

Diskriminierung durch (lernende) Algorithmen

Meike Zehlike

30. September 2021

Humboldt Universität zu Berlin, Zalando

1. Einführung Maschinelles Lernen
2. Bias und Fairness
3. Diskriminierende KI
4. Die Daten sind das Problem
5. Die Algorithmen sind das Problem

Einführung Maschinelles Lernen

Traditional Programming



Machine Learning

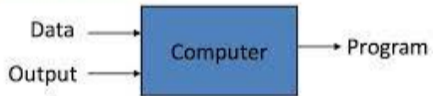


Abbildung 1: Traditionelles Vorgehen vs Machine Learning

MACHINE LEARNING KURZ ERKLÄRT

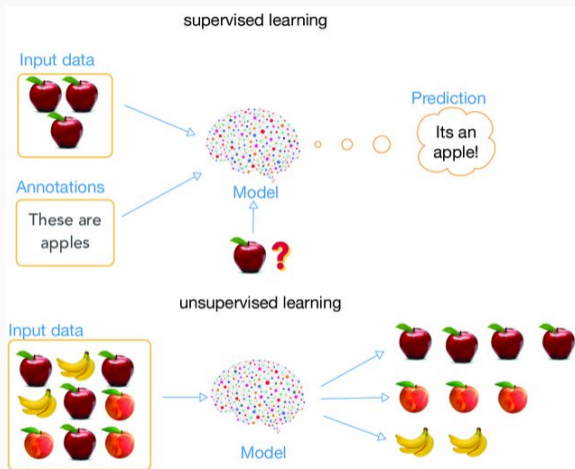


Abbildung 2: Supervised vs Unsupervised Learning

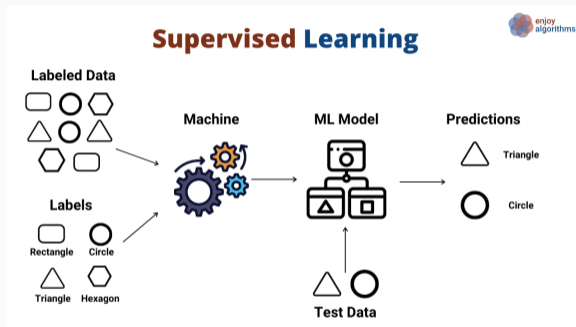


Abbildung 3: Supervised Machine Learning Pipeline

- Um das beste Modell von allen möglichen Modellen zu finden wird ein Optimierungsproblem gelöst:

$$\min \sum_i (\hat{y}_i - y_i)^2$$

- Das beste Modell ist also dasjenige, das die wenigsten Fehler macht

Bias und Fairness

Ein Computer System ist biased, wenn:

Ein Computer System ist biased, wenn:

es bestimmte *Gruppen* oder *Individuen* *systematisch* und *auf unfaire Art und Weise* zum Vorteil anderer diskriminiert.

Ein Computer System ist biased, wenn:

es bestimmte *Gruppen* oder *Individuen* *systematisch* und *auf unfaire Art und Weise* zum Vorteil anderer diskriminiert.

Ein System diskriminiert auf unfaire Art, wenn:

Ein Computer System ist biased, wenn:

es bestimmte *Gruppen* oder *Individuen* *systematisch* und *auf unfaire Art und Weise* zum Vorteil anderer diskriminiert.

Ein System diskriminiert auf unfaire Art, wenn:

es einem *Individuum* oder einer *Gruppe bestimmter Individuen*

Ein Computer System ist biased, wenn:

es bestimmte *Gruppen* oder *Individuen* *systematisch* und *auf unfaire Art und Weise* zum Vorteil anderer diskriminiert.

Ein System diskriminiert auf unfaire Art, wenn:

es einem *Individuum* oder einer *Gruppe bestimmter Individuen* aufgrund *unvernünftiger* oder *unangemessener* Kriterien

Ein Computer System ist biased, wenn:

es bestimmte *Gruppen* oder *Individuen* *systematisch* und *auf unfaire Art und Weise* zum Vorteil anderer diskriminiert.

