

Fair Play for Individuals, Foul Play for Groups?

Auditing Anonymization's Impact on ML Fairness

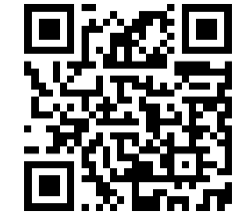
Héber H. Arcolezi[†], Mina Alishahi[‡], Adda-Akram Bendoukha[?], Nesrine Kaaniche[?]

[†]Inria

[‡]Open Universiteit

[?]Télécom SudParis

Paper



GitHub



heber.hwang-arcolezi@inria.fr

“How does **anonymization** impact **fairness metrics** in ML?”

Individual fairness generally **improves**
(smoother, more homogeneous data)

Group fairness generally **worsens**
(up to 4× degradation)

Suppression



Original data



Anonymized data

Generalization

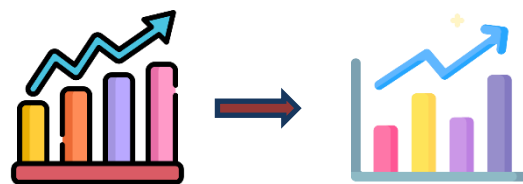


YES!

BUT...

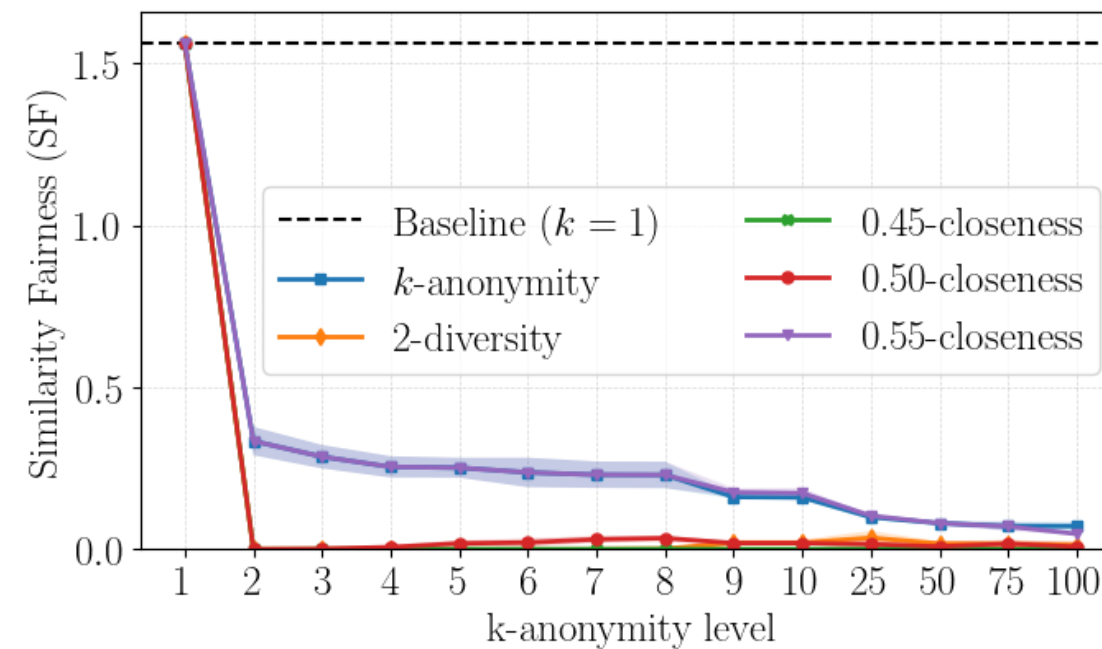


“Safe in the crowd”



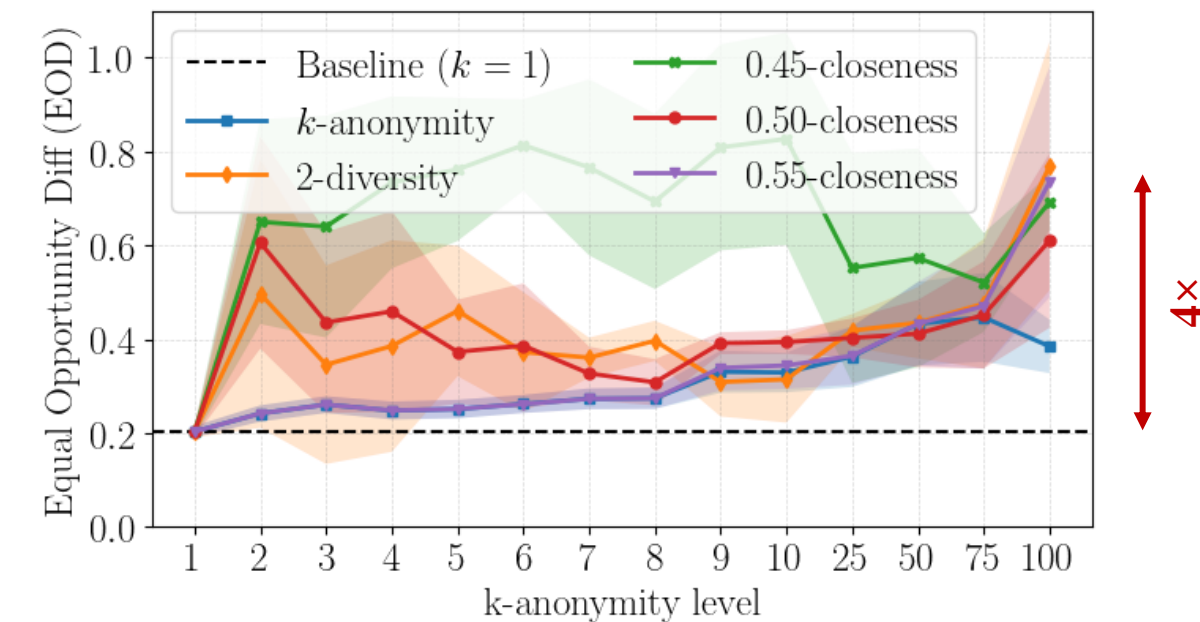
Alters data distributions

“Similar individuals → similar predictions”
(measures *local consistency* of model outputs)



More anonymity

“Equal outcomes across demographic groups”
(measures *parity* between protected groups)



More anonymity

4×