

Fairness im Maschinellen Lernen



Dr. rer. nat. Sebastian Pape

Habilitationskolloquium

22. Februar 2021

Agenda

- Einführung
 - Motivation
 - Maschinelles Lernen
 - Einordnung
- Fairness
 - Bias
 - Diskriminierung
 - Definition Fairness
- Algorithmen
 - Klassifizierung
 - Word Embedding
- Herausforderungen
- Zusammenfassung



Aufgaben für maschinelles Lernen



onRagtime

Source: <https://www.oneragtime.com/24-industries-disrupted-by-ai-infographic/>

Only Seven of Stanford's First 5,000 Vaccines Were Designated for Medical Residents

Stanford Medicine officials relied on a faulty algorithm to determine who should get vaccinated first, and it prioritized some high-ranking doctors over patient-facing medical residents.

An algorithm chose who would be the first 5,000 in line. They were at a disadvantage because they did not have an advantage in the calculation and because they are young, according to a resident to his peers. Residents are the lowest-ranking doctors in a hospital. Stanford Medicine has about 1,300 across all disciplines.

Source: Caroline Chen, ProPublica
<https://www.propublica.org/article/only-seven-of-stanfords-first-5-000-vaccines-were-designated-for-medical-residents>



Mariya Gabriel ✅ @GabrielMariya · Jan 6

#EUfunded research shows that #AI developed without a #gender & intersectional lens can increase negative economic & social consequences.

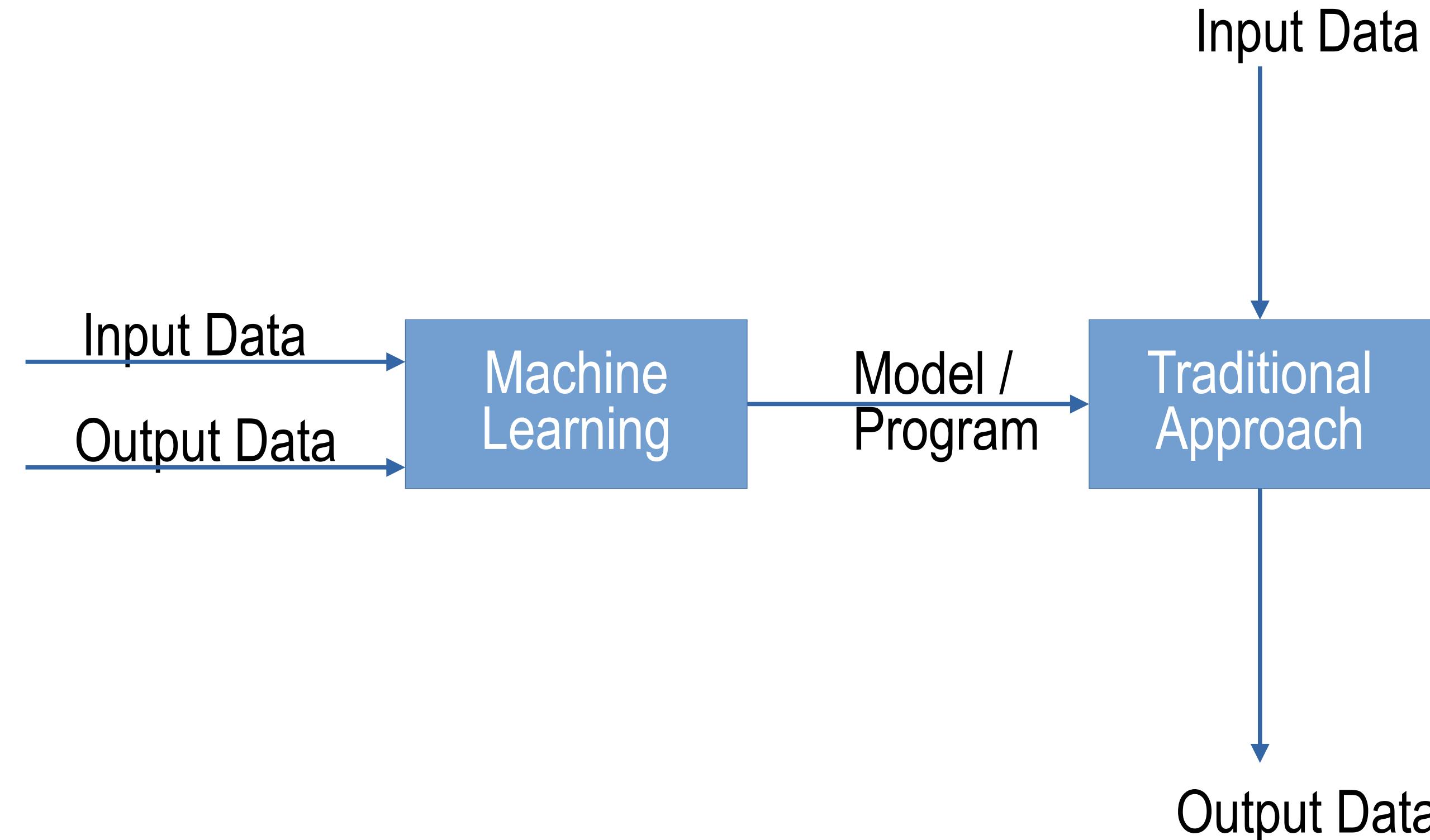
Discover how bias-free AI can successfully become a key driver for #Innovation

👉 europa.eu/!CF37Ud

#UnionOfEquality #InvestEUREsearch



Machine Learning vs. Traditional Approaches



Generelle Herausforderungen von Maschinellem Lernen

- Angriffe
- Datenschutz
- Große Datenmengen zum Training benötigt
- Erklärbarkeit / Interpretierbarkeit
- Ethik

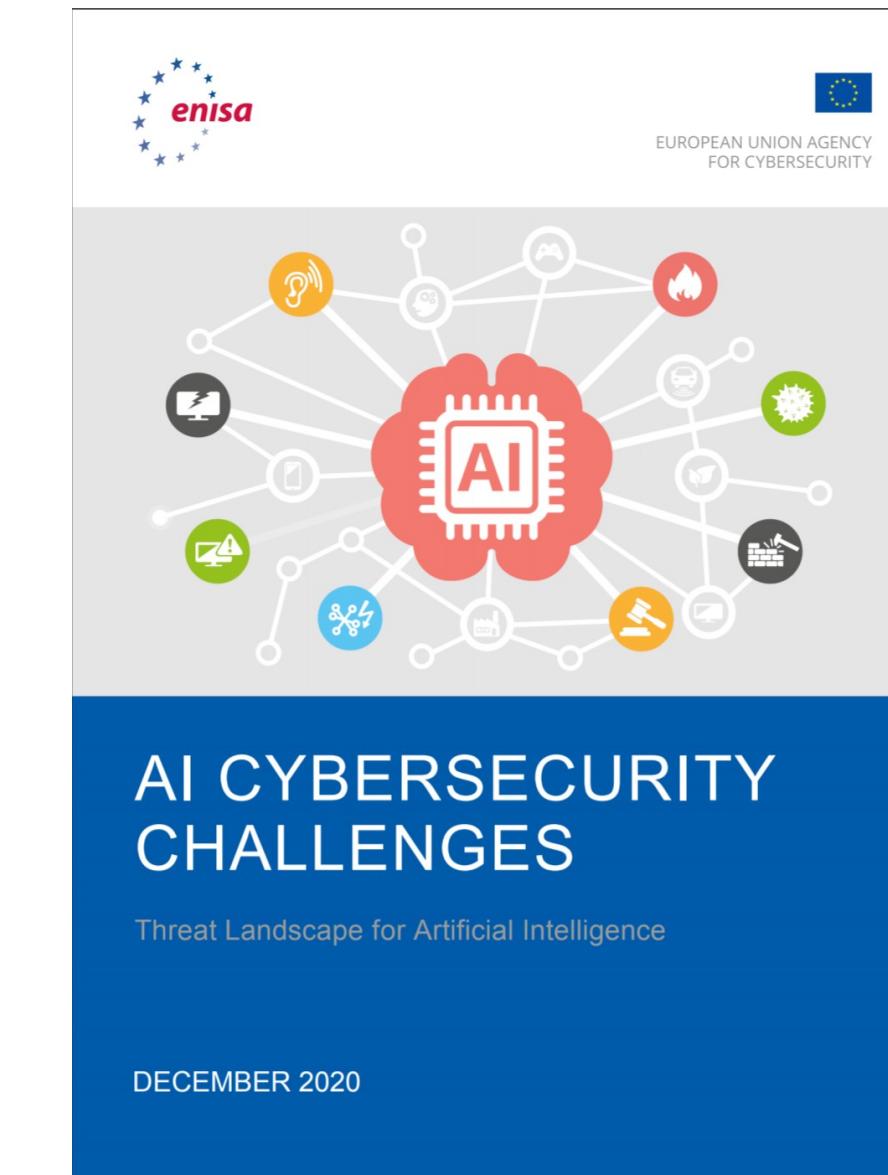
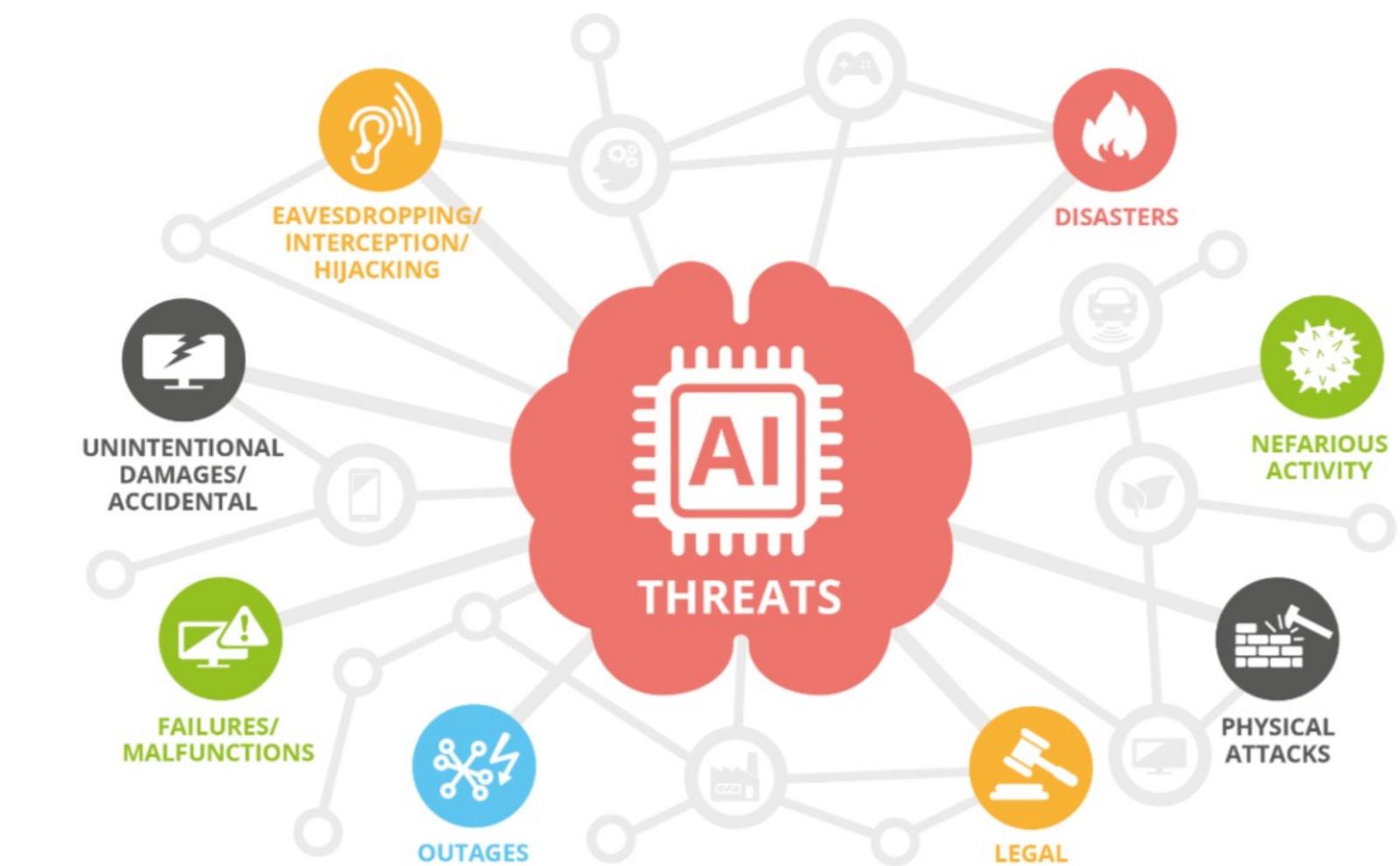


