



Algorithm-Based Clinical Decision Support (ABCDS) Oversight

Duke Health's Approach to Operationalizing & Governing Health AI


Michael Pencina, PhD
Michael Cary, PhD, RN
Nicoleta J. Economou, PhD

September 26, 2023

1

Agenda

- Promise of AI in Healthcare and the Current Landscape
- Algorithm-Based Clinical Decision Support (ABCDS) Oversight
- Bias Mitigation Strategies
- Benefits and Learnings from the Implementation of an Algorithmic Oversight Framework



© 2023 Duke University School of Medicine. All rights reserved.

2

Promise of Artificial Intelligence/Machine Learning in Health Care



© 2023 Duke University School of Medicine. All rights reserved.

3

“Wild West” of Algorithms

“We have a Wild West of algorithms,” said Michael Pencina, coalition [CHAI] co-founder and director of Duke AI Health. There’s so much focus on development and technological progress and not enough attention to its value, quality, ethical principles or health equity implications.”

Politico, April 4, 2023



© 2023 Duke University School of Medicine. All rights reserved.

4

AI/ML Risks

Research

JAMA Internal Medicine | Original Investigation

External Validation of a Widely Implemented Proprietary Sepsis Prediction Model in Hospitalized Patients

Andrew Wong, MD, Erkin Odes, MEng, John P. Donnelly, PhD, Andrew Krumm, PhD, Jeffrey McCullough, PhD, Olivia DeTroyer-Cooley, BSE, Justin Pestrue, MEd, Marie Phillips, BA, Judy Korye, MSN, RN, Carleen Perozo, MHA, RN, Muhammad Ghous, MBBS, Karandeep Singh, MD, MMSc

IMPORTANCE The Epic Sepsis Model (ESM), a proprietary sepsis prediction model, is implemented at hundreds of US hospitals. The ESM's ability to identify patients with sepsis has not been adequately evaluated despite widespread use.

OBJECTIVE To externally validate the ESM in the prediction of sepsis and evaluate its potential clinical value compared with usual care.

DESIGN, SETTING, AND PARTICIPANTS This retrospective cohort study was conducted among 27 697 patients aged 18 years or older admitted to Michigan Medicine, the academic health system of the University of Michigan, Ann Arbor, with 38 455 hospitalizations between December 6, 2018, and October 20, 2019.

EXPOSURE The ESM score, calculated every 15 minutes.

MAIN RESULTS AND MEASURES Sepsis, as defined by a composite of (1) the Centers for Disease Control and Prevention surveillance criteria and (2) International Statistical Classification of Diseases and Related Health Problems, Tenth Revision diagnostic codes accompanied by 2 systemic inflammatory response syndrome criteria and 1 organ...

Editorial page 1040
Multimedia
Supplemental content
CME Quiz at jamacmelookup.com and CME Questions page 1148

RESEARCH

RESEARCH ARTICLE

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

Ziad Obermeyer^{1,2*}, Brian Powers³, Christine Vogels⁴, Sendhil Mullainathan^{2,3}

Health care decisions affect the health of millions of people. But algorithms used to make these decisions are often biased, because they rely on data that may reflect our own biases. Our dataset describes one such typical algorithm. It contains both the algorithm's predictions as well as the data needed to understand its inner workings: that is, the underlying ingredients used to form the algorithm (data, objective function, etc.) and links to a rich set of outcome data. Because we have the inputs, outputs, and eventual outcomes, our data allow us to probe the algorithm's...

“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...”



© 2023 Duke University School of Medicine. All rights reserved.

5

We need to do better

Prediction Models — Development, Evaluation, and Clinical Application

Michael J. Pencina, Ph.D., Benjamin A. Goldstein, Ph.D., and Ralph B. D'Agostino, Ph.D.

“Given the number of emerging prediction models and their diverse applications, no single regulatory agency can review them all. This limitation, however, does not absolve models’ developers and users from applying the utmost scrutiny in demonstrating effectiveness and safety.”

When a prediction model is developed, it is often based on data from a specific population. However, the model's performance may vary when applied to a different population. This is because the model's performance is dependent on the quality and quantity of the data used to train it. Today, we have a large number of prediction models available, but we do not have a standard way to evaluate them. This is a problem because we need to know if a model is accurate and safe before we use it. We need to do better.



