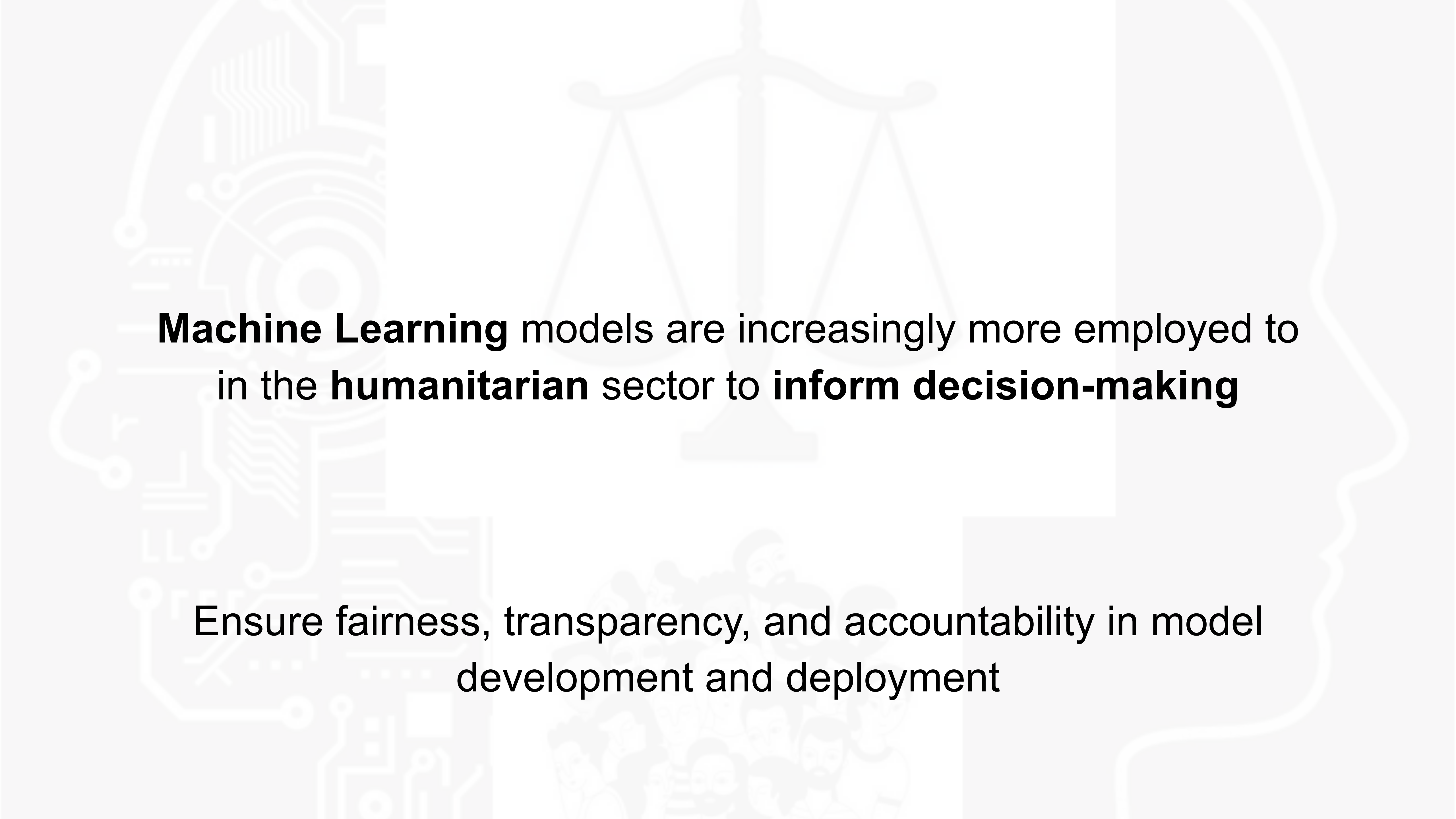


Fairness in vulnerable attribute prediction on social media

Mariano Beiro & Kyriaki Kalimeri

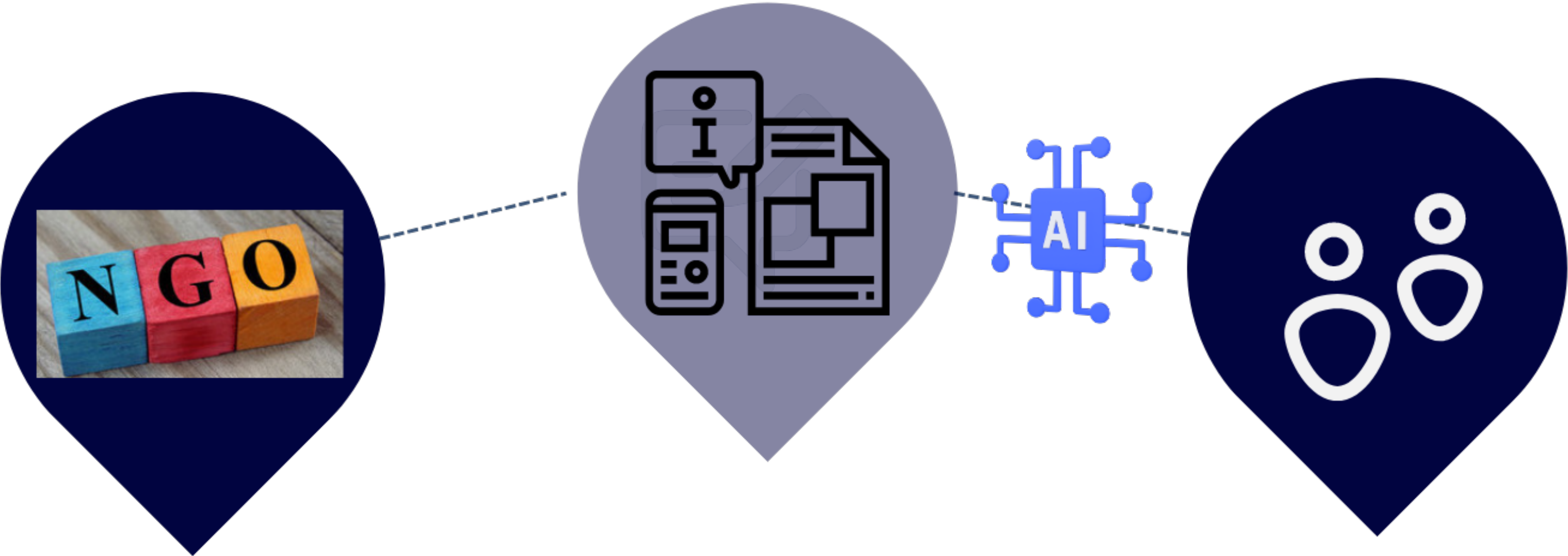




Machine Learning models are increasingly more employed to
in the **humanitarian** sector to **inform decision-making**

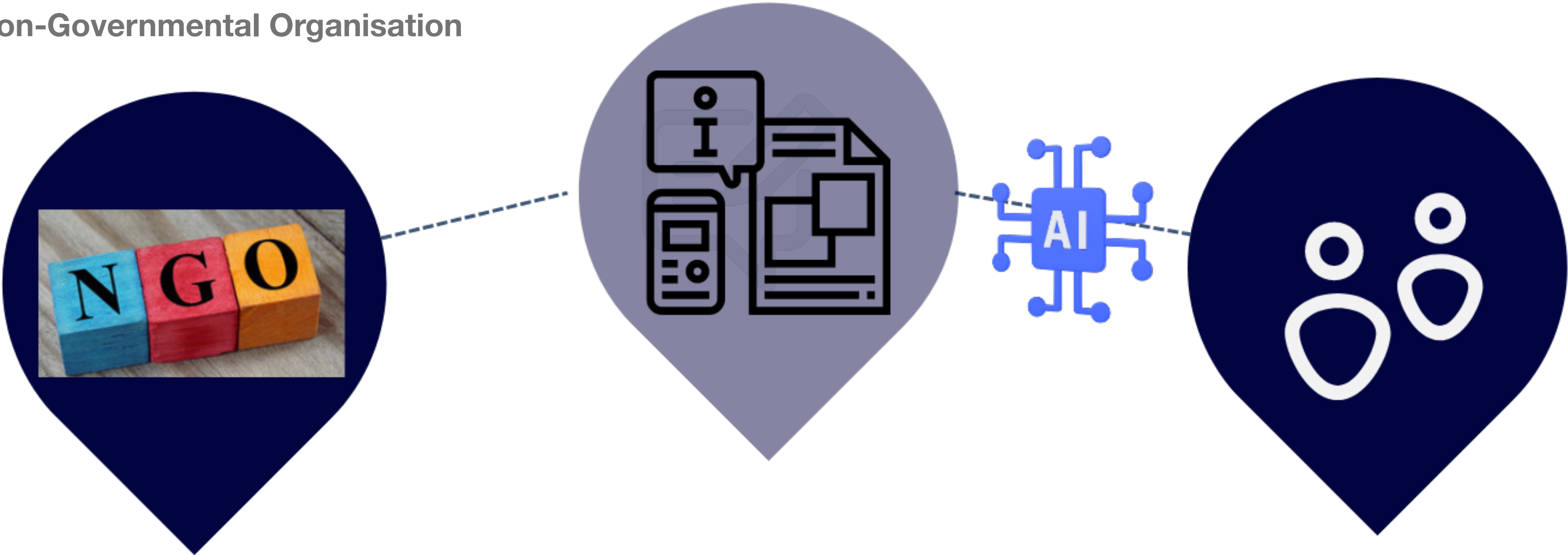
Ensure fairness, transparency, and accountability in model
development and deployment