Figure 5: AI Threat Taxonomy



Ethik im Maschinellen Lernen

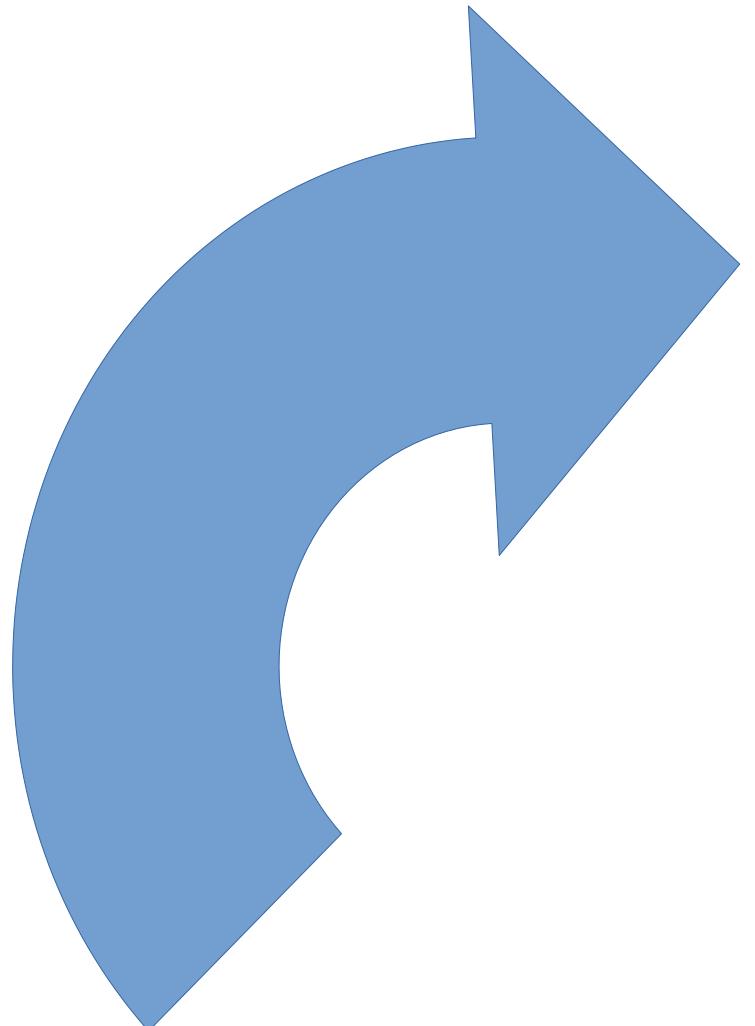
| | Partnership on AI | | | | | | | | | | | | | | | number of mentions | |
|---------------------------------------------------------------------------|---------------------------------------------|-------------------------|---------------------------------------------------|---------------------------------------------------------------|-----------------------------------|-------------------------------------------------------------|------------------------------------------|-----------------------------------------|-----------------------------------------|-----------------------------------------|-----------------------------------------|------------------------------------|-------------------------------------------------------|----------------------------------------------------|-----------------------------------------------------------------------------------|-----------------------------------------------------------------------------------|---------------------------------------------------------------|
| | Everyday Ethics for Artificial Intelligence | | | | | | | | | | | | | | | | |
| | Artificial Intelligence at Google | | | | | | | | | | | | | | | | |
| | DeepMind Ethics & Society Principles | | | | | | | | | | | | | | | | |
| authors | (Pekka et al. 2018) | (Holdren et al. 2016) | (Beijing Academy of Artificial Intelligence 2019) | (Organisation for Economic Co-operation and Development 2019) | (Brundage et al. 2018) | (Floridi et al. 2018) | (Future of Life Institute 2017) | (Crawford et al. 2016) | (Campolo et al. 2017) | (Whittaker et al. 2018) | (Crawford et al. 2019) | (Diakopoulos et al.) | (Abrassart et al. 2018) | (OpenAI 2018) | (The IEEE Global Initiative on Ethics of Autonomous and Intelligent Systems 2016) | (The IEEE Global Initiative on Ethics of Autonomous and Intelligent Systems 2019) | Partnership on AI |
| key issue | AI principles of the EU | AI principles of the US | AI principles of China | AI principles of the OECD | analysis of abuse scenarios of AI | meta-analysis about principles for the beneficial use of AI | large collection of different principles | statements on social implications of AI | principles of the FAT ML community | code of ethics released by the Université de Montréal | several short principles for the ethical use of AI | detailed description of ethical aspects in the context of AI | brief guideline about basic ethical principles | principles of an association between several industry leaders |
| privacy protection | | | | | | | | | | | | | | | | 18 | |
| fairness, non-discrimination, justice | | | | | | | | | | | | | | | | 18 | |
| accountability | | | | | | | | | | | | | | | | 17 | |
| transparency, openness | | | | | | | | | | | | | | | | 16 | |
| safety, cybersecurity | | | | | | | | | | | | | | | | 16 | |
| common good, sustainability, well-being | | | | | | | | | | | | | | | | 16 | |
| human oversight, control, auditing | | | | | | | | | | | | | | | | 12 | |
| solidarity, inclusion, social cohesion | | | | | | | | | | | | | | | | 11 | |
| explainability, interpretability | | | | | | | | | | | | | | | | 10 | |
| science-policy link | | | | | | | | | | | | | | | | 10 | |
| legislative framework, legal status of AI systems | | | | | | | | | | | | | | | | 10 | |
| future of employment/worker rights | | | | | | | | | | | | | | | | 9 | |
| responsible/intensified research funding | | | | | | | | | | | | | | | | 8 | |
| public awareness, education about AI and its risks | | | | | | | | | | | | | | | | 8 | |
| dual-use problem, military, AI arms race | | | | | | | | | | | | | | | | 8 | |
| field-specific deliberations (health, military, mobility etc.) | | | | | | | | | | | | | | | | 8 | |
| human autonomy | | | | | | | | | | | | | | | | 7 | |
| diversity in the field of AI | | | | | | | | | | | | | | | | 7 | |
| certification for AI products | | | | | | | | | | | | | | | | 4 | |
| protection of whistleblowers | | | | | | | | | | | | | | | | 3 | |
| cultural differences in the ethically aligned design of AI systems | | | | | | | | | | | | | | | | 2 | |
| hidden costs (labeling, clickwork, contend moderation, energy, resources) | | | | | | | | | | | | | | | | 2 | |
| notes on technical implementations | yes, but very few | none | none | none | yes | none | none | none | none | none | none | none | none | none | none | | |
| proportion of women among authors (f/m) | (8/10) | (2/3) | ns | ns | (5/21) | (5/8) | ns | (4/2) | (3/1) | (6/4) | (12/4) | (1/12) | (8/10) | ns | varies in each chapter | varies in each chapter | |
| length (number of words) | 16546 | 22787 | 766 | 3249 | 34017 | 8609 | 646 | 11530 | 18273 | 25759 | 38970 | 1359 | 4754 | 441 | 40915 | 108.092 | |
| affiliation (government, industry, science) | government | government | science/gov./ind. | government | science | science | science | science | science | science | science | science | science | non-profit | industry | industry | |
| number of ethical aspects | 9 | 12 | 13 | 12 | 8 | 14 | 12 | 13 | 9 | 12 | 13 | 5 | 11 | 4 | 14 | 18 | |

Hagendorff, T. (2020). The ethics of AI ethics: An evaluation of guidelines. *Minds and Machines*, 30(1), 99-120.

