

UNIVERSITY OF WATERLOO

For Professor Maura Grossman

CS 492

How Biased Clinical Research Data Affects Healthcare

March 14, 2025

Jaskirat Pabla

Samin Sharar Nafi

Sahl Bakshi

Saboor Bakshi

Introduction

Artificial intelligence (AI) is becoming an integral part of modern healthcare, with applications in diagnostics, treatment planning, and patient monitoring. In fact, the technology is expanding so quickly that in 2022, almost 20% of U.S. hospitals integrated some form of AI (Baten & Abdul, 2022). However, AI systems rely on vast datasets, and if these datasets contain biases, the outcomes can be dangerously skewed. A biased AI model can misdiagnose illnesses, delay treatment, or even exacerbate healthcare disparities. In some cases, these biases can be the difference between life and death. Our project explores how biases in healthcare data, particularly racial and gender biases, influence AI decision-making and contribute to systemic inequalities in medical treatment.

To investigate this issue, we are interviewing researchers from SickKids and professors from the University of Waterloo and conducting research. Our goal is to gather expert insights on the impact of biased clinical datasets and explore potential solutions to mitigate these issues. The final project will likely take the form of a website that summarizes our research and interviews. Additionally, we may include an interactive feature (game) where users can engage with real-world case studies and consider how they would respond to the dilemmas and compare their decisions with decisions from AI-driven healthcare. This document, i.e., project update, will outline the research conducted so far, the specific case studies we plan to discuss, and the key questions we will pose during our interviews. Since most of our research has been completed at this stage, the remainder of the term will focus on conducting interviews, developing the website, and potentially implementing the interactive component.

Gender Bias in AI

How Bias Manifests

AI is being used more frequently in the healthcare sector in order to improve the diagnosis of illnesses and to personalize treatments, but gender biases in these systems can lead to unequal and dangerous outcomes.

There are 2 sources of this kind of bias. Firstly, there is the data bias, which is caused by a shortage of women in the training datasets. Datasets often overly represent male patients (Straw & Wu, 2022). For example - in a study for predicting liver disease, of the 583 participants, only 142 were female (Straw & Wu, 2020). The other source is algorithmic bias. That refers to incorrect predictions produced by the AI model that is a direct result of its learning process or design. This kind of bias can result in the model having unequal predictions for different genders (Norori et al, 2021).

Impact on Minority Groups

There was a study involving Parkinson's disease biomarkers, where only 18.6% of the participants were women (Cirillo et al., 2020). The models generated with the data from those

participants were significantly less accurate for female patients as opposed to male patients. Another example of this is that many supervised machine learning models for predicting liver disease had a significantly higher false negative rate for women, because it missed 44% of actual liver disease cases for women as opposed to only 23% in men (University College London, 2022). False negatives for liver disease detection can have significant negative impacts on female patients. This shows how biases can be extremely dangerous.

One study showed that AI models to predict COVID-19 severity performed significantly worse when trained on data with one gender and applied to another (Chung et al., 2021). This shows the importance of training AI models with diverse data. Historically, people have been very negligent toward women's health. 80% of the drugs that were withdrawn were withdrawn because they had unforeseen side effects in females (Joshi, 2024). AI technologies come with the risk of reinforcing these existing biases, and that can be extremely dangerous.

Case Studies

Case 1: Transgender Patient's Emergency Misdiagnosed by Algorithmic Bias

Scenario

A transgender man experienced a life-threatening delay in care due to biases in health records and protocols. The 32-year-old patient went to an emergency room in severe abdominal pain. Despite informing staff that he was transgender, the hospital's intake and record system listed him as "male," and clinicians initially assumed his pain was due to something like obesity – failing to consider pregnancy (Compton, 2019). In fact, he was pregnant and in labour complications. The oversight meant that a pregnancy test and urgent obstetric care were delayed. By the time a pregnancy was confirmed and an emergency C-section was ordered, the situation had worsened: tragically, the baby was stillborn (Ring, 2019).

Impact

This case shows how rigid or biased algorithms in electronic health records (EHRs) and triage can harm transgender and non-binary individuals. The system's binary classification and the providers' biases led to a critical misdiagnosis — treating the patient as a non-pregnant male by default. The NEJM report on this incident noted that a patient identified as a female with similar symptoms "would almost surely have been triaged and evaluated more urgently for pregnancy-related problems" (Ring, 2019). Because the patient was recorded as male, standard alerts or decision support for pregnancy never triggered, and providers' judgment was clouded by gender assumptions. The result was a catastrophic outcome that likely would have been prevented with more inclusive algorithms and training. In essence, the healthcare AI/IT infrastructure did not account for a trans man's reality, illustrating how algorithmic bias and lack of nuance in sex/gender data can lead to incorrect or delayed treatment for transgender patients.