Our Case Study



Our Case Study

A Non-Governmental Organisation

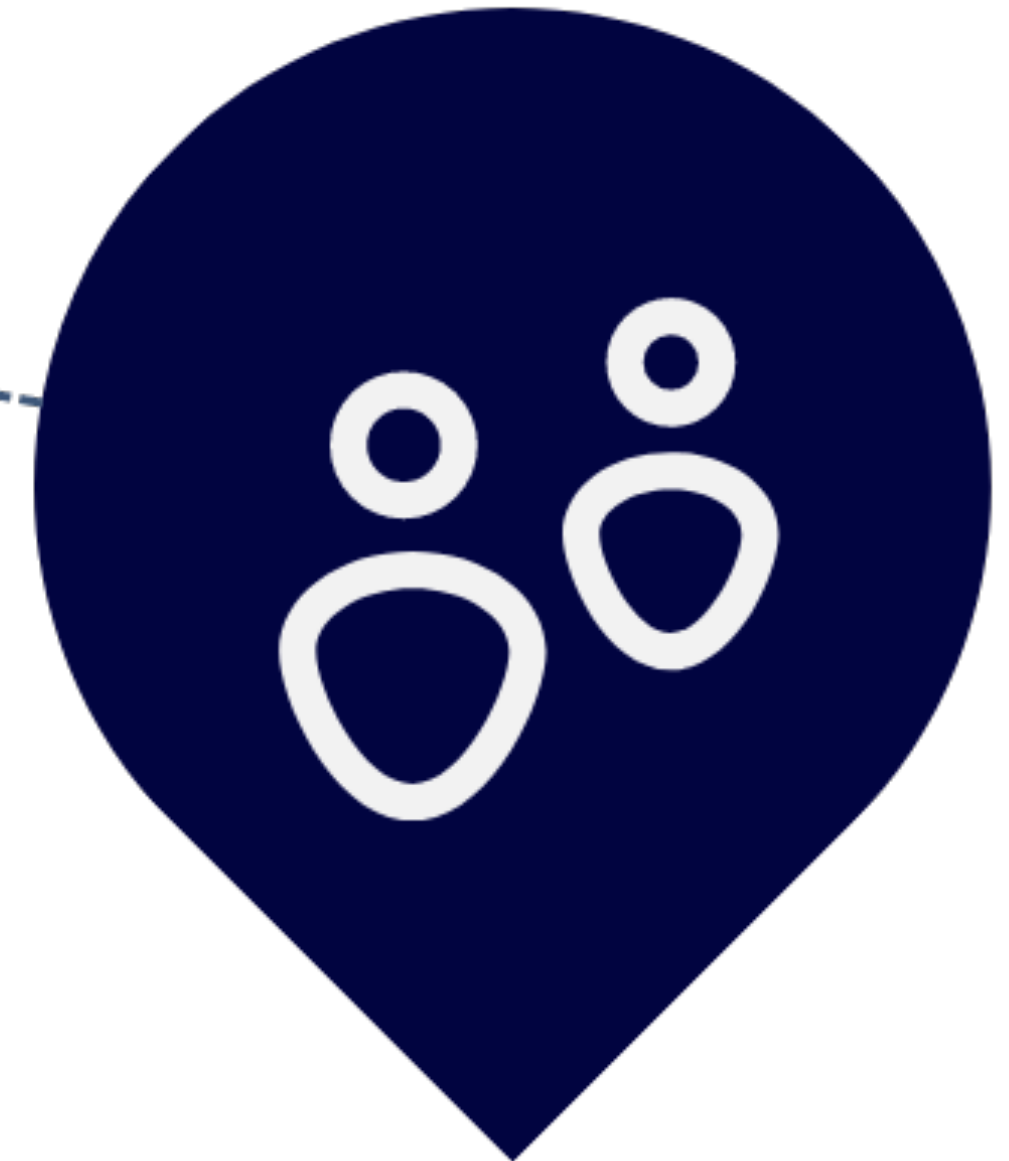
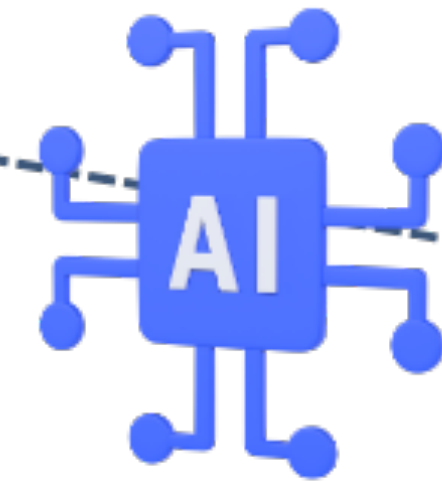


