# Bias in Gender Data: Beyond the Binary

Shale Hunter

2022-06-08

**Last December, I wrote this article on identifying grading bias in an undergraduate course I was TAing at the time. Though I did not discover any *grading* bias in the data, one of the limitations I identified then has stuck with me. I would like to take the time to explore these limitations in more detail in this post.**

## What's Missing?

Here is a summary of the `pronouns` column from the dataset used in my previous analysis:

| Pronouns | Number of Students |
| --- | --- |
| NA | 126 |
| He/Him/His | 62 |
| She/Her/Hers | 38 |
| They/Them/Theirs | 3 |

There are two clear issues highlighted by this table:

- Only 3 of the 229 students in the class reported using the pronouns `They/Them/Theirs`
- More than half of students in the class did not share their pronouns on their page in UCSB's student database

Anyone reading this post as a data scientist will immediately be able to identify the statistical limitations of a dataset with over 50% `NA` values and variable with a category that is almost completely unrepresented. However in this article I would like to take a deeper look at the underlying causes that forced me to use such an inadequate data source in the first place.

## A Review of Potential Causes for Biased Data in SGM Research

While perhaps not quite yet mainstream, research on Sexual and Gender Minorities (SGM) is a field that is growing rapidly as we are discovering the many ways that identifying with a marginalized group can have harmful effects on a person's life because of real or assumed discrimination from broader society. This is exceptionally true in the medical field, but certainly is not limited to it.

At the same time we are trying to understand and combat the unique challenges faced by SGM individuals and communities, these same people are much more likely to mistrust authorities of all kinds. From doctors advocating for conversion therapy to openly transphobic members of UCSB itself, there is a massive deficit of trust between the LGBTQ+ community and institutional power. Lunn et al. (2019) describes this phenomenon in relation to clinical research:

- "With stigmatizing, discriminatory, or dangerous experiences in society, health care, or investigational communities, potential participant interactions with the health care system and traditional research enterprise may be limited out of fear."

This is known to be true not only in SGM research broadly, but also specifically in relation to college students, the target demographic of my own dataset. Women who have sex with women (WSW) have been shown to have significantly higher rates of eating disorders compared to women who have sex with men (WSM) in the university setting (Von Schell et al. 2018). The same study found that the observed rate of eating disorders in WSW was consistent with the observed rate in existing literature on eating disorders - indicative of the historical exclusion of SGM people in the field.

Even in my own graduate-level classes at the Bren School, I have noticed a disturbing affinity for the gender binary in teaching examples: while not explicitly discriminatory, it is the kind of unintentional, unexamined microaggression which quietly invalidates nonbinary students' identity, and with it erodes their trust in the institution our faculty represent.

## Effects of Bias in SGM Data

With this combination of a historical lack of effort on researchers' part to safely and respectfully gather data on minority groups and the ensuing lack of trust between these groups and researchers, it should not be surprising that we can see a clear bias in much of the data that purports to be applicable to SGM people. A review of 43 publicly accessible national, international, and regional data sources on gender and health found that only 14% measured all three dimensions of sexual orientation (identity, behavior, attraction), and no data sources measured transgender-inclusive gender identity (Patterson et al. 2017). Norori et al. (2021) likewise suggest that implicit or explicit biases of research personnel towards SGM individuals can affect data and metadata legitimacy as well as diagnosis accuracy; this is corroborated in Suen et al. (2022), where interviews with SGM focus groups indicated that in-person (as opposed to online or over the phone) is the most uncomfortable, least safe-feeling means of data collection, often because of perceived judgement on the part of the researcher or medical professional.

There also remains substantial misunderstanding of SGM people even among those researchers who are trying to ensure that data collection and use is being done in an ethical and responsible way: Norori et al. (2021) describe non-heterosexual women as "generally higher socioeconomic status than heterosexual women" with hardly a thought for the reality underlying whatever data supports this claim. It may be true that a survey of women (or any gender identity) would find a correlation between sexual orientation and socioeconomic status; but treating this data as fact as Norori does is likely a grave mistake. A more accurate phrasing would be: "women who *describe themselves* as non-heterosexual have higher socioeconomic status" - where the likely explanation is that wealthy and highly educated individuals have the **privilege** of expressing their sexual or gender identity openly in a way that less wealthy or educated individuals may not have because of familial, religious, or other practical restrictions on their ability to live openly (which extends to medical surveys). An extreme version of this is children under 18, most of whom report that they wouldn't participate in SGM research or surveys at all if it required parental consent (Macapagal et al. 2016).

All this is to say that the problem encountered in my data is not unique - but how can we as researchers ensure that the data we *do* collect is unbiased, complete, and respectful of its contributors?

## Ways to Combat Bias in SGM Research

Several approaches have been taken to try and improve SGM data from an ethics and bias perspective. Interviews and focus groups like those described in Suen et al. (2022) offer a means for SGM individuals to express their own concerns directly to the organizers of future research, with the goal of modifying methods to reflect the unique needs of the population. This also helps avoid the problem of non-SGM research personnel making the mistake of thinking they know best when it comes to what makes someone with a marginalized identity feel more or less comfortable in a research setting.

