



Social Norm Bias: Residual Harms of Fairness-Aware Algorithms

Maria De-Arteaga, PhD

Assistant Professor

Information, Risk and Operations Management Department

University of Texas at Austin

Algorithmic group fairness

- Measures disparities across **sensitive groups** on a **measure of interest**.

Examples:

- Gender
- Race
- Ethnicity

Examples:

- Predicted positive
- False positive rate

Algorithmic group fairness

- Measures disparities across **sensitive groups** on a **measure of interest**.

Examples:

- Gender
- Race
- Ethnicity

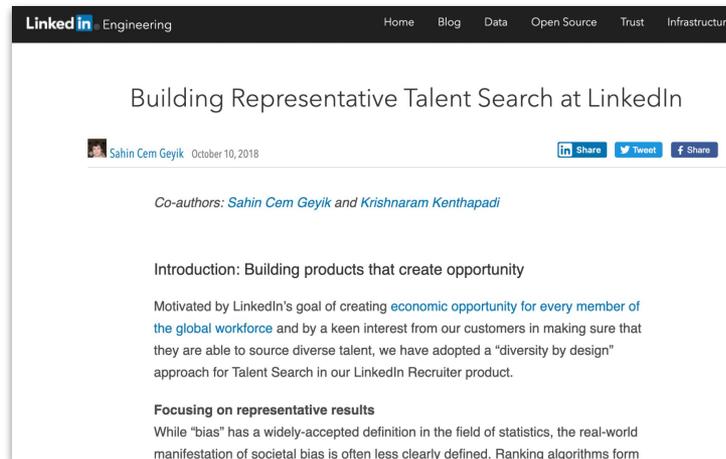
Examples:

- Predicted positive
- False positive rate

- Group fairness constraints enforced via different strategies:
 - pre-processing: transformation of input data
 - in-processing: constraints in optimization objective
 - post-processing: transformation of output scores

Algorithmic group fairness

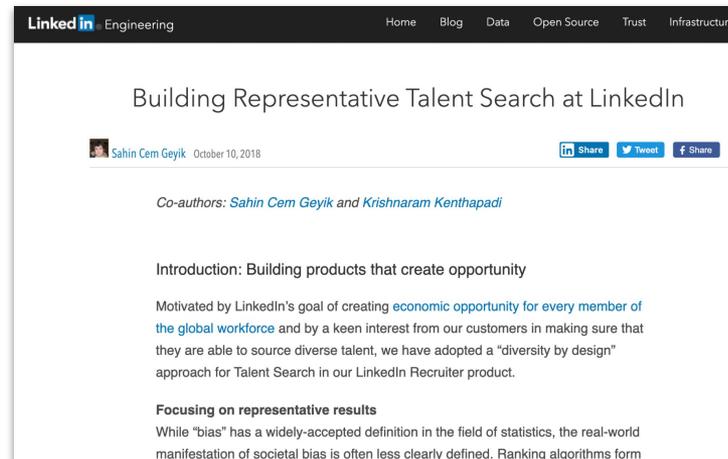
- Very popular due to ease of operationalization, increasingly implemented in ML deployment.
 - Especially true of post-processing (no cost of retraining!)
- Really great that algorithmic fairness research is having an impact! Are we done? Is the problem fixed?



Post-processing

Algorithmic group fairness

- Very popular due to ease of operationalization, increasingly implemented in ML deployment.
 - Especially true of post-processing (no cost of retraining!)
- Really great that algorithmic fairness research is having an impact! Are we done? Is the problem fixed?
 - Not really! Risk of reductive definitions and invisibilized harms.



What's Sex Got to Do With Fair Machine Learning?¹

Lily Hu² & Issa Kohler-Hausmann³

Algorithmic Fairness from a Non-ideal Perspective

Sina Fazelpour & Zachary C. Lipton

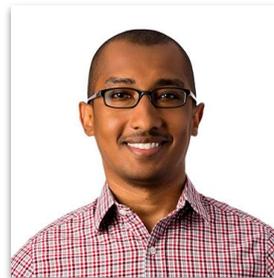
What are some of the invisibilized harms?