Our Case Study

A Non-Governmental Organisation



aimed at providing
educational/training
opportunities

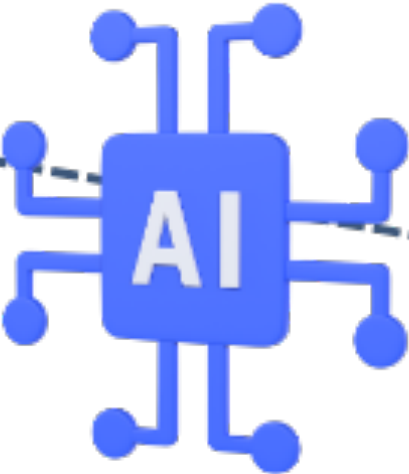


Our Case Study

A Non-Governmental Organisation

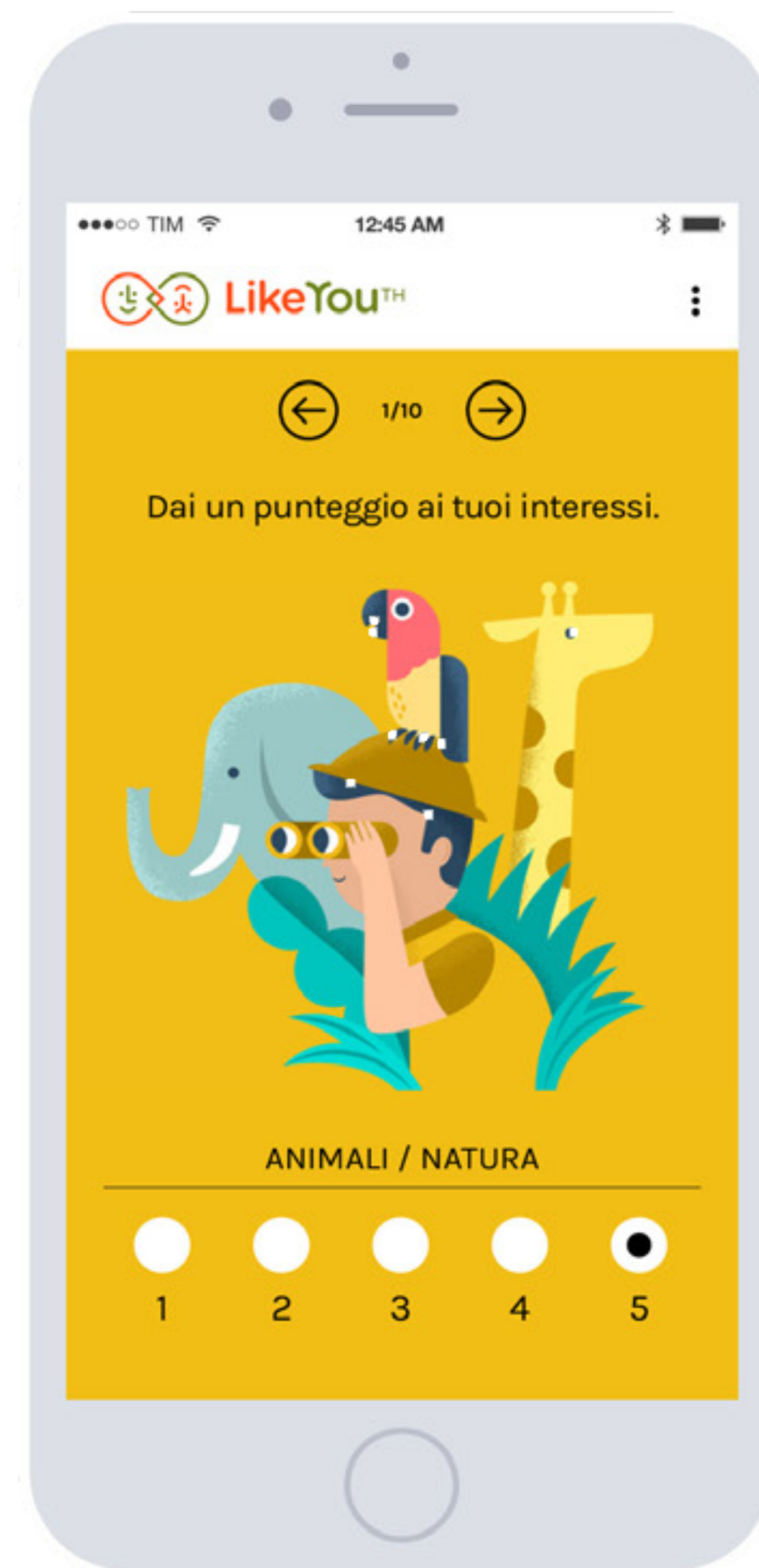


aimed at providing educational/training opportunities



to young unemployed via AI



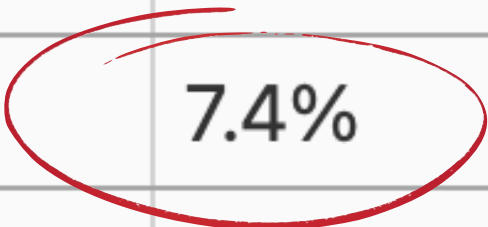


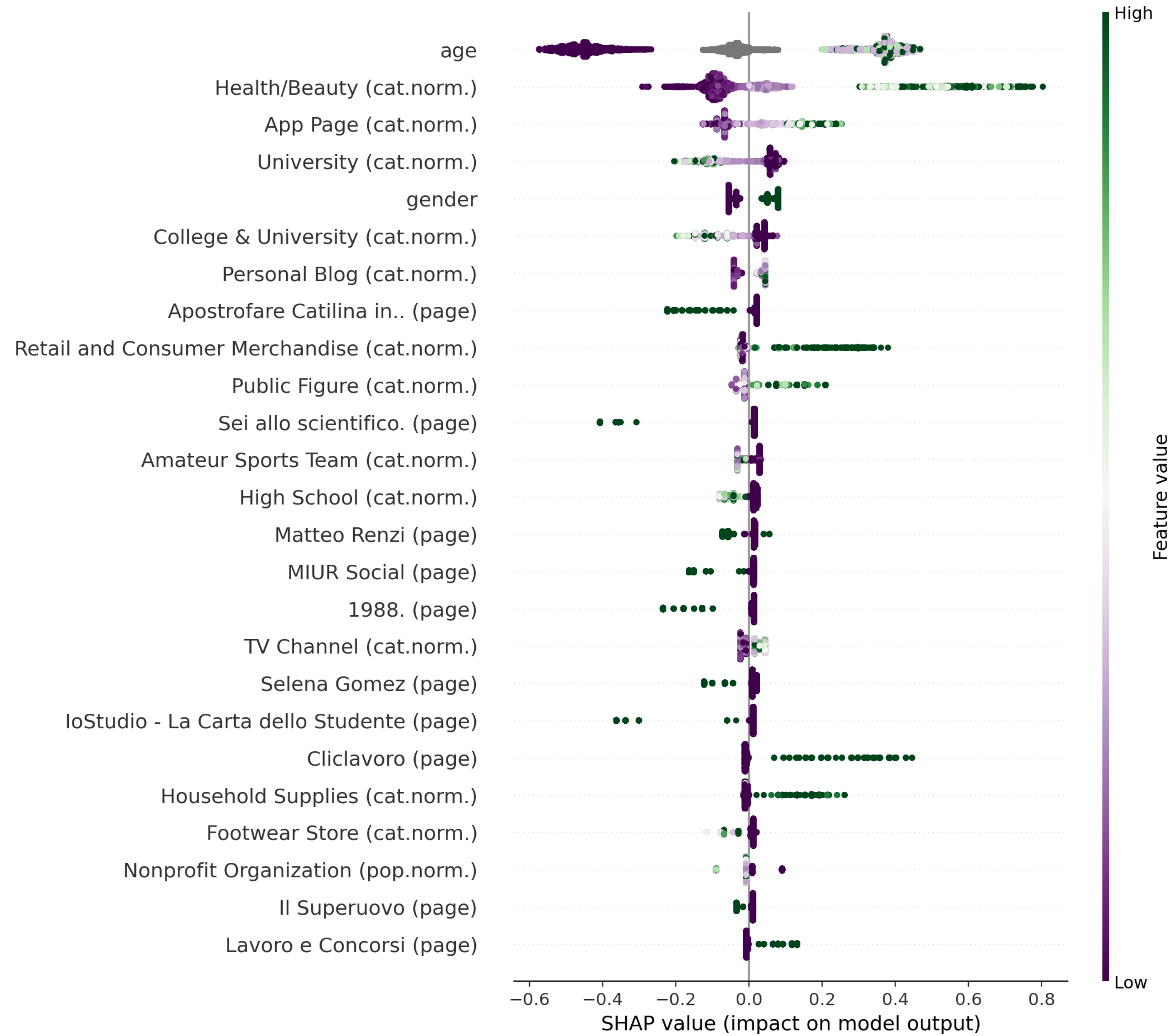
	Census	Dataset
		<i>n</i> = 11,393
<i>Gender</i>		
Female	51.1%	38.1%
Male	48.4%	61.8%
<i>Age</i>		
17–24	7.9%	43.1%
25–34	11.0%	31.2%
35–44	13.8%	13.6%
45–54	16.1%	7.1%
55–64	13.3%	4.5%
65+	24.5%	0.3%
<i>Occupation</i>		
Employed	77%	43.9%
Unemployed	8.7%	7.4%
Student	14.2%	48.5%

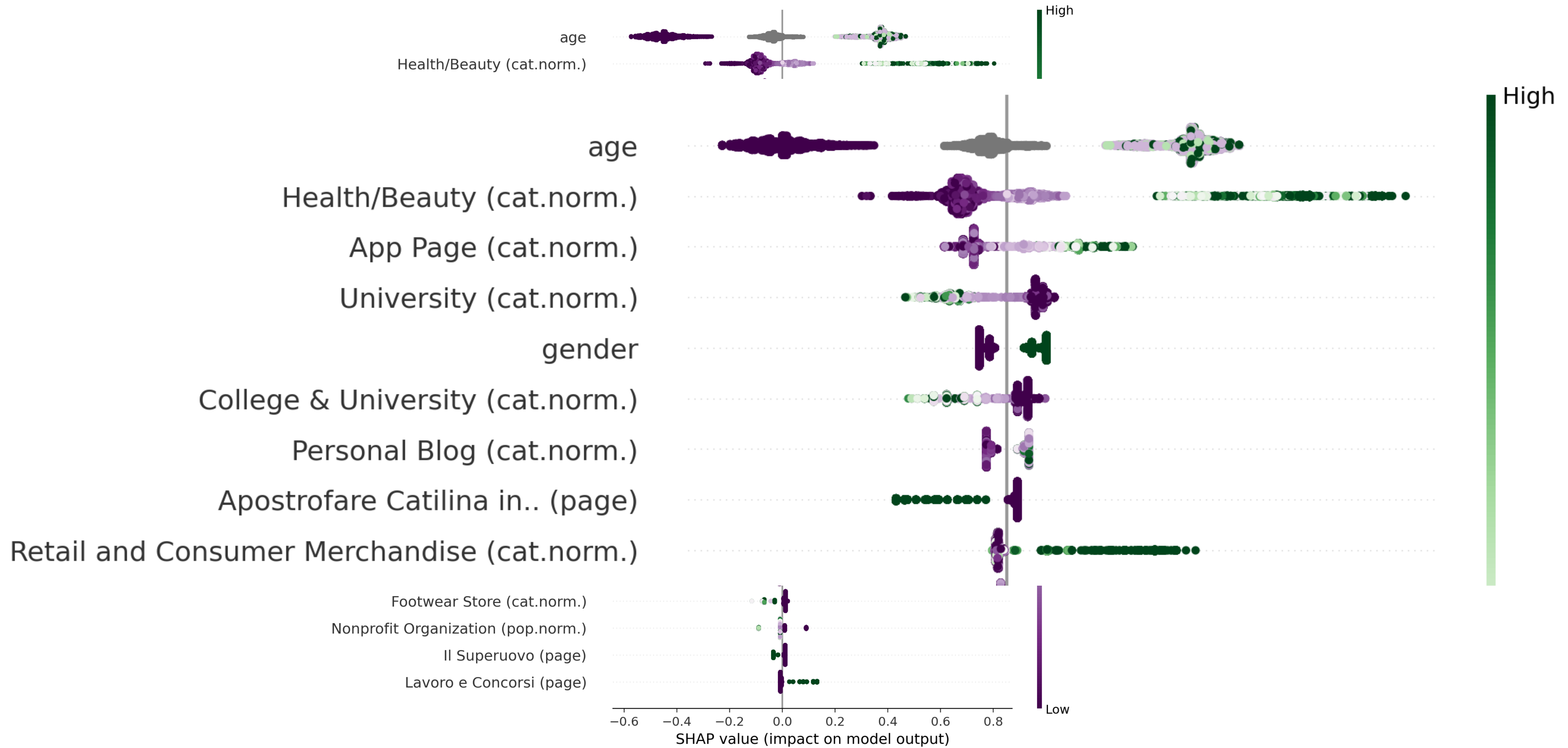
	Census	Dataset
		<i>n</i> = 11,393
<i>Gender</i>		
Female	51.1%	38.1%
Male	48.4%	61.8%
<i>Age</i>		
17–24	7.9%	43.1%
25–34	11.0%	31.2%
35–44	13.8%	13.6%
45–54	16.1%	7.1%
55–64	13.3%	4.5%
65+	24.5%	0.3%
<i>Occupation</i>		
Employed	77%	43.9%
Unemployed	8.7%	7.4%
Student	14.2%	48.5%

	Census	Dataset
		<i>n</i> = 11,393
<i>Gender</i>		
Female	51.1%	38.1%
Male	48.4%	61.8%
<i>Age</i>		
17–24	10.1%	43.1%
25–34	10.1%	31.2%
35–44	10.1%	13.6%
45–54	10.1%	7.1%
55–64	13.3%	4.5%
65+	24.5%	0.3%
<i>Occupation</i>		
Employed	77%	43.9%
Unemployed	8.7%	7.4%
Student	14.2%	48.5%

Unemployment rate per Gender:
 Male: 5.5%
 Female: 9%









Well we can predict unemployment with 74% AUROC. Cool!

