernen tr&#228;gt die Bezeichnung „&#252;berwachtes Lernen“. Im Trainingsprozess sucht das Modell mithilfe mathematischer Optimierungsverfahren eine Funktion, die den unbekannten Zusammenhang zwischen Ein- und Ausgabe so gut wie m&#246;glich abbildet.</p><p>Ein Beispiel f&#252;r eine Klassifizierung w&#228;re eine Sentimentanalyse, die untersuchen soll, ob Tweets positive oder negative Stimmungen (Sentiments) enthalten. In diesem Fall w&#228;re ein Input ein einzelner Tweet, und das dazugeh&#246;rige Label das codierte Sentiment, das f&#252;r diesen Tweet festgelegt wurde (&#8722;1 f&#252;r ein negatives, 1 f&#252;r ein positives Sentiment). Im Trainingsprozess lernt der Algorithmus mit diesen annotierten Trainingsdaten, wie Eingabedaten mit dem Label zusammenh&#228;ngen. Nach dem Training kann der Algorithmus dann neue Tweets selbstst&#228;ndig einer Klasse zuordnen.</p><h3 class=„subheading“ id=„nav\_komplexere5“>Komplexere Komponenten im Machine-Learning-Bereich</h3><p>Somit lernt ein ML-Modell die Entscheidungslogik im Trainingsprozess, statt die Logik mit einer Abfolge von typischen Wenn-Dann-Regeln explizit im Code zu definieren, wie es in der Softwareentwicklung typisch w&#228;re. Dieser grundlegende Unterschied zwischen traditioneller und KI-Software f&#252;hrt dazu, dass sich Methoden des klassischen Softwaretestens nicht direkt auf KI-Systeme &#252;bertragen lassen. Das Testen verkompliziert sich dadurch, dass zus&#228;tzlich zum Code die Daten und das Modell selbst hinzukommen, wobei alle drei Komponenten sich gem&#228; &#223; dem Change-Anything/Change-Everything-Prinzip gegenseitig bedingen (hierzu mehr unter „<a href=„<https://papers.nips.cc/paper/2015/file/86df7dcfd896fcf2674f757a2463eba-Paper.pdf>“ rel=„external noopener“ target=„\_blank“><strong>Hidden Technical Debt in Machine Learning Systems [16]</strong></a>“).</p><p>Unterscheiden sich beispielsweise die Daten im produktiven System von den Daten, mit denen ein Modell trainiert wurde (Distribution Shifts), kommt es zum Leistungsabfall des Modells (Model Decay). In diesem Fall muss ein Modell schnell mit frischen Trainingsdaten trainiert und re-deployed werden. Erschwerend kommt hinzu, dass <a href=„<https://ieeexplore.ieee.org/abstract/document/8920538>“ rel=„external noopener“ target=„\_blank“><strong>das Testen von KI-Software ein noch offenes Forschungsfeld [17]</strong></a> ohne Konsens und ohne Best Practices ist.</p><h3 class=„subheading“ id=„nav\_ethische6“>Ethische Prinzipien als nicht-funktionale Eigenschaften</h3><p>Die relevanten Testaspekte von KI-Software lassen sich in funktionale und nicht-funktionale Eigenschaften einteilen. Correctness als funktionale Eigenschaft l&#228;sst sich <a href=„<https://www.aaai.org/Papers/Workshops/2006/WS-06-06/WS06-06-003.pdf>“ rel=„external noopener“ target=„\_blank“><strong>durch Metriken wie Accuracy und Precision/Recall mathematisch direkt erfassen [18]</strong></a>. Sie geben an, wie hoch die &#222;bereinstimmung zwischen den Vorhersagen des trainierten Modells und den tats&#228;chlichen Predictions ist (Gold Standard). Dazu gibt es etablierte Validierungsverfahren wie die Kreuzvalidierung, die durch Isolation der Testdaten &#252;ber eine Datenstichprobe pr&#252;ft, wie gut das trainierte Modell die richtigen Modellergebnisse (Labels) f&#252;r neue Daten vorhersagt.</p><p><a href=„<https://ieeexplore.ieee.org/document/9000651>“ rel=„external noopener“ target=„\_blank“><strong>Nicht-funktionale Eigenschaften entsprechen den ethischen Prinzipien [19]</strong></a> wie Fairness, Datenschutz, Interpretierbarkeit, Robustheit und