Ein System diskriminiert auf unfaire Art, wenn:

es einem *Individuum* oder einer *Gruppe bestimmter Individuen* aufgrund *unvernünftiger* oder *unangemessener* Kriterien eine *Chance* oder ein *Gut verweigert*, oder wenn

Ein Computer System ist biased, wenn:

es bestimmte *Gruppen* oder *Individuen* *systematisch* und *auf unfaire Art und Weise* zum Vorteil anderer diskriminiert.

Ein System diskriminiert auf unfaire Art, wenn:

es einem *Individuum* oder einer *Gruppe bestimmter Individuen* aufgrund *unvernünftiger* oder *unangemessener* Kriterien eine *Chance* oder ein *Gut verweigert*, oder wenn es denselben ein unerwünschtes Ergebnis zuweist.

Bias ist unfaire Diskriminierung, die *systematisch* auftritt.

Bias ist unfaire Diskriminierung, die *systematisch* auftritt.

Bias ist systematische Diskriminierung, die zusammen mit einem *unfairen Ergebnis* auftritt.

- **Preexisting Bias:** Bias der unabhängig und vor der Schaffung des Systems existiert.

- **Preexisting Bias:** Bias der unabhängig und vor der Schaffung des Systems existiert.
- **Technical Bias:** technische Limitierungen und Zwänge (Hardware und Software Limitierungen, Generation von Pseudozufallszahlen, Formalisierung menschlicher Konstrukte)

- **Preexisting Bias:** Bias der unabhängig und vor der Schaffung des Systems existiert.
- **Technical Bias:** technische Limitierungen und Zwänge (Hardware und Software Limitierungen, Generation von Pseudozufallszahlen, Formalisierung menschlicher Konstrukte)
- **Emergent Bias:** tritt erst nach der Einführung des Systems auf (Veränderung im Nutzerverhalten, Veränderung sozialer Konzepte)

- ein Framework aus der politischen Philosophie

- ein Framework aus der politischen Philosophie
- vier verschiedene Geschmacksrichtungen [1]:

- ein Framework aus der politischen Philosophie
- vier verschiedene Geschmacksrichtungen [1]: libertäre EO,

- ein Framework aus der politischen Philosophie
- vier verschiedene Geschmacksrichtungen [1]: libertäre EO, formelle EO,

- ein Framework aus der politischen Philosophie
- vier verschiedene Geschmacksrichtungen [1]: libertäre EO, formelle EO, substantielle EO

- ein Framework aus der politischen Philosophie
- vier verschiedene Geschmacksrichtungen [1]: libertäre EO, formelle EO, substantielle EO (Rawlsian EO, luck-egalitarian EO)

- ein Framework aus der politischen Philosophie
- vier verschiedene Geschmacksrichtungen [1]: libertäre EO, formelle EO, substantielle EO (Rawlsian EO, **luck-egalitarian EO**)

- Eine jede hat die Freiheit zu tun und zu lassen was sie möchte mit all dem, was ihr legitimerweise gehört (Selbst, Geschäft, etc.),

- Eine jede hat die Freiheit zu tun und zu lassen was sie möchte mit all dem, was ihr legitimerweise gehört (Selbst, Geschäft, etc.), solange dies nicht moralische Rechte anderer Menschen verletzt.

- Eine jede hat die Freiheit zu tun und zu lassen was sie möchte mit all dem, was ihr legitimerweise gehört (Selbst, Geschäft, etc.), solange dies nicht moralische Rechte anderer Menschen verletzt.
- Im Kontext algorithmischer Entscheidungssysteme (ADM) bedeutet dies totale Freiheit jeden beliebigen Algorithmus zu implementieren.

- Eine jede hat die Freiheit zu tun und zu lassen was sie möchte mit all dem, was ihr legitimerweise gehört (Selbst, Geschäft, etc.), solange dies nicht moralische Rechte anderer Menschen verletzt.
- Im Kontext algorithmischer Entscheidungssysteme (ADM) bedeutet dies totale Freiheit jeden beliebigen Algorithmus zu implementieren.
- Dieser kann jede verfügbare Information benutzen, auch solche die eigentlich gesetzlich geschützt sind (Geschlecht, Migrationshintergrund), um statistisch akkuratere Entscheidungen zu treffen.

- Erstrebenswerte Positionen müssen allen zugänglich sein, die für die Erfüllung der Position *notwendige* Eigenschaften mitbringen.