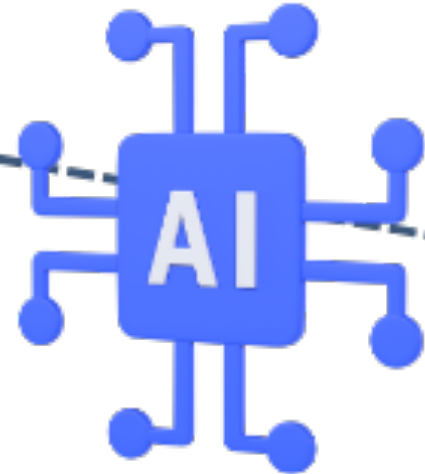


Our Case Study

A Non-Governmental Organisation



aimed at providing educational/training opportunities



to young unemployed via AI



Our Case Study

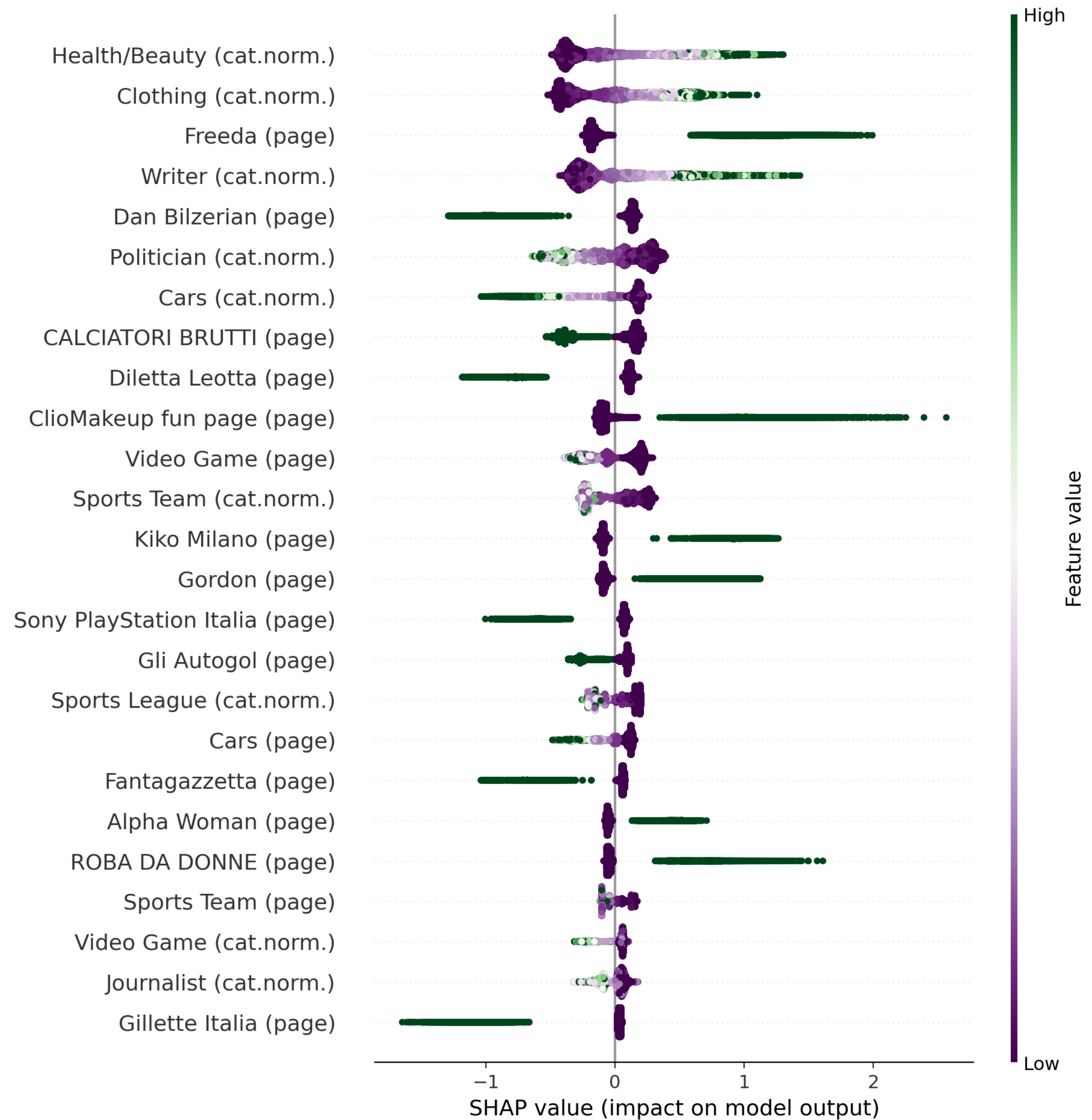
A Non-Governmental Organisation

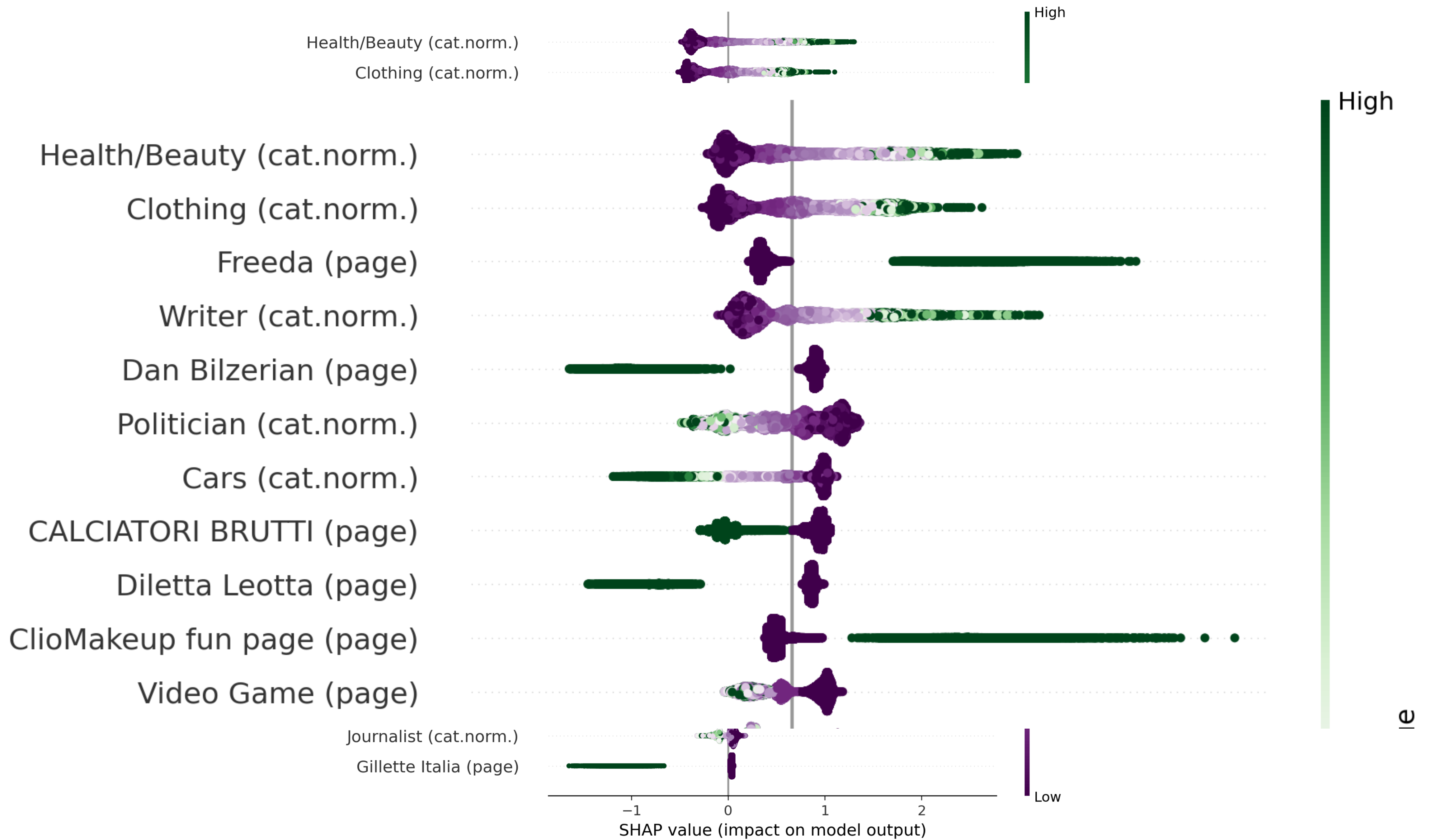
to young unemployed via AI

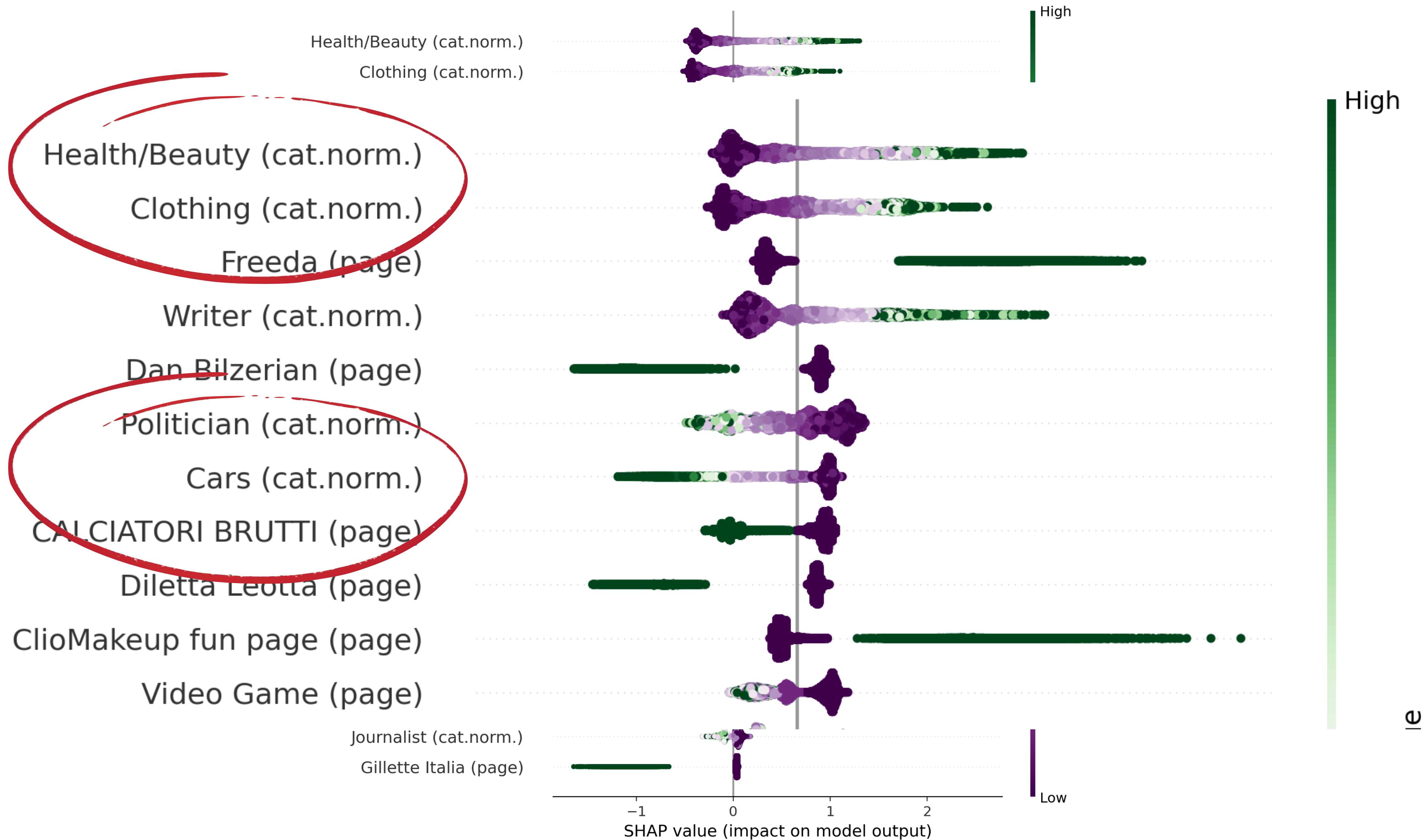
Can we ensure that the machine learning model is fair?

