



# Exploring Fairness and Sociodemographic Bias in Machine Learning

Theodora Chaspari, Ph.D. HUman Bio-Behavioral Signals (HUBBS) Lab Computer Science & Engineering, Texas A&M University

> TAMU Data Science Institute 4/23/2021





# Close your eyes and picture a shoe







## How about a physicist?















## Machine learning







# Machine learning

I STREET hobe check to to the to the 2 - 2 -L \_ = = = = = = 





#### **TECH**

# The problem with AI? Study says it's too white and male, calls for more women, minorities

#### Jessica Guynn USA TODAY

Published 8:00 p.m. ET Apr. 16, 2019 Updated 11:37 a.m. ET Apr. 17, 2019



Two out of four leading face recognition platforms do not reliably detect African American users.



West, S. M., Whittaker, M., & Crawford, K. (2019). Discriminating systems. Al Now. https://www.youtube.com/watch?v=TWWsW1w-BVo&t=45s







Microsoft Computer Vision API

West, S. M., Whittaker, M., & Crawford, K. (2019). Discriminating systems. Al Now. Buolamwini, J., & Gebru, T. (2018, January). Gender shades: Intersectional accuracy disparities in commercial gender classification. In Conference on fairness, accountability and transparency (pp. 77-91). PMLR.





	Is Adult Content Adult Score	False 0.14916780591011047
	Is Racy Content	False
A Dalla	Racy Score	0.12426207214593887
	Categories	[ { "name": "people_swimming", "score": 0.98046875 } ]
	Faces	[ { "age": 28, "gender": "Male", "faceRectangle": { "left": 744, "top": 338, "width": 305, "height": 305 } } ]
Age: 28 Gender: Male	Dominant Color Background	
	Dominant Color Foreground	
Is Adult Content: False	Dominant Colors	
151/tumblr_nnxfl1psuj1u2tjnvo1_500.jpg	Accent Color	#19A482

Microsoft Computer Vision API

8

West, S. M., Whittaker, M., & Crawford, K. (2019). Discriminating systems. Al Now. Buolamwini, J., & Gebru, T. (2018, January). Gender shades: Intersectional accuracy disparities in commercial gender classification. In Conference on fairness, accountability and transparency (pp. 77-91). PMLR.







Line Drawing Type	0 Non-LineDrawing
Black & White Image	False
s Adult Content	False
Adult Score	0.026106031611561775
s Racy Content	False
Racy Score	0.021592045202851295
Categories	[ { "name": "others_", "score": 0.00390625 }, { "name": "people_", "score": 0.5703125 } ]
aces	0
Dominant Color Background	
Dominant Color	

**Microsoft Computer Vision API** 

9

West, S. M., Whittaker, M., & Crawford, K. (2019). Discriminating systems. Al Now. Buolamwini, J., & Gebru, T. (2018, January). Gender shades: Intersectional accuracy disparities in commercial gender classification. In Conference on fairness, accountability and transparency (pp. 77-91). PMLR.







Line Drawing Type	0 Non-LineDrawing
Black & White Image	False
Is Adult Content	False
Adult Score	0.026106031611561775
Is Racy Content	False
Racy Score	0.021592045202851295
Categories	[ { "name": "others_", "score": 0.00390625 }, { "name": "people_", "score": 0.5703125 } ]
Faces	0
Dominant Color Background	
Dominant Color Foreground	

#### Microsoft Computer Vision API

West, S. M., Whittaker, M., & Crawford, K. (2019). Discriminating systems. AI Now. Buolamwini, J., & Gebru, T. (2018, January). Gender shades: Intersectional accuracy disparities in commercial gender classification. In Conference on fairness, accountability and transparency (pp. 77-91). PMLR.





#### MARCH 23, 2020

# Stanford researchers find that automated speech recognition is more likely to misinterpret black speakers

The disparity likely occurs because such technologies are based on machine learning systems that rely heavily on databases of English as spoken by white Americans.

Five leading speech recognition programs make twice as many errors with African American speakers as with Whites



Koenecke, A., Nam, A., Lake, E., Nudell, J., Quartey, M., Mengesha, Z., & Goel, S. (2020). Racial disparities in automated speech recognition. Proceedings of the National Academy of Sciences, 117(14), 7684-7689.





