

¿QUÉ SESGOS IMPORTAN? ENVENENANDO CONJUNTOS DE DATOS PARA CLASIFICACIÓN DE IMÁGENES

Iris Dominguez-Catena

Noviembre 2024

Departamento de Estadística, Informática y Matemáticas,
Universidad Pública de Navarra(UPNA)

upna

Universidad Pública de Navarra
Nafarroako Unibertsitate Publikoa

Fairness and bias

Facial Expression Recognition

Medida de sesgo demográfico en datasets

Análisis de poblaciones

Transferencia de sesgo a modelos

Conclusiones

Fairness and bias

Facial Expression Recognition

Medida de sesgo demográfico en datasets

Análisis de poblaciones

Transferencia de sesgo a modelos

Conclusiones

- **AI ethics:** multidisciplinary field that studies how to optimize AI's beneficial impact while reducing risks and adverse outcomes.
 - **Algorithmic fairness:** Ensuring that algorithms make non-discriminatory decisions.

«Fairness is man's ability to rise above his prejudices.»

Wes Fesler

- **AI ethics:** multidisciplinary field that studies how to optimize AI's beneficial impact while reducing risks and adverse outcomes.
 - **Algorithmic fairness:** Ensuring that algorithms make non-discriminatory decisions.

«Fairness is man's ability to rise above his prejudices.»

Wes Fesler

- **AI ethics:** multidisciplinary field that studies how to optimize AI's beneficial impact while reducing risks and adverse outcomes.
 - **Algorithmic fairness:** Ensuring that algorithms make non-discriminatory decisions.

Algorithmic *machine's* *people's*
y «Fairness is ~~man's~~ ability to rise above ~~his~~ prejudices.»
Wes Fesler

RETAIL OCTOBER 11, 2018 / 1:04 AM / UPDATED 5 YEARS AGO

Amazon scraps secret AI recruiting tool that showed bias against women

By Jeffrey Dastin

8 MIN READ



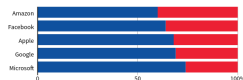
SAN FRANCISCO (Reuters) - Amazon.com Inc's AMZN.O machine-learning specialists uncovered a big problem: their new recruiting engine did not like women.

Dominated by men

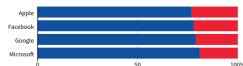
Top U.S. tech companies have yet to close the gender gap in hiring, a disparity most pronounced among technical staff such as software developers where men far outnumber women. Amazon's experimental recruiting engine followed the same pattern, learning to penalize resumes including the word "women's" until the company discovered the problem.

GLOBAL HEADCOUNT

■ Male ■ Female




EMPLOYEES IN TECHNICAL ROLES



Note: Amazon does not disclose the gender breakdown of its technical workforce. Source: Latest data available from the companies, since 2017. By Han Huang | REUTERS GRAPHICS

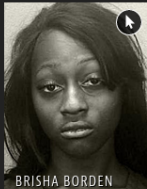
⁰<https://www.reuters.com/article/us-amazon-com-jobs-automation-insight-idUSKCN1MK08G>

Two Petty Theft Arrests



VERNON PRATER

LOW RISK **3**



BRISHA BORDEN

HIGH RISK **8**

Borden was rated high risk for future crime after she and a friend took a kid's bike and scooter that were sitting outside. She did not reoffend.

Two Petty Theft Arrests

VERNON PRATER

Prior Offenses
2 armed robberies, 1 attempted armed robbery

Subsequent Offenses
1 grand theft

LOW RISK **3**

BRISHA BORDEN

Prior Offenses
4 juvenile misdemeanors

Subsequent Offenses
None

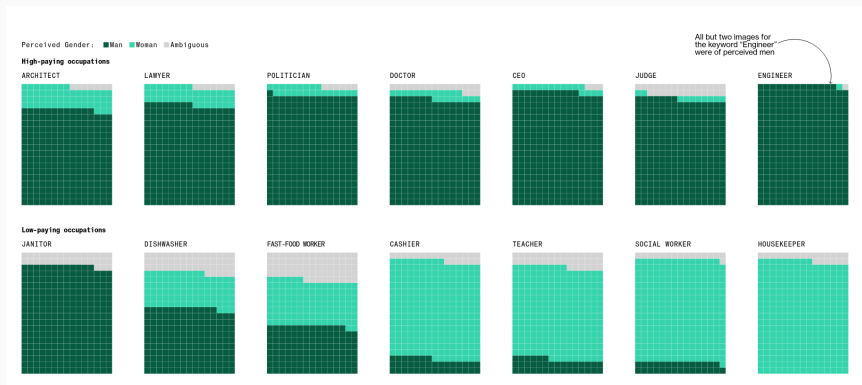
HIGH RISK **8**

Borden was rated high risk for future crime after she and a friend took a kid's bike and scooter that were sitting outside. She did not reoffend.