- Erstrebenswerte Positionen müssen allen zugänglich sein, die für die Erfüllung der Position *notwendige* Eigenschaften mitbringen.
- Die Positionen werden aufgrund der Qualifikation der Kandidatinnen zugewiesen.

- Erstrebenswerte Positionen müssen allen zugänglich sein, die für die Erfüllung der Position *notwendige* Eigenschaften mitbringen.
- Die Positionen werden aufgrund der Qualifikation der Kandidatinnen zugewiesen.
- Im Kontext von ADM fordert formelle Chancengleichheit sensible Informationen aus dem System zu entfernen.

SUBSTANTIELLE CHANCENGLEICHHEIT [3]

- Erstrebenswerte Positionen müssen allen zugänglich sein, die für die Erfüllung der Position **notwendige** Eigenschaften mitbringen.

SUBSTANTIELLE CHANCENGLEICHHEIT [3]

- Erstrebenswerte Positionen müssen allen zugänglich sein, die für die Erfüllung der Position **notwendige** Eigenschaften mitbringen.
- Die Positionen werden aufgrund der Qualifikation der Kandidatinnen zugewiesen.

SUBSTANTIELLE CHANCENGLEICHHEIT [3]

- Erstrebenswerte Positionen müssen allen zugänglich sein, die für die Erfüllung der Position **notwendige** Eigenschaften mitbringen.
- Die Positionen werden aufgrund der Qualifikation der Kandidatinnen zugewiesen.
- zusätzlich braucht es **fairen Zugang** zu solchen Qualifikationen

SUBSTANTIELLE CHANCENGLEICHHEIT [3]

- Erstrebenswerte Positionen müssen allen zugänglich sein, die für die Erfüllung der Position **notwendige** Eigenschaften mitbringen.
- Die Positionen werden aufgrund der Qualifikation der Kandidatinnen zugewiesen.
- zusätzlich braucht es **fairen Zugang** zu solchen Qualifikationen
- **Rawls'sche EO:**

SUBSTANTIELLE CHANCENGLEICHHEIT [3]

- Erstrebenswerte Positionen müssen allen zugänglich sein, die für die Erfüllung der Position **notwendige** Eigenschaften mitbringen.
- Die Positionen werden aufgrund der Qualifikation der Kandidatinnen zugewiesen.
- zusätzlich braucht es **fairen Zugang** zu solchen Qualifikationen
- **Rawls'sche EO:**
 - Jene mit demselben Talent und Bestreben müssen die gleichen Aussichten haben die sozialen Positionen zu erlangen,

SUBSTANTIELLE CHANCENGLEICHHEIT [3]

- Erstrebenswerte Positionen müssen allen zugänglich sein, die für die Erfüllung der Position **notwendige** Eigenschaften mitbringen.
- Die Positionen werden aufgrund der Qualifikation der Kandidatinnen zugewiesen.
- zusätzlich braucht es **fairen Zugang** zu solchen Qualifikationen
- **Rawls'sche EO:**
 - Jene mit demselben Talent und Bestreben müssen die gleichen Aussichten haben die sozialen Positionen zu erlangen, unabhängig von zufälligen Faktoren wie sozio-ökonomischer Status der Eltern.

SUBSTANTIELLE CHANCENGLEICHHEIT [3]

- Erstrebenswerte Positionen müssen allen zugänglich sein, die für die Erfüllung der Position **notwendige** Eigenschaften mitbringen.
- Die Positionen werden aufgrund der Qualifikation der Kandidatinnen zugewiesen.
- zusätzlich braucht es **fairen Zugang** zu solchen Qualifikationen
- **Rawls'sche EO:**
 - Jene mit demselben Talent und Bestreben müssen die gleichen Aussichten haben die sozialen Positionen zu erlangen, unabhängig von zufälligen Faktoren wie sozio-ökonomischer Status der Eltern.
 - Nimmt an, dass Talent und Bestreben sich über alle Individuen hinweg vergleichen lässt.

