

## Einleitung

### Was ist KI

KI (Kurz für Künstliche Intelligenz), beschreibt die Maschine, Menschliche Fähigkeiten zu imitieren. Dadurch ist es möglich das maschinen aufgaben wie z.b logisches denken, kreativität u.v.m nachahmen können.

(Was ist künstliche Intelligenz und wie wird sie genutzt? | Themen | Europäisches Parlament. (o. D.). Themen | Europäisches Parlament.

<https://www.europarl.europa.eu/topics/de/article/20200827STO85804/was-ist-kunstliche-intelligenz-und-wie-wird-sie-genutzt>)

Maschinen können durch Daten (Informationen/Inputs) probleme lösen, dies erfolgt meist signifikant effizienter als wenn diese Probleme von Menschen gelöst werden.

Besonders in der heutigen zeit von "big data", in welcher wir zugriff auf unengen an Datensätzen haben, spielt die KI eine essenzielle rolle in der auswertung und interpretation der großen Datensätze. Zahlreiche industrien nutzen KI bereits um ihre prozesse zu optimieren und Human Kraft zu reduzieren .

(SITNFlash. (2020, 23. April). The History of Artificial Intelligence - Science in the News. Science in The News.

<https://sitn.hms.harvard.edu/flash/2017/history-artificial-intelligence/>)

KI spielt schon heute eine wichtige rolle in unserem leben und ist schon fast ein omnipräsenter faktor im alltag. Es elerichtert uns Menschen und Firmen viel zeit und kraft und bewäätigt aufgaben die sonst als unmöglich erscheinen,

Doch trotz der zahlreichen vorteile, birgt KI auch einige signifikante probleme mit sich, eines der relvantesten ist der Bias in KI systemen.

Bias oder zu deutsch voreingenommenheit bedeutet im allgemeinen, dass das urteil einer Person durch vorurteile nicht objektive ist.

Bias in KI-Systemen beschreibt den fall, bei den Menschliche voreingenommenheit oder Stereotypen (Bias) unbewusst in die planung, entwicklung und nutzung einfließen und dadurch ergebnisse liefern die nicht immer objektive sind. Ein solcher "Bias" kann oftmals in trainingsdaten, dem algorithmus selber und den gegebenen antworten gefunden werden.

Team, Data. & Team, AI IBM (2023, 16. Oktober). *Shedding light on AI bias with real world examples.*

IBM Blog. <https://www.ibm.com/blog/shedding-light-on-ai-bias-with-real-world-examples/>

Trainingsdatensatz-Bias kann beispielsweise durch eine Über- oder Unter Repräsentativität von verschiedenen Personen kommen. Durch ein Ungleichgewicht der Trainingsdaten, im Hinblick auf die Vielfalt (im Sinne von Repräsentierung von verschiedenen Personengruppen) der genutzten Daten, kann dieses Ungleichgewicht in der Repräsentierung von verschiedenen Personengruppen Bias in dem Datensatz entstehen lassen, womit die KI trainiert wird (vgl. IBM, 2023).

Algorithmischer Bias kann durch fehlerhafte Datensätze entstehen, wodurch dann wiederholte Fehler entstehen können und gegebenenfalls der Bias in den Trainingsdaten sogar verstärkt werden kann. Dieser Bias kann auch durch den Programmierungsprozess entstehen. Dies kann eintreten, wenn Entwickler durch eigene Voreingenommenheit bestimmte Aspekte (meist) unbeabsichtigt anders gewichten. Dadurch kann es vorkommen, dass durch die subjektiven Ansichten der Entwickler diese Ansichten sich auch in der Entwicklung widerspiegeln und Ergebnisse hervorbringen, die nicht objektiv sind (vgl. IBM, 2023).

Kognitiver Bias, überschneidet teilweise mit dem algorithmischen Bias. Durch unsere eigenen Erfahrungen und Einstellungen sind wir von Natur aus voreingenommen. Dies kann sich durch die Auswahl von Datensätzen zeigen, welche ausgewählt werden um ein KI-System zu trainieren (vgl. IBM, 2023). Das Nationale Institut für Standards und Technologie aus den Vereinigten Staaten von Amerika gibt an, dass diese Form von Bias häufig auftritt und sich durch einen Mangel von Verständnis und Informationen für verschiedene Bevölkerungsgruppen ergibt. (*There's More to AI Bias Than Biased Data, NIST Report Highlights* | NIST. (2022, 16. März). NIST.