Prediction Fails Differently for Black Defendants

	WHITE	AFRICAN AMERICAN
Labeled Higher Risk, But Didn't Re-Offend	23.5%	44.9%
Labeled Lower Risk, Yet Did Re-Offend	47.7%	28.0%

⁰<https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>



⁰<https://www.bloomberg.com/graphics/2023-generative-ai-bias/>

Stable Diffusion Perpetuates Criminal Stereotypes

Composite average of all images

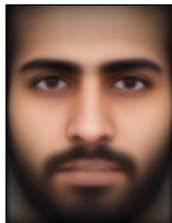
INMATE



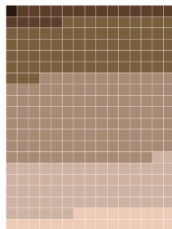
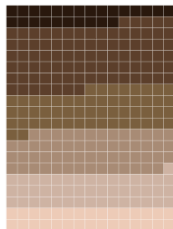
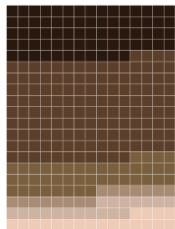
DRUG DEALER



TERRORIST



Distribution of skin tones



⁰<https://www.bloomberg.com/graphics/2023-generative-ai-bias/>

- ChatGPT's political biases¹
- AI agents and bias: Tay.ai ²
- Military AI ³
- Medical AI ⁴
- Social networks and recommendation algorithms ⁵

¹Motoki2023.

²<https://www.theverge.com/2016/3/24/11297050/tay-microsoft-chatbot-racist>

³<https://www.euronews.com/next/2022/10/17/israel-deploys-ai-powered-robot-guns-that-can-track-targets-in-the-west-bank>

⁴<https://www.scientificamerican.com/article/racial-bias-found-in-a-major-health-care-risk-algorithm/>

⁵<https://www.adl.org/resources/report/exposure-alternative-extremist-content-youtube>

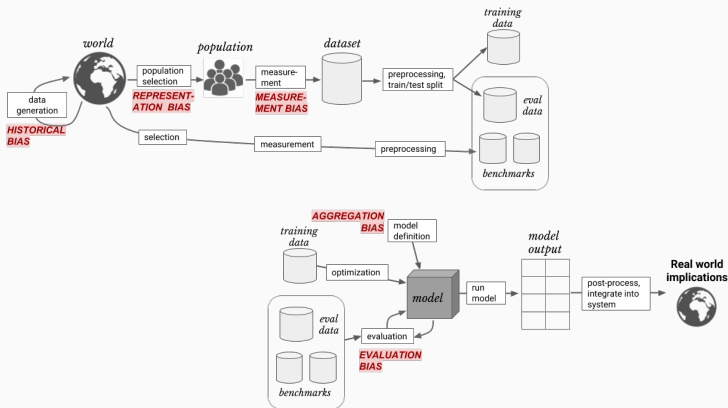


Figure 1: Bias source in the machine learning pipeline⁶

⁶Harini Suresh and John Guttag. "A Framework for Understanding Sources of Harm throughout the Machine Learning Life Cycle". In: *Equity and Access in Algorithms, Mechanisms, and Optimization*. EAAMO '21. 2021, pp. 1–9. 10/34

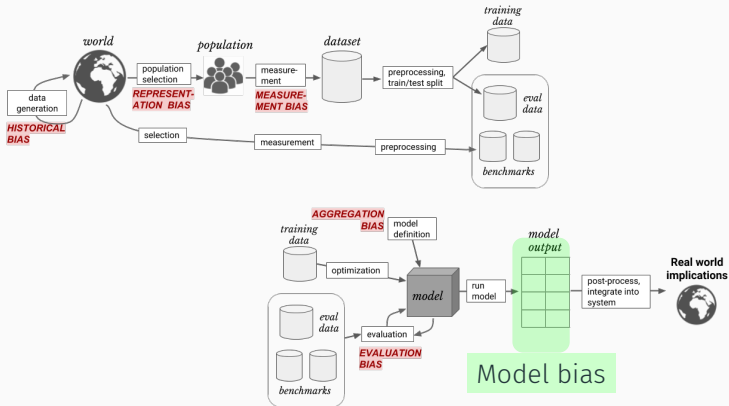


Figure 1: Bias source in the machine learning pipeline⁶

⁶Harini Suresh and John Guttag. "A Framework for Understanding Sources of Harm throughout the Machine Learning Life Cycle". In: *Equity and Access in Algorithms, Mechanisms, and Optimization*. EAAMO '21. 2021, pp. 1–9. 10/34

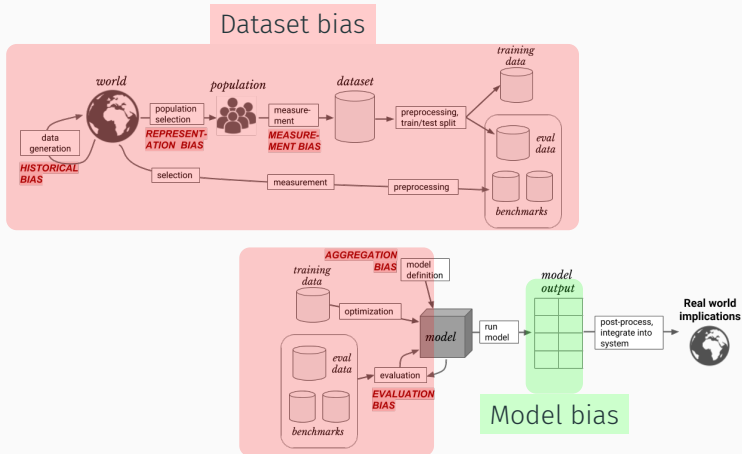


Figure 1: Bias source in the machine learning pipeline⁶

⁶Harini Suresh and John Guttag. "A Framework for Understanding Sources of Harm throughout the Machine Learning Life Cycle". In: *Equity and Access in Algorithms, Mechanisms, and Optimization*. EAAMO '21. 2021, pp. 1–9. 10/34