Sicherheit. Anders als funktionale Eigenschaften können sie nicht auf einen breiten Fundus standardisierter Metriken und Praktiken aus dem Bereich des maschinellen Lernens zurückblicken. Auch hier besteht die Herausforderung darin, dass das Testen nicht-funktionaler Eigenschaften von KI-Software (noch) nicht standardisiert ist. Erschwerend kommen Abweichungen zwischen verschiedenen Eigenschaften hinzu: <a href="https://people.mpi-sws.org/~gummadi/papers/process\_fairness.pdf" rel="external noopener" target="\_blank"><strong>Fairness verringert die Accuracy [20]</strong></a> und umgekehrt.</p><p>Metaphorisch lässt sich KI-Software als Kraftwerk bezeichnen: Ein funktional einwandfreier, reibungsloser Betrieb heißt nicht, dass das Kraftwerk der Umwelt keinen Schaden zufügt. Der funktionsfreie Ablauf entspricht den funktionalen, der Schutz der Umwelt den nicht-funktionalen Kriterien. Die Metapher zeigt, dass es für funktionale und nicht-funktionale Eigenschaften verschiedene Prozesse braucht. Für Erstere sind Best Practices aus dem ML-Fundus anwendbar, für Letztere hingegen braucht es noch Forschungsarbeit. Im weiteren Verlauf soll es hier um den Aspekt der Fairness, eine Beschreibung von Konzepten sowie Teststrategien und um die Hervorhebung von Model Governance gehen, die die Transformation ethischer Prinzipien wie Fairness in der praktischen KI-Softwareentwicklung unterstützen.</p><h3 class="subheading" id="nav\_wie\_entsteht7">Wie entsteht Fairness?</h3><p>Zunächst ist die Frage interessant, <a href="https://papers.ssrn.com/sol3/papers.cfm?abstract\_id=2477899" rel="external noopener" target="\_blank"><strong>wie Ungerechtigkeit (Unfairness) überhaupt entsteht [21]</strong></a>. Die Regel ist dabei einfach: Was die Modelle lernen, manifestiert sich in den Trainingsdaten. Im Übrigen Lernen bestehen Trainingsdaten aus den Eingabedaten und dazugehörigen Labels. Wenn die Datenlabels Bias enthalten, wird das Modell diese Grundeinstellung übernehmen und von Anfang an lernen. Daher ist es wichtig, die Labels ausreichend zu überprüfen. Bias kann sich aber auch innerhalb aus den Daten, nicht nur aus den Labels ergeben: Enthalten die Trainingsdaten an sich bereits Bias, greift der Algorithmus ihn ebenfalls auf. Dieses Problem besteht etwa bei umfangreichen, mit großen Datenmengen aus dem Internet trainierten Sprachmodellen. Es liegt sich nachweisen, dass die Leistungsfähigkeit eines Modells mit der Stärke eines stereotypen Bias korreliert: Mit steigender Präzision nimmt auch der Bias zu.</p><p>Auch ein geringer Stichprobenumfang bei Minderheitsgruppen kann zu einer Homogenisierung des Lernprozesses des Modells zugunsten der Mehrheitsgruppen führen, beispielsweise durch mehr Fotos männlicher als weiblicher Gesichter in den Trainingsdaten. Neben den Daten spielen auch im Trainingsprozess verwendete Merkmale (Features) eine Rolle. Kann das Modell nicht ausreichend viele Merkmale nutzen, erschwert das dem Algorithmus, den Zusammenhang zwischen Ein- und Ausgabe zu lernen. Aus diesem Grund reagierte IBM mit <a href="https://www.ibm.com/blogs/research/2019/01/diversity-in-faces/" rel="external noopener" target="\_blank"><strong>Diversity in Faces als Versuch, die Diversität der Fotos in den Trainingsdaten zu erhöhen [22]</strong></a>. Und schließlich kann es Merkmale „Stellvertreter“ für ausgeschlossene sensible Attribute sein: Auch wenn geschützte Attribute bei der Entscheidungsfindung nicht explizit verwendet werden, können sie implizit beteiligt sein, wenn sie mit den ausgeschlossenen Merkmalen korrelieren.</p><header class="a-boxheader" data-collapse-trigger="">>Fairness in KI-Systemen</header><div class="a-boxtarget a-boxcontent a-inline-textboxcontent a-inline-textboxcontent-container" data-collapse-target="">><p>Das Ziel der Gewährleistung von Fairness ist der Schutz sensibler Attribute wie Geschlecht, Religionszugehörigkeit oder sexueller Orientierung vor unfairer algorithmischer Entscheidungsfindung. Das Recht auf Diskriminierungsfreiheit ist im EU-Rechtsentwurf für KI-Systeme der hohen Risikokategorie explizit verbrieft.</p><p>Während sich die Ungerechtigkeit bei der Gender-Shades-Studie leicht intuitiv erfassen lässt, besteht nun die Herausforderung darin, den abstrakten Begriff der Fairness objektiv, metrikbasiert und möglichst skalierbar zu definieren.</p></div><h3