Pencina MJ, Goldstein BA, D'Agostino RB. *N Engl J Med*. 2020 Apr 23;382(17):1583-1586. doi: 10.1056/NEJMp2000589. © 2023 Duke University School of Medicine. All rights reserved.

6

Considerations for CDS development

- Population at risk
- Outcome of interest
- Time horizon
- Predictors
- Mathematical model
- Model evaluation
- Translation to CDS
- Clinical implementation



Pencina MJ, Goldstein BA, D'Agostino RB. Prediction Models – Development, Evaluation, and Clinical Application. NEJM. 2020;382:1583-1586. doi:10.1056/NEJMp2000589.
© 2023 Duke University School of Medicine. All rights reserved.

7

Principles for Responsible AI

- Ensure that AI technology serves humans
- Define the task we want the AI tool to accomplish
- Describe what the successful use of the AI tool looks like
- Create transparent systems for continuously testing and monitoring AI tools



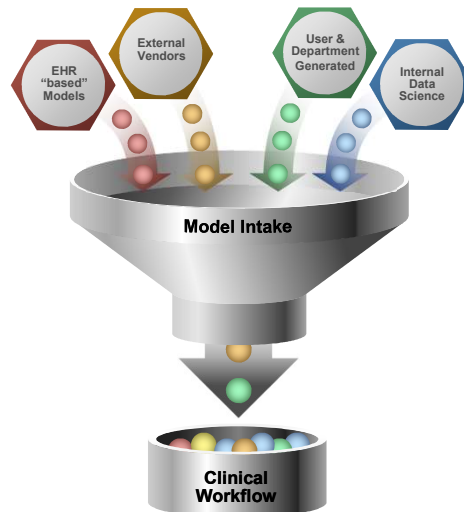
© 2023 Duke University School of Medicine. All rights reserved.

8

Complex environment

Different:

- Skills
- Knowledge bases
- Resources available
- Make up of project teams



© 2023 Duke University School of Medicine. All rights reserved.

9

Formation of the ABCDS Oversight Committee

In recognition of this changing landscape, the Duke Health Chancellor and the Dean of the School of Medicine charged Duke Health leadership to form an oversight framework.



© 2023 Duke University School of Medicine. All rights reserved.

10

Mission Statement

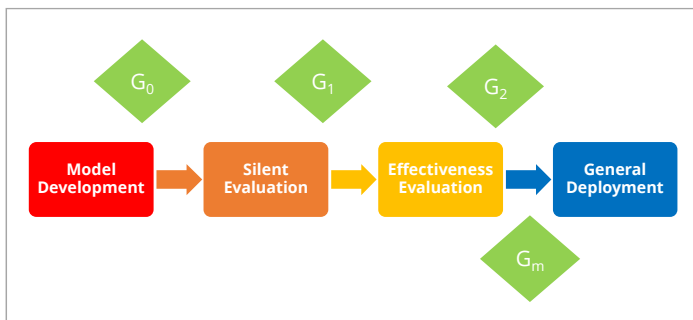
“Out of our primary focus on patient safety and high-quality care, our mission is to guide algorithm-based clinical decision support (ABCDS) tools through their lifecycle by providing governance, evaluation, and monitoring.”



© 2023 Duke University School of Medicine. All rights reserved.

11

ABCDS Lifecycle & Our Framework



'Just-in-time' Check-Points (**G**ates) Help Model Owners Get Ready for What's Ahead



- What are the clinical outcome and performance metrics?
- How has the model been evaluated?
- Who is the Clinical Owner?
- Who will cover maintenance costs in production?
- Will this ABCDS tool be used outside of Duke Health?
- Is this a standard of care model?
- How will the model be used in the clinic and how is it integrated with the workflow?

© 2023 Duke University School of Medicine. All rights reserved.

