

information  
Regulating AI and its inequalities for more  
See Balayn & Grzes (2021): Beyond Debasing.

behind biases.  
more so if we do not question  
societal roots and assumptions  
efforts, bias will persist, and even  
mitigated technically. Despite best  
including discrimination, will be  
because it implies that biases,  
“debasing”, can be misleading  
terminology used, such as  
**solve societal problems**. Also, the  
technical solution might sound  
Even though this promise of a  
tempering, it is **not the right way to**

“unbiased” dataset.  
mitigation methods to gain an  
**fairness metrics** and technical  
treating a biased dataset using  
Debiasing describes the idea of  
problems with discrimination.  
understood as a means to solve  
often used in policy documents and  
called “**debasing**”. It is a term that is  
dealing with biased datasets is  
many entry points. One way of  
and how we deal with them from  
So, biases seep into data, models,  
model learns that the most

recognised by the AI.  
this would not be  
would be great for the job,  
candidate has hobbies that  
learns to focus on that  
in the “career” section. It  
interesting information is  
model learns that the most

**predictors** over others.  
**systematically favour certain**  
models introduce bias into the  
system, causing the **model to**  
The features and labels used in the

up on someone’s CV.  
connections may not end  
networks, ethics or  
qualities such as  
life? Certain important  
PDF summaries a person’s  
kindness, life experiences  
How well can a one-page

distort the model’s view of reality.  
of what we’re really interested in, and  
model are **inaccurate representations**  
The features and labels used in the

## METHODS OF BIASES

### ALGORITHMIC BIAS

#### BIAS

The term **bias** describes a **systematic deviation**. Different research areas each have their own understanding of what exactly “bias” entails. **These different concepts of bias can lead to misunderstandings and misconceptions about how to deal with them.**

In the context of AI, bias was initially used as a term that focussed on technologically induced deviations. The discussion of societal biases in the context of AI is more recent.



Biases seep into AI systems in various places – and there is no single approach to tackle all of them.

Still, Techie wants to know how discriminatory their AI model actually is. How can Techie measure the bias of their model?

Find out more in the next zines!

## WHERE DOES UNFAIRNESS COME FROM?

After implementing the hiring AI, Techie notices discriminatory decisions by the algorithm. So, Techie decides to investigate. Where does the unfairness come from and how did it get into Techie’s hiring algorithm?



## BIAS IN AI CONTEXTS

Biases can influence AI models in **different ways at different stages of development**. Often, a bias is either due to biased data or biased training of the model. AI models then adopt biases, which can have discriminatory effects.

Two examples of biases are explained on the following pages, more can be found on the poster page of this zine!

# FEELED FAIR?

## WHICH BIASES INFLUENCE AI?



How well can a one-page PDF summarise a person's life?

- HISTORICAL BIAS
- Pre-existing societal biases are reflected in the data used to train the model.

- MEASUREMENT BIAS
- The features and labels used in the model are inaccurate representations of what we're really interested in.

- LABEL BIAS
- Binary labels misrepresent our complex reality, e.g. giving the impression that the people hired are perfect and the ones rejected are worthless.

The evaluation is done by C.O.R.P... They care about accuracy and efficiency – not fairness!

- EVALUATION BIAS
- The model is evaluated using criteria that don't accurately reflect its use in the real world.

- FEEDBACK BIAS
- Given our AI's bias towards white men, more white men will be hired and more women of colour will be rejected.

- REPRESENTATION BIAS
- The data used don't accurately reflect the world's diversity and complexity.

- ALGORITHMIC BIAS
- Models introduce bias into the system.

- DEPLOYMENT BIAS
- Occurs when the context in which the model is used differs from the training environment.

- FAIRNESS NOTES:
- This selection of biases is based on van Giffen et al. (2022): Overcoming the pitfalls and perils of algorithms: A classification of machine learning biases and mitigation methods.

- REALISATION BIAS
- Our training dataset spanned the last 20 years. Characteristics that we value in new applicants may not be present in the old ones.

- A BIASED AI IS UNFAIR...



Because of our society's sexist and racist history, biases such as "white men are more successful in leadership positions" are picked up by AI.

- REALISATION BIAS
- If C.O.R.P.'s dataset mainly labels white male applicants as hired, the AI won't learn much about women of colour. If a woman of colour applied to C.O.R.P. in the future, she wouldn't have much of a chance – even if her skills were similar to those hired.

- FAIRNESS NOTES:
- Some of the biases are also known under other names; for details, check the paper "Debiasing" as a solution to biased datasets is mostly mentioned in policy documents. As an investigation from Balayn & Gürses (2021) shows, they found that policy documents did not usually specify what they understood as debiasing.