### Bias related to socioeconomic status (SES)



Visual disparities between communities of different socioeconomic (SES) status





# **Tutorial Overview**

- Sources of sociodemographic bias in machine learning (ML)
- Examples of studies examining sociodemographic bias in ML
  - Health Electronic Health Records
  - Family well-being Speech, language, physiology
- Approaches for mitigating sociodemographic bias in ML
  - Adversarial learning
  - Fairness regularization
  - Explainability





# **Tutorial Overview**

Sources of sociodemographic bias in machine learning (ML)

- Examples of studies examining sociodemographic bias in ML
  - Health Electronic Health Records
  - Family well-being Speech, language, physiology
- Approaches for mitigating sociodemographic bias in ML
  - Adversarial learning
  - Fairness regularization
  - Explainability





#### Bias in training data

- Minority bias: minoritized groups might have insufficient number of samples
- Missing data bias: minoritized groups may have missing data in a non-random fashion (e.g., lower quality sensor devices)
- Confounding factors: socio-demographic factors influencing both input and output variables (e.g., gender influences both resting heart rate and heart disease risk at early age)







#### Bias in model design

- Label bias: the same outcome might not mean the same for all individuals
- Cohort bias: considering traditional groups (e.g., male/female) without considering other protected groups (e.g., LGTBQ) and levels of granularity
- Proprietary algorithms, making it difficult to dissect them







Bias in interaction with experts

- Automation bias: experts are unaware that a model is underperforming for a certain group
- Feedback loops: If the clinician accepts incorrect model outputs, the mistake is propagated next time the model is trained
- Dismissal bias: Desensitization to alerts that are systematically incorrect for a specific group







#### Bias in interaction with users

- Privilege bias: ML models might be unavailable in places where specific groups receive care (e.g., devices with low computational resources, poor internet connectivity)
- Informed mistrust: Users might believe that a model is biased against them due to historical exploitation practices



#### **Privilege bias**



Informed mistrust





# **Tutorial Overview**

- Sources of sociodemographic bias in machine learning (ML)
- Examples of studies examining sociodemographic bias in ML
  - Health Electronic Health Records
  - Family well-being Speech, language, physiology
- Approaches for mitigating sociodemographic bias in ML
  - Adversarial learning
  - Fairness regularization
  - Explainability





### **RESEARCH ARTICLE**

#### ECONOMICS

# Dissecting racial bias in an algorithm used to manage the health of populations

Ziad Obermeyer<sup>1,2</sup>\*, Brian Powers<sup>3</sup>, Christine Vogeli<sup>4</sup>, Sendhil Mullainathan<sup>5</sup>\*†

Health systems rely on commercial prediction algorithms to identify and help patients with complex health needs. We show that a widely used algorithm, typical of this industry-wide approach and affecting millions of patients, exhibits significant racial bias: At a given risk score, Black patients are considerably sicker than White patients, as evidenced by signs of uncontrolled illnesses. Remedying this disparity would increase the percentage of Black patients receiving additional help from 17.7 to 46.5%. The bias arises because the algorithm predicts health care costs rather than illness, but unequal access to care means that we spend less money caring for Black patients than for White patients. Thus, despite health care cost appearing to be an effective proxy for health by some measures of predictive accuracy, large racial biases arise. We suggest that the choice of convenient, seemingly effective proxies for ground truth can be an important source of algorithmic bias in many contexts.





- Commercial risk-prediction algorithm applied to ~ 200 million people in the U.S.
  - Improving patient care by providing additional resources
  - Including greater attention from trained providers
  - Contributing toward ensuring well-coordinated care
- Primary care patients enrolled in risk-based contracts from 2013 to 2015
- 6,079 Black patients and 43,539 White patients
- 71.2% in commercial insurance and 28.8% in Medicare
- 50.9 years old on average and 63% female







Algorithmic outline





• Mean number of chronic conditions by race, plotted against algorithmic score

АМ

• For the same level of algorithmpredicted risk, Black patients found to depict significantly more illness burden compared to White patients



Percentile of Algorithm Risk Score





• For the same level of algorithmpredicted risk, Black patients found to depict significantly more illness burden across significant health markers compared to White patients







- For the same level of algorithm-predicted risk
  - Blacks depict significantly more illness burden compared to Whites
  - Blacks and Whites have (roughly) the same costs the following year



- Substantial disparities in health burden
- Little disparity in costs







• At a given level of health, Black patients generate lower costs than White patients

• Potentially the driving force behind this algorithmic disparity is that Black patients generate less medical expenses, therefore they are considered as lower risk by the algorithm







Examining algorithmic risk predictions with respect to label choice

#### Input features at time (t-1)

- 1.Demographics, excluding race (e.g., biological sex, age)
- 2. Insurance type
- 3.IDC-9 codes (International Statistical Classification of Diseases)
- 4. Prescribed medications
- 5.Medical service encounters (e.g., surgical, radiology),
- 6.Billed amounts, categorized by type (e.g., outpatient specialists, dialysis)

Obermeyer, Z., Powers, B., Vogeli, C., & Mullainathan, S. (2019). Dissecting racial bias in an algorithm used to manage the health of populations. Science, 366(6464), 447-453.