Can we avoid discrimination?

aimed at providing
educational/training
opportunities

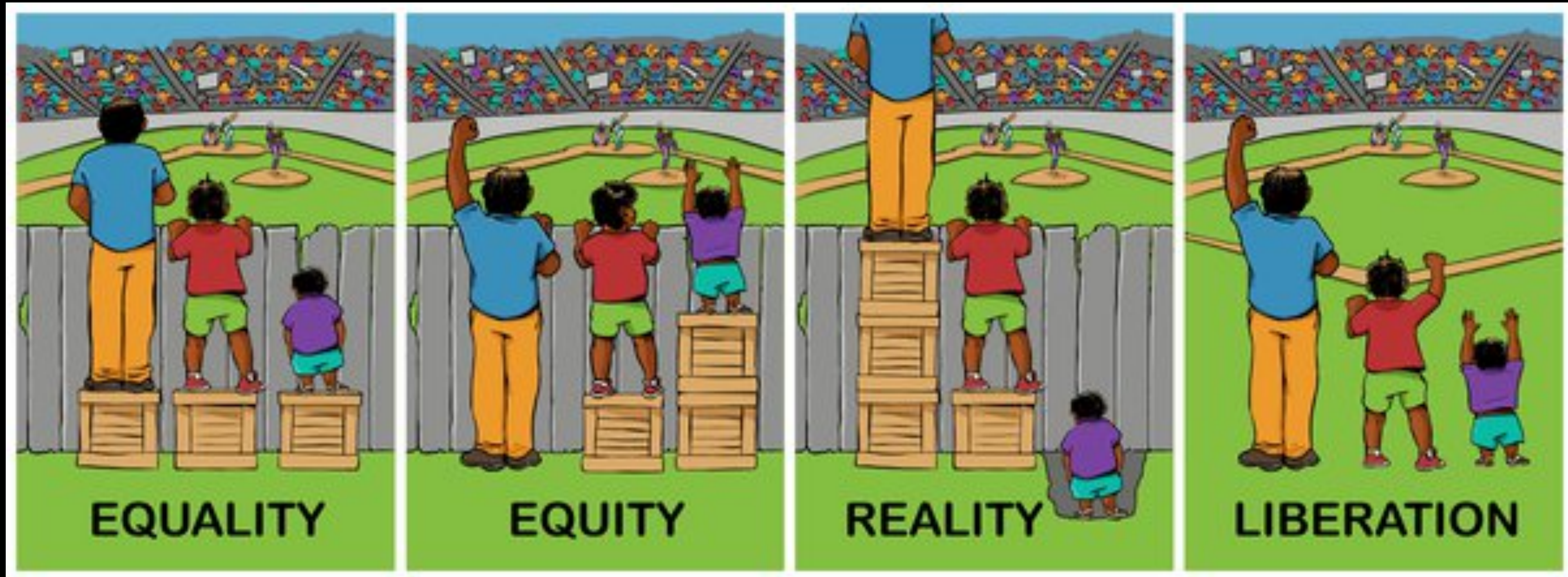




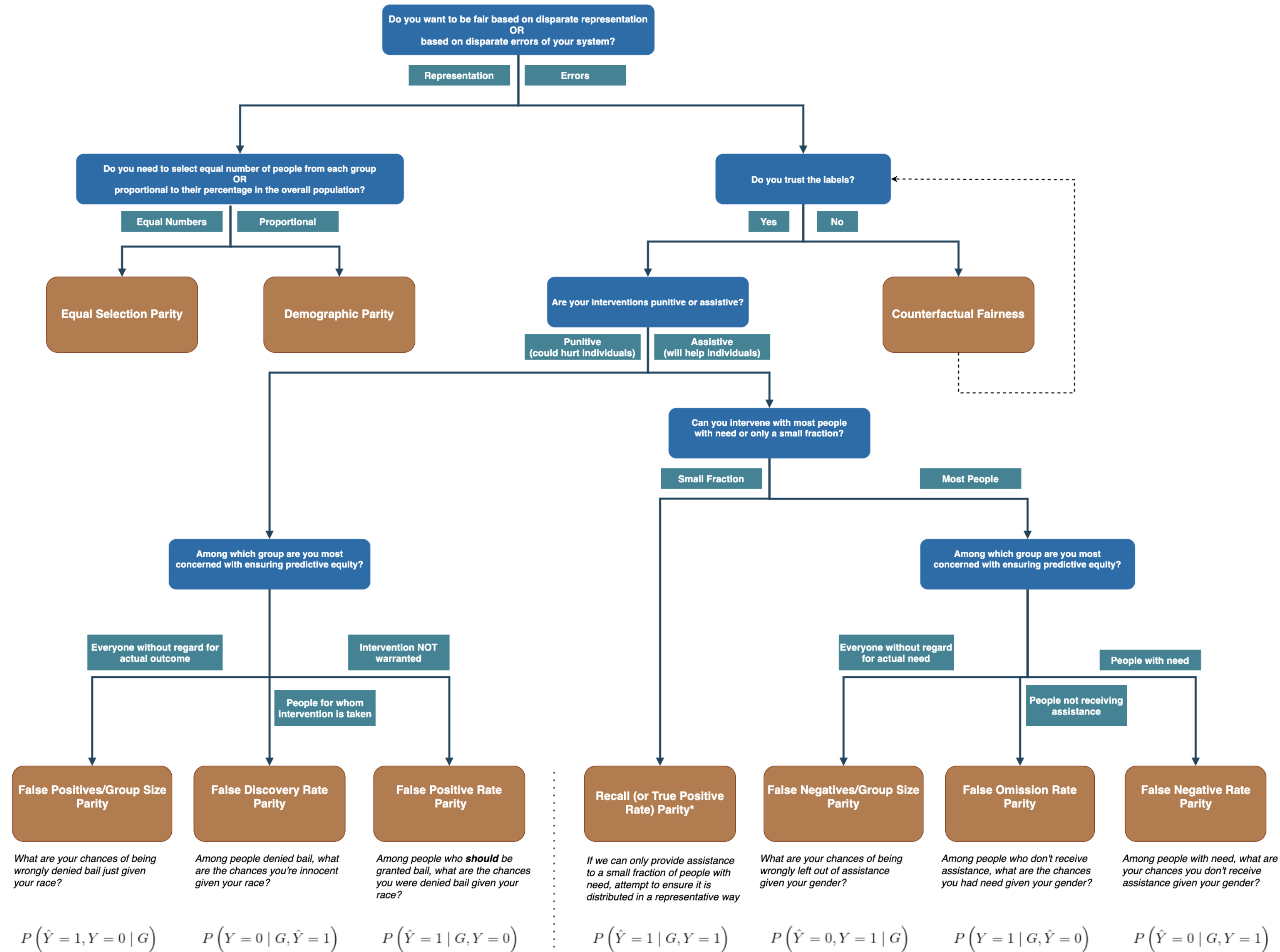


With Fairness Through Unawareness we do not avoid discrimination

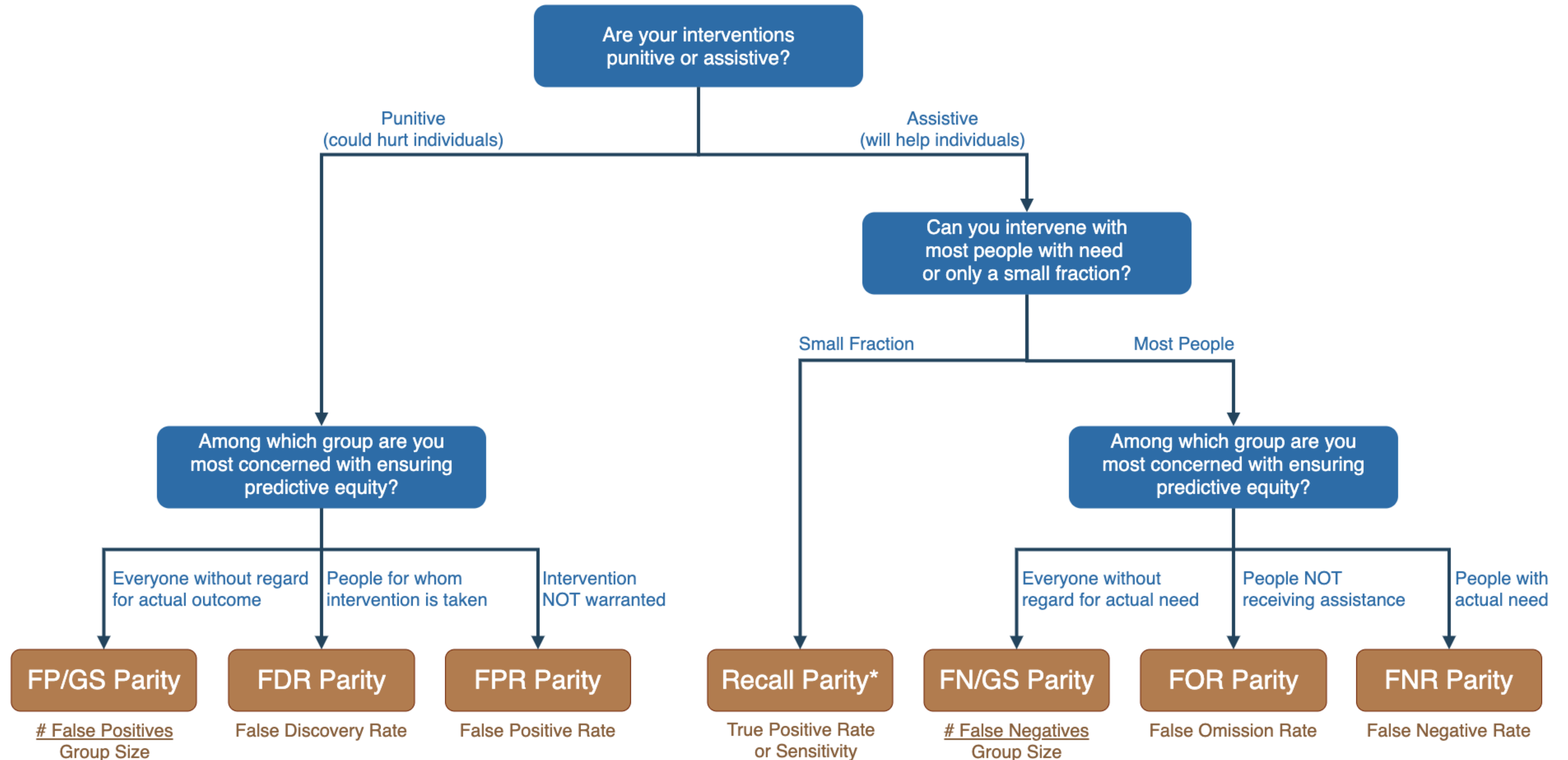
What is "Fairness"?



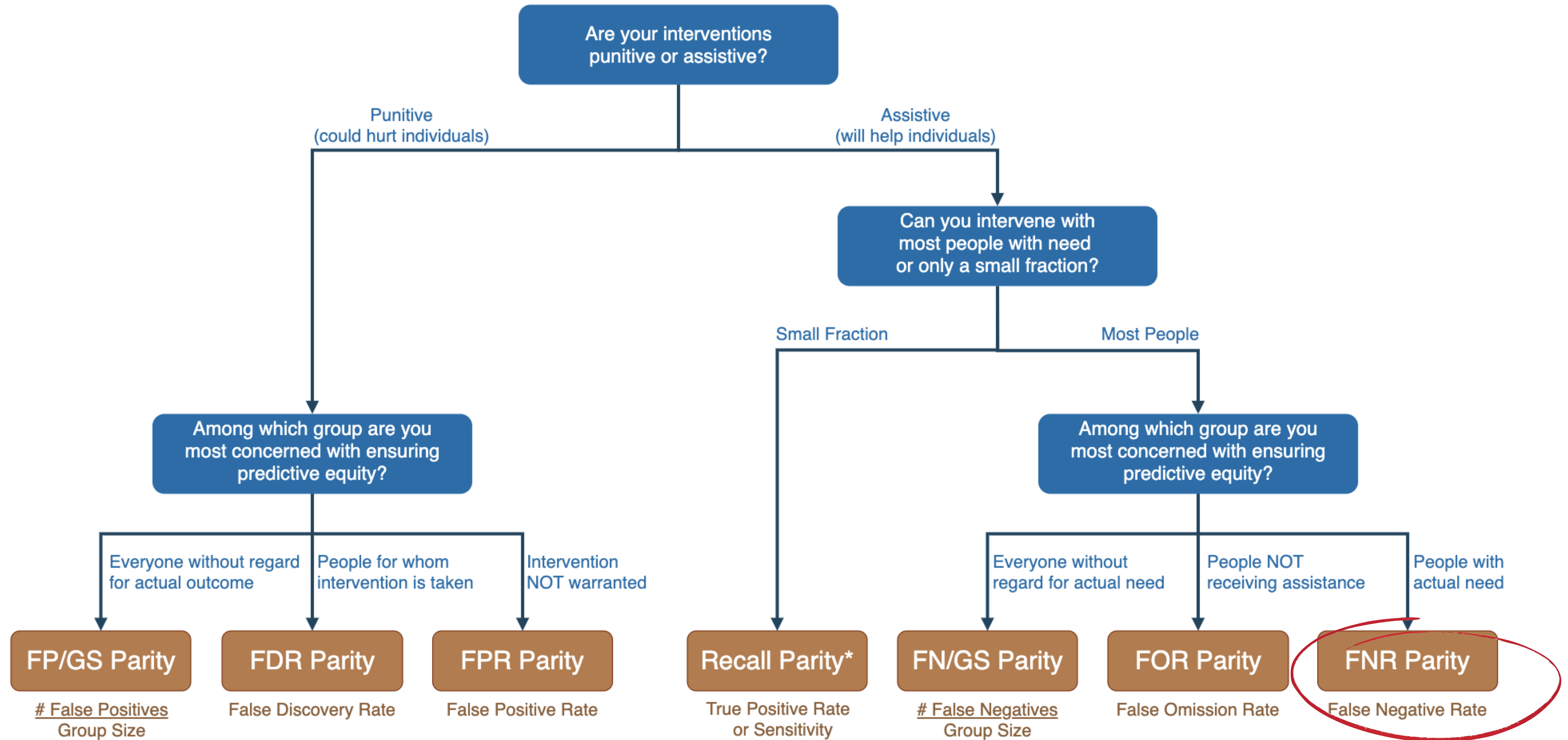
FAIRNESS TREE



FAIRNESS TREE (Zoomed in)



FAIRNESS TREE (Zoomed in)



FAIRNESS TREE (Zoomed in)

Are your interventions punitive or assistive?

Punitive
(could hurt individuals)

Assistive
(will help individuals)

Impossible to satisfy more than one fairness metrics at once
“Fairness and Machine Learning”, S. Barocas, M. Hardt, A. Narayanan, 2022, <https://fairmlbook.org/>

Among which group are you most concerned with ensuring predictive equity?

Everyone without regard for actual outcome

