



What is it? Why is it a concern? How to ensure it?

> Francisca Morgado Data Scientist at NILG.AI

 $\square \square \square$

	•	*	٠		*	*	*	٠	*	*	*	*	٠	*	٠	•	٠	•	*	٠	*	*	*	*	*	*		*	•
		٠	*	*	٠	٠	*	*	*	*	*	*			*	*	*	٠	٠	٠	٠	*	*	*		٠		*	
		•	٠			٠			•	•		·			*	•	٠	•					•			•		•	
,		•				•	*		•	•		•	٠		•	•		•		*		٠	•				,	•	*
,		•	•	*	*	•	*	٠	•	•	٠	•	٠		٠	٠	٠	*		*	•		•	٠	٠	•	·	٠	•
		•		*			+	•	•	-					-	•	•			-	*		•			•		-	
,		•	*	*	٠			•	•	٠		•	•			٠		٠			٠		٠					•	٠
,			*																										
,		•		٠				•	•	•			•		-	•	•	•					•					•	
,		•				•			•	•					•			٠			•	٠	•					٠	
,		٠		٠		٠	*	٠	*	٠	٠	*	٠		٠	٠	٠	٠	٠	٠	٠	٠	٠	٠		٠		٠	
,		٠			۰	٠	٠	٠	٠	٠					٠	٠	٠	٠	٠	۰		٠	٠	٠				•	۰
,		٠	٠	٠	۰	٠	٠	٠	٠	•	*	*			٠	٠	٠	٠	٠	٠	٠	٠	*	٠	·	٠		٠	٠
		٠	٠		•	٠	*		٠	•	٠	•			٠	۰	۰	٠	•	•	٠	٠	٠	٠	*	٠		*	
		٠	*	٠	٠	٠	٠	٠	•	•	٠	٠	•		٠	٠	٠	۰	٠	•	٠	٠	٠	٠	٠	٠		•	•
		٠	٠	*		٠	*		٠	•	٠	*	•		٠	٠	٠	٠	•	٠	•	٠	*	٠	·	•		٠	
		٠	۰	٠	•	٠	٠	•	•	•	٠	٠	•		•	٠	٠	٠		٠	٠	٠	٠	•		•		•	•
		*	*	*	*	*	*	*	*	•	٠	*	*	*	*	*	*	*	*	*	*	*	*	٠	*	٠		٠	٠
		*	*	*	*	*	*	*	٠	•	٠	٠	•	•	٠	٠	٠	*	•	٠	٠	*	٠	٠	٠	*	·	٠	•
		٠	٠	٠	•	٠	•		٠	•	·	•	•		٠	٠	٠	٠	٠	٠	•	٠	٠	٠	•	•	·	•	•
,	•	•	*	*		•		•		•	٠	•	*		•	•	•	•	*	•	•	*	•	*	•			•	
		•	•	•	•	•	*		•	•	•	•	•		•	•	•	•	•	•	•	•	•	•	·	•		•	•
,		٠	*	*	٠	٠	*	٠	•	•	•	*	٠		٠	٠	٠	*	*	*	٠	٠	•	•	•	•	,	٠	•
,		٠		•	۰	٠	٠	٠	٠	•	٠	•	٠		•	6	۰	٠	•	•	•	٠	٠	٠	•	•		٠	
,		•	٠	٠	۰	•	*	٠	*	*	٠	*	٠		٠	٠		*	*	*	•	٠	*	٠		٠		٠	۰
,	•	٠	٠	٠	۰	٠	٠	٠	•	٠	٠	٠	٠	*	٠	٠	٠	٠	۰	۰	٠	٠	٠	٠	•	٠		٠	•
,		•	۰	٠		•	*	۰	*	*	٠	*	•		•		•	•	•	٠	۰	٠	*	٠	•	•		٠	
		۰	۰	۰	۰	8	٠	•	٠	•	٠	٠	•		٠	۰	٠	۰		۰	٠	٠	*	٠	•	٠		٠	•
		۰	۰		۰	•	۰	*	•	*	*	٠	•		۰	۰	۰	۰	۰	۰		۰	٠	*		۰		٠	
,	•	٠	0	۰	•	٠	*	۰	•	*	۰	*	٠		۰	۰	۰	٠	٠	۰	۰	٠	*	۰	•	٠	•	٠	•
		٠	٠	*	*	٠	*	*	*	•	*	*	*		*	*	*	*	*	*	*	*	*	*		•		*	*
		*	۰	*	*	•	*	•	*	•	•	*	*		•	۰	*	*	*	*	•	٠	٠	•	•	•	•	•	•
	•	•	*	*	*	•	*	•	•	•	•	•	•	•	•	۰	•	•	•	•	•	•	٠	•		•	•	*	•
	•	٠	٠	٠	٠	*	٠	۰	*	•	•	•	٠	•	٠	۰	٠	۰	٠	٠	٠	٠	٠	٠	•	•	'	•	•
									•	-		•											•			•	,	•	
	•	٠	٠	٠		٠	٠	٠	٠	٠		*	٠		٠		٠		٠				٠			٠			

Why Foirness? Types of bias and bias loop

· ·

. .