SUBSTANTIELLE CHANCENGLEICHHEIT [3]

- Erstrebenswerte Positionen müssen allen zugänglich sein, die für die Erfüllung der Position **notwendige** Eigenschaften mitbringen.
- Die Positionen werden aufgrund der Qualifikation der Kandidatinnen zugewiesen.
- zusätzlich braucht es **fairen Zugang** zu solchen Qualifikationen
- **Rawls'sche EO:**
 - Jene mit demselben Talent und Bestreben müssen die gleichen Aussichten haben die sozialen Positionen zu erlangen, unabhängig von zufälligen Faktoren wie sozio-ökonomischer Status der Eltern.
 - Nimmt an, dass Talent und Bestreben sich über alle Individuen hinweg vergleichen lässt.
- **Luck-egalitarian EO:**

SUBSTANTIELLE CHANCENGLEICHHEIT [3]

- Erstrebenswerte Positionen müssen allen zugänglich sein, die für die Erfüllung der Position **notwendige** Eigenschaften mitbringen.
- Die Positionen werden aufgrund der Qualifikation der Kandidatinnen zugewiesen.
- zusätzlich braucht es **fairen Zugang** zu solchen Qualifikationen
- **Rawls'sche EO:**
 - Jene mit demselben Talent und Bestreben müssen die gleichen Aussichten haben die sozialen Positionen zu erlangen, unabhängig von zufälligen Faktoren wie sozio-ökonomischer Status der Eltern.
 - Nimmt an, dass Talent und Bestreben sich über alle Individuen hinweg vergleichen lässt.
- **Luck-egalitarian EO:**
 - Bietet eine relative Sichtweise auf Talent und Bestreben und erlaubt, dass diese vom sozio-ökonomischen Hintergrund beeinflusst werden.

SUBSTANTIELLE CHANCENGLEICHHEIT [3]

- Erstrebenswerte Positionen müssen allen zugänglich sein, die für die Erfüllung der Position **notwendige** Eigenschaften mitbringen.
- Die Positionen werden aufgrund der Qualifikation der Kandidatinnen zugewiesen.
- zusätzlich braucht es **fairen Zugang** zu solchen Qualifikationen
- **Rawls'sche EO:**
 - Jene mit demselben Talent und Bestreben müssen die gleichen Aussichten haben die sozialen Positionen zu erlangen, unabhängig von zufälligen Faktoren wie sozio-ökonomischer Status der Eltern.
 - Nimmt an, dass Talent und Bestreben sich über alle Individuen hinweg vergleichen lässt.
- **Luck-egalitarian EO:**
 - Bietet eine relative Sichtweise auf Talent und Bestreben und erlaubt, dass diese vom sozio-ökonomischen Hintergrund beeinflusst werden.
 - Nimmt an, dass Talent und Bestreben nur innerhalb einer Gruppe verglichen werden darf.

Diskriminierende KI

- Apples Credit Scoring Algorithmus bewertet Frauen schlechter

BEISPIELE

- Apples Credit Scoring Algorithmus bewertet Frauen schlechter
- Google Ad Service suggeriert Schwarze seien kriminell

- Apples Credit Scoring Algorithmus bewertet Frauen schlechter
- Google Ad Service suggeriert Schwarze seien kriminell
- Google Bilderkennung erkennt auf Bildern mit schwarzen Menschen Gorillas

- Apples Credit Scoring Algorithmus bewertet Frauen schlechter
- Google Ad Service suggeriert Schwarze seien kriminell
- Google Bilderkennung erkennt auf Bildern mit schwarzen Menschen Gorillas
- Siri und Alexa reagieren bei sexueller Belästigung unangemessen

- Apples Credit Scoring Algorithmus bewertet Frauen schlechter
- Google Ad Service suggeriert Schwarze seien kriminell
- Google Bilderkennung erkennt auf Bildern mit schwarzen Menschen Gorillas
- Siri und Alexa reagieren bei sexueller Belästigung unangemessen
- Tay (Windows Chatbot) wird zum Rassisten

- Apples Credit Scoring Algorithmus bewertet Frauen schlechter
- Google Ad Service suggeriert Schwarze seien kriminell
- Google Bilderkennung erkennt auf Bildern mit schwarzen Menschen Gorillas
- Siri und Alexa reagieren bei sexueller Belästigung unangemessen
- Tay (Windows Chatbot) wird zum Rassisten
- Amazon's Hiring Tool sortiert Lebensläufe von Frauen aus