12

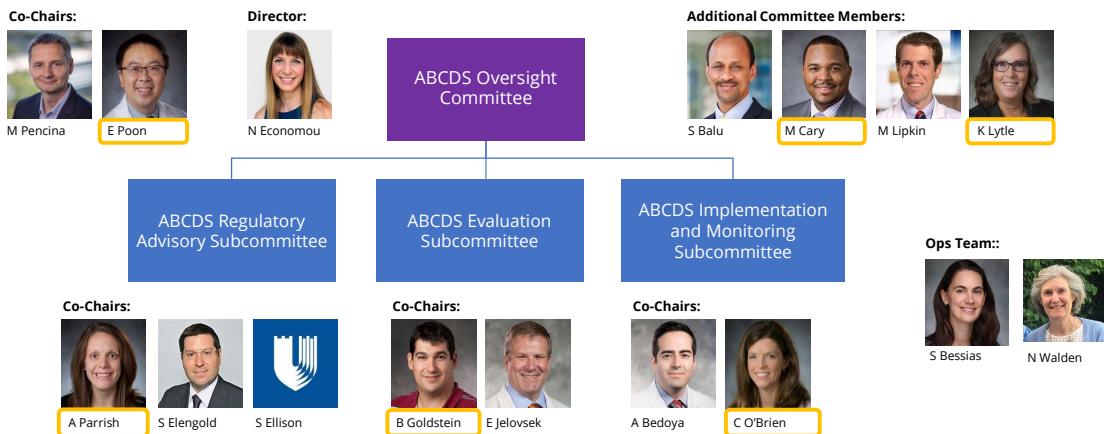
People



© 2023 Duke University School of Medicine. All rights reserved.

13

People: ABCDS Oversight Committee



© 2023 Duke University School of Medicine. All rights reserved.

14

Process



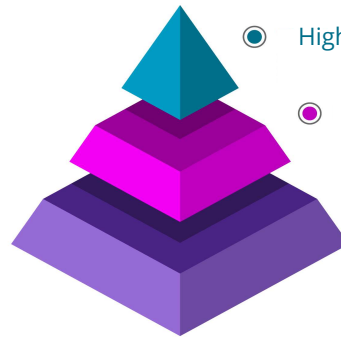
© 2023 Duke University School of Medicine. All rights reserved.

15

Scope of ABCDS Oversight Framework

ABCDS Tool = Algorithm(s) + Interface Algorithms Are Presented In

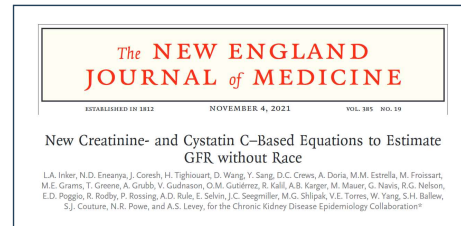
All electronic algorithms that could impact patient care at Duke Health fall within the scope of the ABCDS Oversight Committee and must undergo registration.



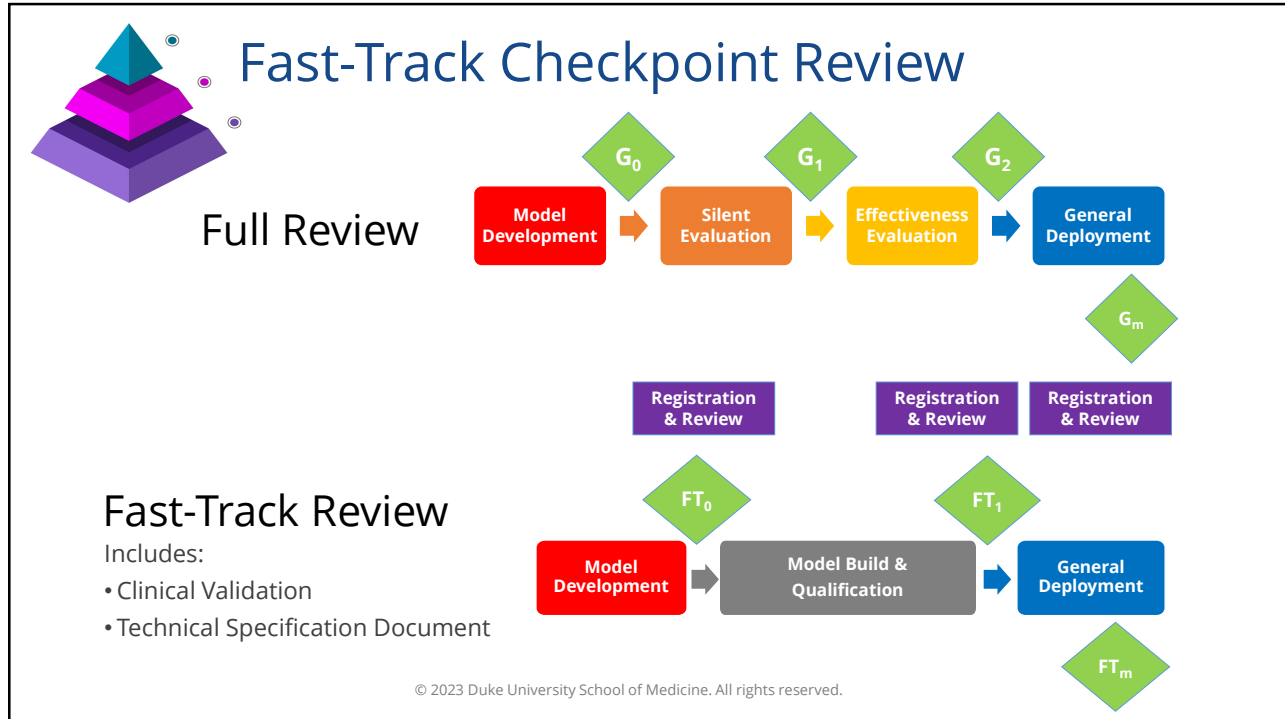
- High Risk: Data-Derived
- Medium Risk (e.g., Clinical Consensus)
- Low Risk: Standard of Care