class=„subheading“ id=„nav\_definitionen8“>Definitionen f&#252;r Fairness und Ableitungen von Teststrategien</h3><p>Welche Audits und welche Metriken bieten sich an, um Fairness zu testen? Die bereits bekannte Konsensl&#252;cke klafft auch f&#252;r die <a href=„<https://arxiv.org/pdf/1906.10742.pdf>“ rel=„external noopener“ target=„\_blank“><strong>Definition von Fairness [23]</strong></a> auseinander. Erschwerend kommt hinzu, dass die Vielf&#228;ltigkeit der verschiedenen Ursachen f&#252;r Fairness zeigt, dass sich Fairness nicht mit einer simplen Metrik oder Teststrategie herstellen l&#228;sst &#8211; Fairness-Audits m&#252;ssen Teil der Model-Governance-Prozesse sein, die die Qualit&#228;tssicherung der Trainingsdaten und des Modells sicherstellen. Dazu kommt, dass die verschiedenen Anwendungsf&#228;lle f&#252;r KI zu vielf&#228;ltig sind, als dass es eine gut generalisierende One-Size-Fits-All-L&#246;sung geben k&#246;nnte. Die Frage, wie Fairness gemessen und nachgewiesen werden kann, l&#228;sst sich also nicht nur an einer simplen Metrik festmachen. Dennoch soll es zun&#228;chst um konkrete M&#246;glichkeiten gehen, Fairness quantitativ zu erfassen, bevor diese Audits in das Model-Governance-Framework eingebettet werden.</p><p>Statistische Ans&#228;tze bieten die am leichtesten messbaren Definitionen von Fairness, und sie bilden gleichzeitig die Grundlage f&#252;r weiterf&#252;hrende Ans&#228;tze. Zur Quantifizierung von Fairness lassen sich statistische Metriken nutzen. Von diesen Messgr&#246;ßen leiten sich Definitionen ab, die sich auf die Ausgabe von Modellen konzentrieren. Fairness l&#228;sst sich aufgrund &#228;hnlicher <a href=„[https://people.mpi-sws.org/~gummadi/papers/process\\_fairness.pdf](https://people.mpi-sws.org/~gummadi/papers/process_fairness.pdf)“ rel=„external noopener“ target=„\_blank“><strong>Fehlerquoten der Ausgaben f&#252;r unterschiedlich sensible demografische Gruppen [24]</strong></a> definieren. Entsprechend ist ein Algorithmus dann fair, wenn Gruppen, die auf der Grundlage sensibler Attribute ausgew&#228;hlt werden, die gleiche Wahrscheinlichkeit von vorteilhaften Entscheidungsergebnissen haben („<a href=„[https://people.ece.ubc.ca/mjulia/publications/Fairness\\_Definitions\\_Explained\\_2018.pdf](https://people.ece.ubc.ca/mjulia/publications/Fairness_Definitions_Explained_2018.pdf)“ rel=„external noopener“ target=„\_blank“><strong>Group Fairness [25]</strong></a>“).