- Apples Credit Scoring Algorithmus bewertet Frauen schlechter
- Google Ad Service suggeriert Schwarze seien kriminell
- Google Bilderkennung erkennt auf Bildern mit schwarzen Menschen Gorillas
- Siri und Alexa reagieren bei sexueller Belästigung unangemessen
- Tay (Windows Chatbot) wird zum Rassisten
- Amazon's Hiring Tool sortiert Lebensläufe von Frauen aus
- Google Ads zeigt Werbung für gutbezahlte Jobs überwiegend Männern

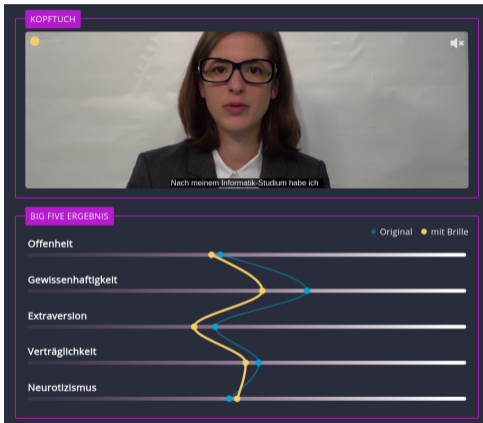
- Apples Credit Scoring Algorithmus bewertet Frauen schlechter
- Google Ad Service suggeriert Schwarze seien kriminell
- Google Bilderkennung erkennt auf Bildern mit schwarzen Menschen Gorillas
- Siri und Alexa reagieren bei sexueller Belästigung unangemessen
- Tay (Windows Chatbot) wird zum Rassisten
- Amazon's Hiring Tool sortiert Lebensläufe von Frauen aus
- Google Ads zeigt Werbung für gutbezahlte Jobs überwiegend Männern
- Gesichtserkennung zum Entsperren von iPhones funktioniert für schwarze Menschen wesentlich schlechter als für Männer

- Apples Credit Scoring Algorithmus bewertet Frauen schlechter
- Google Ad Service suggeriert Schwarze seien kriminell
- Google Bilderkennung erkennt auf Bildern mit schwarzen Menschen Gorillas
- Siri und Alexa reagieren bei sexueller Belästigung unangemessen
- Tay (Windows Chatbot) wird zum Rassisten
- Amazon's Hiring Tool sortiert Lebensläufe von Frauen aus
- Google Ads zeigt Werbung für gutbezahlte Jobs überwiegend Männern
- Gesichtserkennung zum Entsperren von iPhones funktioniert für schwarze Menschen wesentlich schlechter als für Männer
- Ein Algorithmus zur Bettenbelegung auf einer Intensivstation diskriminiert gegen Schwarze

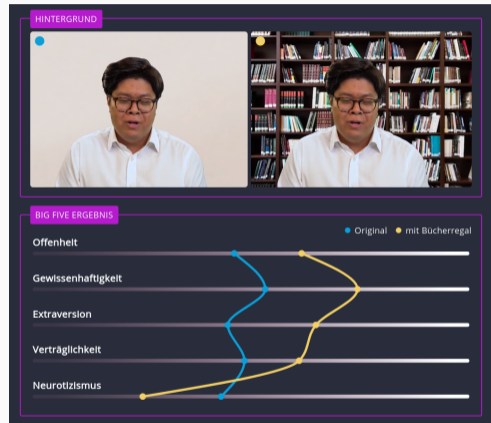
- Apples Credit Scoring Algorithmus bewertet Frauen schlechter
- Google Ad Service suggeriert Schwarze seien kriminell
- Google Bilderkennung erkennt auf Bildern mit schwarzen Menschen Gorillas
- Siri und Alexa reagieren bei sexueller Belästigung unangemessen
- Tay (Windows Chatbot) wird zum Rassisten
- Amazon's Hiring Tool sortiert Lebensläufe von Frauen aus
- Google Ads zeigt Werbung für gutbezahlte Jobs überwiegend Männern
- Gesichtserkennung zum Entsperren von iPhones funktioniert für schwarze Menschen wesentlich schlechter als für Männer
- Ein Algorithmus zur Bettenbelegung auf einer Intensivstation diskriminiert gegen Schwarze
- Facebook erkennt Profile von Native-Americans als Fake-Profile