What are some of the invisibilized harms?

Social Norm Bias: Residual Harms of Fairness-Aware Algorithms

Myra Cheng, Maria De-Arteaga, Lester Mackey, Adam Tauman Kalai

ICML Machine Learning for Data Workshop, ICML Socially Responsible Machine Learning Workshop



What are some of the invisibilized harms?

Social Norm Bias: Residual Harms of Fairness-Aware Algorithms

Myra Cheng, Maria De-Arteaga, Lester Mackey, Adam Tauman Kalai

ICML Machine Learning for Data Workshop, ICML Socially Responsible Machine Learning Workshop



Social Norm Bias (SNoB): The associations between an **algorithm's predictions** and **adherence to social norms**

Description and Prescription: How Gender Stereotypes Prevent Women's Ascent Up the Organizational Ladder

Madeline E. Heilman*

New York University

How Women Engineers Do and Undo Gender: Consequences for Gender Equality

Abigail Powell ✉, Barbara Bagilhole, Andrew Dainty

Gender stereotypes and workplace bias

Madeline E. Heilman

Description and Prescription: How Gender Stereotypes Prevent Women's Ascent Up the Organizational Ladder

Madeline E. Heilman*
New York University

How Women Engineers Do and Undo Gender: Consequences for Gender Equality

Abigail Powell ✉ Barbara Bagilhole, Andrew Dainty

Gender stereotypes and workplace bias

Madeline E. Heilman

Do algorithms exhibit social norm bias?

Is this bias mitigated by group fairness approaches?

In other words, am I more likely to benefit from group fairness approaches if I sound/act "like a man"?

Experiments using task/dataset *Bias in Bios*

Enter the bio

She is a fifth year PhD student in the joint Machine Learning and Public Policy program at Carnegie Mellon University's Machine Learning Department and Heinz College. She is co-advised by Prof. Artur Dubrawski and Prof. Alexandra Chouldechova, and she is part of the Auton Lab.

Currently, her main focus is algorithmic fairness, studying how to measure and prevent bias and discrimination that may arise when using machine learning for decision support. She is passionate about developing novel machine learning algorithms that are

PREDICT TITLE

SHE

HE

she is a fifth year phd student in the joint machine learning and public policy program at carnegie mellon university <unk> s machine learning department and heinz college . she is co-advised by prof. artur <unk> and prof. alexandra chouldechova , and she is part of the auton lab . currently , her main focus is algorithmic fairness , studying how to measure and prevent bias and discrimination that may arise when using machine learning for decision support . she is passionate about developing novel machine learning algorithms that are motivated by existing policy problems , and understanding how machine learning can better help us overcome important societal challenges . prior to graduate school she received her b.sc . in mathematics from universidad nacional de colombia and worked as a journalist for one of colombia <unk> s main news magazine , semana . she is the recipient of a microsoft