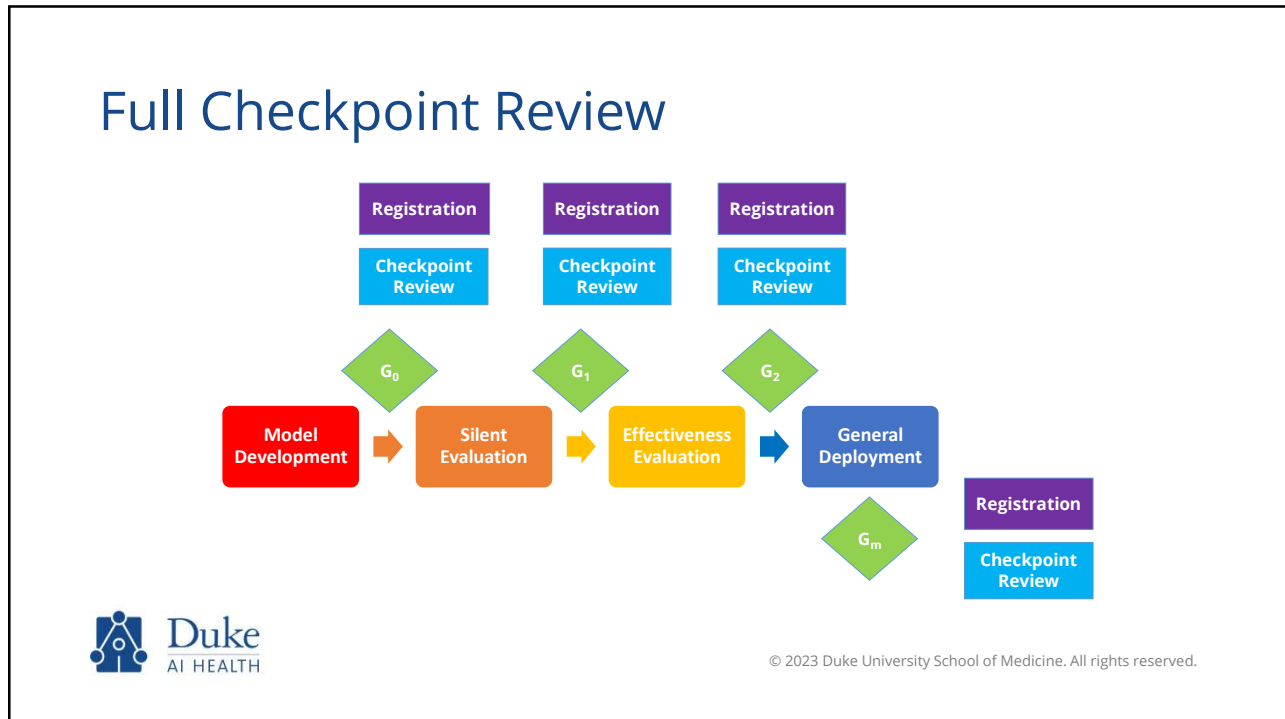
© 2023 Duke University School of Medicine. All rights reserved.