People for whom intervention is taken

Intervention NOT warranted

FP/GS Parity

$\frac{\# \text{ False Positives}}{\text{Group Size}}$

FDR Parity

False Discovery Rate

FPR Parity

False Positive Rate

Recall Parity*

True Positive Rate or Sensitivity

Among which group are you most concerned with ensuring predictive equity?

Everyone without regard for actual need

People NOT receiving assistance

People with actual need

FN/GS Parity

$\frac{\# \text{ False Negatives}}{\text{Group Size}}$

FOR Parity

False Omission Rate

FNR Parity

False Negative Rate

Parity of Opportunity

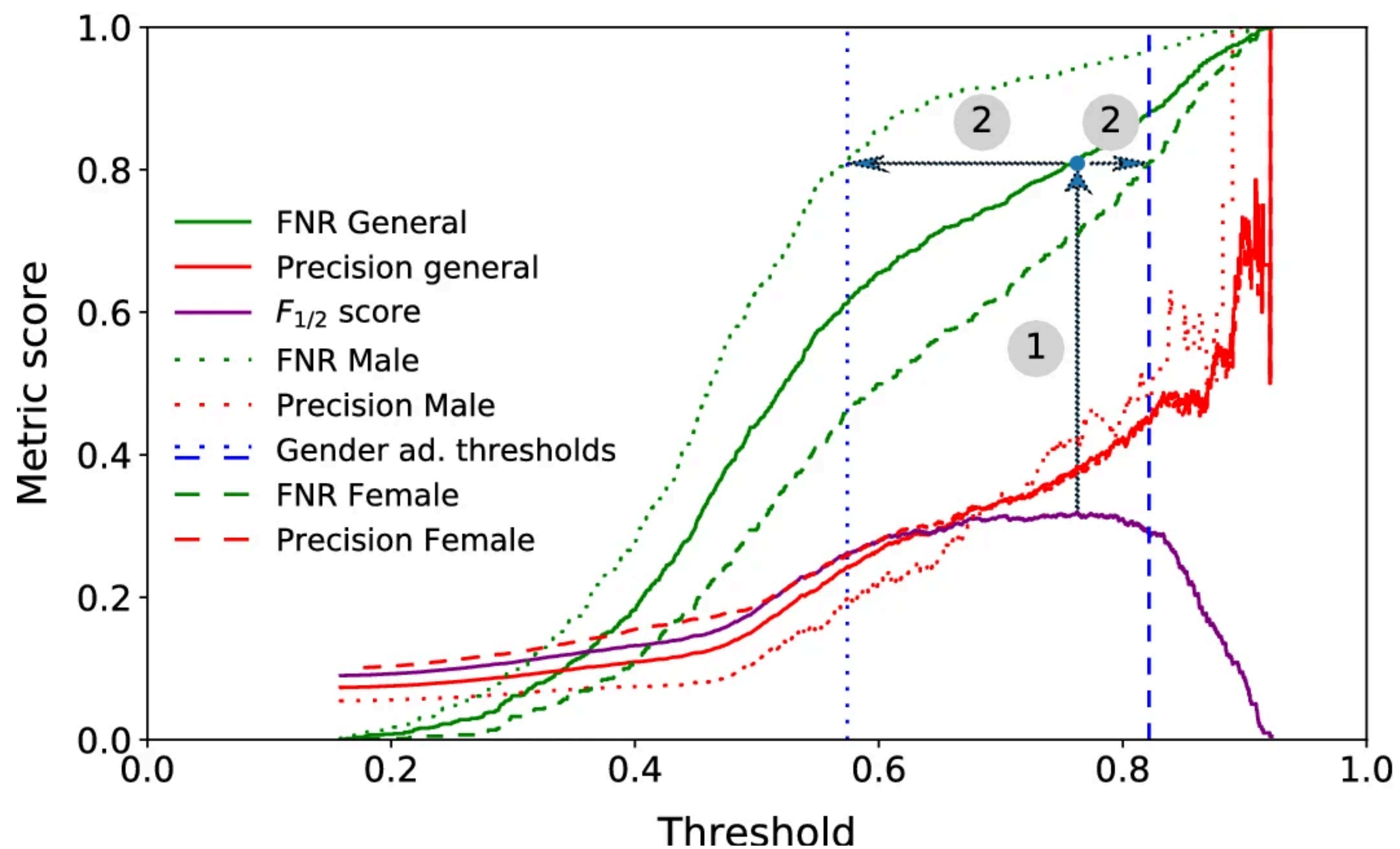
$$\begin{aligned}\text{FNR}_g \text{ disp.} &= \frac{\text{FNR}_g}{\text{FNR}_{ref.group}} \\ &= \frac{\text{Pr}[\hat{Y}=0|Y=1 \wedge G=g]}{\text{Pr}[\hat{Y}=0|Y=1 \wedge G=ref.group]}\end{aligned}$$