Another approach is described in Lunn et al. (2019): the PRIDE Study (Population Research in Identity and Disparities for Equality) is a longitudinal SGM research project based out of Stanford that takes advantage of the well-documented preference in the SGM community for remote interaction with researchers by developing a mobile application to complete surveys and give feedback on research topics and methodologies. Using

this approach allowed 16,394 SGM respondents (98% sexual minority, 15% gender minority) from a range of geographic and demographic groups in the US to safely and comfortably provide 3,544 suggested important SGM health topics, complete 24,022 surveys, and use forum and voting features to share opinions on the methods and functionality of the app itself.

Finally, our responsibility in relation to the biases discovered in SGM data is not just as researchers, but also as educators. The data science and R communities have made progress in addressing other types of latent injustices through the replacement of problematic datasets with better ones. There is certainly room for such a solution in SGM data education: with the limitations of current data sources clearly highlighted by reviews such as Patterson et al. (2017), new teaching data files with complete gender information would be an invaluable tool both in service of the demarginalization of diverse gender identities and education of the broader population.

## Final Thoughts

In the end, the only way to ensure researchers are collecting unbiased and complete data about our SGM community is by ensuring that the community feels safe providing their data as well as that the data is collected in a way that benefits that community. More than 25% of SGM respondents to one survey said of this kind of research: "I don't need anything in return. I will keep participating to improve LGBTQ health" (Lunn et al. 2019).

Even so, we have a long way to go before this is broadly achieved in the US, let alone the world. Even in a highly educated liberal institution like USCB there remain barriers to unbiased SGM data collection, as evidenced by my own experience with data collection:

| Pronouns | Number of Students |
|---|---:|
| NA | 126 |
| He/Him/His | 62 |
| She/Her/Hers | 38 |
| They/Them/Theirs | 3 |

Even here there appears to be some level of apathy or lack of respect for the importance of gender identity in the student body as a whole. And it remains difficult to ascertain exactly how much peer, faculty, or institutional judgement is affecting SGM individuals' willingness to share their identity on school platforms such as the one I used to collect data for my project last winter. Spizzirri et al. (2021) report global estimations of gender-diverse individuals between 0.1 to 2% of the population, which my proportion of 3/229 students (1.3%) would be consistent with. However, Lunn et al. (2019) report over 15% gender minority in their sample - of course this is a sample that is specifically targeting the LGBTQ+ community and would be expected to be greater than the overall population. But there remains a concern that the true proportion is higher than the 2% upper bound of observed gender minorities because of real and persistent stigmas relating to non-standard gender identities.

Regardless of what the true proportion is, it is important that research methods take into account the many ways in which SGM populations have been mistreated in the past, and how that mistreatment affects current relationships between data collector and data contributor. Hopefully then I will be able to **really** get a complete and unbiased estimation of my objectivity as a grader!

## References

- Lunn, M. R., Capriotti, M. R., Flentje, A., Bibbins-Domingo, K., Pletcher, M. J., Triano, A. J., Sooksaman, C., Frazier, J., & Obedin-Maliver, J. (2019). Using mobile technology to engage sexual and gender minorities in clinical research. PLoS ONE, 14(5), e0216282–e0216282. https://doi.org/10.1371/journal.pone.0216282

- Macapagal, K., Coventry, R., Arbeit, M. R., Fisher, C. B., & Mustanski, B. (2016). "I Won't Out Myself Just to Do a Survey": Sexual and Gender Minority Adolescents' Perspectives on the Risks and Benefits of Sex Research. Archives of Sexual Behavior, 46(5), 1393–1409. https://doi.org/10.1007/s10508-016-0784-5
- May, A., Wachs, J., & Hannák, A. (2019). Gender differences in participation and reward on Stack Overflow. Empirical Software Engineering : An International Journal, 24(4), 1997–2019. https://doi.org/10.1007/s10664-019-09685-x
- 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 (New York, N.Y.), 2(10), 100347–100347. https://doi.org/10.1016/j.patter.2021.100347
- Patterson, J. G., Jabson, J. M., & Bowen, D. J. (2017). Measuring Sexual and Gender Minority Populations in Health Surveillance. LGBT Health, 4(2), 82–105. https://doi.org/10.1089/lgbt.2016.0026
- Spizzirri, G., Eufrásio, R., Lima, M.C.P. et al. Proportion of people identified as transgender and non-binary gender in Brazil. Sci Rep 11, 2240 (2021). https://doi.org/10.1038/s41598-021-81411-4
- Suen, L. W., Lunn, M. R., Sevelius, J. M., Flentje, A., Capriotti, M. R., Lubensky, M. E., Hunt, C., Weber, S., Bahati, M., Rescate, A., Dastur, Z., & Obedin-Maliver, J. (2022). Do Ask, Tell, and Show: Contextual Factors Affecting Sexual Orientation and Gender Identity Disclosure for Sexual and Gender Minority People. LGBT Health, 9(2), 73–80. https://doi.org/10.1089/lgbt.2021.0159
- Von Schell, A., Ohrt, T. K., Bruening, A. B., Perez, M. (2018). Rates of disordered eating behaviors across sexual minority undergraduate men and women. Psychology of Sexual Orientation and Gender Diversity, 5(3), 352–359. https://doi.org/10.1037/sgd0000278