Fairness and bias

Facial Expression Recognition

Medida de sesgo demográfico en datasets

Análisis de poblaciones

Transferencia de sesgo a modelos

Conclusiones

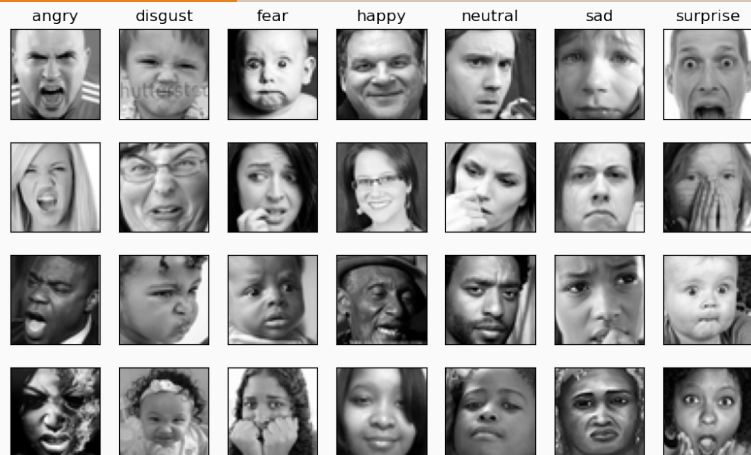






Figure 2: A sample of FER2013/FER+, a popular FER dataset⁷.

⁷Emad Barsoum et al. "Training Deep Networks for Facial Expression Recognition with Crowd-Sourced Label Distribution". In: *Proc. 18th ACM International Conference on Multimodal Interaction*. 2016, pp. 279–283.

Modalities

- **Image** or video
- **RGB**, IR, Depth...
- **Discrete** (Ekman's basic emotions ⁸) or continuous (NRC-VAD) labeling...

Applications

-  Interactive multimedia
 - Emotional Films
-  Healthcare ⁹
-  Assistive robotics ¹⁰
-  Public safety ¹¹

⁹Paul Ekman and Wallace V. Friesen. "Constants across Cultures in the Face and Emotion.". In: *Journal of Personality and Social Psychology* 17.2 (1971), pp. 124–129

¹⁰Philipp Werner et al. "Automatic Recognition Methods Supporting Pain Assessment: A Survey". In: *IEEE Trans. on Affective Computing* 13.1 (2022), pp. 530–552

¹¹Ritvik Nimmagadda, Kritika Arora, and Miguel Vargas Martin. "Emotion Recognition Models for Companion Robots". In: *The Journal of Supercomputing* (2022)

¹²Luntian Mou et al. "Isotropic Self-Supervised Learning for Driver Drowsiness Detection With Attention-Based Multimodal Fusion". In: *IEEE Trans. on Multimedia* 25 (2023), pp. 529–542

- Age, race and gender biases:
 - In research models¹³.
 - In commercial systems^{14,15}.

¹³Tian Xu et al. “Investigating Bias and Fairness in Facial Expression Recognition”. In: *Computer Vision – ECCV 2020 Workshops*. 2020, pp. 506–523.

¹⁴Khurshid Ahmad et al. “Comparing the Performance of Facial Emotion Recognition Systems on Real-Life Videos: Gender, Ethnicity and Age”. In: *Proc. Future Technologies Conference (FTC) 2021, Volume 1*. Vol. 358. 2022, pp. 193–210.

¹⁵Eugenia Kim et al. “Age Bias in Emotion Detection: An Analysis of Facial Emotion Recognition Performance on Young, Middle-Aged, and Older Adults”. In: *Proc. 2021 AAAI/ACM Conference on AI, Ethics, and Society*. 2021, pp. 638–644.

- Deep Learning approaches: CNN¹⁶ and Transformers¹⁷.
 - They require large amounts of data!
- Shift to large datasets gathered from the Internet¹⁸.
 - Datasets with little to no demographic metadata.

¹⁶Shan Li and Weihong Deng. "Deep Facial Expression Recognition: A Survey". In: *IEEE Trans. on Affective Computing* (2020), pp. 1-1.

¹⁷Alexey Dosovitskiy et al. *An Image Is Worth 16x16 Words: Transformers for Image Recognition at Scale*. 2021. arXiv: 2010.11929 [cs].

¹⁸Emily Denton et al. "On the Genealogy of Machine Learning Datasets: A Critical History of ImageNet". In: *Big Data & Society* 8.2 (2021).

Table 1: 20 selected datasets.

Abbreviation	Year	Collection	Images	Videos	Subjects
JAFFE	1998	LAB	213	—	10
KDEF	1998	LAB	4,900	—	70
CK	2000	LAB	8,795	486	97
Oulu-CASIA	2008	LAB	66,000	480	80
CK+	2010	LAB	10,727	593	123
GEMEP	2010	LAB	2,817	1,260	10
MUG	2010	LAB	70,654	—	52
SFEW	2011	ITW-M	1,766	—	330
FER2013	2013	ITW	32,298	—	—
WSEFEP	2014	LAB	210	—	30
ADFES	2016	LAB	—	648	22
FERplus	2016	ITW	32,298	—	—
AffectNet	2017	ITW	291,652	—	—
ExpW	2017	ITW	91,793	—	—
RAF-DB	2017	ITW	29,672	—	—
CAER-S	2019	ITW-M	70,000	—	—
LIRIS-CSE	2019	LAB	26,000	208	12
iSAFE	2020	LAB	—	395	44
MMAFEDB	2020	ITW	128,000	—	—
NHFIER	2020	ITW	5,558	—	—

- Lab: Laboratory-gathered.
- ITW-I: From Internet.
- ITW-M: From motion pictures.