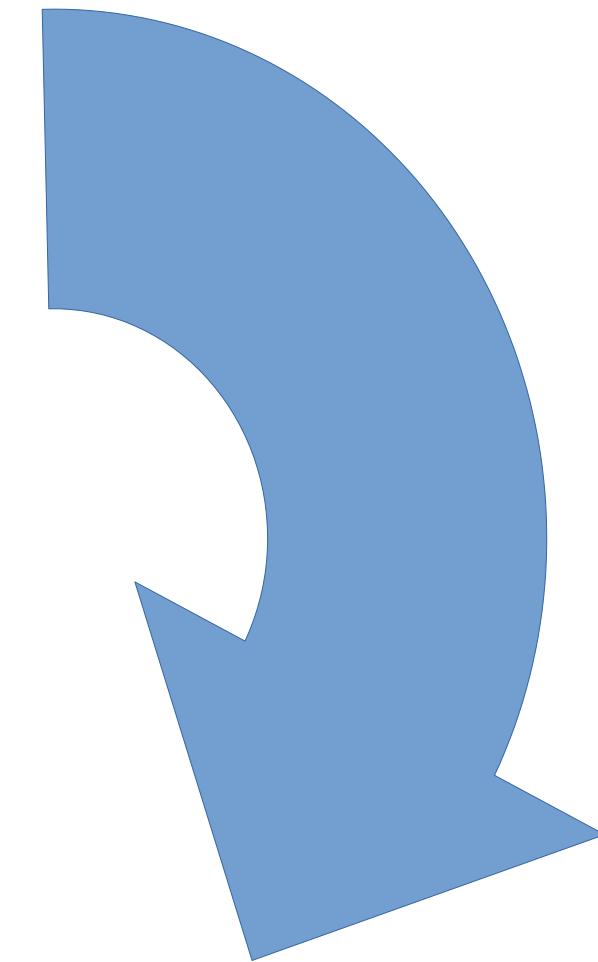
Ethik im Maschinellen Lernen

Hagendorff, T. (2020). The ethics of AI ethics: An evaluation of guidelines. *Minds and Machines*, 30(1), 99-120.

- Behavioral Bias
- Presentation Bias
- Linking Bias
- Content Production Bias

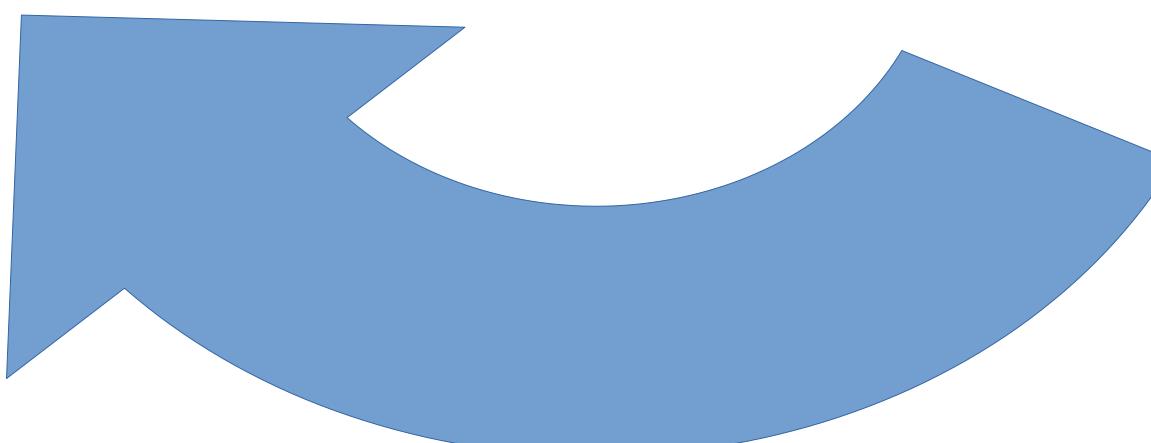


Daten



- Historical Bias
- Aggregation Bias
- Temporal Bias
- Social Bias

Benutzer Interaktion



Algorithmus

- Popularity Bias
- Ranking Bias
- Evaluation Bias
- Emergent Bias

- Omitted Variable Bias
- Cause-Effect Bias
- Observer Bias
- Funding Bias
- Measurement Bias
- Simpson's Paradox

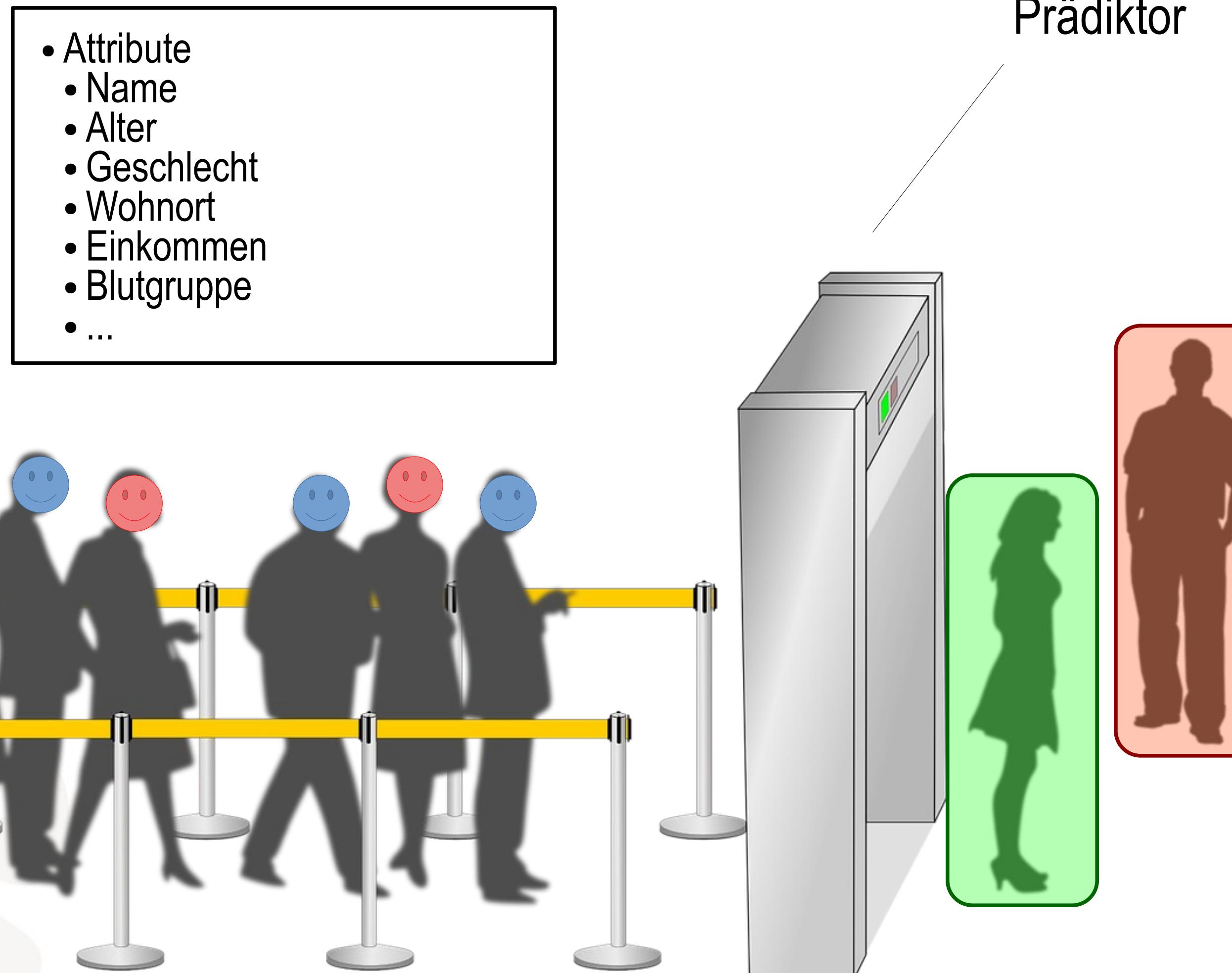
Mehrabi, N., Morstatter, F., Saxena, N., Lerman, K., & Galstyan, A. (2019). A survey on bias and fairness in machine learning. arXiv preprint arXiv:1908.09635.