16



17



18

What to Expect: ABCDS Checkpoint Review

Registration

- Pre-Registration
- Triage

Full Review (Asynchronous)

- Preview Letter
- Review Meeting w/ Committee (Optional)

Outcomes:

- Approval
- Approval with Contingencies
- Re-review
- Denial

Outcome

- Outcome Letter

© 2023 Duke University School of Medicine. All rights reserved.

19

Portfolio Metrics (April 2023)

ABCDS Model Registration & Review	Total
Number of registered tools	52
Number of evaluated tools	31

Lifecycle Phase	Count
Model Development	10
Silent Evaluation	12
Effectiveness Evaluation	8
General Deployment	14
Out of Scope	8

© 2023 Duke University School of Medicine. All rights reserved.

20

Implementing Quality & Ethics with Our Framework

Transparency & Accountability

Impact & Safety

Fairness & Equity

Usability & Adoption

Regulatory Compliance

© 2023 Duke University School of Medicine. All rights reserved.

21

Implementing Quality & Ethics with Our Framework

Principle	Criteria	Submission Materials
<p>Transparency & Accountability</p> <p style="color: green;">Impact & Safety</p> <p>Fairness & Equity</p> <p>Usability & Adoption</p> <p>Regulatory Compliance</p>	<p>Clinical Impact & Safety</p> <p>The ABCDS software, in comparison to current state, stands to improve clinical care.</p> <p>Plans for Silent Evaluation will inform the decision to proceed with pilot implementation in clinic.</p>	<ul style="list-style-type: none"> ✓ Evidence that the tool has potential to impact clinical outcomes or processes ✓ List of key impact metrics (clinical outcomes and/or process improvement) with definitions, following TRIPOD guidelines⁵ ✓ List of core performance metrics (e.g. sensitivity, PPV, etc.) and results from development ✓ Calibration curves, threshold selections and justification if applicable <p>Silent Evaluation Plan, including:</p> <ul style="list-style-type: none"> ✓ Summary of benefits you expect to demonstrate and criteria to proceed into Effectiveness Evaluation ✓ Study design, including in/exclusion criteria, timeframe and sample size considerations ✓ Core performance metrics with shell tables ✓ Data analysis plan ✓ Data quality evaluation plan

Sample evaluation criteria supporting the principle of clinical impact and safety at the G₀ Checkpoint evaluation between pilot implementation and general deployment

© 2023 Duke University School of Medicine. All rights reserved.

(Unpublished work)

22

Implementing Quality & Ethics with Our Framework

Transparency & Accountability

Impact & Safety

Fairness & Equity

Usability & Adoption

Regulatory Compliance

Principle	Criteria	Submission Materials
Fairness & Equity	The principles of fairness and equity are reflected in the development process.	<ul style="list-style-type: none"> ✓ Summary of fairness & equity considerations during development ✓ Strategy for subgroup analysis with list of key stratification variables (e.g. race, age group, etc.) and rationale ✓ Subgroup analysis of chosen impact metrics ✓ Subgroup analysis of chosen model performance metrics ✓ Interpretation of findings and recommendations

Sample evaluation criteria supporting the principle of clinical impact and safety at the G_0 Checkpoint evaluation between pilot implementation and general deployment

