AMA Journal of Ethics[®]

January 2025, Volume 27, Number 1: E14-20

MEDICAL EDUCATION: PEER-REVIEWED ARTICLE What Should Health Professions Students Learn About Data Bias?

Douglas Shenson, MD, MPH, MA, MS, Beverley J. Sheares, MD, MS, and Chelesa Fearce

Abstract

In epidemiology, bias is defined as systematic deviation from the truth, and it can arise at different stages of scientific investigation (eg, data collection, methodological application, and outcomes analysis). Epidemiological bias can appear as a consequence of data bias (usually categorized as selection bias or information bias) or social bias (prejudice). Such forms of bias may occur separately or together. This article explores what health professions students should learn about the relationship between data bias and social bias—generated by racial, ethnic, gender, or other kinds of prejudice, singly or in combination—as a source of ethical and clinical concern in health care practices and policies that influence patient care and community health.

The American Medical Association designates this journal-based CME activity for a maximum of 1 AMA PRA Category 1 Credit[™] available through the AMA Ed Hub[™]. Physicians should claim only the credit commensurate with the extent of their participation in the activity.

Thinking Critically About Evidence

Bias is defined as the presence of systematic error in a study, and its adverse impact has significant ethical importance: a cascade of diagnoses or broader interventions that are erroneous, leading to treatment plans that harm patients and populations. The recognition of data bias is a foundational clinical medicine skill because evidence-based practice depends on accurate information. Students tend to consider the answer to the question, "Are the numbers in this table correct?" to be straightforward. But accuracy cannot be assessed by scrutiny of numerical data alone; only a close look at the methods that produced those values can reveal bias. In short, data bias emerges from a *process*. A reader in the health professions must understand the steps that generate the numbers and the assumptions made by investigators about their data sources. It is important to know whether bias stems from the availability of information, its collection, its methodological manipulation, or the analysis of findings.

Broader judgment comes into play because epidemiological bias, which includes data bias, does not arise from methodological errors alone. It can also result from socially discriminatory choices that inform data selection, classification, and analysis. Health professions students are often surprised that, as used in epidemiology, the term *bias* is

putatively unconnected from its everyday meaning: prejudice or devaluation of a particular social group. Here, we argue that the everyday and scientific meanings of the term are at times closely interrelated. A focus on implicit bias in the clinical realm enhances this awareness,¹ and recent attention to the need for more equity in public health data² reinforces the importance of the issue. This article canvasses concepts pertaining to epidemiological bias that health professions students should understand, regardless of whether bias is generated by epidemiological errors or by racial, ethnic, gender, or other prejudices. Such sources of inaccuracy are of ethical and clinical concern because they can influence patient care and community health.

Importance of Accurate Data

The goal of epidemiological and clinical research is to produce accurate data that are useful. The presence of bias in a study implies that there are systematic errors in the data. Nonetheless, bias is not a dichotomous concept: it can exist to a greater or lesser degree and may distort a true association in one direction or the other.³ Where it is present, bias will influence the validity of data for the population under study (internal validity) and for populations for whom results are assumed to be relevant (external validity). For example, in a drug trial, internal validity represents the extent to which observed outcomes can be ascribed to the treatment regimen, allowing for causal inference. There can be no external validity (broader effectiveness) without internal validity, although the presence of the second does not guarantee the first.⁴

Sources of Data Bias

Data bias is a capacious concept. Its largest categories are selection bias and information bias, which in turn encompass numerous subvarieties.⁵ Notably, certain study designs are structurally vulnerable to data bias. For example, retrospective cohort studies are particularly prone to selection bias. In such chart-based studies, the investigator identifies a cohort that has been assembled in the past, identifies potential predictor variables from measurements made in the past, and evaluates outcome variables. Since data will likely not have been collected for research, some charts might be excluded due to missing but crucial information.⁶ Interviewer bias can occur in case-control studies if investigators question patients who are "cases" more intensively about exposures that are already known to be associated with the disease.^{3,7} Even randomized controlled studies are vulnerable to bias resulting from misallocation of participants, insufficient data blinding, or loss of subjects to follow-up.⁸

Selection bias. In most studies, only a sample of the target population is chosen for observation or intervention. Consequently, studies are susceptible to selection bias, that is, to the recruitment of a nonrepresentative assemblage of subjects.⁹ Individuals within the sample may systematically differ with respect to social and economic status, educational level, age, or other consequential characteristics. Such errors can obscure causal associations between an exposure, such as a treatment, and a health-related outcome.¹⁰ Biases in which errors of inclusion or exclusion play a role often have their own designation or eponym. This inventory of biases includes nonresponse bias, volunteer bias, Berksonian bias, attrition bias, incidence-prevalence bias, confounding by indication, surveillance bias, and other named biases.^{5,9}