- Apples Credit Scoring Algorithmus bewertet Frauen schlechter
- Google Ad Service suggeriert Schwarze seien kriminell
- Google Bilderkennung erkennt auf Bildern mit schwarzen Menschen Gorillas
- Siri und Alexa reagieren bei sexueller Belästigung unangemessen
- Tay (Windows Chatbot) wird zum Rassisten
- Amazon's Hiring Tool sortiert Lebensläufe von Frauen aus
- Google Ads zeigt Werbung für gutbezahlte Jobs überwiegend Männern
- Gesichtserkennung zum Entsperren von iPhones funktioniert für schwarze Menschen wesentlich schlechter als für Männer
- Ein Algorithmus zur Bettenbelegung auf einer Intensivstation diskriminiert gegen Schwarze
- Facebook erkennt Profile von Native-Americans als Fake-Profile
- ...

MEHR BEISPIELE



(a) Mit vs. Ohne Brille



(b) Mit Bücherregal

Abbildung 4: <https://web.br.de/interaktiv/ki-bewerbung/>

Die Daten sind das Problem

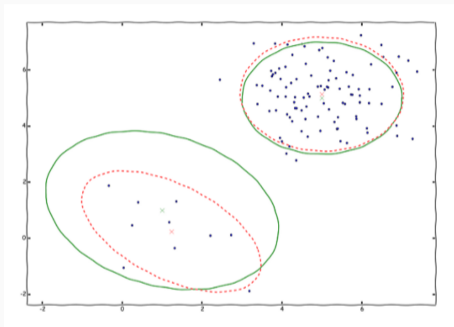


Abbildung 5: Modelle funktionieren mit weniger Trainingsdaten schlechter (grün – ground truth, rot – Modell)

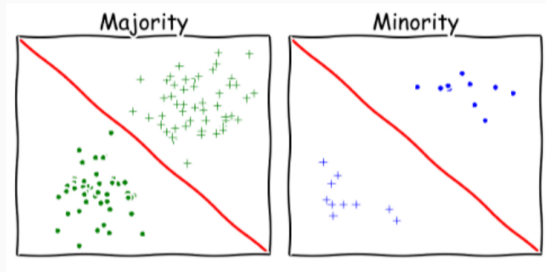


Abbildung 6: Innerhalb der Population kann es große statistische Unterschiede geben

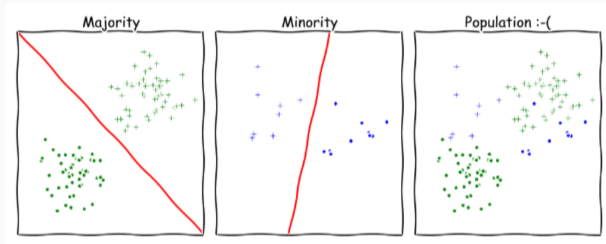


Abbildung 7: Was nun?

SIND TRAININGSDATEN NEUTRAL?



(a) Was ist hier zu sehen?

SIND TRAININGSDATEN NEUTRAL?



(a) Was ist hier zu sehen?



(b) Was ist hier zu sehen?

SIND TRAININGSDATEN NEUTRAL?



(a) Was ist hier zu sehen?



(b) Was ist hier zu sehen?

Abbildung 8: Beschreibung und Parameter werden als Norm und Abweichung von der Norm dargestellt

Beispiel Gender Pay Gap:

- Frauen verdienen 80€/h, Männer verdienen 100€/h

Beispiel Gender Pay Gap:

- Frauen verdienen 80€/h, Männer verdienen 100€/h
- “Noch immer verdienen Frauen 20% weniger als Männer”

Beispiel Gender Pay Gap:

- Frauen verdienen 80€/h, Männer verdienen 100€/h
- “Noch immer verdienen Frauen 20% weniger als Männer”
- Aus Sicht der Frauen: “Noch immer verdienen Frauen 25% weniger als Männer”

Was ist die Kategorie? Eigenschaft oder Wahl?

Was ist die Kategorie? Eigenschaft oder Wahl?

Konstruierte Kategorien? Biologisch oder sozial konstruiert?

Was ist die Kategorie? Eigenschaft oder Wahl?

Konstruierte Kategorien? Biologisch oder sozial konstruiert?

Modularität Modelle sind Systeme mit reduzierter Komplexität

Was ist die Kategorie? Eigenschaft oder Wahl?