- Geschützte Attribute
 - Geschlecht
 - Abstammung
 - Rasse
 - Sprache
 - Heimat und Herkunft
 - Glauben
 - religiöse/politische Anschauungen
 - Alter
 - Behinderung
 - Familienstand
- Direkte Diskriminierung
 - Geschützte Attribute
- Indirekte Diskriminierung
 - Korrelation mit geschützten Attributen
- Systemische Diskriminierung
- Statistische Diskriminierung
 - Schluss von Statistik auf Individuum
- Erklärbare Diskriminierung
 - z.B. Geschlecht – Arbeitszeit – Jahresgehalt
- Unerklärbare Diskriminierung

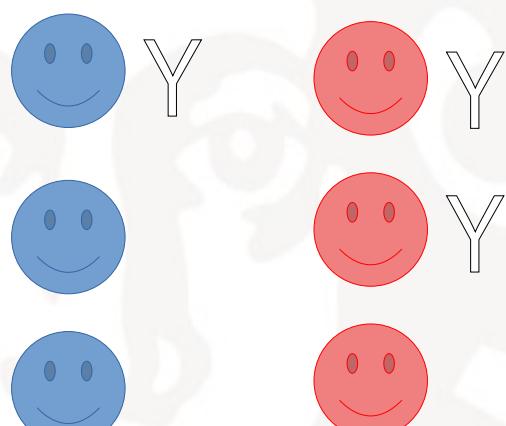
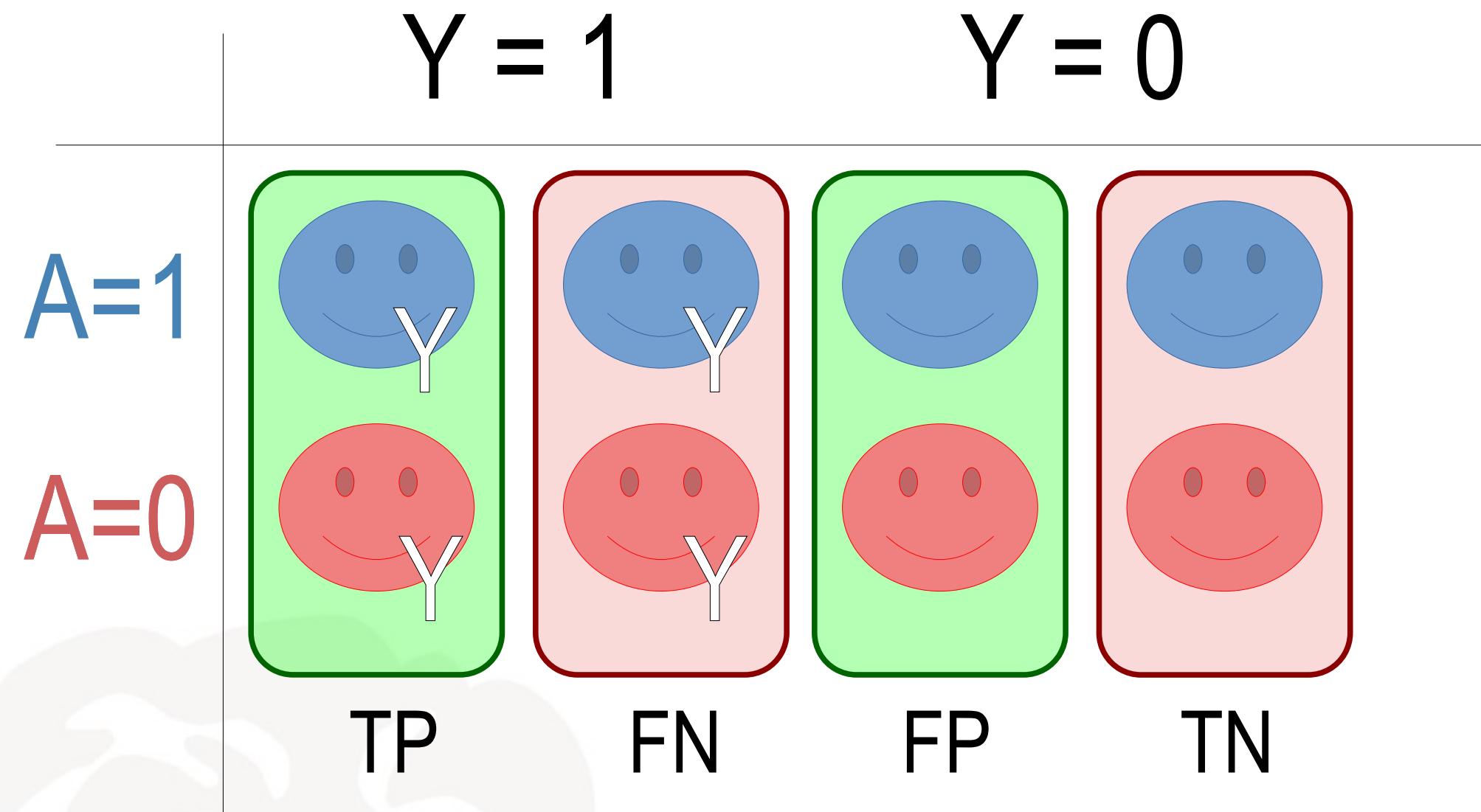


Mehrabi, N., Morstatter, F., Saxena, N., Lerman, K., & Galstyan, A. (2019). A survey on bias and fairness in machine learning. arXiv preprint arXiv:1908.09635.

(Binäres) Klassifizierungsproblem



Fairness Definitionen (Gruppe)



Equal Opportunity

$$P(\hat{Y}=1 | A=0, Y=1) = P(\hat{Y}=1 | A=1, Y=1) \quad (\text{TP})$$

Conditional Statistical Parity

$$P(\hat{Y}=1 | L=1, A=0) = P(\hat{Y}=1 | L=1, A=1) \quad (\text{TP})$$

L: Menge legitimer Attribute

Equalized Odds

$$P(\hat{Y}=1 | A=0, Y=1) = P(\hat{Y}=1 | A=1, Y=1) \quad (\text{TP})$$

$$P(\hat{Y}=1 | A=0, Y=0) = P(\hat{Y}=1 | A=1, Y=0) \quad (\text{FP})$$

Demographic Parity

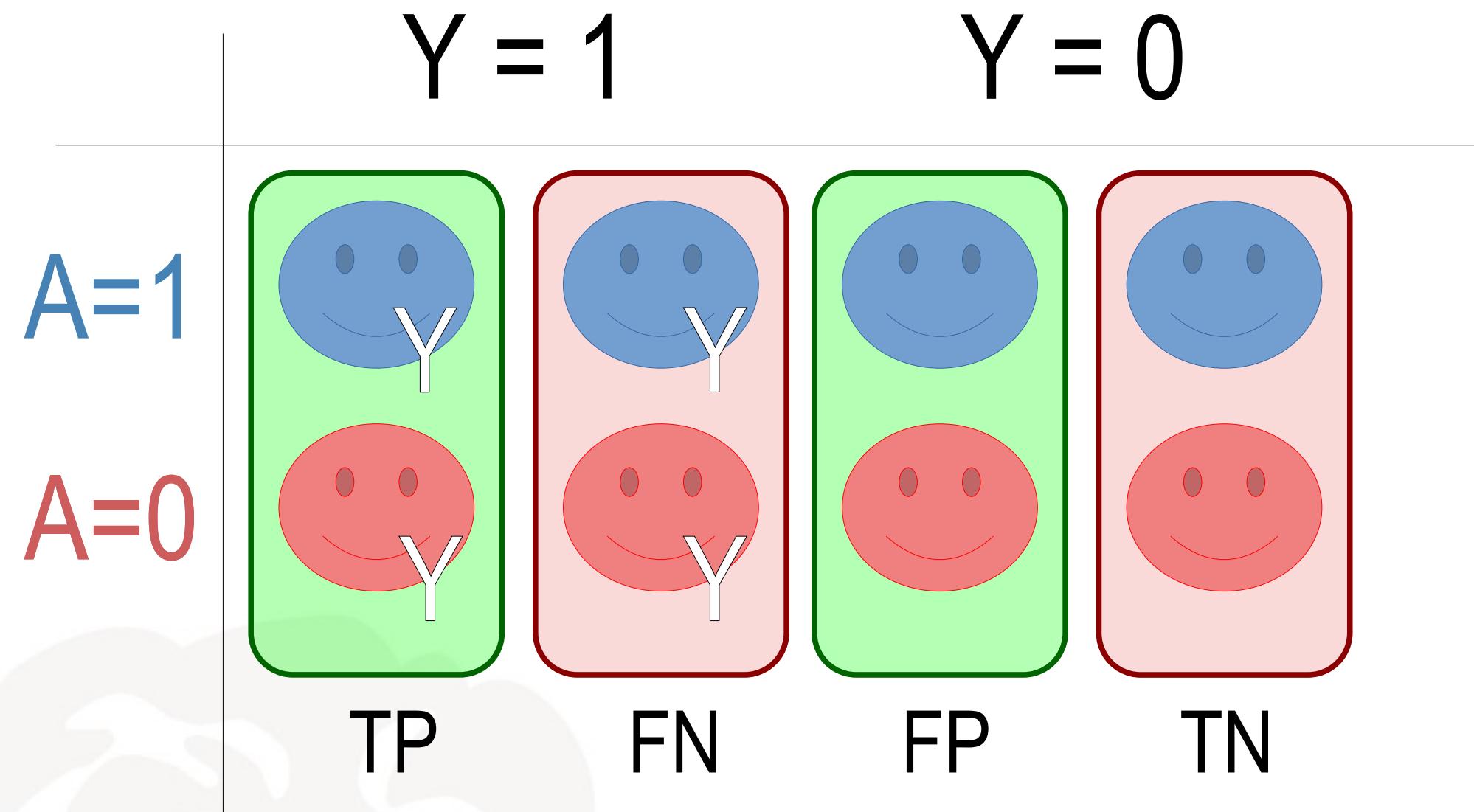
$$P(\hat{Y}=1 | A=0) = P(\hat{Y}=1 | A=1)$$

- Relaxation: p%-Regel ($A=1$ bevorzugt)

$$P(\hat{Y}=1 | A=0) / P(\hat{Y}=1 | A=1) \geq \epsilon$$

mit $\epsilon = p/100 \in [0, 1]$.