Interview Questions

If you were a doctor and a patient had severe abdominal pain, but no outside wounds, what would you think the problem is?

What if you knew the person was a man?

What if you knew the person was a woman?

What if you knew the person was a transgender?

Knowing now that the person was transgender and pregnant, how would you solve the issue so that the AI would not make this mistake again?

Case 2: Symptom Checker Underestimates a Woman's Heart Attack

Scenario

An AI-driven symptom-checker app provided vastly different recommendations for a man and a woman with identical health inputs, revealing a dangerous gender bias. In an analysis by an NHS doctor, two hypothetical patients – one male, one female – both 59-year-old smokers with sudden chest pain and nausea, queried a popular health chatbot. The only difference was the gender selected. The female patient was told her symptoms might be due to depression or a panic attack, with no urgent action needed beyond maybe a GP visit (Trendall, 2024). In stark contrast, the male patient with the same profile was warned that it could be gastritis or even serious heart problems like unstable angina or heart attack – in which case he should seek emergency care or call an ambulance (Trendall, 2024). In other words, the AI did not even consider a cardiac emergency for the woman, whereas it did for the man, solely because of gender. This chatbot (used in the UK's "GP at Hand" service by Babylon Health) was purportedly basing its advice on statistical evidence that women's chest pain is less likely to be heart attack – but in doing so, it risked missing a real heart attack in a female patient (Trendall, 2024).

Impact

A female patient following this AI advice could have delayed going to the ER for a true heart attack, with potentially fatal consequences. Heart disease in women often goes underdiagnosed precisely because symptoms can present differently and biases lead to attributing them to anxiety or other causes. Here the AI essentially mirrored and amplified that bias. The public outcry around this example was significant – observers were "deeply concerned" that the program failed to even raise the possibility of a heart attack in the woman's case (Trendall, 2024). This case highlights how AI triage tools, if not carefully designed, can perpetuate harmful stereotypes (e.g. "women are hysterical, men have real heart attacks"), thus providing suboptimal or dangerous guidance. Babylon Health defended the system as operating as intended, citing medical data differences (Trendall, 2024). However, even if statistically fewer women present with classic heart attacks, many do – and an AI that dismisses women's cardiac symptoms can lead to delayed treatment, poorer outcomes, or even preventable death for female patients. It underscores the need for AI in healthcare to be rigorously tested for gender bias and for algorithms to err on the side of caution with life-threatening possibilities for all patients.

Interview Questions

If a 59-year-old patient, regardless of gender, presents with sudden chest pain and nausea, what would be your first thoughts on possible diagnoses?

What if I told you this patient was a man and women's chest pain is less likely to be heart attack?

If this was a woman, and knowing that women's chest pain is less likely to be heart attack, would you still take the case seriously, have precautionary backup options to the patient, just in case?

How do you think an AI symptom-checker should handle cases where statistical likelihoods differ between genders? Should it always list all serious possibilities?

Do you think this AI's response was justified based on medical data, or do you believe it failed in its duty to provide safe recommendations? Why?

Case 3: Male Breast Cancer Patient Denied Treatment by Gendered Algorithm

Scenario

Raymond Johnson, a 26-year-old man in South Carolina, faced a life-threatening algorithmic bias after being diagnosed with breast cancer. When he applied to a federal Medicaid program for breast cancer treatment, he was denied coverage solely because he was a male (Park, 2022). The program (created by the Breast and Cervical Cancer Prevention and Treatment Act) was designed to cover cancer care for patients diagnosed via federal screening programs – but those programs only screened women, so by policy only women qualified for treatment coverage (Park, 2022). In Raymond's case, the insurance algorithm automatically excluded men, rendering him ineligible for the chemotherapy and surgery he urgently needed simply due to his gender (Park, 2022).

Impact

This is a clear example of gender bias embedded in a treatment plan algorithm. Raymond was left to struggle for access to life-saving care because the system failed to consider that men can get breast cancer. Denying benefits based solely on gender meant a potentially deadly delay or enormous out-of-pocket costs for his treatment. Advocacy groups intervened; the ACLU condemned the policy, noting that refusing cancer coverage to a patient "simply because they are men" is a blatant violation of law and basic fairness (Park, 2022). Raymond's case not only illustrates bias against a male patient with a so-called "women's disease," but also led to calls for policy change so that diagnostic and coverage algorithms include all patients who need care (Park, 2022). It underlines how men can also be harmed when medical algorithms or guidelines incorrectly treat certain serious conditions as "female-only," resulting in suboptimal or delayed care for male patients.

