

Fairness bedeutet, dass sich Fehler des KI-Systems auf beide Gruppen in gleicher Weise auswirken. Das beruht immer noch auf einer ungerütteten „Ground Truth“. Selbst wenn die Fehlerrquoten gleich sind, könnten die unterschiedlich sehr unterschiedlich sein – und ungerüttet!



Was ist mit Qualität fizzierten, aber weniger selbstbewussten Gruppen? Es gäbe weniger Bewerbungen von ihnen, aber der Anteil der Qualifizierteren könnte höher sein, und die würden abgelehnt ...

Damit können nur ähnliche Individuen verglichen werden, nicht alle. Damit kann nur ähnliche Individualität unterscheiden, was überhaupt entscheidend ist. Sind Entschiedungen Unkenntnis ist unmöglich. Vorherzusagen. Völligkeitsgraden, sensible Merkmale keiten, sensible Merkmale bei der Entschiedungsfindung nicht explizit berücksichtigt werden.



Ahnlichkeit bedeutet, was überhaupt fair ist. Was darf darüber entscheiden. Wer darf überprüfen, ob es Mögliche eingeoordnet werden. Von ihrer Gruppe ähnlich ahnliche Personen unabhängig Fairness bedeutet, dass

TREATMENT EQUALITY

DEMOGRAPHIC PARITY
DER RÜCKSEITEN! Deep dive auf

FAIRNESS-METRIKEN SIND NICHT GENUG!

Fühlt sich das fair an? Jede der Metriken mag fair klingen, aber es ist nicht einfach, Fairness zu erreichen. Die Metriken liefern Informationen für Maßnahmen. Doch sie alle beruhen auf binären Unterscheidungen und widersprechen sich teilweise! Kritik an den Metriken findest du auf der Posterseite des nächsten Zines. Was kann Techie sonst noch tun, um faire(re) KI zu entwickeln? Erforsche fairere Prozesse auf den Seiten des nächsten Zines!

WIE WIRD FAIRNESS GEMESSEN?

Unter Fairness versteht man die Gleichbehandlung verschiedener Gruppen. Doch wie lässt sich dieses fluide Konzept in KI-Kontexten umsetzen? Auf den folgenden Seiten erfahrt ihr mehr über Fairness-Metriken. Das sind Ansätze, um Fairness messbar zu machen.



FAIRNESS DURCH UNAWARENESS

SENSIBLE MERKMALE

Wenn es bei Fairness um die Gleichbehandlung verschiedener Gruppen geht – wer sind dann diese Gruppen? Bei der Definition davon werden **sensible Merkmale**, binäre Unterteilungen auf der Grundlage von Privilegien, hinzugezogen:



Privilegiert



z. B. Menschen, die als nicht-männlich, nicht-weiß, queer gelesen werden

Nicht-privilegiert

Mit Hilfe der Testdaten kann herausgefunden werden, ob sich das KI-Modell für verschiedene Gruppen unterschiedlich verhält. Dafür werden **Fairness-Metriken** verwendet. Davon gibt es viele verschiedene und jede einzelne ist eine spezifische 1 2 Definition von Fairness. Auf den folgenden Seiten gibt es einige Beispiele.

FEE LS FAIR?

FAIRNESS MESSSEN?

Fairness-Metriken basieren auf dem Vergleich von Prognosen für Einzelpersonen oder Gruppen. Aber wie werden sie verglichen? Es gibt zwei Ansätze:

Individuenbezogene Fairness-Metriken ...

Gruppenbezogene Fairness-Metriken ...

... beruhen auf dem Vergleich der Vorhersagen von Einzelpersonen, z.B. die Metrik Fairness durch awareness oder unawareness (auf den Zine-Seiten). Hier wird z.B. die Vorhersage für ein Individuum mit den Ergebnissen für ein ähnliches oder anderes Individuum verglichen.

... beruhen auf dem Vergleich der Vorhersagen des KI-Modells mit den gelabelten Daten für verschiedene Gruppen. Siehe z.B. Demographic Parity oder Treatment Equality (auf den Zine-Seiten). Der Vergleich der Vorhersagen pro Gruppe basiert auf einer Tabelle, die einen Überblick darüber gibt, wo und für welche Gruppe das KI-System richtig vorhergesagt hat:

		KI prognostiziert: ablehnen	KI prognostiziert: einstellen
Label von C.O.R.P.: ablehnen	WÄHR POSITIV für Gruppe A UND Gruppe B	FÄLSCH NEGATIV für Gruppe A UND Gruppe B	
	FÄLSCH POSITIV für Gruppe A UND Gruppe B	WÄHR NEGATIV für Gruppe A UND Gruppe B	

Jede Fairness-Metrik ist eine bestimmte Definition von Fairness. Mathematisch ausgedrückt, könnte so eine Definition so aussehen:

Das KI-System ist fair, wenn

$$\text{Gruppe A: } \frac{\text{Alle Positiven}}{\text{Alle}} \approx \text{Gruppe B: } \frac{\text{Alle Positiven}}{\text{Alle}}$$

Das kann auch so ausgedrückt werden:

Das Verhältnis zwischen den voraussichtlichen Einstellungen und der Gesamtzahl der Bewerbungen ist in allen Gruppen ähnlich. 40 % der Bewerber*innen sind weiblich – 40 % der eingestellten Mitarbeiter*innen sollten weiblich sein!
Diese Metrik definiert Fairness als die gleiche Wahrscheinlichkeit für beide Gruppen, eingestellt zu werden. Das nennt man demografische Parität – und diese Definition hast du bereits auf den Zine-Seiten kennengelernt!

MERO NOTIZEN:

Fairness zu finden – diese Unschärfe ist aber auch eine Stärke von Fairness. Wir haben hier eine Metrik im Detail dargestellt, um zu erklären, wie die Metriken funktionieren. Wenn man andere Metriken genauer vergleicht, wird deutlich, dass sie sich manchmal widersprechen. Oder die Ergebnisse einer Metrik

würden ein KI-System als fair beschreiben, während eine andere Metrik das nicht tun würde.
Die Fairness-Metriken sind eine Auswahl der in Mehrabi et al. (2022) vorgestellten Fairness-Metriken: A survey on bias and fairness in Machine Learning; und Verma & Rubin (2018): Fairness definitions explained.