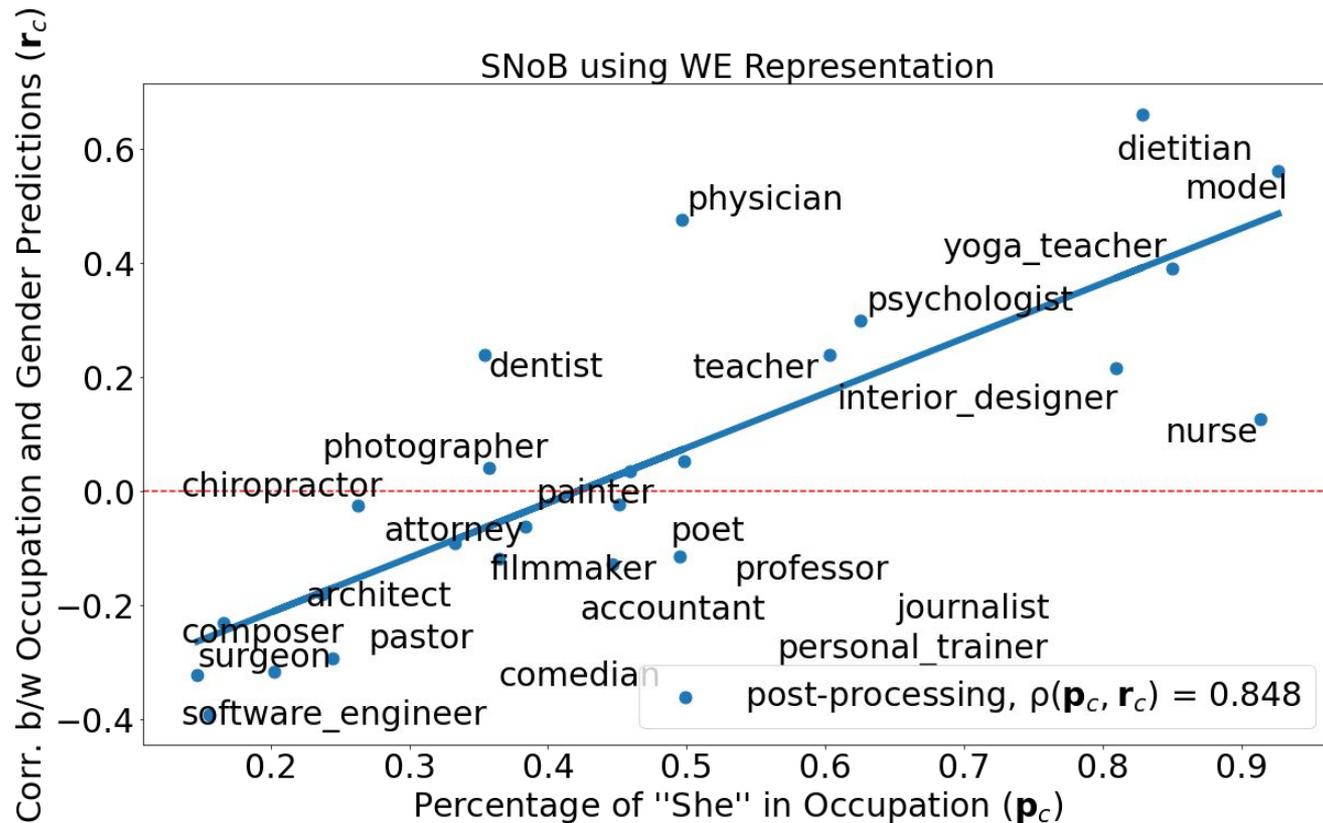
teacher

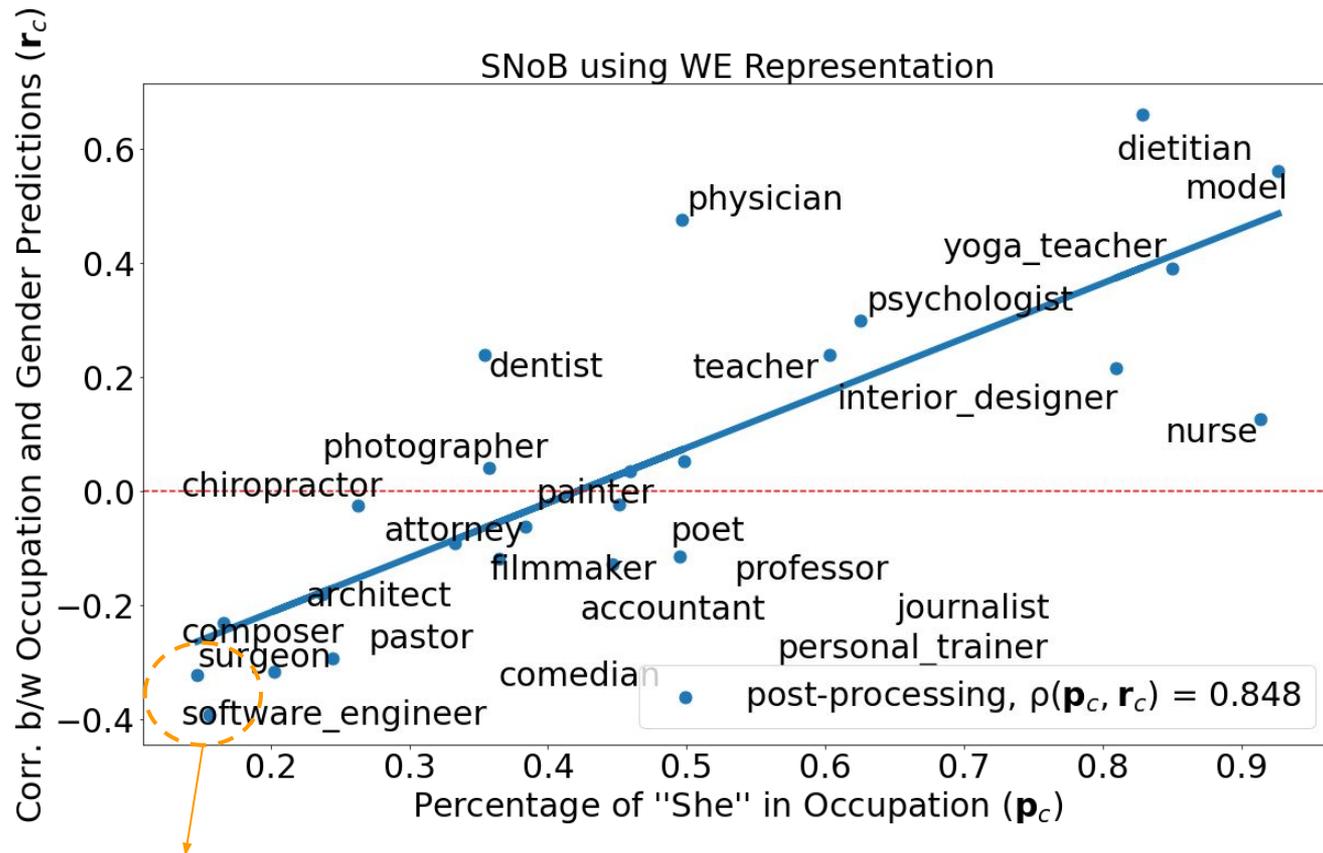
Bias in Bios: A Case Study of Semantic Representation Bias in a High-Stakes Setting (FAcCT 2019)