Interview Questions

If a young patient presents with a lump in their chest, what factors would you consider when assessing their condition?

If they were a female, what would your decision lean more towards?

What if they were male?

Does the fact that breast cancer is rarer in men impact your decision to recommend testing for it? Why or why not?

How do you think the assumption that breast cancer is a "women's disease" influenced the Medicaid program's algorithmic decision to deny coverage?

How could the algorithm and policy be modified to ensure that men with breast cancer receive equal access to treatment?

Do you think medical AI and insurance programs should be based strictly on statistical data trends, or should they incorporate flexibility for outlier cases? Why?

Racial Bias in AI

How Bias Manifests

Racial bias in medical AI occurs when algorithms aren't equally effective spanning different racial or ethnic groups. That often comes from the inequality in the data that is used to train the models or in the presumptions that were incorporated into the clinical algorithms. For instance, patients suffering from skin diseases and AI diagnostic tools for dermatology have demonstrated poor accuracy towards individuals who have darker skin due to the predominantly lighter-skinned image datasets used to train the models (Nicholls, 2022). That type of diagnostic bias entails AI that was not designed with black or brown patients in mind and would fail to recognize or in worse cases, misdiagnose them.

Biased Algorithms

Racial bias discrimination can also be seen in recommendations regarding medical treatment as well as the classification of risk factors and categories. A well-known example is an algorithm that was created in hospitals to flag patients who may be eligible for care management and monitoring (Manke, 2019). A particular study done in 2019 showed that the software would consistently favour white patients over more ill black patients because it utilized the proxy of healthcare spending as a substitute for health needs.

Historically Black patients did not have equal access to care therefore, in the eyes of the algorithm, the Black patients incurred less medical spending leading the algorithm to not appreciate black patients costing more due to having lower access to medical assistance. An algorithm that was used to aid an assessment of kidney disease also incorporated a correction with bias concerning race in which the function was assumed to be better than it was in coloured people. As a result, many black patients were inappropriately delayed for specialty referrals or consideration for transplant operations. The need for these patients was higher but the assumption made it more difficult for those who were in need leading to discrimination. Structural racism that is embedded in a data set because of poor integration through healthcare systems artificially creates inequities that allow the power of these systems to implement deeper discriminatory boundaries.

Biased Medical Devices

Racial biases don't just exist in automated software; they also extend to AI-based medical devices and sensors. For example, light-based pulse oximeters (devices that estimate blood oxygen levels) have been shown to overestimate oxygen saturation in patients who have darker skin (Department of Epidemiology & Biostatistics, UCSF, 2022). Research from the COVID-19 pandemic showed that pulse oximeters were three times more likely to not detect dangerously low oxygen levels in black patients when compared to white patients. The reason for this bias lies in how light absorption differs across people with varying skin tones (Allen, 2024). Another case is the blood oxygen level monitoring function in commercial smartwatches and fitness trackers: marketed wearable pulse oximeters were found to be much less precise in estimating the oxygen saturation level in darker skin. In the listed algorithms, tools and devices, prejudice may result in minority patients being overdiagnosed or receiving the wrong treatment recommendations.

Impact on Minority Groups

Black, Hispanic, and Indigenous communities tend to be impacted from the AI healthcare bias the most. These communities have had to deal with systemic inequalities and discrimination not just in education and employment, but healthcare, which includes access to services, representation in clinical studies, and discrimination (Penn Medicine, 2024). AI systems that carry forward biases from historical data risk worsening these disparities. Take, for instance, the biased risk algorithm described earlier that resulted in Black patients not receiving enough preventive care. With regards to kidney disease, the racial bias adjustment in kidney function score calculations resulted in Black patients being added to the transplant waiting list much slower than white patients who were equally qualified. This meant that Black patients were forced to wait years to receive transplants, while it is well known that Black Americans with end-stage kidney disease dramatically outnumber whites in their need for these organ donations.

The way technology suffers from bias against certain groups of people truly stood out in the case of the minority of patients during the COVID-19 pandemic. The use of pulse oximeters on patients with darker skin resulted in tragically poor treatment for patients severely suffering from COVID-19. One study estimated that such errors could have caused an average delay of 4.5 hours in Black patients receiving COVID treatment (Allen, 2024). The overestimation of blood oxygen levels which were used as a key metric to decide hospitalization and therapy were over before Black and other non-white patients life-saving oxygen or medications. The treatment given after these prolonged periods was insufficient and did not help avoid increased mortality rates in affected communities.