Percentage of Black patients in group with highest predicted risk







- At every level of algorithm-predicted risk, Blacks and Whites have (roughly) the same costs the following year
- we find substantial disparities in health conditional on risk but little disparity in costs.







# **Tutorial Overview**

- Sources of sociodemographic bias in machine learning (ML)
- Examples of studies examining sociodemographic bias in ML
  - Health Electronic Health Records
  - Family well-being Speech, language, physiology
- Approaches for mitigating sociodemographic bias in ML
  - Adversarial learning
  - Fairness regularization
  - Explainability







✓ quality of interactions



electrodermal activity (2)
✓ skin conductance level
✓ skin conductance response



context and interaction (11)

- ✓ GPS
- ✓ activity count
- ✓ body temperature
- ✓ alcohol/caffeine/drugs

THE USC COUPLE MOBILE SENSING PROJECT

http://homedata.github.io/



physiological synchrony (2)
✓ joint sparse representation
✓ multiple time scales





#### language use

- ✓ linguistic constructs (25)
- ✓ psychological factors (32)
- ✓ personal concern (7)
- ✓ paralinguistic (3)



acoustic analysis (8) ✓ pitch (F0) ✓ intensity







- 50 couples, 1438 samples (548 conflict)
- 5-fold stratified cross-validation
- 90.7% precision for non-conflict, 19.5% precision of conflict, 55.15% balanced accuracy





Participants' demographic distribution







Adherence of self-reports of interpersonal conflict with respect to race







Adherence of self-reports of interpersonal conflict with respect to age







Discrepancies between self-reported and algorithm-detected conflict with respect to race







Discrepancies between self-reported and algorithm-detected conflict with respect to age







# **Tutorial Overview**

- Sources of sociodemographic bias in machine learning (ML)
- Examples of studies examining sociodemographic bias in ML
  - Health Electronic Health Records
  - Family well-being Speech, language, physiology
- Approaches for mitigating sociodemographic bias in ML
  - Adversarial learning
  - Fairness regularization
  - Explainability





# **Tutorial Overview**

- Sources of sociodemographic bias in machine learning (ML)
- Examples of studies examining sociodemographic bias in ML
  - Health Electronic Health Records
  - Family well-being Speech, language, physiology
- Approaches for mitigating sociodemographic bias in ML
  - Adversarial learning
  - Fairness regularization
  - Explainability





#### Overarching research questions

- Can we characterize user re-identification risk in human behavior recognition models?
- Can we learn anonymized signal transformations that preserve behavioral information?

#### Case study

• Anonymized models of facial emotion recognition

#### Challenges

- Necessary to capture the subtlety of emotional expression
- Images captured in close proximity to user's face
  - Iterative adversarial learning with alternate training between minimizing emotion classification cost and maximizing user classification cost



Arora & Chaspari, ACM ICMI 2018 Narula & Chaspari, ACM ICMI, 2020





#### Japanese Female Facial Expression (JAFFE)

- 213 images
- 7 emotions

#### **YALE Face Dataset**

- 60 images
- 4 emotions

#### IEMOCAP

- 32,902 images
- 4 emotions



• Data augmentation through rotation, horizontal flip, and noise for JAFFE and YALE





#### RQ1: Can we characterize user re-identification risk?

• Emotion: CNN trained only on emotion classification



• Emotion to Face: CNN trained on emotion classification and fine tuned on user classification







#### **RQ1: Can we characterize user re-identification risk?**

- User classification: CNN trained on user classification only
- Emotion to user classification: CNN trained on emotion classification and fine tuned on user classification



Significant amount of user-dependent information preserved in facial emotion classification models (even after blurring)





# **RQ2: Can we learn anonymized signal transformations that preserve behavioral information?**

- Anonymizing the convolutional transformation so that user re-identification is not possible
- Alternate training between emotion and user classification losses