where Y and \hat{Y} represent the real and predicted target values respectively (1 represents the ‘unemployed’, 0 the employed)

Disparity threshold 80% with respect to a reference group

Adaptive Threshold

$$F_{\beta} = (1 + \beta^2) \cdot \frac{\text{precision} \cdot \text{recall}}{(\beta^2 \cdot \text{precision}) + \text{recall}} \quad \text{with } \beta \text{ at } 0.5 \text{ to favour precision}$$



	Demo	NoDemo	Demo+AT.	NoDemo+AT.
<i>Global Accuracy (Metric: AUC(std))</i>				
Baseline	.50	.50	.50	.50
State of the Art	–	.61(.01) (*)	–	–
Our Approach	.74(.02)	.71(.02)	.74(.02)	.71(.02)
<i>Precision and Recall</i>				
Precision	.16(.02)	.18(.01)	.26(.05)	.25(.03)
Recall	.56(.05)	.48(.02)	.21(.05)	.22(.04)
<i>Demographic accuracy (Metric: AUC(std))</i>				
Gender (M)	.66(.05)	.64(.04)	.66(.05)	.64(.04)
Gender (F)	.78(.02)	.76(.02)	.78(.02)	.76(.02)
Age (17-24)	.70(.08)	.69(.08)	.70(.08)	.69(.08)
Age (25-34)	.66(.05)	.65(.05)	.66(.05)	.65(.05)
Age (35-44)	.74(.09)	.73(.08)	.74(.09)	.73(.08)
Age (45-54)	.61(.17)	.54(.16)	.61(.17)	.54(.16)
Age (55+)	.46(.31)	.46(.29)	.46(.31)	.46(.29)
<i>Fairness (Metric: $\frac{FNR}{FNR_{ref}}$)</i>				
Gender (ref.class: Male)				
Female	.47(.11)	.58(.14)	1.0(.07)	1.02(.14)
Age (ref.class: 17–24)				
25-34	.35(.08)	.62(.12)	.75(.09)	.80(.08)
35-44	.26(.12)	.49(.2)	.71(.09)	.73(.1)
45-54	.41(.24)	.82(.35)	.82(.17)	.84(.19)
55+	.59(.36)	.99(.42)	.82(.18)	.91(.19)

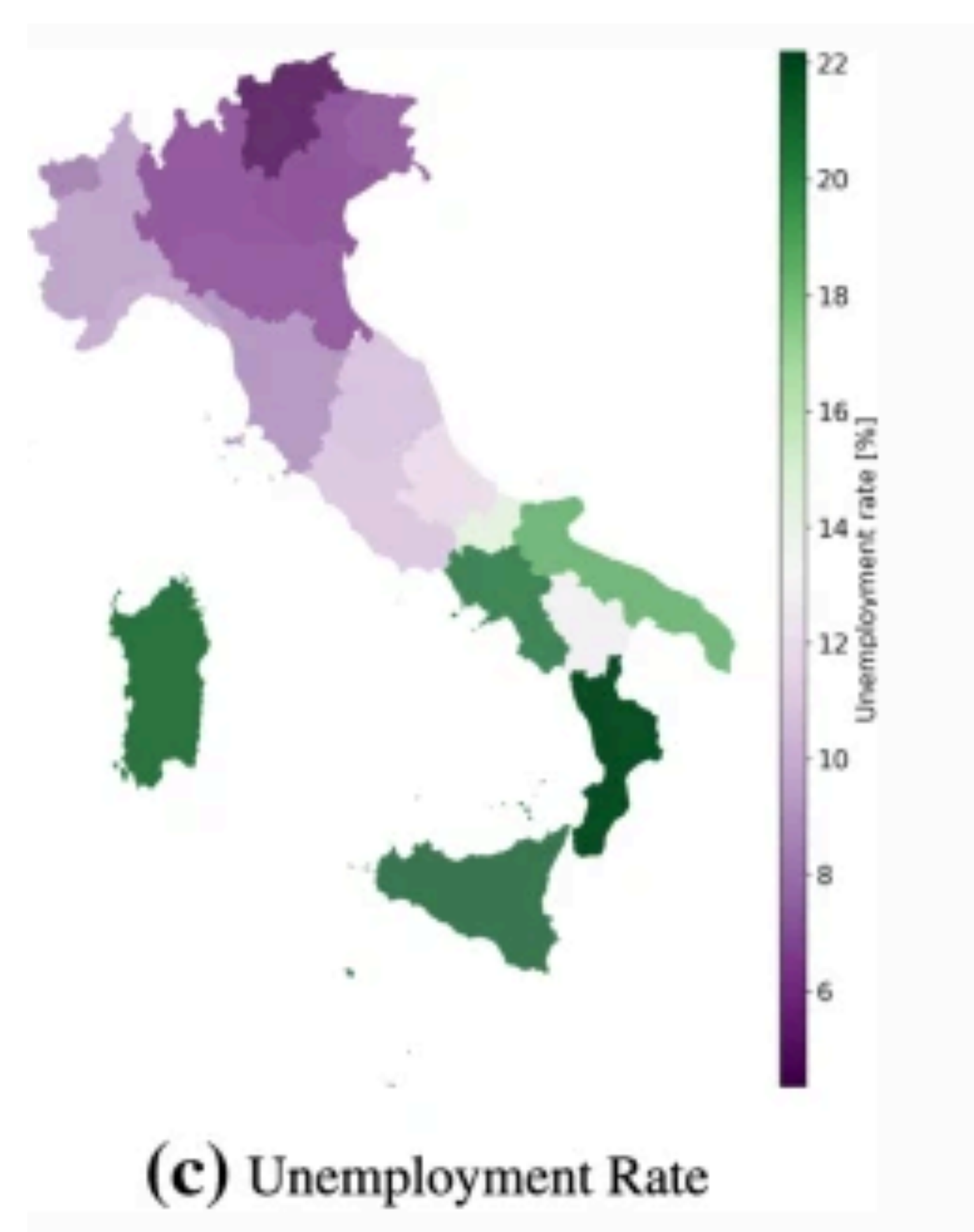
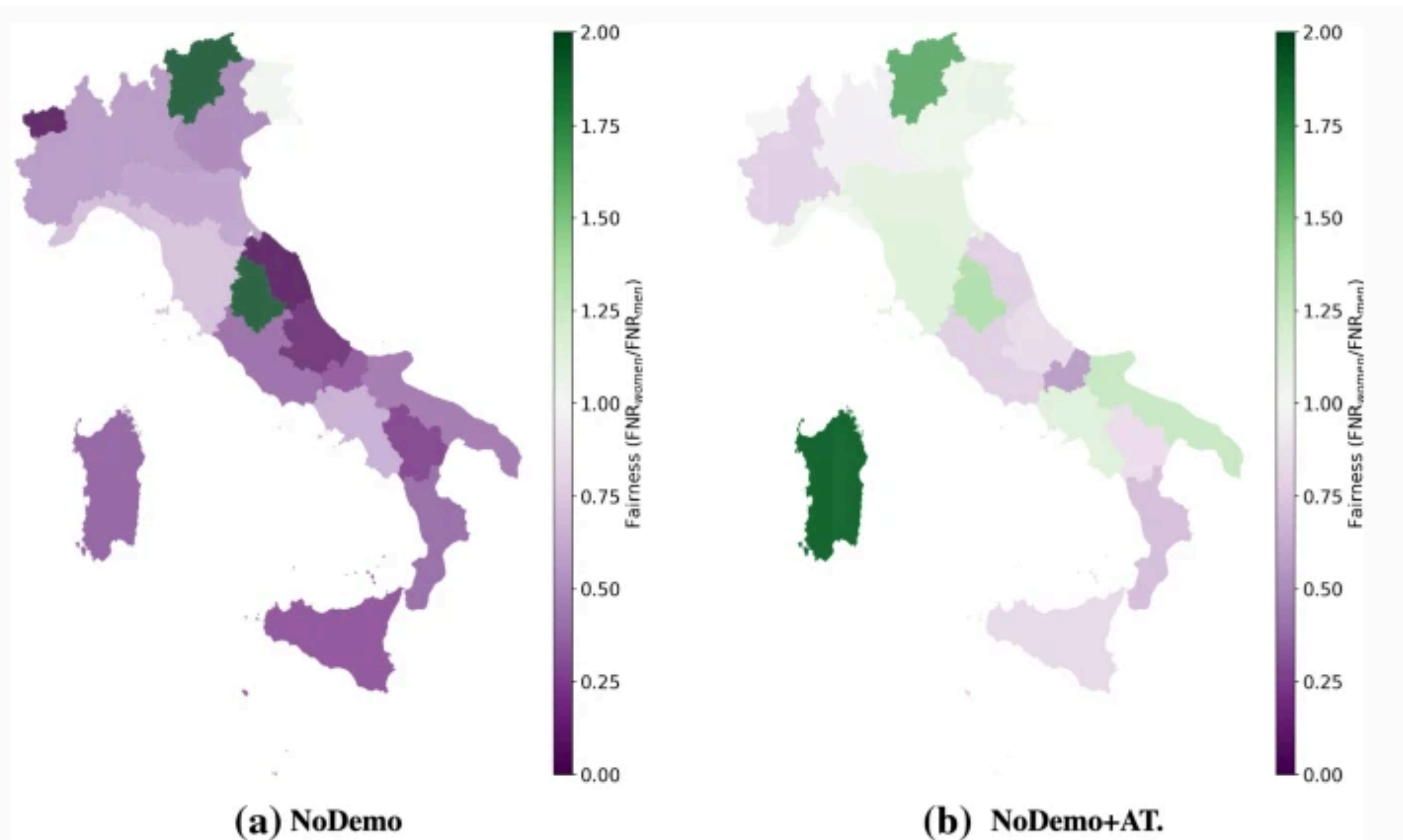
	Demo	NoDemo	Demo+AT.	NoDemo+AT.
<i>Global Accuracy (Metric: AUC(std))</i>				
Baseline	.50	.50	.50	.50
State of the Art	–	.61(.01) (*)	–	–
Our Approach	.74(.02)	.71(.02)	.74(.02)	.71(.02)
<i>Precision and Recall</i>				
Precision	.16(.02)	.18(.01)	.26(.05)	.25(.03)
Recall	.56(.05)	.48(.02)	.21(.05)	.22(.04)
<i>Demographic accuracy (Metric: AUC(std))</i>				
Gender (M)	.66(.05)	.64(.04)	.66(.05)	.64(.04)
Gender (F)	.78(.02)	.76(.02)	.78(.02)	.76(.02)