Fairness and bias

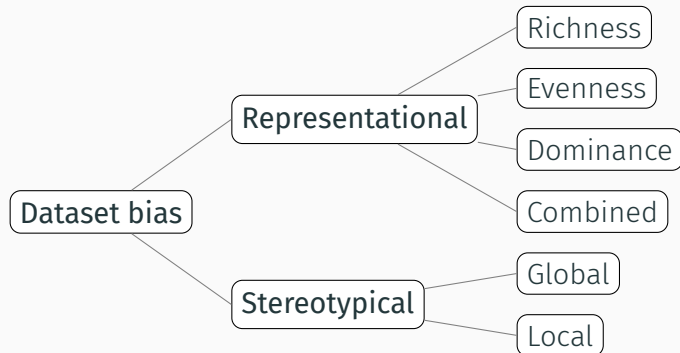
Facial Expression Recognition

Medida de sesgo demográfico en datasets

Análisis de poblaciones

Transferencia de sesgo a modelos

Conclusiones



Representational bias

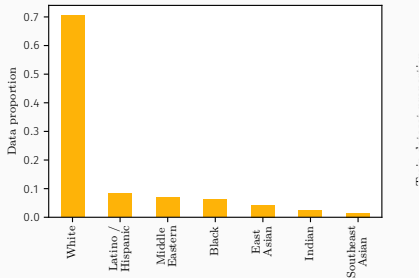


Figure 3: Apparent race distribution in FER+.

Representational bias

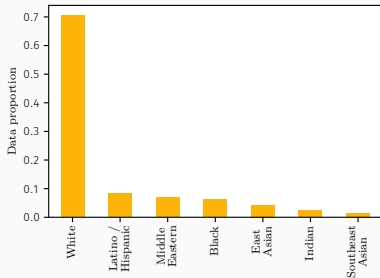


Figure 3: Apparent race distribution in FER+.

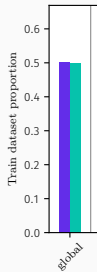


Figure 4: Apparent *per-label* gender distribution in FER+.

Representational bias

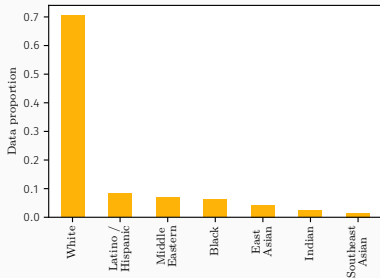


Figure 3: Apparent race distribution in FER+.

Stereotypical bias

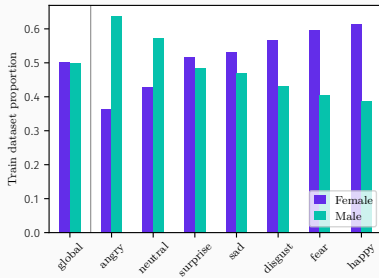


Figure 4: Apparent *per-label* gender distribution in FER+.

	Laboratory	ITW-M	ITW-I	
Age	9 – ENS	7.067 ± 0.932	5.551 ± 0.990	3.334 ± 0.286
	1 – SEI	0.409 ± 0.200	0.437 ± 0.154	0.211 ± 0.024
	ϕ_C	0.075 ± 0.063	0.084 ± 0.031	0.104 ± 0.027
Race	7 – ENS	5.168 ± 0.634	5.070 ± 0.080	3.724 ± 0.321
	1 – SEI	0.384 ± 0.151	0.663 ± 0.021	0.393 ± 0.050
	ϕ_C	0.092 ± 0.083	0.058 ± 0.018	0.063 ± 0.018
Gender	2 – ENS	0.139 ± 0.280	0.074 ± 0.055	0.005 ± 0.005
	1 – SEI	0.039 ± 0.052	0.055 ± 0.041	0.004 ± 0.004
	ϕ_C	0.067 ± 0.090	0.199 ± 0.062	0.167 ± 0.018

Figure 5: Average representational bias (ENS), evenness (SEI) and stereotypical bias (ϕ_C) of Lab, ITW-M and ITW-I datasets.

		Laboratory	ITW-M	ITW-I
Age	9 – ENS	7.067 ± 0.932	5.551 ± 0.990	3.334 ± 0.286
	1 – SEI	0.409 ± 0.200	0.437 ± 0.154	0.211 ± 0.024
	ϕ_C	0.075 ± 0.063	0.084 ± 0.031	0.104 ± 0.027
Race	7 – ENS	5.168 ± 0.634	5.070 ± 0.080	3.724 ± 0.321
	1 – SEI	0.384 ± 0.151	0.663 ± 0.021	0.393 ± 0.050
	ϕ_C	0.092 ± 0.083	0.058 ± 0.018	0.063 ± 0.018
Gender	2 – ENS	0.139 ± 0.280	0.074 ± 0.055	0.005 ± 0.005
	1 – SEI	0.039 ± 0.052	0.055 ± 0.041	0.004 ± 0.004
	ϕ_C	0.067 ± 0.090	0.199 ± 0.062	0.167 ± 0.018

Figure 5: Average representational bias (ENS), evenness (SEI) and stereotypical bias (ϕ_C) of Lab, ITW-M and ITW-I datasets.