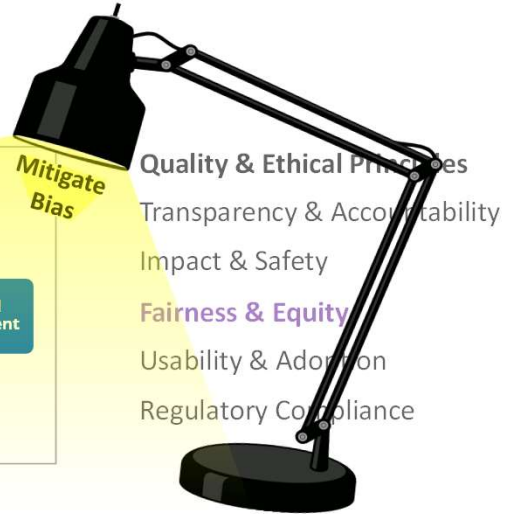
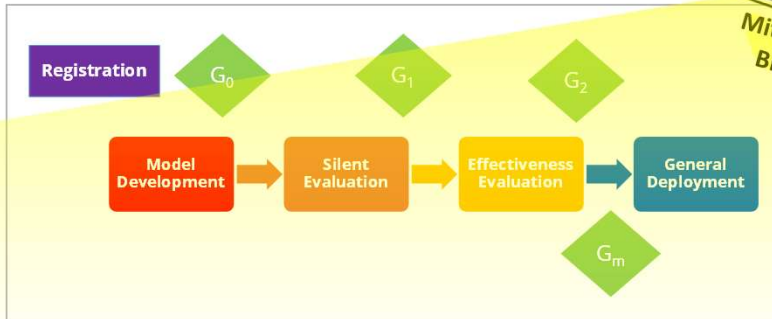


© 2023 Duke University School of Medicine. All rights reserved.
(Unpublished work)

23

Mitigating Bias Through Algorithmic Oversight

ABCDS Oversight process for the governance, evaluation and monitoring of algorithms to be deployed at Duke Health

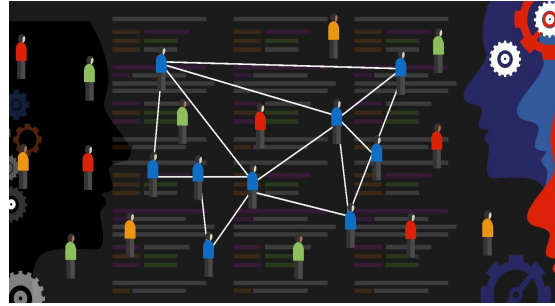


© 2023 Duke University School of Medicine. All rights reserved.

24

What is Bias in Clinical Algorithms?

Bias refers to the difference in how one or more subgroups is treated, represented or perceived, resulting in unfair/unjust outcomes.



Accessed on July 25, 2023, <https://nvlpubs.nist.gov/nistpubs/ai/NIST.AI.100-1.pdf>.
© 2023 Duke University School of Medicine. All rights reserved.

25

Societal Bias

Bias Type	Example	Assessment	Mitigation Strategy
Societal Bias Bias due to training data shaped by present and historical inequities and their fundamental causes	Predictive policing algorithms ¹ are trained on data that reflects structural racism and criminalization of, e.g., homelessness and poverty. Groups that are more likely to interact with the police are more likely to be identified by policing algorithms as “at risk” for future offense.	<i>Please discuss the real-world inequities reflected in your training data and how they inform the problem formulation and intended purpose of your model.</i>	<ul style="list-style-type: none"> • <i>Restriction to particular settings or use cases</i> • <i>Human-in-the-loop deployment design</i> • <i>Multi-stakeholder engagement</i>

- Label Bias
- Aggregation Bias
- Learning Bias
- Representation Bias
- Evaluation Bias
- Human Use Bias



¹ Angwin, J. Larson, S. Mattu, L. Kirchner, “Machine bias: There’s software used across the country to predict future criminals. And it’s biased against blacks,” *ProPublica*, 23 May 2016; www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing.
© 2023 Duke University School of Medicine. All rights reserved.

26

Label Bias

Bias Type	Example	Assessment	Mitigation Strategy
Label Bias Use of a biased proxy target variable in place of the ideal prediction target.	An algorithm ¹ used to identify patients for high-risk care management services predict healthcare costs as a proxy for healthcare <i>need</i> . Despite having greater health needs, Black patients have lower average healthcare spending (due to structural barriers in access to care) and are thus less likely to be recognized by the algorithm as 'high risk.'	<i>Please discuss any proxies used as inputs or outputs. Provide a rationale and describe implications for use.</i>	<ul style="list-style-type: none"> Eliminating proxies (where possible) or choosing a proxy as close as possible to the intended idea or concept



¹Obermeyer Z, Powers B, Vogeli C, Mullainathan S. Dissecting racial bias in an algorithm used to manage the health of populations. *Science*. 2019 Oct 25;366(6464):447-453. doi: 10.1126/science.aax2342. © 2023 Duke University School of Medicine. All rights reserved.