. .

· ·

. .

• •

· ·

. .

• •

· ·

. .

• •

• •

2



01

What is Fairness? Fairness definitions



How to ensure Foirness in ML? Bias mitigation and fair algorithms



Where to apply it? Use cases on Fairness

01 Why Fairness?

																																				×	+	\leftarrow
																																	×	\uparrow		40	¥.	
																														•	1	1		•			· /	$\overline{\neg}$
																														• /			\succ	4			1	• `\
																						~							• /	1.		Ń		•		X		• •
																				•	1								×	<.				• /	\mapsto		K	
																							×	R				×			1		>	•		\mathbf{z}^{\dagger}	• \}	\mathbf{N}
																•	5		7	Ś				$ \rightarrow $	-	~				· >*					. ¥	→	•	\mathbf{a}
																		. ~	γ.		Y			6			Y		1		- 64		- 0 1					<u> </u>

Why Fairness?

2	2
	5



An algorithm used in USA courts, attributes a higher criminal risk on black defendants. [1]

Automatic sentiment analysis rates sentences as "I'm homosexual" with a negative score. [2] Automatic translations have gender bias. When translating "She works in an Hospital, my friend is a doctor" from English to Portuguese, it assumes the doctor is a man. [3] SE

E-commerce website was displaying different prices depending on people localization and distance from a rival store. [4]



 [1] https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing
 [2] https://mashable.com/2017/10/25/google-machine-learning-bias/?europe=true
 [3] https://tanslate.google.com/#view=home&op=translate&sl=en&tl=pt&text=She%20works%20in%20an%20Hospital%2C%20my%20friend%20is%20a%20nurse.
 202

 [4] https://tanslate.google.com/#view=home&op=translate&sl=en&tl=pt&text=She%20works%20in%20an%20Hospital%2C%20my%20friend%20is%20a%20nurse.
 202

2020 nilg.ai. All rights reserved

Types of Bias

Historical Bias

Bias present in our history that still affect the current systems. e.g., gender and racial bias.





Observer Bias

It happens when the observer projects his/her expectations onto the problem, picking the solution that best suits the initial goal, instead of suiting the real data.





Population Bias

This bias results when the statistics, demographics, representatives, and user characteristics represented in the training dataset differs from the population on the test set.

Measurement Bias

It occurs when unfair features are used as proxy labels. e.g., use family criminal records as a proxy to measure the level of "individual crime risk".





Bias Feedback loop



7



What is Fairness?



What is Fairness?

Def 1. Equalized Odds

The protected and unprotected groups should have equal rates for true positives and false positives

Def 2. Equal Opportunity

The protected and unprotected groups should have equal true positive rates

Def 3. Fairness Through Awareness

Any two individuals who are similar with respect to a similarity (inverse distance) metric defined for a particular task should receive a similar outcome

Def 4. Demographic Parity Def 5. Fairness Through Unawareness Def 6. Treatment Equality

Def 7. Test Fairness

Def 8. Counterfactual Fairness Def 9. Fairness in Relational Domains



Mehrabi, Ninareh, et al. "A survey on bias and fairness in machine learning." arXiv preprint arXiv:1908.09635 (2019).

03 How to ensure Fairness in ML?



Mitigating Unfairness in the AI framework





Mitigating Unfairness in the AI framework



Data Collection

Rethink your data collection process, aiming representative datasets and avoiding unfair biased features.



Mitigating Unfairness in the AI framework



Data Collection

Rethink your data collection process, aiming representative datasets and avoiding unfair biased features.

Pre-processing

Try to transform the data, removing underlying discrimination and unbalanced representations.



13

Mitigating Unfairness in the AI framework



Rethink your data collection process, aiming representative datasets and avoiding unfair biased features. Try to transform the data, removing underlying discrimination unbalanced representations. Modify the algorithms in order to remove discrimination during the training process. E.g. Adding Fairness regularization.



Mehrabi, Ninareh, et al. "A survey on bias and fairness in machine learning." arXiv preprint arXiv:1908.09635 (2019).

Mitigating Unfairness in the AI framework



Data Collection

Rethink your data collection process, aiming representative datasets and avoiding unfair biased features.

Pre-processing

Try to transform the data, removing underlying discrimination unbalanced representations.

In-processing

Modify the algorithms in order to remove discrimination during the training process. E.g. Adding Fairness regularization.