</p><h3 class=„subheading“ id=„nav\_gleichheit\_der9“>Gleichheit der Gesamtgenauigkeit</h3><p>Zudem l&#228;sst sich untersuchen, ob die Genauigkeit des Modells f&#252;r verschiedene Subgruppen gleich ist (Gleichheit der Gesamtgenauigkeit). Am Beispiel einer Kreditw&#252;rdigkeitspr&#252;fung w&#228;re diese Definition von Fairness dann erf&#252;llt, wenn die Wahrscheinlichkeit f&#252;r Personen m&#228;nnlichen und weiblichen Geschlechts gleich ist, dass Antragsstellende mit einem tats&#228;chlich guten Kreditscore als kreditw&#252;rdig eingestuft werden und dass solchen mit einem schlechten Kreditscore die Kreditw&#252;rdigkeit abgesprochen wird, ohne Ansehen der Geschlechtszugeh&#246;rigkeit.</p><p>F&#252;r das Testen statistischer Ans&#228;tze sind bereits erste L&#246;sungen verf&#252;gbar: <a href=„<https://github.com/tensorflow/fairness-indicators>“ rel=„external noopener“ target=„\_blank“><strong>Fairness Indicators von TensorFlow [26]</strong></a> ist eine Bibliothek, die die Berechnung h&#228;ufig identifizierter Fairness-Metriken mit verbesserter Skalierbarkeit auf gro&#223;en Datens&#228;tzen und Modellen bietet. Dar&#252;ber hinaus unterst&#252;tzt Fairness Indicators die Auswertung der Verteilung von Datens&#228;tzen und der Modellleistung &#252;ber verschiedene Benutzergruppen sowie die Berechnung statistisch signifikanter Unterschiede auf der Basis von Konfidenzintervallen.</p><h3 class=„subheading“ id=„nav\_statistische10“>Statistische Ans&#228;tze und Counterfactual Fairness</h3><p>Zwar sind statistische Ans&#228;tze gut messbar, sie k&#246;nnen jedoch zu kurz greifen. Fairness l&#228;sst sich nicht allein durch &#228;hnliche Fehlklassifizierungsquoten erkl&#228;ren, insbesondere, wenn alle anderen Attribute mit Ausnahme des sensiblen Attributs ignoriert werden. Beispielsweise k&#246;nnte ein KI-System zur Kreditw&#252;rdigkeitspr&#252;fung denselben Anteil m&#228;nnlicher und weiblicher Bewerber eine positive Bewertung zuweisen &#8211; statistische Ans&#228;tze w&#252;rden das Modell dann als gerecht beurteilen. Doch wenn die m&#228;nnlichen Bewerber zuf&#228;llig ausgew&#228;hlt wurden, w&#228;rend weibliche Bewerber schlicht diejenigen sind, die die meisten Ersparnisse haben, w&#228;re Fairness nicht