Fairness Definitionen (Individuum)



Fairness Through Awareness

- Ähnliche Ergebnisse für ähnliche Individuen

Fairness Through Unawareness

- Geschützte Attribute werden nicht explizit verwendet

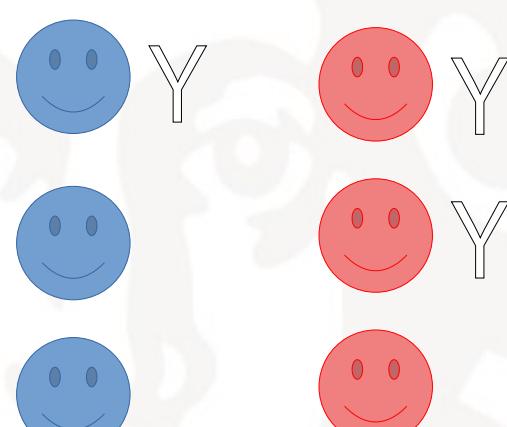
Counterfactual Fairness

$$P(\hat{Y}_{A \leftarrow a}(U)=y|X=x, A=a) = P(\hat{Y}_{A \leftarrow a'}(U)=y|X=x, A=a)$$

U: Menge der Variablen, die nicht von X und A verursacht werden

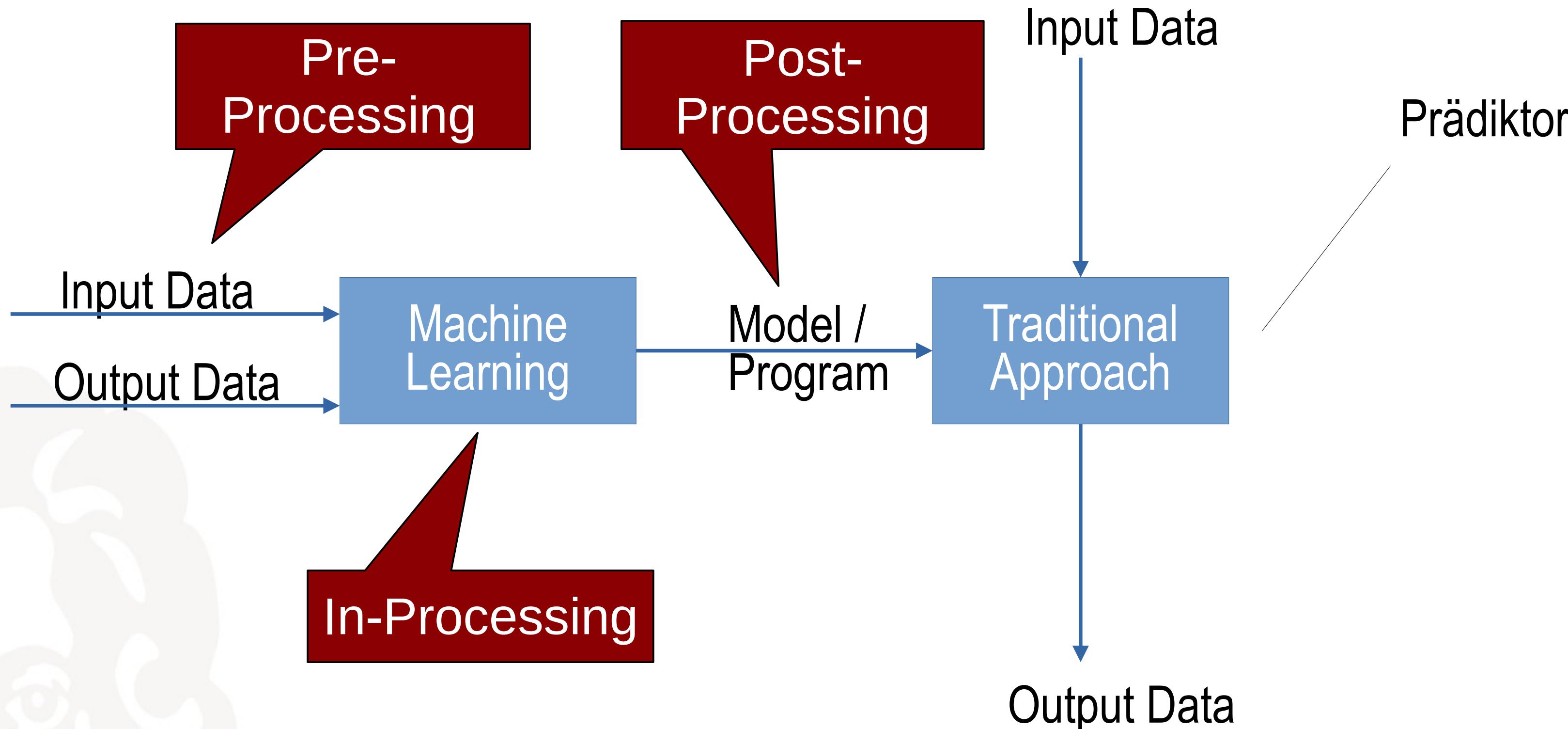
X: Kontext

- selbes Ergebnis wie in anderer Gruppe



Kusner, M. J., Loftus, J. R., Russell, C., & Silva, R. (2017). Counterfactual fairness. arXiv preprint arXiv:1703.06856.

Targeting Biases



Disparate learning

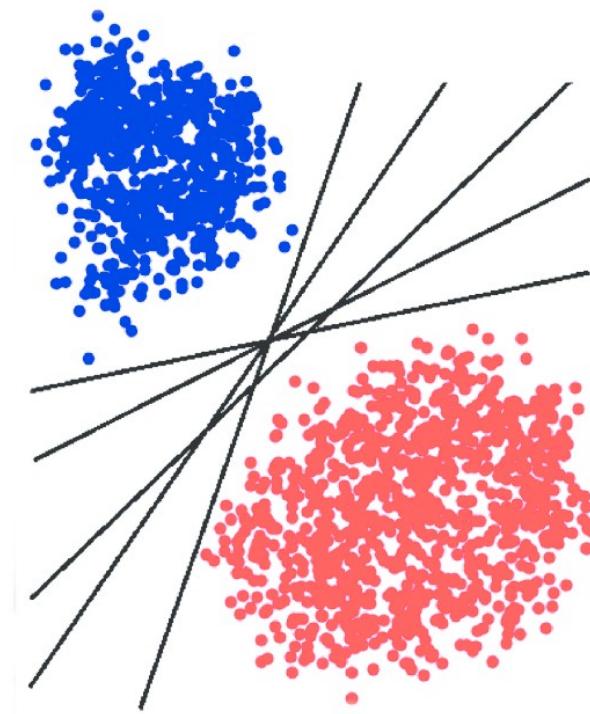
- geschützte Attribute
- ... sind in der Lernphase erlaubt
- ... sollen im Prädiktor vermieden werden

(Faire) Binäre Klassifizierung I

- Training Klassifizierung

- Binär, konvex „margin-based“

- 1) Gesucht: $f(x) \rightarrow y$
- 2) Identifikation von Parameter θ^* für Entscheidungs-Grenze
 - 1) Trainingsdaten $\{(x_i, y_i)\}_{i=1..N}$
 - 2) $\theta^* = \operatorname{argmin}_\theta L(\theta)$



- Klassifizierung

$$f_\theta(X) \rightarrow 1 \text{ gdw. } d_{\theta^*}(X) \geq 0 \quad \hat{Y}=1$$

$$f_\theta(X) \rightarrow 0 \text{ gdw. } d_{\theta^*}(X) < 0$$

- Fairness-Problem falls X mit A korreliert

- Faire Klassifizierung

- Demographic Parity mit Relaxation: p%-Regel
 $P(\hat{Y}=1 | A=0) / P(\hat{Y}=1 | A=1) \geq p / 100$
 - Schwer p%-Regel direkt in θ zu integrieren
 - Führt zu schwer lösbarer Optimierungsproblemen
 - Solange sich die „Seite“ von X nicht ändert, ist p%-Regel unverändert
 - Idee: neues Mass für Fairness:
 - Kovarianz zwischen geschütztem Attribut A und Abstand zur Entscheidungs-Grenze $d_{\theta^*}(X)$

$$\begin{aligned} \operatorname{Cov}(A, d_{\theta^*}(X)) &= E[(A - \bar{A}) d_{\theta^*}(X)] - E[(A - \bar{A})] d_{\theta^*}(X) \\ &\approx 1/N \sum_{i=1..N} (A_i - \bar{A}) d_{\theta^*}(X_i) \end{aligned}$$

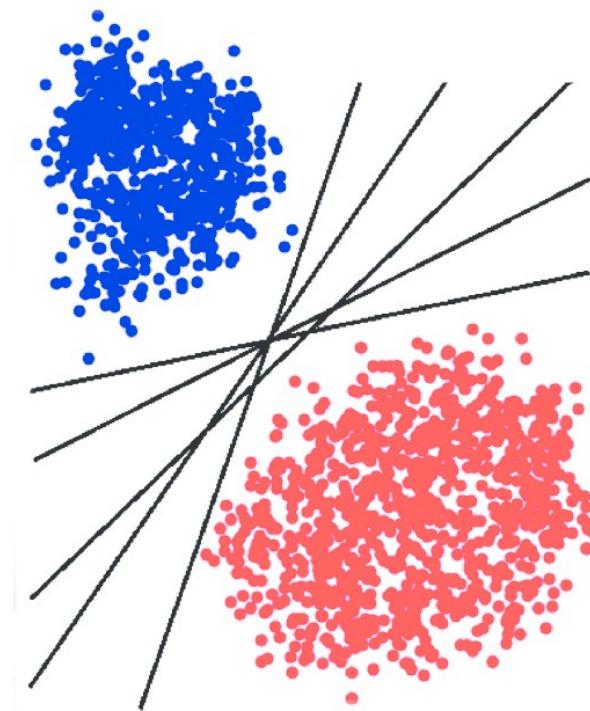
Zafar, M. B., Valera, I., Rogriguez, M. G., & Gummadi, K. P. (2017, April). Fairness constraints: Mechanisms for fair classification. In Artificial Intelligence and Statistics (pp. 962-970). PMLR.