M. De-Arteaga, A. Romanov, H. Wallach, J. Chayes, C. Borgs, A. Chouldechova, S. Geyik, K. Kenthapadi, A. Kalai

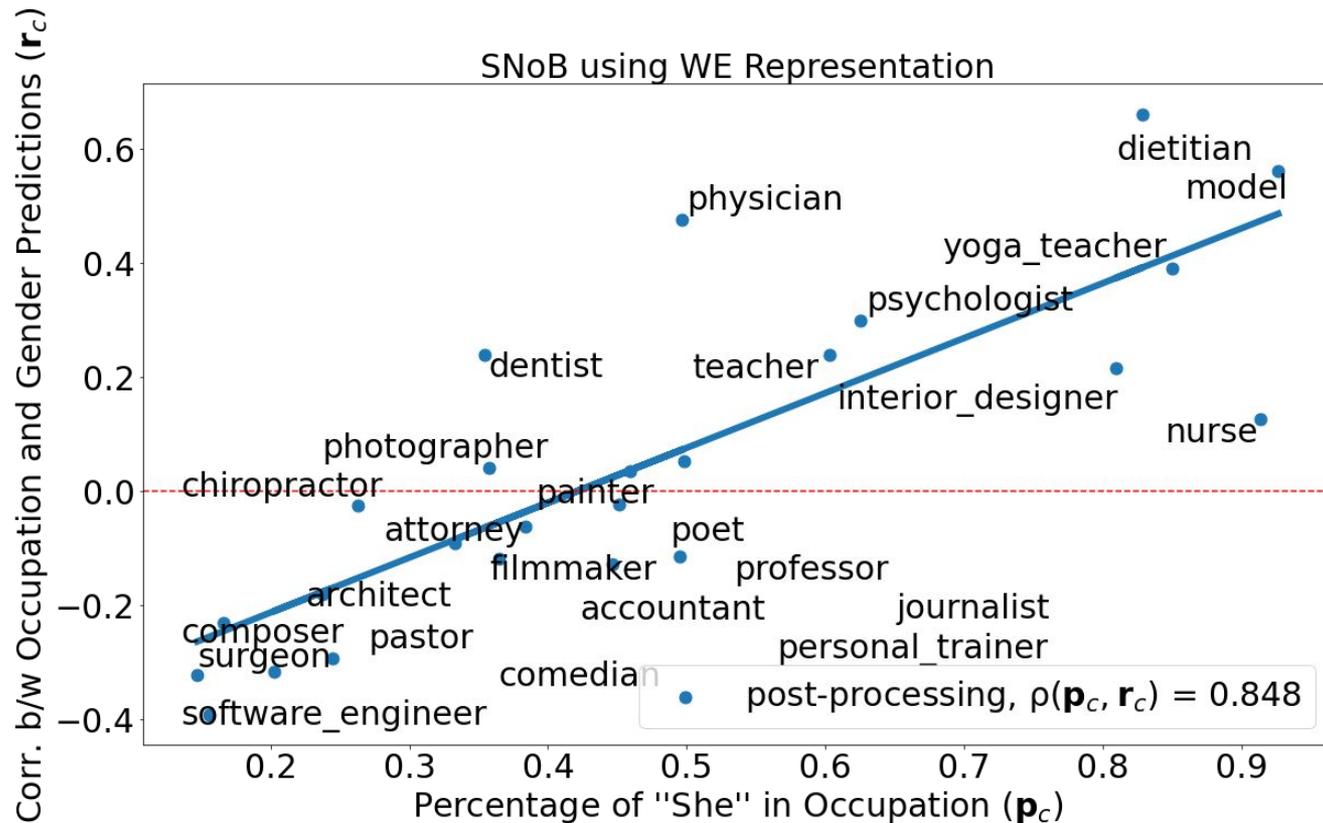
Measuring Inferred Social Norms

- Trained a classifier **G(x)** to predict group membership (e.g. “he”/“she”) from bio x
 - Gender-balanced within occupation
- Use probability of “he”/“she” as measure of aligning to inferred masculine/feminine gender norms
 - Does the biography align with what the algorithm associates as masculine/feminine?
 - Validation to ensure associations learned by algorithms align with human’s gender norm associations