gegeben.</p><p>Similarity-based Measures stellen nicht die Modellergebnisse und Fehlklassifizierungsquoten, sondern den Prozess der Entscheidungsfindung sowie die Verwendung von Merkmalen im Trainingsprozess in den Vordergrund. Daraus lässt sich „Fairness through Unawareness“ als Konzept für Fairness ableiten: Algorithmen können als fair gelten, wenn geschätzte Attribute aus den Trainingsdaten ausgeschlossen worden sind. In unserem Beispiel bedeutet dies, dass geschlechtsspezifische Merkmale nicht für das Training des Modells verwendet werden, sodass Entscheidungen nicht auf diesen Merkmalen beruhen können. Doch auch dieser Ansatz hat Einschränkungen: Das Ausschließen geschätzter Attribute reicht nicht aus, da andere, ungeschätzte Attribute Informationen enthalten können, die mit den ausgeschlossenen geschätzten Attributen korrelieren („[Counterfactual Fairness \[27\]](https://papers.nips.cc/paper/2017/file/a486cd07e4ac3d270571622f4f316ec5-Paper.pdf)“). In diesem Fall wäre das ursprünglich ausgeschlossene Attribut implizit in anderen Attributen enthalten und würde den Entscheidungsprozess indirekt beeinflussen (siehe auch Kasten „Kontrafaktisch testen“).</p><header class="a-boxheader" data-collapse-trigger="Kontrafaktisch testen"></header><div class="a-boxtarget a-boxcontent a-inline-textboxcontent a-inline-textboxcontent-container" data-collapse-target=""><p>Causal-Reasoning-Ansätze stützen sich auf Werkzeuge der Kausalinferenz. Die Definition der kontrafaktischen Fairness basiert auf der Intuition, dass eine Entscheidung gegenüber einer Person dann fair ist, wenn sie in der tatsächlichen Welt und in einer kontrafaktischen Welt, in der die Person einer anderen demografischen Gruppe angehört, gleich ist.</p><p>Damit ist Counterfactual Fairness dann gegeben, wenn sich eine Prediction nicht ändert, obwohl das geschätzte Attribut in das kontrafaktische Gegenteil verkehrt wird. Beispielsweise muss die Entscheidung für oder gegen die Kreditwürdigkeit einer Person gleich ausfallen, wenn das Attribut von männlich in weiblich; verändert wird.</p></div><h3 class="subheading" id="nav\_unfairness\_mit11">Unfairness mit manipulierten Daten aufdecken</h3><p>Adversariales Testen ist eine gängige Strategie, die zur Aufdeckung von Schwachstellen einen bewilligen Angriff auf ein System simuliert. Beim adversarialen Testen erhält das Modell Eingabedaten, die mit kleinen, absichtlichen Merkmalsmanipulationen versehen sind. Auf diese Weise wird getestet, <a href="https://christophm.github.io/interpretable-ml-book/adversarial.html" rel="external noopener" target="\_blank"><strong>ob das Modell für speziell zugeschnittene Eingabedaten unerwartete Vorhersagen trifft [28]</strong></a>. Die Manipulation der Eingabedaten ist domänenunspezifisch und kann von Analysen algorithmischer Unfairness inspiriert sein. Die Idee, die Reaktion eines Modells auf Eingabedaten zu testen und auf diese Art von Bias zu quantifizieren, findet sich mittlerweile auch in Frameworks.</p><p>So lässt sich <a href="https://stereoset.mit.edu/" rel="external noopener" target="\_blank"><strong>mit dem Benchmarking-Datensatz StereoSet prüfen [29]</strong></a>, ob und wie stark die Vorurteile von Sprachmodellen bezüglich Geschlecht, Ethnie, Religion und Profession ausgeprägt sind: Entwicklerinnen und Entwickler können fertig trainierte Sprachmodelle einreichen, um diskriminierende Entscheidungsfindung in Sprachmodellen zu messen und gleichzeitig die Sprachmodellierungsleistung zu berücksichtigen. StereoSet betrachtet die Gesamtleistung des Modells als gut, wenn das Modell in der Lage ist, den Zielkonflikt zwischen Genauigkeit und Fairness abzuschwächen und so ein genaues Verständnis natürlicher Sprache bei gleichzeitiger Minimierung von Verzerrungen zu gewährleisten. Frameworks wie diese können zwar eine gute Leitlinie sein, ersetzen aber nicht das individuelle Testen, das fest in einem Prüfungsprozess eingebettet sein muss.</p><h3 class="subheading" id="nav\_warum\_audits12">Warum Audits alleine nicht ausreichen</h3><p>Die verschiedenen Prüfprozesse sind wichtig; sie allein genügen aber nicht, um Fairness zu gewährleisten. Vielmehr kommen Audits als eine der ersten Möglichkeiten in Betracht,