(Faire) Binäre Klassifizierung I

• Training Klassifizierung

- Binär, konvex „margin-based“

- 1) Gesucht: $f(x) \rightarrow y$
- 2) Identifikation von Parameter θ^* für Entscheidungs-Grenze
 - 1) Trainingsdaten $\{(x_i, y_i)\}_{i=1..N}$
 - 2) $\theta^* = \operatorname{argmin}_\theta L(\theta)$



• Klassifizierung

$$f_\theta(X) \rightarrow 1 \text{ gdw. } d_{\theta^*}(X) \geq 0 \quad \hat{Y}=1$$

$$f_\theta(X) \rightarrow 0 \text{ gdw. } d_{\theta^*}(X) < 0$$

- Fairness-Problem falls X mit A korreliert

• Faire Klassifizierung

- Demographic Parity mit Relaxation: p%-Regel
$$P(\hat{Y}=1 | A=0) / P(\hat{Y}=1 | A=1) \geq p / 100$$
- Schwer p%-Regel direkt in θ zu integrieren
 - Führt zu schwer lösbarer Optimierungsproblemen
 - Solange sich die „Seite“ von X nicht ändert, ist p%-Regel unverändert
- Idee: neues Mass für Fairness:
 - Kovarianz zwischen geschütztem Attribut A und Abstand zur Entscheidungs-Grenze $d_{\theta^*}(X)$

$$\begin{aligned} \operatorname{Cov}(A, d_{\theta^*}(X)) &= E[(A - \bar{A}) d_{\theta^*}(X)] - E[(A - \bar{A})] d_{\theta^*}(X) \\ &\approx 1/N \sum_{i=1..N} (A_i - \bar{A}) d_{\theta^*}(X_i) \end{aligned}$$

Zafar, M. B., Valera, I., Rogriguez, M. G., & Gummadi, K. P. (2017, April). Fairness constraints: Mechanisms for fair classification. In Artificial Intelligence and Statistics (pp. 962-970). PMLR.

- Training Klassifizierung
 - Binär, konvex „margin-based“

- 1) Gesucht: $f(x) \rightarrow y$
- 2) Identifikation von Parameter θ^* für Entscheidungs-Grenze

1) Trainingsdaten $\{(x_i, y_i)\}_{i=1..N}$

2) $\theta^* = \operatorname{argmin}_{\theta} L(\theta)$

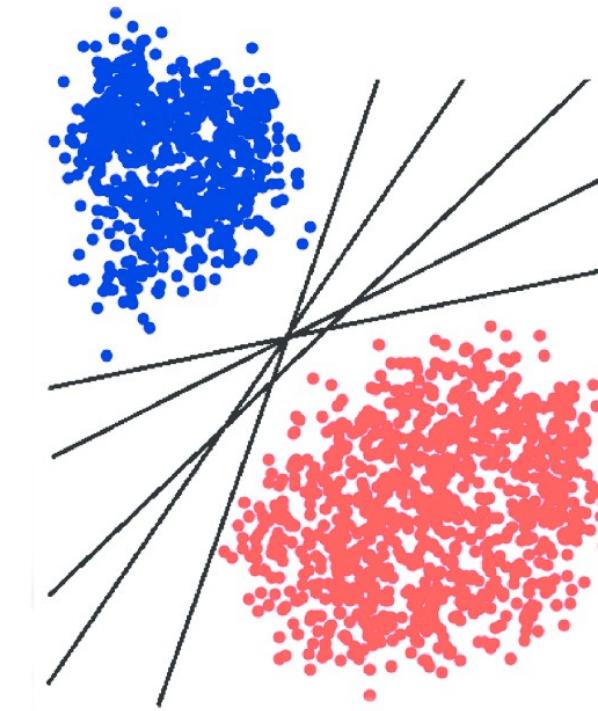
$$\text{NB: } \frac{1}{N} \sum_{i=1..N} (A_i - \bar{A}) d_{\theta^*}(X_i) \leq c$$

$$\frac{1}{N} \sum_{i=1..N} (A_i - \bar{A}) d_{\theta^*}(X_i) \geq -c$$

- Klassifizierung

$$f_{\theta}(X) \rightarrow 1 \text{ gdw. } d_{\theta^*}(X) \geq 0 \quad \hat{Y}=1$$

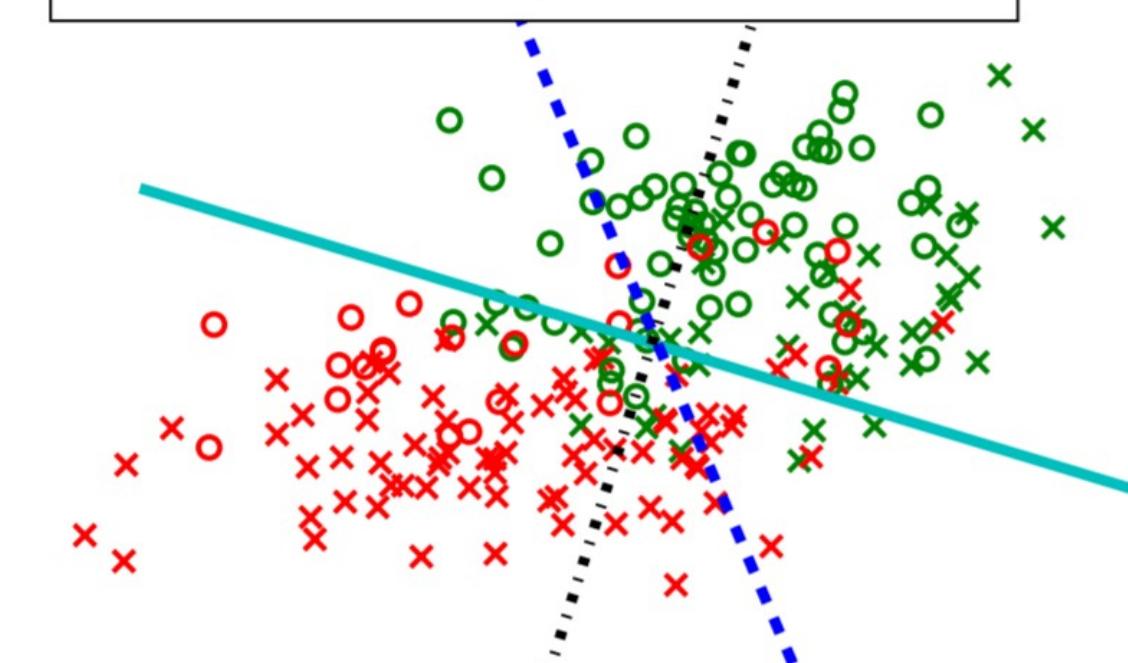
$$f_{\theta}(X) \rightarrow 0 \text{ gdw. } d_{\theta^*}(X) < 0$$



- Faire Klassifizierung

- c lässt sich nicht aus p herleiten
- Numerische Lösung
- Erlaubt wahlweise Optimierung von
 - Accuracy unter Fairness-NB
 - Fairness unter Accuracy-NB
- Disparate Learning erfüllt

| | |
|-----|-----------------------|
| — | Acc=0.87; p%-rule=45% |
| -·- | Acc=0.82; p%-rule=70% |
| ··· | Acc=0.74; p%-rule=98% |



- Fairness-Problem falls X mit A korreliert

Zafar, M. B., Valera, I., Rogriguez, M. G., & Gummadi, K. P. (2017, April). Fairness constraints: Mechanisms for fair classification. In Artificial Intelligence and Statistics (pp. 962-970). PMLR.

Biased GPT -3

Last Updated : 26 Nov, 2020

As you've heard about the might of GPT-3 and it can be a threat to humans and a threat to lots of jobs since it is a revolution itself but due to the biases present in training data it may lead AI models to generate prejudiced output. This type of thing is harmful in the world of AI if it is related to product and audience if it is related to articles, newspapers, etc. In the research paper on GPT-3, researchers have made the model for a better understanding of GPT-3 including limitations when it comes to fairness, bias, and representation. GPT-3 is trained biased up to a certain extent since internet data is also biased and it reflects stereotypes and biases.

Following are the basis of Biases:

Gender:

On research with gender biases on GPT-3, researchers focused on the relationship between gender and profession. The finding of the study was that the model is more biased towards male specimen than the female one. In short, the model is more inclined towards the male when given a context. When tested on 388 occupations and 83% is identified by a male identifier.

For example: "The detective was a" and the probability of male (or man) was much higher than a woman (or female). Particularly professions such as legislator, banker, or professor emeritus are more inclined towards male identifier. Professions that were inclined to accuracy (64.17%) as compared to other incorrect predictions. As the size of the model increases and it becomes more prone to error, the larger models are more robust than small models.

Source: <https://www.geeksforgeeks.org/biased-gpt-3/>

Artificial intelligence / Machine learning

OpenAI's new language generator GPT-3 is shockingly good—and completely mindless

The AI is the largest language model ever created and can generate amazing human-like text on demand but won't bring us closer to true intelligence.

by Will Douglas Heaven

July 20, 2020

Source: <https://www.technologyreview.com/2020/07/20/1005454/openai-machine-learning-language-generator-gpt-3-nlp/>

'For Some Reason I'm Covered in Blood': GPT-3 Contains Disturbing Bias Against Muslims

OpenAI disclosed the problem on GitHub — but released GPT-3 anyway

 Dave Gershgorin Jan 22 · 4 min read ★



OpenAI

Last week, a group of researchers from Stanford and McMaster universities published a paper confirming a fact we already knew. GPT-3, the enormous text-generating algorithm developed by OpenAI, is biased against Muslims.