	Laboratory	ITW-M	ITW-I	
Age	9 – ENS	7.067 ± 0.932	5.551 ± 0.990	3.334 ± 0.286
	1 – SEI	0.409 ± 0.200	0.437 ± 0.154	0.211 ± 0.024
	ϕ_C	0.075 ± 0.063	0.084 ± 0.031	0.104 ± 0.027
Race	7 – ENS	5.168 ± 0.634	5.070 ± 0.080	3.724 ± 0.321
	1 – SEI	0.384 ± 0.151	0.663 ± 0.021	0.393 ± 0.050
	ϕ_C	0.092 ± 0.083	0.058 ± 0.018	0.063 ± 0.018
Gender	2 – ENS	0.139 ± 0.280	0.074 ± 0.055	0.005 ± 0.005
	1 – SEI	0.039 ± 0.052	0.055 ± 0.041	0.004 ± 0.004
	ϕ_C	0.067 ± 0.090	0.199 ± 0.062	0.167 ± 0.018

Figure 5: Average representational bias (ENS), evenness (SEI) and stereotypical bias (ϕ_C) of Lab, ITW-M and ITW-I datasets.

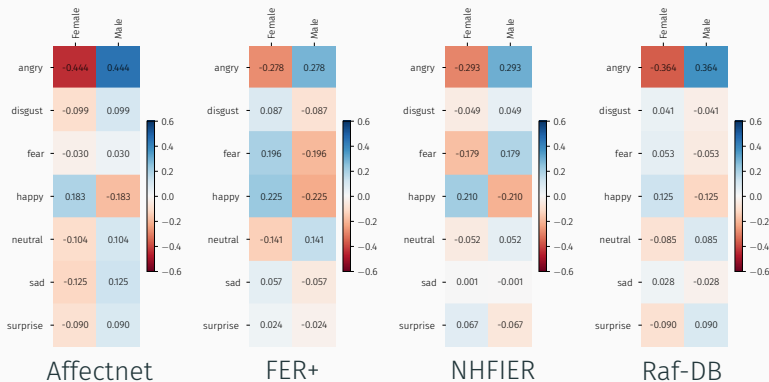


Figure 6: Local stereotypical bias (Ducher's Z) for some ITW-I datasets.

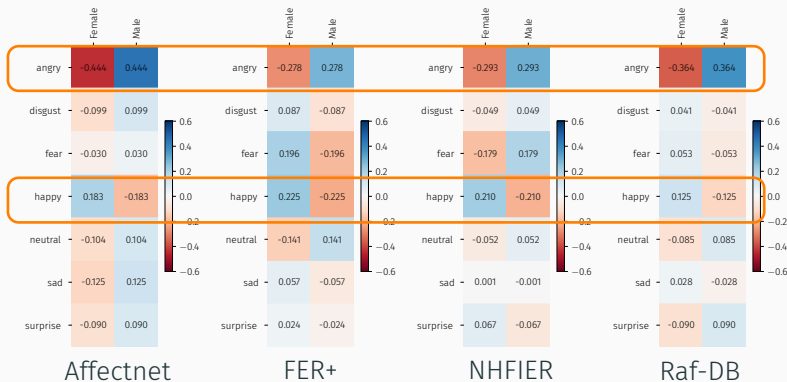


Figure 6: Local stereotypical bias (Ducher's Z) for some ITW-I datasets.

Fairness and bias

Facial Expression Recognition

Medida de sesgo demográfico en datasets

Análisis de poblaciones

Transferencia de sesgo a modelos

Conclusiones

- Dataset bias metrics still lack **interpretability**.
 - Issues of comparability and range disparities.
 - Lack of clear meaning.
- Most datasets lack demographic information.
 - Especially modern ITW datasets.
- There is no way to study the evolution and changes in bias across datasets.
 - Do equally biased dataset represent the same populations?
 - Analogous problems in archaeology¹⁹ and ecology^{20,21}.

¹⁹W. S. Robinson. "A Method for Chronologically Ordering Archaeological Deposits", in: *American Antiquity* 16A (1951), pp. 293-301.

²⁰M. V. Wilson and A. Shmida. "Measuring Beta Diversity with Presence-Absence Data", in: *Journal of Ecology* 72.3 (1984), pp. 1053-1064. JSTOR: 2359551.

²¹C. Ricotta and I. Podani. "On Some Properties of the Bray-Curtis Dissimilarity and Their Ecological Meaning", in: *Ecological Complexity* 31 (2017), pp. 201-205.

- Dataset bias metrics still lack **interpretability**.
 - Issues of comparability and range disparities.
 - Lack of clear meaning.
- **Most datasets lack demographic information.**
 - **Especially modern ITW datasets.**
- There is no way to study the evolution and changes in bias across datasets.
 - Do equally biased dataset represent the same populations?
 - Analogous problems in archaeology¹⁹ and ecology^{20,21}.

¹⁹W. S. Robinson. "A Method for Chronologically Ordering Archaeological Deposits", in: *American Antiquity* 16A (1951), pp. 293-301.

²⁰M. V. Wilson and A. Shmida. "Measuring Beta Diversity with Presence-Absence Data", in: *Journal of Ecology* 72.3 (1984), pp. 1053-1064. JSTOR: 2359551.