um Probleme zu identifizieren. Sie müssen Teil des Model-Governance-Frameworks sein und sollten es ergänzen. Für sich allein besitzen sie hingegen keine Aussagekraft; nur ein ganzheitlicher Ansatz kann alle Aspekte berücksichtigen, die für Fairness eine Rolle spielen. Neben der Validierung funktionaler und nicht-funktionaler Anforderungen, die die hier beschriebenen Auditprozesse zum Testen von Fairness enthalten kann, ist sauberes Dokumentieren eine weitere wichtige Komponente im Model-Governance-Framework.

Dokumentationen sollten bereits in der ersten Phase des ML-Lebenszyklus, der Entwicklung, begonnen werden (mehr dazu im „[Practitioners' Guide to MLOps \[30\]](https://services.google.com/fh/files/misc/practitioners_guide_to_mlops_whitepaper.pdf)“). In der Development-Phase geht es um den Aufbau einer robusten und reproduzierbaren Trainingsprozedur, die aus Datenverarbeitungs- und Modellaufbauschritten besteht. Dieser Aufbauprozess ist experimentell und iterativ, wobei wichtige Informationen über Daten (Auswahl und Definition von Features, Aufteilung der Trainings-, Validierungs- und Testdaten, Schema und Statistiken), Modelle (verschiedene getestete Modelle) und Parameter festzuhalten sind (Experimental Tracking).

**Nach dem Training: KI-Modelle evaluieren**

Nach dem Aufbau der Trainingsprozedur gilt es, entwickelte Modelle hinsichtlich funktionaler und nicht-funktionaler Eigenschaften zu evaluieren (an dieser Stelle sind die beschriebenen Auditstrategien zum Testen von Fairness relevant). Die Ergebnisse der Evaluation und alle Informationen über den Aufbauprozess der Trainingsprozedur fließen in die Dokumentation ein, die zusätzlich eine Erklärung des Use-Case-Kontextes, eine High-Level-Erklärung des Algorithmus, Modellparameter, Anweisungen zur Reproduktion des Modells und Beispiele für das Training der Algorithmen sowie Beispiele für das Treffen von Predictions durch den Algorithmus enthalten sollte.

Die Dokumentation lässt sich durch Toolkits wie Model Cards und Data Sheets praktisch unterstützen. Data Sheets halten fest, welche Mechanismen oder Verfahren für die Datenerhebung verwendet wurden oder ob ethische Überprüfungsverfahren (Audits) stattgefunden haben.

[Model Cards \[31\]](https://arxiv.org/pdf/1810.03993.pdf) ergänzen [Data Sheets \[32\]](https://arxiv.org/pdf/1803.09010.pdf) und informieren über die Art der Modellerstellung, die bei der Entwicklung getroffenen Annahmen sowie bezüglich des Modellverhaltens bei verschiedenen kulturellen, demografischen oder physischen Gruppen.

**Unternehmenspolitik als Schlüssel für ethische KI**

Vollständige Dokumentation schafft Reproduzierbarkeit und Transparenz nach außen. Nach dem Deployment

[muss diese Sichtbarkeit \(Observability\) im produktiven System gegeben sein \[33\]](https://www.oreilly.com/library/view/the-framework-for/9781098100483/ch01.html). Hier spielt zum einen die Versionierung von Modellen und Datensätzen eine wichtige Rolle. Die Versionierung dient der Wahrung des Unveränderlichkeitsgrundsatzes der Modelle, sodass alle Modelle sich ohne Datenverluste und Veränderung reproduzieren lassen. Damit ist auch gewährleistet, dass eine Model Prediction der Modellversion, die sie produziert hat, zugeordnet werden kann.