 $\min_{\{\mathbf{U}_{\mathbf{c}},\mathbf{U}_{\mathbf{e}},\mathbf{U}_{\mathbf{i}}\}}\{L_{e}\left(g_{\mathbf{U}_{\mathbf{c}}}(\mathbf{x})\right),y_{e}\right)-\alpha L_{i}\left(g_{\mathbf{U}_{\mathbf{i}}}(g_{\mathbf{U}_{\mathbf{c}}}(\mathbf{x})),y_{i}\right)\}$ 





# **RQ2: Can we learn anonymized signal transformations that preserve behavioral information?**

- Original: Original images
- Emotion classification: Images transformed by the emotion classification models
- **Proposed:** Images transformed by the proposed anonymization approach







### **RQ2:** Can we learn anonymized signal transformations that preserve behavioral information?

- Baseline: Adversarial learning without alternate training
- Proposed: Adversarial learning with alternate training







# **Tutorial Overview**

- Sources of sociodemographic bias in machine learning (ML)
- Examples of studies examining sociodemographic bias in ML
  - Health Electronic Health Records
  - Family well-being Speech, language, physiology
- Approaches for mitigating sociodemographic bias in ML
  - Adversarial learning
  - Fairness regularization
  - Explainability
- General recommendations for fair machine learning





### Fairness-aware learning through regularization

• Discrimination score

$$CV = Pr[Y = HighRisk | S = SensitiveGroup] - Pr[Y = HighRisk | S = NonSensitiveGroup]$$

$$\uparrow$$
Predicted risk of sensitive group
Predicted risk of non-sensitive group

• Fairness-aware regularization



Kamishima, T., Akaho, S., & Sakuma, J. (2011, December). Fairness-aware learning through regularization approach. In 2011 IEEE 11th International Conference on Data Mining Workshops (pp. 643-650). IEEE.





#### Comparison between adversarial learning and fairness regularization

- Gender de-biasing in speech emotion recognition
- The non-adversarial approach maintains equality of odds (CCC) levels similar to no bias mitigation
- The adversarial approach yields lower CCC



Gorrostieta, C., Lotfian, R., Taylor, K., Brutti, R., & Kane, J. (2019). Gender De-Biasing in Speech Emotion Recognition. In INTERSPEECH (pp. 2823-2827).





#### Comparison between adversarial learning and fairness regularization

• The non-adversarial approach, on the other hand, achieves much better consistency with all metrics



Gorrostieta, C., Lotfian, R., Taylor, K., Brutti, R., & Kane, J. (2019). Gender De-Biasing in Speech Emotion Recognition. In INTERSPEECH (pp. 2823-2827).





# **Tutorial Overview**

- Sources of sociodemographic bias in machine learning (ML)
- Examples of studies examining sociodemographic bias in ML
  - Health Electronic Health Records
  - Family well-being Speech, language, physiology
- Approaches for mitigating sociodemographic bias in ML
  - Adversarial learning
  - Fairness regularization
  - Explainability
- General recommendations for fair machine learning





# Improving fairness via explainability

• Understanding inner mechanisms of the model via explainable methods





# Improving fairness via explainability

#### Local Interpretable Model-agnostic Explanations (LIME)

Ā M

- Presenting artifacts that provide qualitative understanding of the relationship between the instance's components (e.g. words in text, patches in an image) and the model's prediction via explainable methods
- Quantify trust in individual predictions and entire model







# **Tutorial Overview**

- Sources of sociodemographic bias in machine learning (ML)
- Examples of studies examining sociodemographic bias in ML
  - Health Electronic Health Records
  - Family well-being Speech, language, physiology
- Approaches for mitigating sociodemographic bias in ML
  - Adversarial learning
  - Fairness regularization
  - Explainability
- General recommendations for fair machine learning





### General recommendations for fair machine learning

#### Design

- Define the goal of the machine learning model and review with diverse stakeholders
- Discuss ethical concerns of how the model could be used and what are the protective groups, also informed by historical data







### General recommendations for fair machine learning

#### Data collection

- Collect and document training data
- Ensure that participants in the protected group can be identified
- Make sure that the protected group is adequately represented in terms of numbers and features





### General recommendations for fair machine learning

#### Training & Evaluation

A M

- Train the model to take into account fairness goals
- Measure algorithmic output differences between sensitive and non-sensitive groups
- Assess model output with diverse stakeholders

#### Deployment

- Systematically review data and continuously evaluate metrics
- Collect feedback from participants and stakeholders