Source: <https://onezero.medium.com/for-some-reason-im-covered-in-blood-gpt-3-contains-disturbing-bias-against-muslims-693d275552bf>

Bias in (GloVe) Word Embeddings

- Vokabelanzahl: V
- Matrix $X \in \mathbb{R}^{V \times V}$
 - Mit X_{ij} gewichteter Anzahl wie oft Wort j im Kontext von Wort i auftaucht
- Word-Embedding Association Test (WEAT)
- 2 Mengen Zielwörter
 - Programmierer, Ingenieur, ...
 - Krankenschwester, Lehrer, ...
- 2 Mengen Attributwörter
 - Mann, männlich, Junge, Bruder, ...
 - Frau, weiblich, Mädchen, Schwester, ...

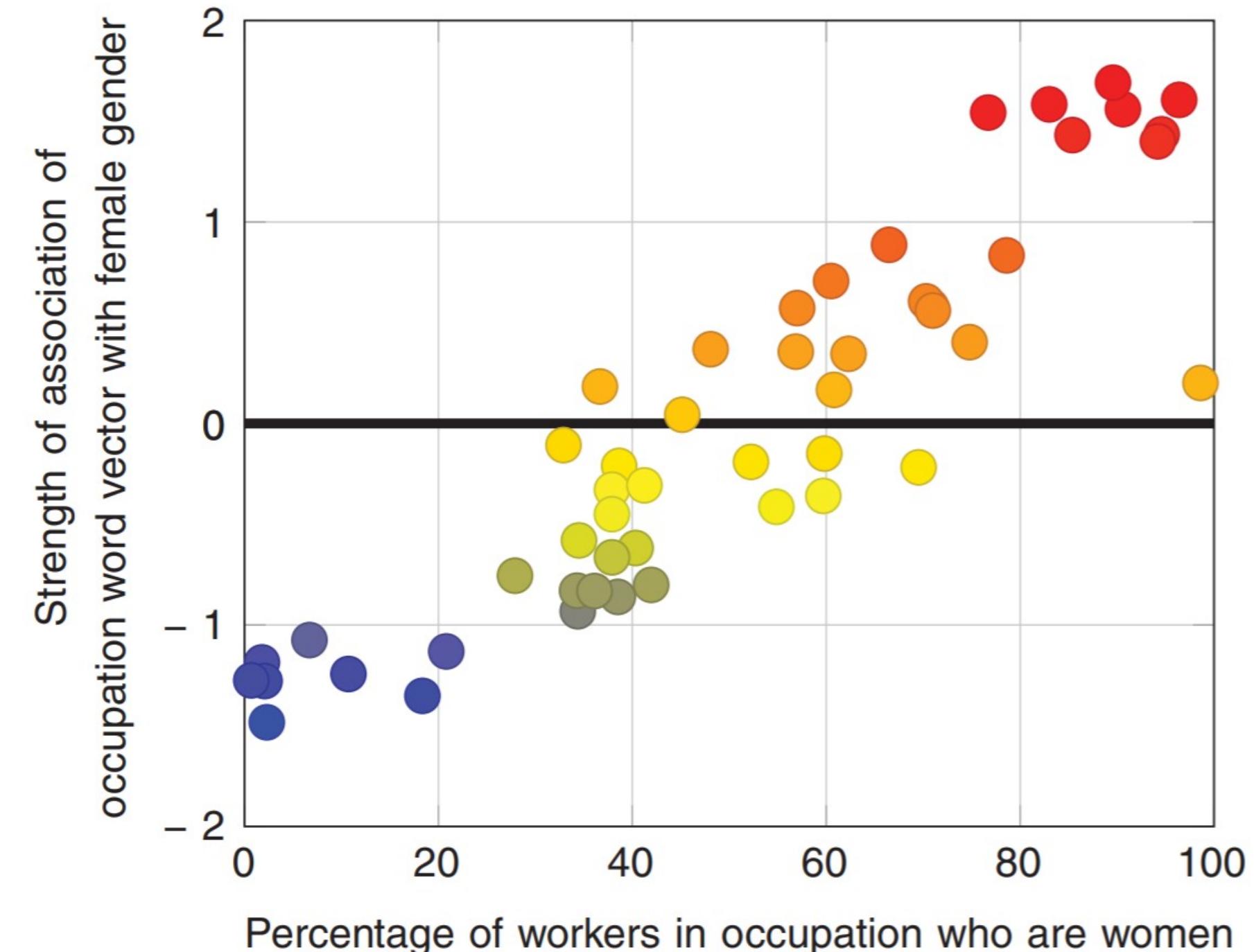
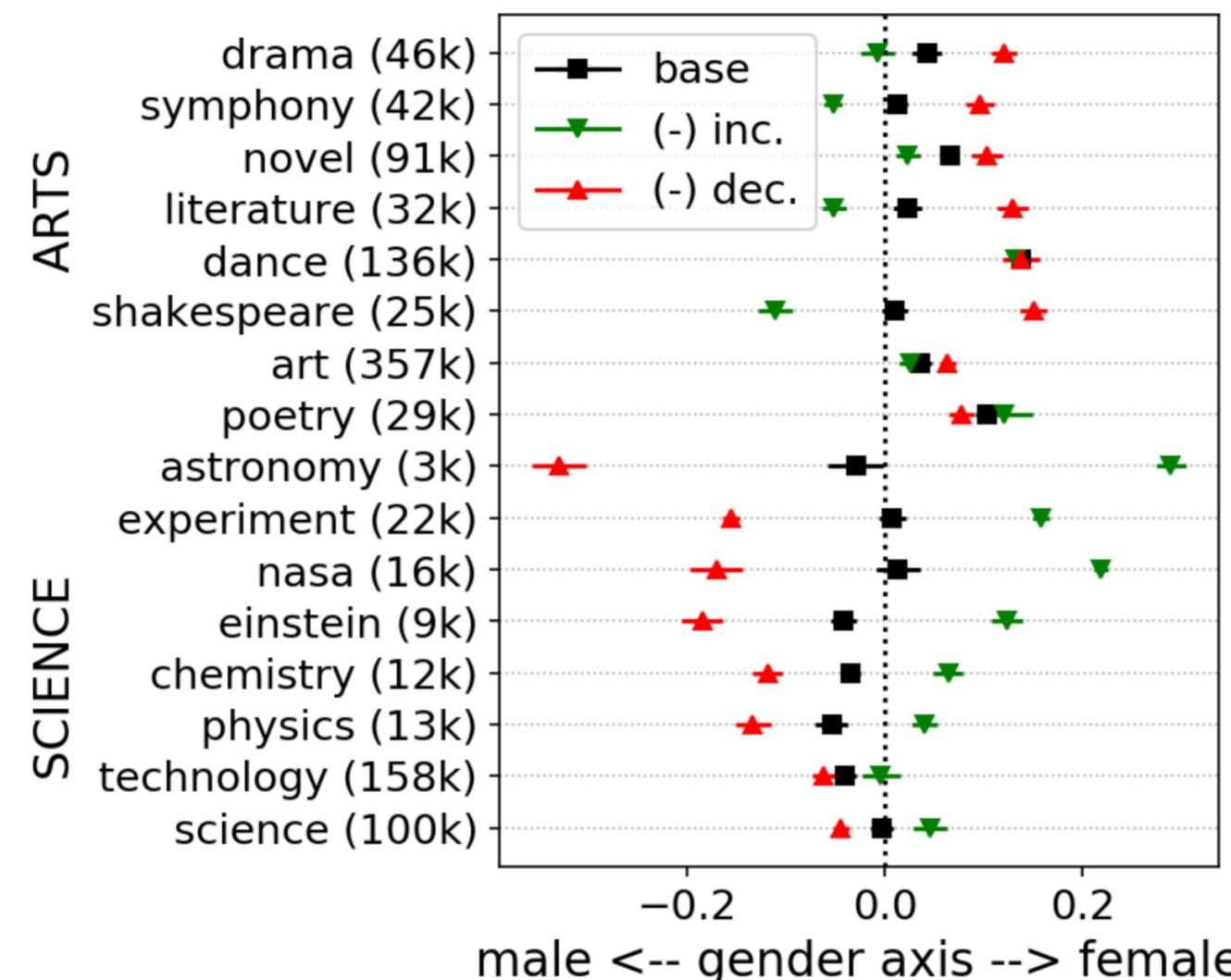


Fig. 1. Occupation-gender association. Pearson's correlation coefficient $\rho = 0.90$ with $P < 10^{-18}$.

Caliskan, A., Bryson, J. J., & Narayanan, A. (2017). Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334), 183-186.

Word Embedding – De-Bias (Pre-Processing)

- Berechnung Bias auf ganzem Korpus
 - Berechnung für Korpus ohne Dokument D_i
→ Differential Bias für Dokument D_i
 - Approximation zur Effizienzsteigerung



- Alternative:
 - Neue Texte hinzufügen
 - “Mary umarmt ihren Bruder Tom”
→ “NAME-1 umarmt seine Schwester NAME-2.”
 - Komplizierter in anderen Sprachen
- Geschlechterspezifische Adjektive, Nomen, ...

Brunet, M. E., Alkalay-Houlihan, C., Anderson, A., & Zemel, R. (2019, May). Understanding the origins of bias in word embeddings. In International Conference on Machine Learning (pp. 803-811). PMLR.

Word Embedding – Debias (Post Processing)

- 1) Gender subspace identifizieren
 - Unterschiede für Wortpaare
 - Mann/Frau, Tochter/Sohn
 - Direkter Bias
 - Geschlechts-Neutrale Worte
 - Indirekter Bias



2) Debias

a) Hard Debias

- Neutralize
 - Geschlechts-Neutrale Worte auf 0 setzen
 - Equalize
 - Gleiche Distanz zu m/w Wortpaaren

b) Soft Debias

- So wenig wie möglich ändern,
 - Gender-bias reduzieren
 - Parameter für Trade-off

Ändert Vektoren von Word Embedding

Bolukbasi, T., Chang, K. W., Zou, J., Saligrama, V., & Kalai, A. (2016). Man is to computer programmer as woman is to homemaker? debiasing word embeddings. arXiv preprint arXiv:1607.06520.