Information bias. Information bias can arise during or after data collection and refers to systematic errors in the measurement of variables or classification of subjects. Errors of measurement can occur because of faulty instrumentation or discernment, the latter of which includes recall bias, interview bias, observer bias, or confirmation bias.^{9,11} As a

rule, understanding the relationship between an exposure and an outcome requires subjects to be classified into categories, such as "exposed" or "non-exposed," and to isolate variables responsible for differing outcomes.^{12,13} Misclassifications commonly arise in observational studies but can be present in randomized controlled studies.^{12,13} Nondifferential and differential misclassification bias refer to whether measurement error due to misclassification of subjects is symmetrically or asymmetrically distributed between the intervention and comparison groups. For example, in a study of the impact of a drug on obesity, the scales used to weigh patients may not all be accurate. Depending on whether those inaccuracies are similar (say, 5% higher for all patients) or dissimilar, this error will have a divergent bearing on the results.

Bias should be differentiated from other problems of accuracy, particularly from confounding. Confounding describes an association between 2 variables—an exposure and outcome—that appears causal but that surfaces only due to influence from a hidden yet consequential variable. A well-cited example is the association between heavy coffee drinking and cancer of the pancreas.⁷ This mirage is present only because heavy coffee drinkers are more likely than light or non-coffee drinkers to smoke cigarettes, an action responsible for the elevated risk of pancreatic cancer. Confounding represents a distortion of the relationship between exposure and outcome due to the presence of one or more extraneous variables, and, like data bias, it can lead to incorrect inferences about causality. Typically, it is not possible to correct for data bias, whereas if a confounding variable is known and measured, the real effect of the exposure on the outcome can be obtained by adjustment for this factor. In sum, confounding produces errors of interpretation despite the accuracy of the measurement.⁴

Random error, in contrast to bias, is nonsystematic and affects the precision rather than the validity of research findings. This lack of exactness results from sampling variability, producing errors that are unsystematic. Data can be both biased and imprecise, but, unlike bias, lack of precision is best addressed by increasing a study's sample size.

Social Bias and Data Bias

In the epidemiological literature, data biases are implicitly considered oversights, mistakes, or unavoidable failures in research protocols. Indeed, epidemiologists distinguish sharply between data bias and social bias: "Bias undermines the internal validity of research. Unlike the conventional meaning of bias—i.e., prejudice—bias in research denotes deviation from the truth."¹⁴ In short, data bias is an operational error and social bias (prejudice) is a disposition of judgment. This distinction is not always clear, however. The consequences of social bias can lead researchers to deviate from the truth, and clinicians can collect biased data by using measurement tools that have social biases structured into them.

An important example of overlap between data bias and social bias is found in research on risk factors for cardiovascular disease in young men identified as Black. The available data are fraught with selection bias. Given the enormity of the population with a history of incarceration and the disproportionate incarceration of Black men,^{15,16} the exclusion of incarcerated persons from household-based surveys poses a large obstacle to obtaining unbiased samples. Examples of surveys that exclude people who are currently incarcerated include the National Health and Nutrition Examination Survey and the National Health Interview Survey.¹⁶

Information bias arising from a legacy of medical racism continues to affect diagnostic and eligibility criteria. Indeed, race is embedded in clinical algorithms and decisionmaking tools across many medical specialties.^{17,18} For example, in 2019 researchers revealed algorithmic bias in a widely used medical artificial intelligence tool that incorporates health care costs into the prediction of clinical risk, with deleterious consequences for Black patients. Since the health care system spends more money, on average, on White patients than on Black patients, the tool returns higher risk scores for White patients than for Black patients. Use of this tool might have led to more referrals for White patients to specialty services, perpetuating both spending discrepancies and race bias in health care.¹⁸ Moreover, in pulmonary medicine, the observation of differences in lung capacity between a population characterized in the 1800s as "Full Blacks" and White soldiers was attributed to a biological difference associated with race rather than the effect of enslavement and environmental exposures now known to alter lung function. Subsequently, a "race correction" was built into equations used in spirometry.¹⁹ Whether in the assessment of occupational lung diseases such as asbestosis or "objective" eligibility for lung transplantations, the incorporation of biased reference standards for lung function can lead to worse outcomes for Black patients.