²¹C. Ricotta and J. Podani. "On Some Properties of the Bray-Curtis Dissimilarity and Their Ecological Meaning", in: *Ecological Complexity* 31 (2017), pp. 201-205.

- Dataset bias metrics still lack **interpretability**.
 - Issues of comparability and range disparities.
 - Lack of clear meaning.
- Most datasets lack demographic information.
 - Especially modern ITW datasets.
- There is no way to study the evolution and changes in bias across datasets.
 - Do equally biased dataset represent the same **populations?**
 - Analogous problems in archaeology¹⁹ and ecology^{20,21}.

¹⁹W. S. Robinson. "A Method for Chronologically Ordering Archaeological Deposits". In: *American Antiquity* 16.4 (1951), pp. 293–301.

²⁰M. V. Wilson and A. Shmida. "Measuring Beta Diversity with Presence-Absence Data". In: *Journal of Ecology* 72.3 (1984), pp. 1055–1064. JSTOR: 2259551.

²¹C. Ricotta and J. Podani. "On Some Properties of the Bray-Curtis Dissimilarity and Their Ecological Meaning". In: *Ecological Complexity* 31 (2017), pp. 201–205.

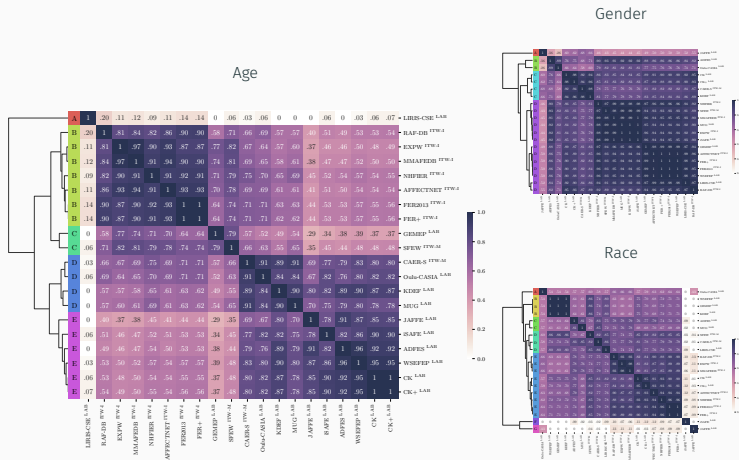


Figure 7: DSAP based comparison of datasets (age, gender and race axis).

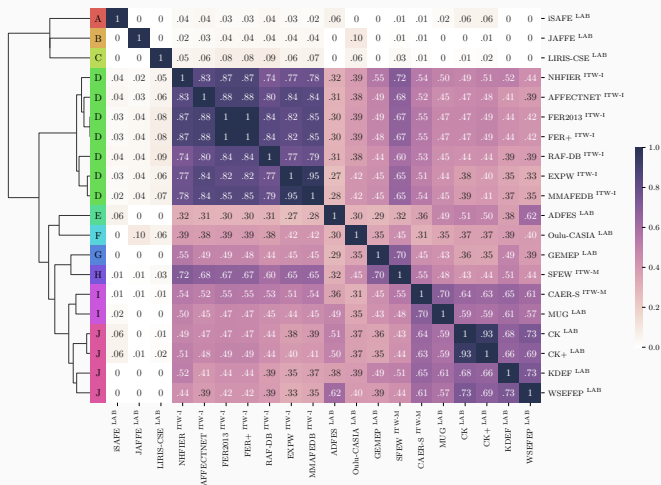


Figure 8: DSAP based comparison of datasets (combination axis, 126 subgroups).

Fairness and bias

Facial Expression Recognition

Medida de sesgo demográfico en datasets

Análisis de poblaciones

Transferencia de sesgo a modelos

Conclusiones

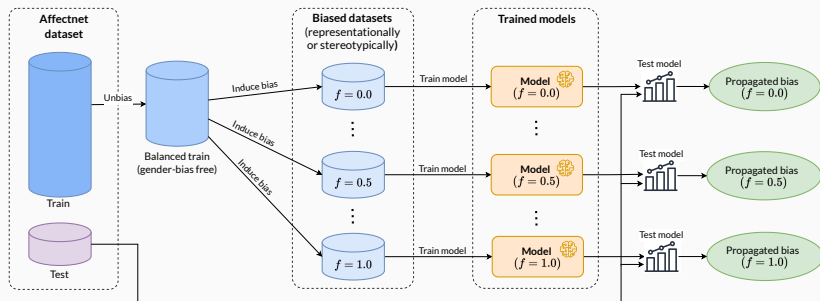


Figure 9: Summary of the methodology

Original

Gender	Female	Male	% F
angry	6,962	17,803	28.11%
disgust	1,733	2,054	45.76%
fear	3,117	3,222	49.17%
happy	79,797	54,241	59.53%
neutral	33,708	40,858	45.21%
sad	11,175	14,156	44.12%
surprise	6,439	7,588	45.90%
Total	142,931	139,922	50.53%

Balanced

Gender	Female	Male	% F
angry	6,962	6,962	50%
disgust	1,733	1,733	50%
fear	3,117	3,117	50%
happy	54,241	54,241	50%
neutral	33,708	33,708	50%
sad	11,175	11,175	50%
surprise	6,439	6,439	50%
Total	117,375	117,375	50%