27

Why is it Important to Identify Racial/Ethnic Bias in Health Algorithms?

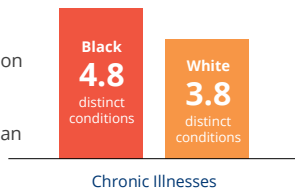
Algorithms are used to identify patients with complex health needs in order to provide more comprehensive care management. However, these algorithms can exhibit significant racial bias.

A 2019 study of one such algorithm found:



Black patients who are considerably sicker than White patients are given the same risk score

At the risk level that would result in automatic identification for the care management program, Black patients had **26%** more chronic illnesses than White patients.

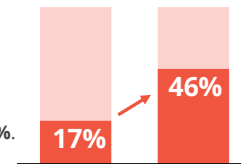


Why is this?



This algorithm assigned risk scores based on past health care spending. Black patients have lower spending than White patients for a given level of health.

If this bias was eliminated, the percentage of Black patients automatically enrolled in the program would rise from **17%** to **46%**.



¹Obermeyer Z, Powers B, Vogeli C, Mullainathan S. Dissecting racial bias in an algorithm used to manage the health of populations. *Science*. 2019 Oct 25;366(6464):447-453. doi: 10.1126/science.aax2342. © 2023 Duke University School of Medicine. All rights reserved.

28

Aggregation Bias

Bias Type	Example	Assessment	Mitigation Strategy
<p>Aggregation Bias</p> <p>Bias due to use of a one-size-fits-all model for data in which there are underlying groups or types of examples.</p>	<p>A natural language processing (NLP) model developed to scan clinical notes and suggest medication review is used across hospitals in a large health system in which documentation practices differ between locations, leading to poor performance in recently-acquired rural hospitals switching EHR systems.</p>	<p><i>Please discuss the ways that the data used to train your model may be observed differently across subgroups.</i></p>	<ul style="list-style-type: none"> • Use of subpopulation-specific models instead of or in addition to one-size-fits-all models • Use of subgroup-specific thresholds in a one-size-fits-all model • Imputation or other strategies to improve mapping from inputs to labels across subgroups



© 2023 Duke University School of Medicine. All rights reserved.

29

Learning Bias

Bias Type	Example	Assessment	Mitigation Strategy
<p>Learning Bias</p> <p>Bias due to modeling choices that amplify performance disparities across subgroups.</p>	<p>A development team is working on prediction of asthma exacerbation and uses a variety of methods to generate candidate models. The final model is selected by ranking the candidates on a single performance metric, AUROC. The focus on a single summary metric conceals large performance differences by race leading to poor prediction in the demographic most exposed to environmental asthma triggers.</p>	<p><i>Please describe how the model was optimized and the performance metrics used among candidate models.</i></p>	<ul style="list-style-type: none"> • Penalized optimization methods • Subgroup analysis to inform model selection



© 2023 Duke University School of Medicine. All rights reserved.

30

Representation Bias

Bias Type	Example	Assessment	Mitigation Strategy
<p>Representation Bias</p> <p>Bias emerging from non-representative training data which can lead to poor performance in subsets of the deployment population.</p>	<p>A melanoma detection model¹ achieved accuracy parity with a board-certified dermatologist; however, the model was trained primarily on light-colored skin. As such, the algorithm is likely to underperform for patients with dark skin.</p>	<p><i>Please discuss the quality and representativeness of your training data.</i></p> <p><i>If your model is adaptive, please discuss how you will ensure representativeness of the training data on an ongoing basis.</i></p>	<ul style="list-style-type: none"> • Integration with data from other sources • Supplementation with synthetic data • Up- or down-sampling approaches • Acknowledgement of limitations in model brief or other training materials • Refitting an out-of-the-box model to the local population



¹Wang HE, et al. A bias evaluation checklist for predictive models and its pilot application for 30-day hospital readmission models. *J Am Med Inform Assoc.* 2022 Jul 12;29(8):1323-1333. doi: 10.1093/jamia/ocac065. © 2023 Duke University School of Medicine. All rights reserved.