Gender (ref.class: Male)

Female	.47(.11)	.58(.14)	1.0(.07)	1.02(.14)
--------	-----------------	-----------------	-----------------	------------------

Age (18–24)	.51(.17)	.51(.16)	.51(.17)	.51(.16)
Age (55+)	.46(.31)	.46(.29)	.46(.31)	.46(.29)
<i>Fairness (Metric: $\frac{FNR}{FNR_{ref}}$)</i>				
Gender (ref.class: Male)				
Female	.47(.11)	.58(.14)	1.0(.07)	1.02(.14)
Age (ref.class: 17–24)				
25–34	.35(.08)	.62(.12)	.75(.09)	.80(.08)
35–44	.26(.12)	.49(.2)	.71(.09)	.73(.1)
45–54	.41(.24)	.82(.35)	.82(.17)	.84(.19)
55+	.59(.36)	.99(.42)	.82(.18)	.91(.19)

	Demo	NoDemo	Demo+AT.	NoDemo+AT.
<i>Global Accuracy (Metric: AUC(std))</i>				
Baseline	.50	.50	.50	.50
State of the Art	–	.61(.01) (*)	–	–
Our Approach	.74(.02)	.71(.02)	.74(.02)	.71(.02)
<i>Precision and Recall</i>				
Precision	.16(.02)	.18(.01)	.26(.05)	.25(.03)
Recall	.56(.05)	.48(.02)	.21(.05)	.22(.04)
<i>Demographic accuracy (Metric: AUC(std))</i>				
Gender (M)	.66(.05)	.64(.04)	.66(.05)	.64(.04)
Our Approach	.74(.02)	.71(.02)	.74(.02)	.71(.02)
<i>Precision and Recall</i>				
Precision	.16(.02)	.18(.01)	.26(.05)	.25(.03)
Recall	.56(.05)	.48(.02)	.21(.05)	.22(.04)
Age (55+)	.46(.31)	.46(.29)	.46(.31)	.46(.29)
<i>Fairness (Metric: $\frac{FNR}{FNR_{ref}}$)</i>				
Gender (ref.class: Male)				
Female	.47(.11)	.58(.14)	1.0(.07)	1.02(.14)
Age (ref.class: 17–24)				
25-34	.35(.08)	.62(.12)	.75(.09)	.80(.08)
35-44	.26(.12)	.49(.2)	.71(.09)	.73(.1)
45-54	.41(.24)	.82(.35)	.82(.17)	.84(.19)
55+	.59(.36)	.99(.42)	.82(.18)	.91(.19)




Gender fairness per region in the NoDemo (left) and NoDemo+Thresh. (right) models. Gender fairness is computed as the FNR of females in relation to that of males. The color extremities are both unfair (Color figure online)

Conclusions

- Fairness through unawareness does not suffice
- Models with **good overall accuracy aren't always efficient** in humanitarian domain
- Easily **generalisable approach** to any **fairness metrics, demographic feature, and digital data**

Link to the paper



Kyriaki Kalimeri
kyriaki.kalimeri@isi.it
 @KyriakiKalimeri

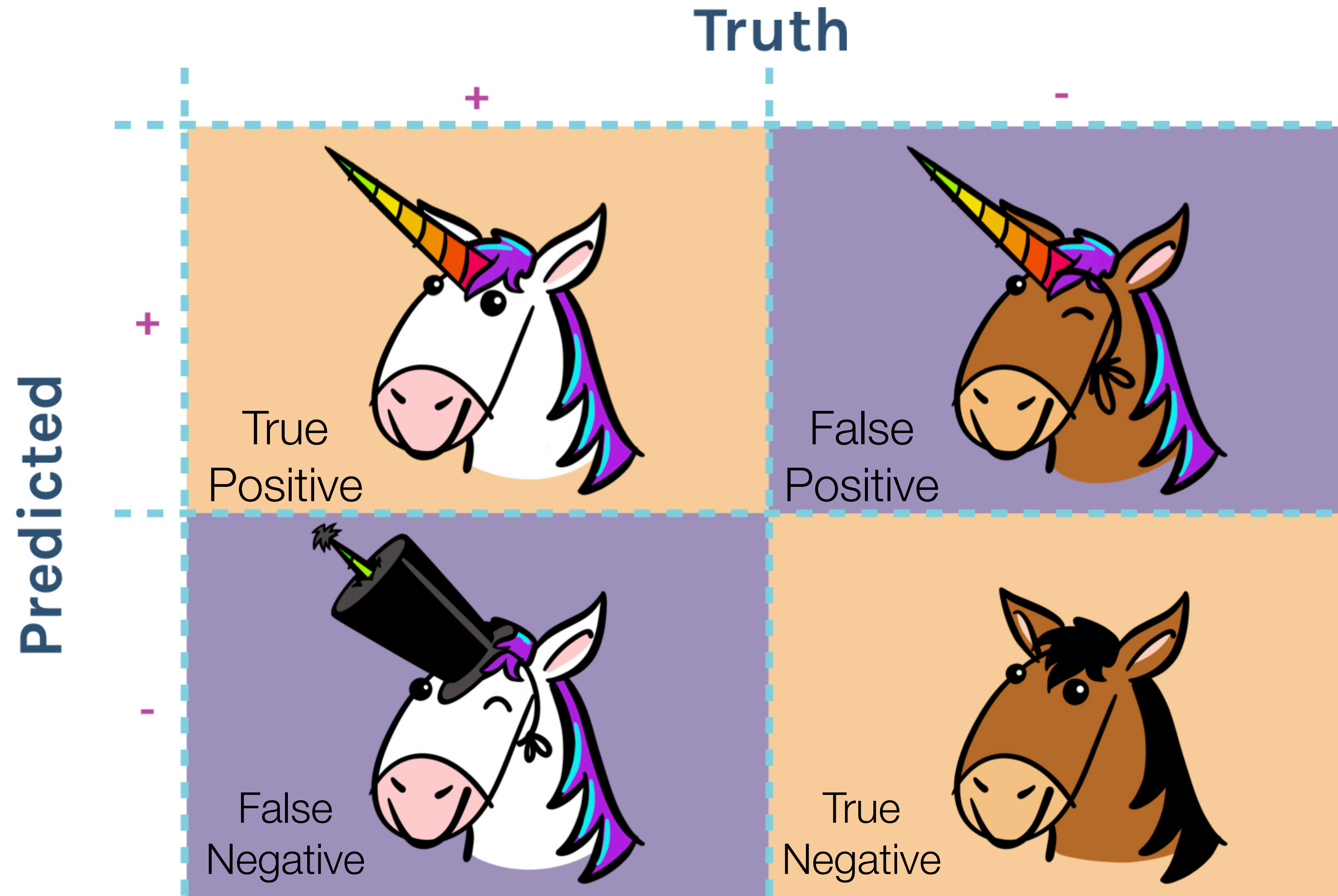


Description	# Features
Liked Pages	2,063,944
Liked Pages per Category To express how much a participant is interested in different categories of pages, we compute the number of pages he gave like to inside each Category.	1.553
Normalised Categories: As the participants' activity can greatly vary, we normalise the <i>Liked Pages per Category</i> to have sum 1.	1.553
Median Page Popularity This index shows how much a participant likes popular pages. The <i>popularity</i> of a Page is the number of users that gave like to it, as reported by Facebook in the Page profile.	1
Standard Deviation of Page Popularity	1
Median Category Popularity This index shows how much a participant likes popular categories.	1
Total number of Page likes One feature containing the total number of pages liked by the participant.	1
Total number of liked Categories One feature containing the total number of categories with pages liked by the participant.	1

Machine Learning and Vulnerable Populations

Precision = Of all the positives how many are truly positives?

Recall = Of the real positives, how many are predicted correctly?



LightGBM

10-fold stratified x validation
balanced weighting