Representational bias

Gender	Female	Male	% F
angry	6,962	0	100%
disgust	1,733	0	100%
fear	3,117	0	100%
happy	54,241	0	100%
neutral	33,708	0	100%
sad	11,175	0	100%
surprise	6,439	0	100%
Total	117,375	0	100%

Stereotypical bias

Gender	Female	Male	% F
angry	0	6,962	0%
disgust	866	866	50%
fear	1,558	1,558	50%
happy	27,120	27,120	50%
neutral	16,854	16,854	50%
sad	5,587	5,587	50%
surprise	3,219	3,219	50%
Total	55,204	62,166	47.4%

Original

Gender	Female	Male	% F
angry	6,962	17,803	28.11%
disgust	1,733	2,054	45.76%
fear	3,117	3,222	49.17%
happy	79,797	54,241	59.53%
neutral	33,708	40,858	45.21%
sad	11,175	14,156	44.12%
surprise	6,439	7,588	45.90%
Total	142,931	139,922	50.53%

Balanced

Gender	Female	Male	% F
angry	6,962	6,962	50%
disgust	1,733	1,733	50%
fear	3,117	3,117	50%
happy	54,241	54,241	50%
neutral	33,708	33,708	50%
sad	11,175	11,175	50%
surprise	6,439	6,439	50%
Total	117,375	117,375	50%

Representational bias

Gender	Female	Male	% F
angry	6,962	0	100%
disgust	1,733	0	100%
fear	3,117	0	100%
happy	54,241	0	100%
neutral	33,708	0	100%
sad	11,175	0	100%
surprise	6,439	0	100%
Total	117,375	0	100%

Stereotypical bias

Gender	Female	Male	% F
angry	0	6,962	0%
disgust	866	866	50%
fear	1,558	1,558	50%
happy	27,120	27,120	50%
neutral	16,854	16,854	50%
sad	5,587	5,587	50%
surprise	3,219	3,219	50%
Total	55,204	62,166	47.4%

Original

Gender	Female	Male	% F
angry	6,962	17,803	28.11%
disgust	1,733	2,054	45.76%
fear	3,117	3,222	49.17%
happy	79,797	54,241	59.53%
neutral	33,708	40,858	45.21%
sad	11,175	14,156	44.12%
surprise	6,439	7,588	45.90%
Total	142,931	139,922	50.53%

Balanced

Gender	Female	Male	% F
angry	6,962	6,962	50%
disgust	1,733	1,733	50%
fear	3,117	3,117	50%
happy	54,241	54,241	50%
neutral	33,708	33,708	50%
sad	11,175	11,175	50%
surprise	6,439	6,439	50%
Total	117,375	117,375	50%

Representational bias

Gender	Female	Male	% F
angry	6,962	0	100%
disgust	1,733	0	100%
fear	3,117	0	100%
happy	54,241	0	100%
neutral	33,708	0	100%
sad	11,175	0	100%
surprise	6,439	0	100%
Total	117,375	0	100%

Stereotypical bias

Gender	Female	Male	% F
angry	0	6,962	0%
disgust	866	866	50%
fear	1,558	1,558	50%
happy	27,120	27,120	50%
neutral	16,854	16,854	50%
sad	5,587	5,587	50%
surprise	3,219	3,219	50%
Total	55,204	62,166	47.4%

Original

Gender	Female	Male	% F
angry	6,962	17,803	28.11%
disgust	1,733	2,054	45.76%
fear	3,117	3,222	49.17%
happy	79,797	54,241	59.53%
neutral	33,708	40,858	45.21%
sad	11,175	14,156	44.12%
surprise	6,439	7,588	45.90%
Total	142,931	139,922	50.53%

Balanced

Gender	Female	Male	% F
angry	6,962	6,962	50%
disgust	1,733	1,733	50%
fear	3,117	3,117	50%
happy	54,241	54,241	50%
neutral	33,708	33,708	50%
sad	11,175	11,175	50%
surprise	6,439	6,439	50%
Total	117,375	117,375	50%

Representational bias

Gender	Female	Male	% F
angry	6,962	0	100%
disgust	1,733	0	100%
fear	3,117	0	100%
happy	54,241	0	100%
neutral	33,708	0	100%
sad	11,175	0	100%
surprise	6,439	0	100%
Total	117,375	0	100%

Stereotypical bias

Gender	Female	Male	% F
angry	0	6,962	0%
disgust	866	866	50%
fear	1,558	1,558	50%
happy	27,120	27,120	50%
neutral	16,854	16,854	50%
sad	5,587	5,587	50%
surprise	3,219	3,219	50%
Total	55,204	62,166	47.4%

Original

Gender	Female	Male	% F
angry	6,962	17,803	28.11%
disgust	1,733	2,054	45.76%
fear	3,117	3,222	49.17%
happy	79,797	54,241	59.53%
neutral	33,708	40,858	45.21%
sad	11,175	14,156	44.12%
surprise	6,439	7,588	45.90%
Total	142,931	139,922	50.53%

Balanced

Gender	Female	Male	% F
angry	6,962	6,962	50%
disgust	1,733	1,733	50%
fear	3,117	3,117	50%
happy	54,241	54,241	50%
neutral	33,708	33,708	50%
sad	11,175	11,175	50%
surprise	6,439	6,439	50%
Total	117,375	117,375	50%

Representational bias

Gender	Female	Male	% F
angry	1,392	5,570	20%
disgust	347	1,386	20%
fear	623	2,494	20%
happy	10,848	43,393	20%
neutral	6,742	26,966	20%
sad	2,235	8,940	20%
surprise	1,288	5,151	20%
Total	23,475	93,900	20%

Stereotypical bias

Gender	Female	Male	% F
angry	1,392	5,570	20%
disgust	866	866	50%
fear	1,558	1,558	50%
happy	27,120	27,120	50%
neutral	16,854	16,854	50%
sad	5,587	5,587	50%
surprise	3,219	3,219	50%
Total	55,596	60,774	47.7%

- **Dataset original:** Affectnet²².
- **Modelos:** ResNet50²³ y ViT-Base²⁴.
- **Configuraciones de sesgo:** 1 representacional, 7 estereotípicas.
 - **Proportions:** [0%, 10%, . . . , 100%].