Links between social and data biases are also evident in biomedical research. For example, the evolving field of precision medicine is driven by lab-based sequencing of the genetic code, creating large databases that are curated and organized to extract clinically relevant information. The underrepresentation of non-European populations in genomic databases, like all selection bias, is problematic for clinical care because the exclusion of such data limits their generalizability.^{20,21} Moreover, while at the cellular level racial identity is nonexistent, once a pathophysiological process is understood and given a label, the resulting diagnostic category can take on racialized associations, leading to information bias. There are many examples of such racialized diagnostic categories, including sickle cell disease, sarcoidosis, gallstones, and cystic fibrosis. In the clinical setting, this racial "essentialism" leads to assumed or missed diagnoses, misclassification through confirmation bias, and harmful consequences.²¹

Social biases that contribute to data bias are not limited to race. A body of public health research documents gaps in national survey data of sexual orientation and gender identity.^{22,23,24} Although data collection procedures have evolved, prior to 2016, biological sex in Behavioral Risk Factor Surveillance System telephone surveys could be assigned based on a respondent's "vocal timbre," a practice vulnerable to confirmation bias.²⁵ Research documents that this approach has resulted in substantial misclassification of answers, especially those of persons who identify as transgender or gender diverse. ^{25,26,27} Moreover, the collection of data on sex and gender, where explicitly sought, does not by itself guarantee validity, as there is widespread misunderstanding of the meaning of sex and gender.²⁴

Conclusion

Data and social biases are oblique to one another: they are separate frames, but at times they interlock; together, they contribute to epidemiological bias. Data bias refers to systematic errors in a sequence of tasks that produces data; social bias refers to actions and attitudes that can shape those operations. And when these frames coincide, it is not always clear which is a subset of the other. The exclusion of a group from a survey or study can reflect selection bias, but this exclusion may more accurately be ascribed to prejudice.

While there is a path to identifying data bias, there are no shortcuts. Some degree of bias is always present in a published study, so the challenge of bias recognition is ongoing.³ Awareness of the nature and types of bias in research studies allows for a more meaningful scrutiny of results and conclusions. As researchers, careful planning is needed in each step of research design, and, when presenting results, a full acknowledgment of any sources of bias is essential.²⁸ The health professional's commitment to a close examination of evidence must remain steadfast, as the presence of bias—whether of epidemiological or social origin—undermines the provision of effective and acceptable clinical care.

References

- 1. Gopal DP, Chetty U, O'Donnell P, Gajria C, Blackadder-Weinstein J. Implicit bias in healthcare: clinical practice, research and decision making. *Future Healthc J*. 2021;8(1):40-48.
- 2. Ponce NA, Lau DT. Toward more equitable public health data: an AJPH special section. *Am J Public Health*. 2023;113(12):1276-1277.
- 3. Pannucci CJ, Wilkins EG. Identifying and avoiding bias in research. *Plast Reconstr Surg.* 2010;126(2):619-662.
- 4. Rothman KJ. Modern Epidemiology. Little Brown & Co; 1986.
- 5. Delgado-Rodríguez M, Llorca J. Bias. *J Epidemiol Community Health*. 2004;58(8):635-641.
- 6. Talari K, Goyal M. Retrospective studies—utility and caveats. *J R Coll Physicians Edinb*. 2020;50(4):398-402.
- 7. Gordis L. Epidemiology. 5th ed. Elsevier/Saunders; 2014.
- Higgins JPT, Savović J, Page MJ, Elbers RG, Sterne JAC. Assessing risk of bias in a randomized trial. In: Higgins J, Thomas J, Chandler J, et al, eds. Cochrane Handbook for Systematic Reviews of Interventions. Cochrane Training; 2024:chap 8. Accessed June 24, 2024. https://training.cochrane.org/handbook/current/chapter-08
- 9. Jager KJ, Tripepi G, Chesnaye NC, Dekker FW, Zoccali C, Stel VS. Where to look for the most frequent biases? *Nephrology (Carlton)*. 2020;25(6):435-441.
- 10. Hsu JL, Banerjee D, Kuschner WG. Understanding and identifying bias and confounding in the medical literature. *South Med J.* 2008;101(12):1240-1245.
- 11. Glick M. Believing is seeing: confirmation bias. *J Am Dent Assoc.* 2017;148(3):131-132.
- 12. Moseley AM, Pinheiro MB. Research note: evaluating risk of bias in randomised controlled trials. *J Physiother*. 2022;68(2):148-150.
- 13. Bosdriesz JR, Stel VS, van Diepen M, et al. Evidence-based medicine—when observational studies are better than randomized controlled trials. *Nephrology* (*Carlton*). 2020;25(10):737-743.
- 14. Grimes DA, Schulz KF. Bias and causal associations in observational research. *Lancet*. 2002;359(9302):248-252.
- 15. Robey JP, Massoglia M, Light MT. A generational shift: race and the declining lifetime risk of imprisonment. *Demography*. 2023;60(4):977-1003.
- 16. Wang EA, Redmond N, Dennison Himmelfarb CR, et al. Cardiovascular disease in incarcerated populations. *J Am Coll Cardiol*. 2017;69(24):2967-2976.
- 17. Cerdeña JP, Plaisime MV, Tsai J. From race-based to race-conscious medicine: how anti-racist uprisings call us to act. *Lancet*. 2020;396(10257):1125-1128.
- 18. Vyas DA, Eisenstein LG, Jones DS. Hidden in plain sight—reconsidering the use of race correction in clinical algorithms. *N Engl J Med*. 2020;383(9):874-882.