Konstruierte Kategorien? Biologisch oder sozial konstruiert?

Modularität Modelle sind Systeme mit reduzierter Komplexität

Ontologische Instabilität Bedeutung eines Begriffs verändert sich über die Zeit

Was ist die Kategorie? Eigenschaft oder Wahl?

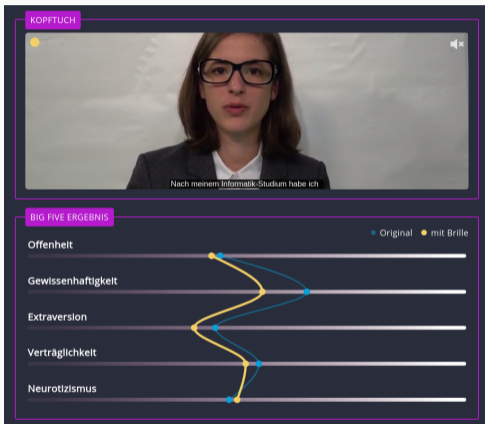
Konstruierte Kategorien? Biologisch oder sozial konstruiert?

Modularität Modelle sind Systeme mit reduzierter Komplexität

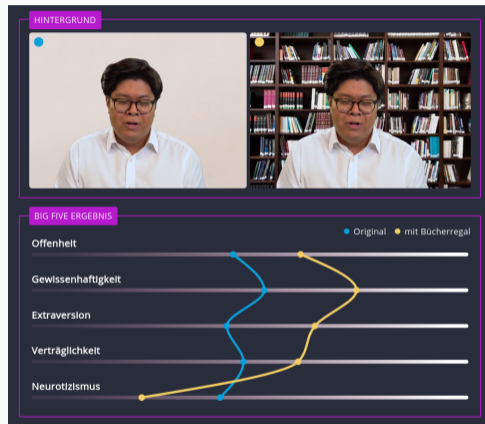
Ontologische Instabilität Bedeutung eines Begriffs verändert sich über die Zeit
Zuschreibungen führen zu veränderten Verhaltensweisen

Die Algorithmen sind das Problem

WELCHES PROBLEM WIRD ÜBERHAUPT GELÖST?



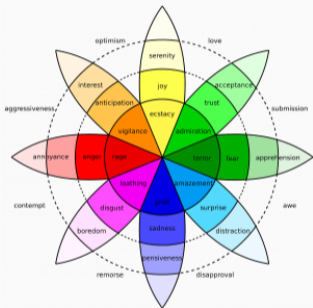
(a) Mit vs. Ohne Brille



(b) Mit Bücherregal

Abbildung 9: <https://web.br.de/interaktiv/ki-bewerbung/>

ABSURDE VEREINFACHUNG KOMPLEXER KONZEPTE



(a) Emotionsrad nach Plutchik



(b) Beispiel für Trainingsdaten

Utilitaristische Ethik:

Utilitaristische Ethik: $\min \sum_i (\hat{y}_i - y_i)^2$

Utilitaristische Ethik: $\min \sum_i (\hat{y}_i - y_i)^2$

Rawls'sche Fairness:

Utilitaristische Ethik: $\min \sum_i (\hat{y}_i - y_i)^2$

Rawls'sche Fairness: $\min \max_i (\hat{y}_i - y_i)^2$




WIE BAUT MAN EINEN FAIREN ALGORITHMUS?



WIE ÜBERSETZT MAN GERECHTIGKEIT IN MATHEMATIK?

$$\begin{aligned} \min \sum_i (\hat{y}_i - y_i)^2 \quad \text{subject to} \quad & P(\hat{y}_i \neq y_i | \text{protected} = 1) \\ & = P(\hat{y}_i \neq y_i | \text{protected} = 0) \end{aligned}$$

Fragen?

-  R. Arneson.
Equality of opportunity.
In E. N. Zalta, editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, summer 2015 edition, 2015.
-  B. Friedman and H. Nissenbaum.
Bias in computer systems.
ACM Transactions on Information Systems, 14(3):330–347, 1996.
-  H. Heidari, M. Loi, K. P. Gummadi, and A. Krause.
A moral framework for understanding fair ml through economic models of equality of opportunity.
In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pages 181–190. ACM, 2019.