²²Ali Mollahosseini, Behzad Hasani, and Mohammad H. Mahoor. "AffectNet: A Database for Facial Expression, Valence, and Arousal Computing in the Wild". In: *IEEE Trans. on Affective Computing* 10.1 (2019), pp. 18–31. arXiv: 1708.03985.

²³Kaiming He et al. *Deep Residual Learning for Image Recognition*. 2015. arXiv: 1512.03385 [cs].

²⁴Alexey Dosovitskiy et al. *An Image Is Worth 16x16 Words: Transformers for Image Recognition at Scale*. 2021. arXiv: 2010.11929 [cs].

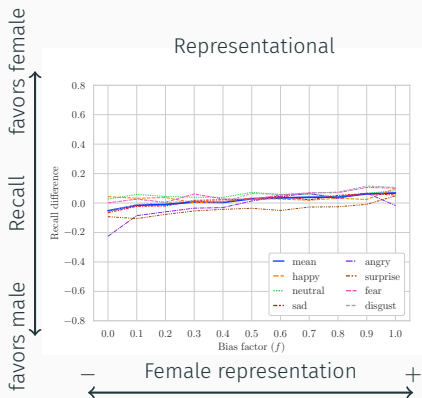


Figure 10: Recall difference (female recall minus male recall).

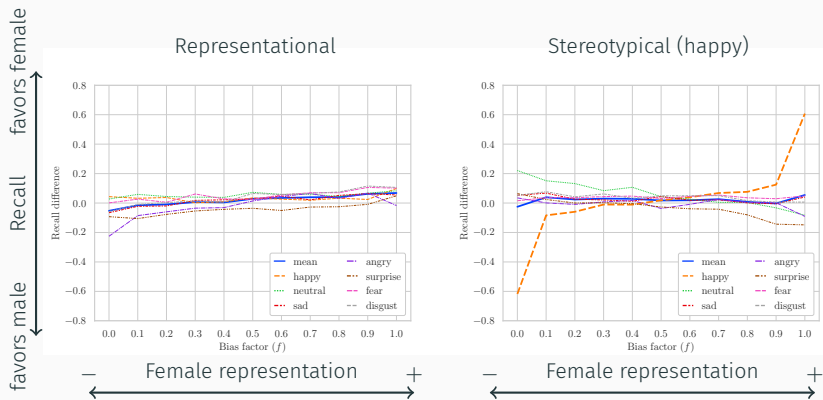


Figure 10: Recall difference (female recall minus male recall).

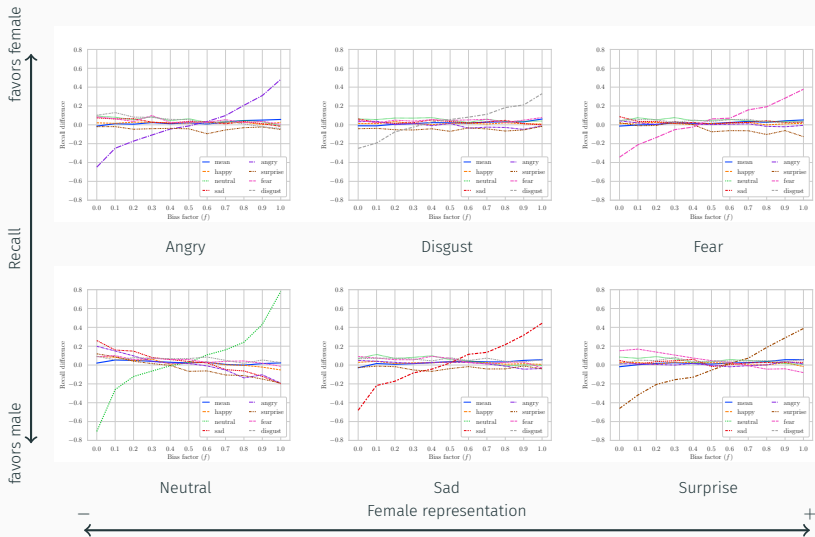


Figure 11: Effect of stereotypical bias.

Fairness and bias

Facial Expression Recognition

Medida de sesgo demográfico en datasets

Análisis de poblaciones

Transferencia de sesgo a modelos

Conclusiones

- Los nuevos datasets ITW-I han cambiado de perfil de sesgo. Predomina el sesgo estereotípico.
- Los datasets ITW-I tienden a ser extremadamente homogéneos.
- Los modelos entrenados son mucho más sensibles a sesgo estereotípico.

THANK YOU FOR THE ATTENTION.

¿QUESTIONS?

✉ IRIS.DOMINGUEZ@UNAVARRA.ES



<https://irisai.neocities.org>