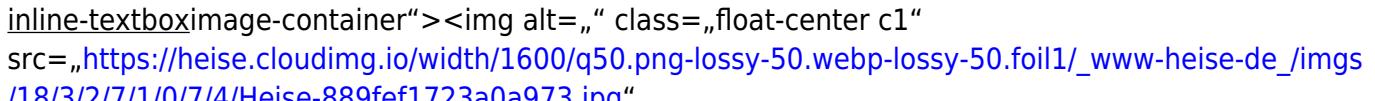
Zum anderen muss ein Monitoring-System die Leistung des produktiven Modells kontinuierlich überwachen und relevante Metriken in einem Report zusammenfassen sowie visualisieren. Diese Werte aus dem Model-Logging sollten in Metriken aufbereitet und in Dashboards zu Protokollierungs-, Analyse- und Kommunikationszwecken visualisierbar sein. Wird im Monitoring der Leistungsabfall eines Modells (Model Decay) festgestellt, muss das Modell mit neuen Trainingsdaten trainiert und dann re-deployed werden.

**Empfehlung: Audits vor jedem neuen Deployment**

Vor jedem neuen Deployment sollten erneut Audits stattfinden, um ethische, rechtliche oder geschäftliche Risiken zu kontrollieren. Ethische Prinzipien wie Fairness sind auf

jeder Ebene der Softwareentwicklung zu berücksichtigen, unter anderem bereits bei der Datenbeschaffung. Fairness lässt sich nicht mit einer simplen Metrik oder Teststrategie herstellen: Es braucht eine entsprechend ausgerichtete Unternehmenspolitik, die das anerkennt. Ohne Model Governance sind KI-Systeme hinsichtlich der Einhaltung gesetzlicher Vorgaben sowie der Minderung des unternehmerischen Risikos unkalkulierbar.

Young Professionals schreiben darüber Young Professionals



Isabel B. studiert Data Engineering am Hasso-Plattner-Institut und arbeitet als Werkstudentin bei INNOQ. Sie beschäftigt sich mit Fragen rund um den langfristig erfolgreichen Einsatz von Künstlicher Intelligenz (KI), wozu insbesondere MLOps und das Implementieren von Model Governance gehören.

Links in diesem Artikel:

[1] <https://www.heise.de/hintergrund/Alpine-js-Das-Schweizer-Taschenmesser-fuer-dynamische-Weboberflaechen-6177628.html>

[2] <https://www.heise.de/hintergrund/Developer-Experience-Glueckliche-Entwickler-schreiben-besseren-Code-6150890.html>

[3] <https://www.heise.de/hintergrund/Ethik-und-Kuenstliche-Intelligenz-ein-neuer-Umgang-mit-KI-Systemen-6056059.html>

[4] <https://www.heise.de/developer/young-professionals-6065678.html>

[5] <https://algorithmia.com/blog/new-report-discover-the-top-10-trends-in-enterprise-machine-learning-for-2021>

[6] <https://digital-strategy.ec.europa.eu/en/policies/european-approach-artificial-intelligence>

[7] [https://ec.europa.eu/info/strategy/priorities-2019-2024/europe-fit-digital-age/excellence-trust-artificial-intelligence\\_de](https://ec.europa.eu/info/strategy/priorities-2019-2024/europe-fit-digital-age/excellence-trust-artificial-intelligence_de)

[8] <https://algorithmia.com/blog/model-governance>

[9] [https://germany.representation.ec.europa.eu/news/fuer-vertrauenswurdige-kunstliche-intelligenz-eu-kommission-legt-weltweit-ersten-rechtsrahmen-vor-2021-04-21\\_de](https://germany.representation.ec.europa.eu/news/fuer-vertrauenswurdige-kunstliche-intelligenz-eu-kommission-legt-weltweit-ersten-rechtsrahmen-vor-2021-04-21_de)

[10] <https://planit.legal/das-ki-gesetz-der-eu-entwurf-und-diskussionsstand>

[11] <https://algorithmia.com/blog/model-governance>

[12] <https://algorithmia.com/blog/new-report-discover-the-top-10-trends-in-enterprise-machine-learning-for-2021>

[13] <https://dl.acm.org/doi/pdf/10.1145/3351095.3372873>

[14] <https://www.heise.de/hintergrund/Ethik-und-Kuenstliche-Intelligenz-ein-neuer-Umgang-mit-KI-Systemen-6056059.html>