Post-processing

Use Fair metrics to choose the best model; Use predictions from a "black-box" model to feed a Fair algorithm; Use the predictions from the model and apply fairness in the decision process.

Mehrabi, Ninareh, et al. "A survey on bias and fairness in machine learning." arXiv preprint arXiv:1908.09635 (2019).

Adversarial Debiasing Network



- 1. Given X predicts Y
- 2. The dense layer is passed to an adversary network
- 3. Adversary network predicts Z (protected variable) given Y and \hat{Y}
- 4. Goal: minimizing the accuracy of the Adversary while maximizing Predictor accuracy









Zhang, Brian Hu, Blake Lemoine, and Margaret Mitchell. "Mitigating unwanted biases with adversarial learning." Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society. 2018.

FairGAN



- Use Generative Adversarial Networks (GAN) to generate artificial data, PG, adding noise to real data;
- 2. Use D1 Discriminator to predict if the data is real or fake;
- 3. Use D2 Discriminator to predict if the data is "protected" or not;
- 4. Train the GAN in order to fool the discriminators;
- 5. Use the generated debiased data to train your model







Depeng, et al. "Fairgan: Fairness-aware generative adversarial networks." 2018 IEEE International Conference on Big Data (Big Data). IEEE, 2018..

17

Adversarial SHAP - Methods

- 1. Use SHapley Additive Explanations (SHAP) to measure the relevance of the "protected" features, Z, with respect to the label, Y;
- 2. Define "Fairness by Explicability" metrics:
 - a. FE: difference in mean attribution of Z between the "protected" and "unprotected" groups
 - b. SFE: total attribution of Z across the population
- 3. Train the model using Fairness regularization in the Loss function









Adversarial SHAP - Results

From λ =0 to λ =0.9, AUC and accuracy drop 0.04







Where to apply it?



Literature Use Case

Gender Bias in Education Recommendation





Yao, Sirui, and Bert Huang, "Beyond parity: Fairness objectives for collaborative filtering." Advances in Neural Information Processing Systems. 2017...

Literature Use Case

Gender Bias in Education Recommendation

Value Unfairness

Measures inconsistency in signed prediction error across the user type.

Absolute Unfairness

Similar to Value Unfairness but does not consider the direction of the error

Underestimation Unfairness

Measures inconsistency in how much the predictions underestimate the true ratings.

Overestimation Unfairness

Measures inconsistency in how much the predictions overestimate the true ratings.

Non-parity unfairness

Measures the absolute difference between the overall average ratings of disadvantaged users and advantaged ones.



Yao, Sirui, and Bert Huang. "Beyond parity: Fairness objectives for collaborative filtering." Advances in Neural Information Processing Systems. 2017..

Literature Use Case

Gender Bias in Education Recommendation



Conclusions:

- Regularizing on one of the metrics tends to decrease the other fairness metrics;
- Decreasing fairness metrics (decreasing unfairness) does minimized the prediction error;
- Regularization on the Value Unfairness was the most effective;

ao, Sirui, and Bert Huang. "Beyond parity: Fairness objectives for collaborative filtering." Advances in Neural Information Processing Systems. 2017...

NILG.AI Use Cases

Fairness by Construction:



- Project Goal:
 - a. AI models that promote fair access to high-quality healthcare
 - b. Data quality to prevent improper conclusions when decisions are provided by entry-level personnel
 - c. See more at: <u>https://tinyurl.com/modtsympbio</u>



- Data collection:
 - a. Data collection that ensures representation of main protected groups:
 - i. Gender, race
 - ii. Others: staff skills, tatoos, room conditions
 - b. See more on "Innovation in Medical Device Decontamination" at: <u>https://tristelopenday.com/</u>



NILG.AI Use Cases

Group invariance:

race, gender, country, deep-learning-framework-preference



I NANK YOU TOF YOUF ATTENTION!	
I NANK YOU FOR YOUR ATTENTION!	
I NANK YOU FOR YOUR ATTENTION!	
I NANK YOU FOR YOUR ATTENTION!	
I NANK YOU FOR YOUR ATTENTION!	
I NANK YOU FOR YOUR ATTENTION!	
I NANK YOU FOR YOUR ATTENTION!	
I NANK YOU FOR YOUR ATTENTION!	
I NANK YOU FOR YOUR ATTENTION!	
I NANK YOU FOR YOUR ATTENTION!	
I NANK YOU FOR YOUR ATTENTION!	
I NANK YOU FOR YOUR ATTENTION!	
I NANK YOU FOR YOUR ATTENTION!	
I NANK YOU TOF YOUF Attention!	
I NANK YOU FOR YOUR ATTENTION!	
I NANK YOU FOR YOUR ATTENTION!	
I NANK YOU TOF YOUF ATTENTION!	