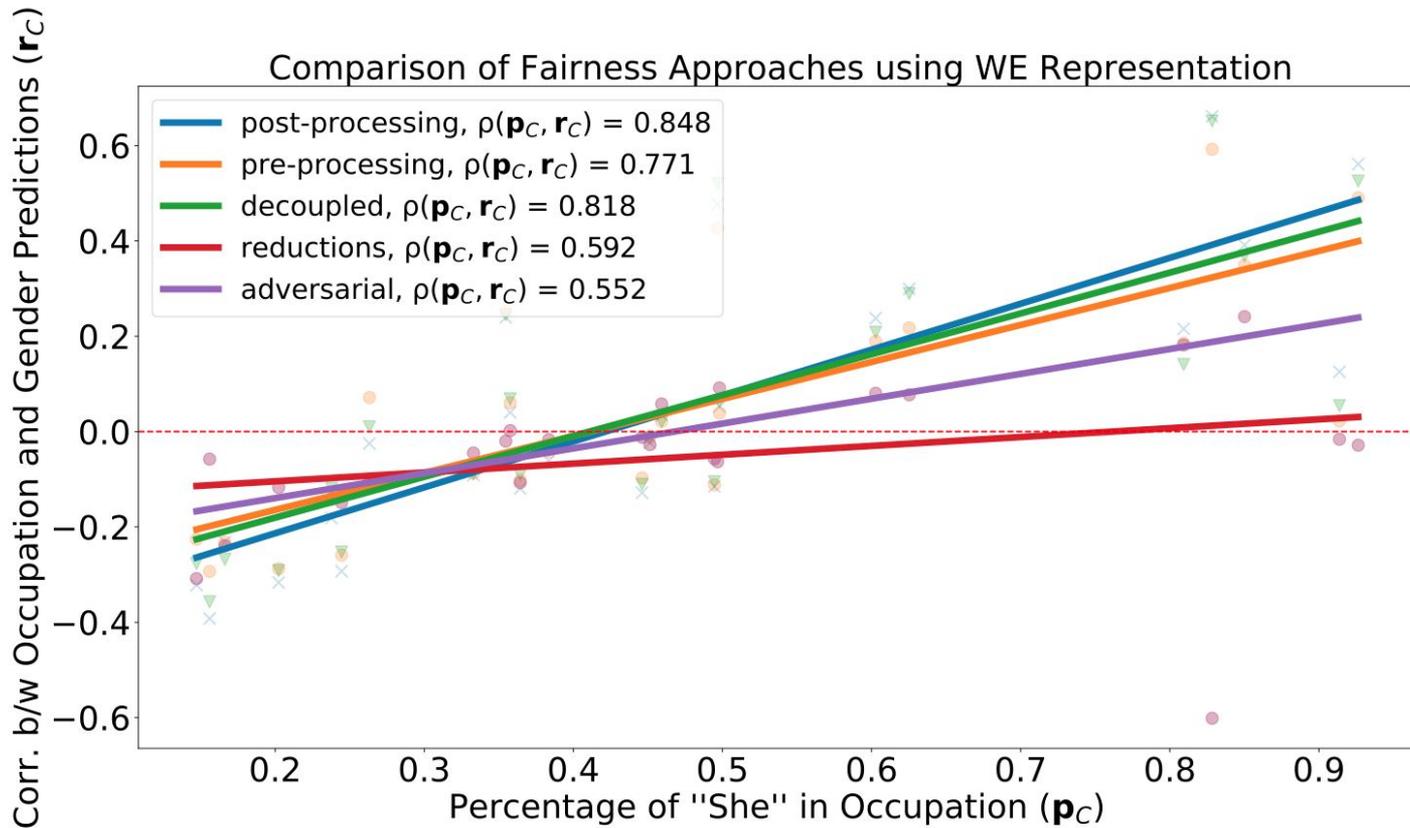




women surgeons and software engineers are less likely to be predicted as such if their bio "sound like a woman's"



Strength of correlation increases with occupation's gender imbalance



Gender correlations persist using “fair” interventions

— post-processing, $\rho(\mathbf{p}_c, \mathbf{r}_c) = 0.848$

Invisibilized harms of group fairness

- Algorithms may display **social norm bias**: occupation classification associated to individuals' adherence to social norms.
- Group fairness approaches may **veil harm** to individuals, e.g. fewer opportunities for people perceived as feminine in predominantly masculine occupations.
- Fairness-aware algorithms are not all equal: post-processing that preserves within-group order **does not mitigate** this harm at all.

Percentage of "She" in Occupation (\mathbf{p}_c)

— post-processing, $\rho(\mathbf{p}_c, \mathbf{r}_c) = 0.848$

Invisibilized harms of group fairness

- Algorithms may display **social norm bias**: occupation classification associated to individuals' adherence to social norms.
- Group fairness approaches may **veil harm** to individuals, e.g. fewer opportunities for people perceived as feminine in predominantly masculine occupations.
- Fairness-aware algorithms are not all equal: post-processing that preserves within-group order **does not mitigate** this harm at all.

Thank you!