/><small><code><strong>[15]</strong>&#160;<https://christophm.github.io/interpretable-ml-book/what-is-machine-learning.html></code></small><br>/><small><code><strong>[16]</strong>&#160;<https://papers.nips.cc/paper/2015/file/86df7dcfd896fcacf2674f757a2463eba-Paper.pdf></code></small><br>/><small><code><strong>[17]</strong>&#160;<https://ieeexplore.ieee.org/abstract/document/8920538></code></small><br>/><small><code><strong>[18]</strong>&#160;<https://www.aaai.org/Papers/Workshops/2006/WS-06-06/WS06-06-003.pdf></code></small><br>/><small><code><strong>[19]</strong>&#160;<https://ieeexplore.ieee.org/document/9000651></code></small><br>/><small><code><strong>[20]</strong>&#160;[https://people.mpi-sws.org/~gummadi/papers/process\\_fairness.pdf](https://people.mpi-sws.org/~gummadi/papers/process_fairness.pdf)</code></small><br>/><small><code><strong>[21]</strong>&#160;[https://papers.ssrn.com/sol3/papers.cfm?abstract\\_id=2477899](https://papers.ssrn.com/sol3/papers.cfm?abstract_id=2477899)</code></small><br>/><small><code><strong>[22]</strong>&#160;<https://www.ibm.com/blogs/research/2019/01/diversity-in-faces/></code></small><br>/><small><code><strong>[23]</strong>&#160;<https://arxiv.org/pdf/1906.10742.pdf></code></small><br>/><small><code><strong>[24]</strong>&#160;[https://people.mpi-sws.org/~gummadi/papers/process\\_fairness.pdf](https://people.mpi-sws.org/~gummadi/papers/process_fairness.pdf)</code></small><br>/><small><code><strong>[25]</strong>&#160;[https://people.ece.ubc.ca/mjulia/publications/Fairness\\_Definitions\\_Explained\\_2018.pdf](https://people.ece.ubc.ca/mjulia/publications/Fairness_Definitions_Explained_2018.pdf)</code></small><br>/><small><code><strong>[26]</strong>&#160;<https://github.com/tensorflow/fairness-indicators></code></small><br>/><small><code><strong>[27]</strong>&#160;<https://papers.nips.cc/paper/2017/file/a486cd07e4ac3d270571622f4f316ec5-Paper.pdf></code></small><br>/><small><code><strong>[28]</strong>&#160;<https://christophm.github.io/interpretable-ml-book/adversarial.html></code></small><br>/><small><code><strong>[29]</strong>&#160;<https://stereoset.mit.edu/></code></small><br>/><small><code><strong>[30]</strong>&#160;[https://services.google.com/fh/files/misc/practitioners\\_guide\\_to\\_mlops\\_whitepaper.pdf](https://services.google.com/fh/files/misc/practitioners_guide_to_mlops_whitepaper.pdf)</code></small><br>/><small><code><strong>[31]</strong>&#160;<https://arxiv.org/pdf/1810.03993.pdf></code></small><br>/><small><code><strong>[32]</strong>&#160;<https://arxiv.org/pdf/1803.09010.pdf></code></small><br>/><small><code><strong>[33]</strong>&#160;<https://www.oreilly.com/library/view/the-framework-for/9781098100483/ch01.html></code></small><br>/><small><code><strong>[34]</strong>&#160;<https://hpi.de/></code></small><br>/><small><code><strong>[35]</strong>&#160;<https://www.innoq.com/de/></code></small><br>/><small><code><strong>[36]</strong>&#160;<mailto:sih@ix.de></code></small><br /></p><p>class=„printversioncopyright“><em>Copyright &#169; 2022 Heise Medien</em></p> </html>

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

Permanent link:  
[https://schnipsl.qgelm.de/doku.php?id=wallabag:wb2fairness-und-knstliche-intelligenz\\_-warum-metriken-nicht-ausreichen](https://schnipsl.qgelm.de/doku.php?id=wallabag:wb2fairness-und-knstliche-intelligenz_-warum-metriken-nicht-ausreichen)

Last update: 2025/06/27 11:17