Discriminatory technology can deteriorate healthcare trust among marginalized communities. It's no secret that medicine has a history of being racist- from the Tuskegee syphilis experiment to current inequalities- and the application of AI technology that systematically ignores Black or Indigenous people serves to make them more cynical. Indigenous populations and other minority groups need to worry that algorithms based on

mostly white or urban populations will not capture their distinctive health characteristics and subsequently get diagnosed or treated as if they do not exist. In other words, the existing prejudice AI systems tend to have will likely widen the gap and worsen the care that minorities receive by delaying treatment, providing inaccurate diagnoses, and ignoring advanced care options.

Case Studies

Anthony Randall and Kidney Transplant Algorithm Bias

Anthony Randall is a Black man from Los Angeles who was on dialysis, waiting for a kidney transplant for over five years (TheGrio, 2023). What he did not know is that an algorithm from the transplant system incorporated a race-based “modifier” that made Black patients’ kidney scores seem better than they were. This modifier caused Randall’s kidney disease to be less severe than it truly was, leading to his placement on the national transplant waiting list being significantly delayed. In mid-2023, he filed a case against his hospital (Cedars Sinai Medical Center) and the United Network for Organ Sharing alleging that he was unfairly deprived of a fair chance to get the transplant because of the racially biased formula. It was no secret that the algorithm had a bias.

The board of the transplant system understood the modifier was resulting in Black patients’ illnesses being severely underestimated. By early 2023, all hospitals were directed to stop the usage of race adjustment and Black patients’ waiting times were to be changed to reflect the postponement. Randall claims that had these changes come sooner; he could have already had the kidney that he desperately needs. His case highlights how the goal of the clinical algorithm was good, but the execution was not due to the insertion of race which caused Black patients to not receive quality care in a timely manner.

Dr. Noha Aboelata and Pulse Oximeter Bias During COVID-19

During the pandemic, race-related biases in medical technology confronted Dr. Noha Aboelata, a family physician and the Chief Executive Officer of Roots Community Health Center based in Oakland. In late 2020, one of her patients was an elderly African American gentleman who suffered from chronic lung illness (Department of Epidemiology & Biostatistics, UCSF, 2022). One of the checks done previously, the pulse oxygenation check, revealed that his oxygen saturation levels were high. Even though the device showed a relatively normal oxygen level, Dr. Aboelata’s clinical instinct indicated the patient was much more distressed. She conducted an arterial blood gas test which confirmed her worst fears, the oxygen content in the patient’s blood was too low and he needed oxygen.

Sometime later, she came across an article in the New England Journal of Medicine that confirmed her hunch; the oximeters were unable to register low oxygen levels in dark-skinned patients as compared to white patients. Her and her colleagues were outraged by a device that was supposed to help their patients but was grossly inaccurate for the Black population. Finally,

her clinic participated in a class-action lawsuit against easier manufacturers and sellers of pulse oximeters for more detailed warnings and up-to-date devices. She did not stand idle while demanding the FDA take pulse oximeter discrimination towards races very seriously.

Alex Morales and Smartwatch Blood-Oxygen Reading Bias

Alex Morales, a New York resident, brought attention to a case of possible racial bias in the consumer health device, Apple Watch (Stempel, 2023). Morales, who has a darker complexion, bought an Apple Watch with the expectation that its blood oxygen sensor would accurately log his oxygen levels for fitness and health purposes. To his surprise, he later found out that the device's oximeter may not work well with people from his demographic.

In late 2022, Morales initiated a class action lawsuit against Apple for allegedly containing a blood oxygen app that was racially discriminatory and did not function as promised for non-white customers. Supported in part by complaints of other studies claiming that more advanced pulse oximetry devices are “massively” less useful on people with darker skin, Morales asserted that Apple owed the public an explanation. After all, paying smartphone users assumed that the device would be equal for all users, which is not the case. While the judge dismissed the case in 2023, it did start an important sole discussion regarding Apple products. Their case showed the world that there are, in fact, indirect biases in medical-grade equipment. This shows us that Alex Morales and the rest of the community are still subjected to discrimination based on race even in the technology they choose to use. Moreover, it demonstrates the keen eye for responsibility the tech industry has in these situations.

Interview Questions (to be done)

We will be developing interview questions for the above scenarios based on racial bias similar to ones created in the examples covered in gender biases.