31

Evaluation Bias

Bias Type	Example	Assessment	Mitigation Strategy
<p>Evaluation Bias</p> <p>Bias emerging from a validation dataset that is not reflective of the deployment population and/or the training population.</p>	<p>A health system implements a new vendor model to predict in-hospital deterioration after receiving a validation report showing strong performance in other health systems that share the same EHR. Once the model is connected to the local data source, it produces an unexpected number of false alerts.</p>	<p><i>Briefly summarize plans for local validation.</i></p>	<ul style="list-style-type: none"> • Local validation (required) • Re-fitting the model on development sample that better represents the deployment population • Post-deployment monitoring with chart review (required)



© 2023 Duke University School of Medicine. All rights reserved.

32

Human Use Bias

Bias Type	Example	Assessment	Mitigation Strategy
<p>Human Use Bias</p> <p>Inconsistent user response to algorithm outputs for different subgroups.</p>	<p>A machine learning algorithm¹ developed to help pathologists differentiate liver cancer types did not improve every pathologist's accuracy despite the model's high rate of correct classification. Instead, pathologists' accuracy was improved when the model's prediction was correct but decreased when the model's prediction was incorrect.</p>	<p><i>Briefly describe how your algorithm fits into the clinical workflow. If it will replace an existing model or process, please include a comparison to baseline.</i></p>	<ul style="list-style-type: none"> • Workflow design solutions • End user training • Post-deployment monitoring with chart review (required) • Collection of end user feedback and metrics of adoption

- Label Bias
- Aggregation Bias
- Learning Bias
- Representation Bias
- Evaluation Bias
- Human Use Bias



Wang HE, et al. A bias evaluation checklist for predictive models and its pilot application for 30-day hospital readmission models. *J Am Med Inform Assoc.* 2022 Jul 12;29(8):1323-1333. doi: 10.1093/jamia/ocac065. © 2023 Duke University School of Medicine. All rights reserved.

33

Impacting How We Deliver Patient Care



© 2023 Duke University School of Medicine. All rights reserved.

34

Lessons Learned

- **Successful AI Governance is a Team Sport**
 - ✓ Lots of skillsets, perspectives and languages to bring together
- **Culture Shift is Hard**
 - ✓ Show Teams how to succeed by addressing gaps in their knowledge, skillsets, and/or bandwidth
 - ✓ Governance's role is Coach and Facilitator (not Punisher)
 - ✓ There is no such thing as over-communication in a complex system
- **Benefits of Centralized Governance**
 - ✓ Transparency of process expectations
 - ✓ Institutional visibility into all the 'skeletons in the closet'
- **Conscious Decision (thus far) Not to Regulate Who Gets to Build AI Models**



© 2023 Duke University School of Medicine. All rights reserved.

35

Future Directions

- QMS
 - SOPs
 - Centralized model monitoring / safety surveillance
- Incorporating the patient voice

36

The screenshot shows two sections of the CHAI website. The top section is titled "Providing guidelines for the responsible use of AI in healthcare" and features a colorful brain graphic. Below it is a dark blue section titled "Our Purpose" with text describing the coalition's mission. The bottom section is titled "Workshop Papers" and includes a "Call for review" button and a "View Paper" button. Navigation links for "Learn More", "Insights", and "Join Us" are visible in the top right of both sections.

<https://www.coalitionforhealthai.org/>


<https://www.coalitionforhealthai.org/insights/>

© 2023 Duke University School of Medicine. All rights reserved.

37

Learn More...

<https://aihealth.duke.edu/algorithm-based-clinical-decision-support-abcds/>





What is ABCDS?

Algorithm-Based Clinical Decision Support (ABCDS) Oversight is a "people-process-technology" framework for the governance and evaluation of clinical algorithms created for use at Duke Health. This framework fosters innovative, safe, equitable, and high-quality patient care by introducing checkpoints throughout the development lifecycle as well as after deployment to ensure that transparency, quality, and ownership are maintained for ABCDS algorithms and tools. The ABCDS Oversight is a collaborative effort between the Duke University School of Medicine and the Duke University Health System.

Bedoya AD, et al. A framework for the oversight and local deployment of safe and high-quality prediction models. *J Am Med Inform Assoc.* 2022 Aug 16;29(9):1631-1636. doi: 10.1093/jamia/ocac078. PMID: 35641123; PMCID: PMC9382367.

Contact us at abcds@duke.edu or nicoleta.economou@duke.edu

38

Thank you



© 2023 Duke University School of Medicine. All rights reserved.