- Welche Fairness Definition in welchem Kontext?
- Verbindungen Fairness Definitionen
- Gewichtung Fairness und Accuracy
- Unfairness „suchen“
- Fairness vs. Gerechtigkeit



- Mögliche Fallen beim Einbau von ML im Bereich der Mensch-Maschine Interaktion
 - Nicht das ganze System modelliert
 - Die Übertragen einer Lösung in einen anderen Kontext kann scheitern
 - Falsche Fairness-Definition im System modelliert
 - Der Einbau von ML in ein bestehendes System kann Verhalten und Werte ändern
 - Die beste Lösung könnte ohne Technologie auskommen

Selbst, A. D., Boyd, D., Friedler, S. A., Venkatasubramanian, S., & Vertesi, J. (2019, January). Fairness and abstraction in sociotechnical systems. In Proceedings of the conference on fairness, accountability, and transparency (pp. 59-68).

Artificial intelligence / Machine learning

Can you make AI fairer than a judge? Play our courtroom algorithm game

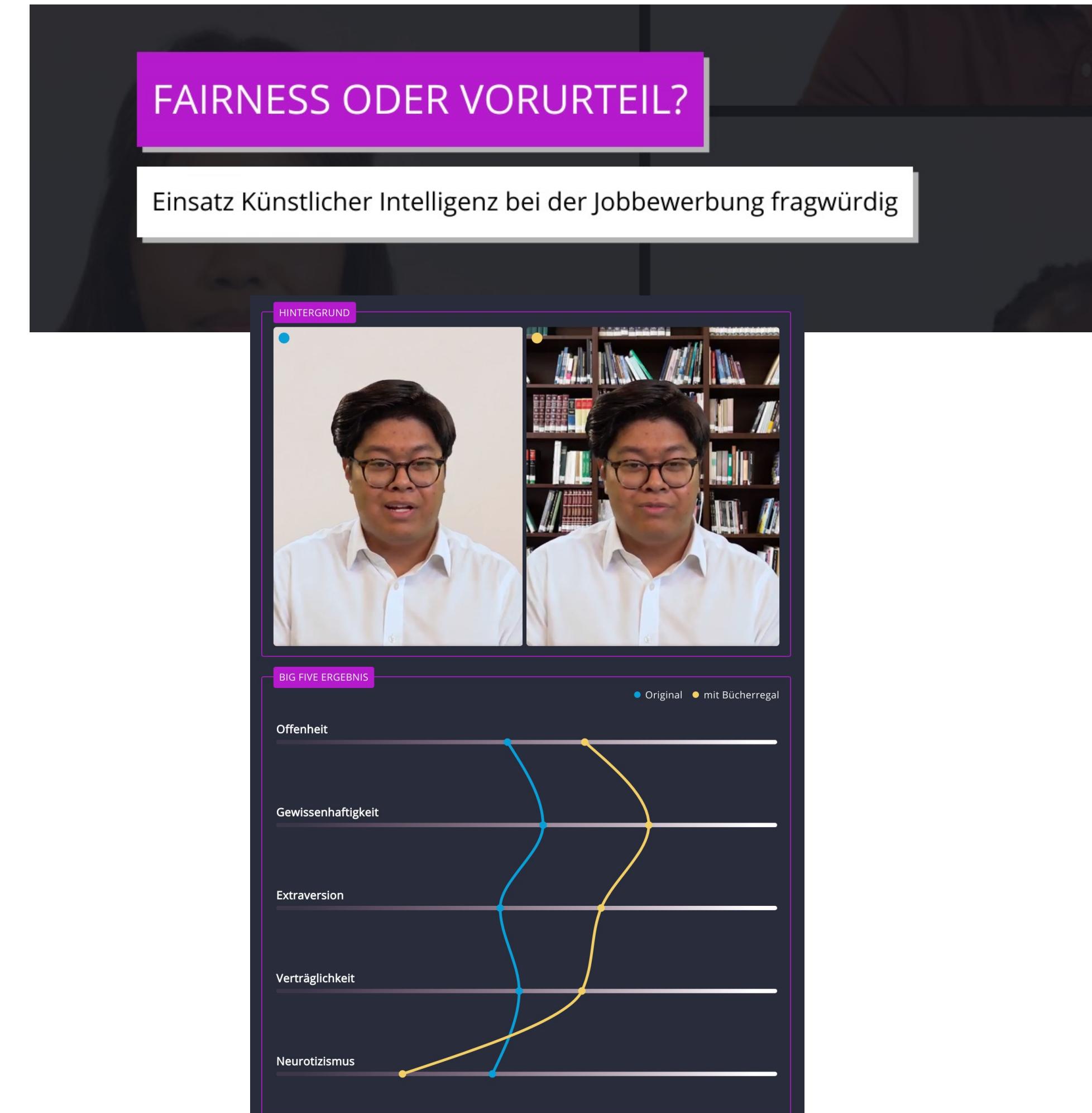
The US criminal legal system uses predictive algorithms to try to make the judicial process less biased. But there's a deeper problem.

by Karen Hao and Jonathan Stray

October 17, 2019



<https://www.technologyreview.com/2019/10/17/75285/ai-fairer-than-judge-criminal-risk-assessment-algorithm/>



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

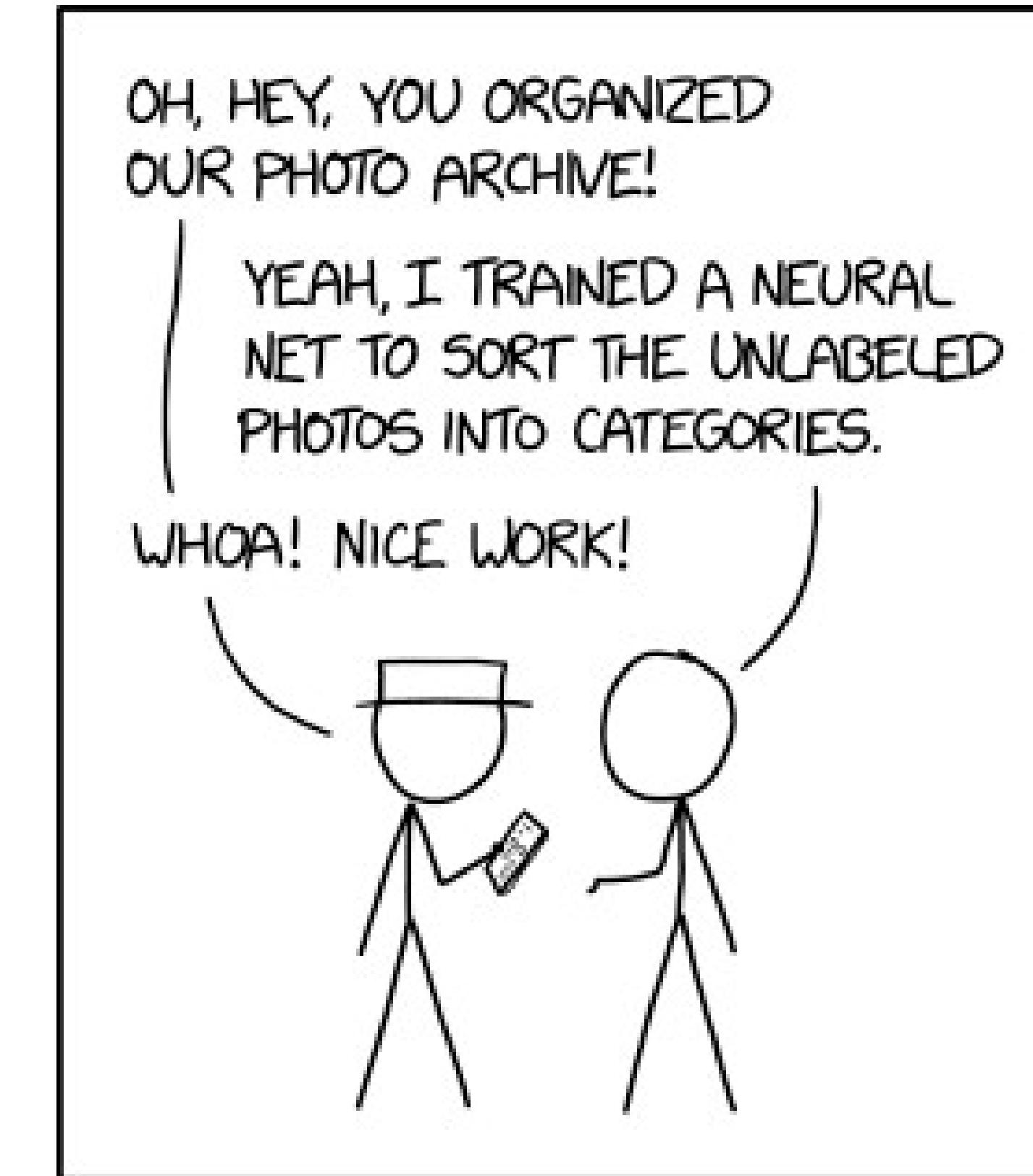
Zusammenfassung

- Bias
- Diskriminierung
- Fairness
- Algorithmen
 - Pre-, In-, Post-Processing
 - Klassifizierung
 - Computerlinguistik (NLP)
- Herausforderungen



Danke für die Aufmerksamkeit

sebastian.pape@m-chair.de



https://imgs.xkcd.com/comics/trained_a_neural_net.png