<https://www.nist.gov/news-events/news/2022/03/theres-more-ai-bias-biased-data-nist-report-highlights>)

KI birgt zahlreiche Vorteile, aber auch die Nachteile von unausgefilterter/vorgezogener KI sind gravierend und können in vielen Fällen das Leben von verschiedenen Personen ins Negative verändern, ohne dass diese etwas dafür können.

Ein Bereich in welcher Bias KI Personen benachteiligt hat ist der Recruiting Bereich. Amazon entwickelte ein Recruiting Tool, welches dabei helfen sollte den Bewerbungsprozess bei Amazon zu verbessern und effizienter zu gestalten. Dieses Programm sollte die Lebensläufe der Bewerber scannen und einer Vorauswahl treffen, welche dann weitergegeben wurde für die finale Entscheidung. Die Problematik entstand hierbei durch den genutzten Datensatz, denn dieser basierte auf einem historischen Datensatz welcher mehr Männer beinhaltete als Frauen. Durch diesen Datensatz trainierte sich der Algorithmus auf eine Art und Weise in welcher Männer in der Vorauswahl bevorzugt wurden und Schlüsselwörter die auf Frauen zurückzuführen sind wurden schlechter gewertet. Durch diese Ungleichheit im Datensatz, hatten weibliche Bewerberinnen weniger Chancen vorgeschlagen zu werden als Männer. (Walmsley, J. (2020). Artificial intelligence and the value of transparency. *AI & SOCIETY*, 36(2), 585–595. <https://doi.org/10.1007/s00146-020-01066-z>)

In einer Studie von Obermeyer et al., wurde ein Branchentypischer Algorithmus, aus den USA, untersucht, welcher in im Gesundheitswesen genutzt wird um vorherzusagen für Patienten mit komplexen Krankheitsbildern besser zu behandeln. Hierbei wird ein Risikowert ermittelt welcher genutzt wird um weitere Behandlungsschritte zu determinieren. Bei der Untersuchung des Algorithmus wurde jedoch festgestellt, dass Schwarze Patienten bei einem bestimmten Risikowert, wesentlich Kränker sind als weiße Patienten. Durch diesen Bias im Algorithmus kann es also vorkommen, dass Schwarze Patienten nicht die nötige Hilfe bekommen, welche sie eigentlich benötigen. Dieser Bias entstand nach Obermeyer et al. dadurch, dass der Algorithmus die Behandlungskosten vorhersagt und keine Krankheitsbilder, durch ungleichen Zugang zu für medizinische Behandlung für Schwarze Patienten, im Vergleich zu weißen

Patienten, es wird also weniger Geld für die Behandlung von schwarzen Patienten ausgegeben. Bei dieser KI wurde Behandlungskosten als Indikator für die Vorhersagen genutzt, jedoch durch die Ungleichheit in den Behandlungskosten wurden dadurch verzerrte Risikowerte ausgegeben (Obermeyer et al. (2019))

Obermeyer, Z., Powers, B., Vogeli, C. & Mullainathan, S. (2019). Dissecting racial bias in an algorithm used to manage the health of populations. *Science*, 366(6464), 447–453.  
<https://doi.org/10.1126/science.aax2342>

Es gibt noch weitere Beispiele für die Auswirkungen für Bias in KI-System, welche das Leben der betroffenen Personen beeinflussen.

Aufgrund dieser Problematik befassen wir uns mit der Thematik, welche Herausforderung bei der Entwicklung von verantwortungsvoller KI entstehen. Durch unsere Erkenntnisse wollen wir verstehen, welche Problematiken bei der Entwicklung von solchen Systemen aufkommen, sodass man diese frühzeitig erkennt und entsprechend handhabt.

Hierfür analysierten wir die Erkenntnisse aus zahlreichen wissenschaftlichen Arbeiten, welche von uns im Hinblick auf unsere Fragestellung analysiert wurden und interpretieren diese Ergebnisse auch und versuchen dadurch unsere Fragestellung zu beantworten und die allgemeine Tätigkeit besser zu verstehen.