- 19. Braun L. Race correction and spirometry: why history matters. *Chest.* 2021;159(4):1670-1675.
- 20. Landry LG, Ali N, Williams DR, Rehm HL, Bonham VL. Lack of diversity in genomic databases is a barrier to translating precision medicine research into practice. *Health Aff (Millwood)*. 2018;37(5):780-785.
- 21. Tsai J. How should educators and publishers eliminate racial essentialism? *AMA J Ethics*. 2022;24(3):E201-E211.
- 22. Baker KE, Streed CG Jr, Durso LE. Ensuring that LGBTQI+ people count collecting data on sexual orientation, gender identity, and intersex status. *N Engl J Med*. 2021;384(13):1184-1186.
- 23. Patterson CJ, Sepúlveda MJ, White J, eds. Understanding the Well-Being of LGBTQI+ Populations. National Academies Press; 2020.
- 24. Jacobs JW, Bibb LA, Shelton KM, Booth GS. Assessment of the use of sex and gender terminology in US federal, state, and local databases. *JAMA Intern Med.* 2022;182(8):878-879.
- Riley NC, Blosnich JR, Bear TM, Reisner SL. Vocal timbre and the classification of respondent sex in US phone-based surveys. *Am J Public Health*. 2017;107(8):1290-1294.
- 26. Tordoff D, Andrasik M, Hajat A. Misclassification of sex assigned at birth in the behavioral risk factor surveillance system and transgender reproductive health: a quantitative bias analysis. *Epidemiology*. 2019;30(5):669-678.
- 27. Gonzales G, Tran NM, Bennett MA. State policies and health disparities between transgender and cisgender adults: considerations and challenges using population-based survey data. *J Health Polit Policy Law*. 2022;47(5):555-581.
- 28. Althubaiti A. Information bias in health research: definition, pitfalls, and adjustment methods. *J Multidiscip Healthc*. 2016;9:211-217.

Douglas Shenson, MD, MPH, MA, MS is an adjunct associate professor in the Section of General Internal Medicine at the Yale School of Medicine (YSM) in New Haven, Connecticut, where he is deputy leader of YSM's Health Equity Thread. Dr Shenson is also director of YSM's required preclinical course: "Populations & Methods: The Application of Epidemiology and Biostatistics to Public Health."

Beverley J. Sheares, MD, MS is the inaugural leader of the Health Equity Thread at the Yale School of Medicine in New Haven, Connecticut, where she is an associate professor of pediatrics in the Pulmonary, Allergy, Immunology, and Sleep Medicine Section. Dr Sheares' clinical, teaching, and research experiences all coalesce around reducing health disparities and promoting equity.

Chelesa Fearce is a MD/PhD student studying chemistry at Yale University in New Haven, Connecticut, who is interested in drug development for psychiatric disorders. After graduating from Spelman College in 2017, she spent 2 years at the National Institutes of Health studying dopamine receptor signaling. Chelesa plans to make health equity an integral part of her career as a physician-scientist.

Citation

AMA J Ethics. 2025;27(1):E14-20.

DOI 10.1001/amajethics.2025.14.

Conflict of Interest Disclosure

Authors disclosed no conflicts of interest.

The viewpoints expressed in this article are those of the author(s) and do not necessarily reflect the views and policies of the AMA.

Copyright 2025 American Medical Association. All rights reserved. ISSN 2376-6980