References

Chung, H., Park, C., Kang, W. S., & Lee, J. (2021). Gender bias in artificial intelligence: Severity prediction at an early stage of COVID-19. *Frontiers in Physiology*, 12, Article 778720. <https://doi.org/10.3389/fphys.2021.778720>

Cirillo, D., Catuara-Solarz, S., Morey, C., Guney, E., Subirats, L., Mellino, S., Gigante, A., Valencia, A., Rementeria, M. J., Santuccione Chadha, A., & Mavridis, N. (2020). Sex and gender differences and biases in artificial intelligence for biomedicine and healthcare. *npj Digital Medicine*, 3, Article 81. <https://doi.org/10.1038/s41746-020-0288-5>

Compton, J. (2019, May 22). Trans dads tell doctors: “you can be a man and have a baby.” NBC News. <https://www.nbcnews.com/feature/nbc-out/trans-dads-tell-doctors-you-can-be-man-have-baby-n1006906>

Joshi, A. (2024). Big data and AI for gender equality in health: Bias is a big challenge. *Frontiers in Big Data*, 7, Article 1436019. <https://doi.org/10.3389/fdata.2024.1436019>

Norori, N., Hu, Q., Aellen, F. M., Faraci, F. D., & Tzovara, A. (2021). Addressing bias in big data and AI for health care: A call for open science. *Patterns*, 2(10), 100347. <https://doi.org/10.1016/j.patter.2021.100347> accuray.com

Park, S. (2022, September 8). Breast cancer doesn't discriminate against men. American Civil Liberties Union. <https://www.aclu.org/news/womens-rights/breast-cancer-doesnt-discriminate-against-men>

Ring, T. (2019, May 16). A pregnant trans man's misdiagnosis leads to Stillbirth. Advocate. <https://www.advocate.com/transgender/2019/5/16/pregnant-trans-mans-misdiagnosis-leads-stillbirth>

Straw, I., & Wu, H. (2022). Investigating for bias in healthcare algorithms: A sex-stratified analysis of supervised machine learning models in liver disease prediction. *BMJ Health & Care Informatics*, 29(1), e100457. <https://doi.org/10.1136/bmjhci-2021-100457>

Trendall, S. (2024, July 29). Gender bias concerns raised over GP app. Public Technology. <https://www.publictechnology.net/2019/09/13/health-and-social-care/gender-bias-concerns-raised-over-gp-app/>

University College London. (2022, July 11). Gender bias revealed in AI tools screening for liver disease. UCL News. <https://www.ucl.ac.uk/news/2022/jul/gender-bias-revealed-ai-tools-screening-liver-diseases-ucl.ac.uk>

Baten, R. B. A. (2024, October). How are US hospitals adopting artificial intelligence? Early evidence from 2022. *Health Affairs Scholar*, 2(10), qxae123.
<https://pubmed.ncbi.nlm.nih.gov/39403132/>

Nicholls, M. (2022, January 10). AI in skin cancer detection: Darker skin, inferior results? *Healthcare in Europe*.
<https://healthcare-in-europe.com/en/news/ai-in-skin-cancer-detection-darker-skin-inferior-results.html>

Manke, K. (2019, October 24). Widely used health care prediction algorithm biased against Black people. *Berkeley News*.
<https://news.berkeley.edu/2019/10/24/widely-used-health-care-prediction-algorithm-biased-against-black-people/>.

TheGrio Staff. (2023, April 11). Black man awaiting kidney transplant alleges racial bias. TheGrio.
<https://www.washingtonpost.com/health/2023/04/10/lawsuit-unos-kidney-transplant-race-discrimination/>.

Allen, A. (2024, October 7). FDA's promised guidance on pulse oximeters unlikely to end decades of racial bias. *KFF Health News*.
<https://kffhealthnews.org/news/article/pulse-oximeters-racial-bias-fda-rules/>.

Department of Epidemiology & Biostatistics, UCSF. (2022, September 29). Pulse oximeters don't work as well on darker skin, leading to flawed COVID care. *UCSF News (Epi/Biostatistics)*.
<https://epibiostat.ucsf.edu/news/pulse-oximeters-dont-work-well-darker-skin-leading-flawed-covid-care>

Stempel, J. (2023, August 21). Lawsuit claiming Apple Watch sensor exhibits "racial bias" is dismissed. *Reuters*.
<https://www.reuters.com/legal/lawsuit-claiming-apple-watch-sensor-exhibits-racial-bias-is-dismissed-2023-08-21>

Penn Medicine. (2024, March 25). Friend or foe: A closer look at the role of health care algorithms in racial and ethnic disparities. *Penn Medicine News Release*.
<https://www.pennmedicine.org/news/news-releases/2024/march/the-role-of-health-care-algorithms-in-racial-bias-inequity>