

Letting Data Speaking for Themselves*

Xiaoyu (Alice) Miao

The notion of letting “data speak for themselves” is appealing in its simplicity and purported objectivity. However, this perspective, when scrutinized through the lenses provided by Jordan (2019), D’Ignazio and Klein (2020), and Au (2020), reveals a landscape where data cannot and should not be expected to stand alone as arbiters of truth. Their arguments collectively underscore the complexities involved in data generation, processing, and analysis, challenging the feasibility and desirability of adopting a hands-off approach to data interpretation. This essay explores the extent to which we should allow data to speak for themselves, drawing on these authors’ insights and considering broader implications for data science and ethics.

The Myth of Unmediated Data

At the heart of the discussion is the misconception that data, in their raw form, are neutral and unbiased reflections of reality. Jordan (2019) dismantles this notion by highlighting the entanglement of data with human decisions, biases, and societal structures. The very act of data collection involves choices about what to measure, how to measure it, and which data to record or discard. These decisions, often invisible in the final dataset, imbue the data with the perspectives and prejudices of those who collect and process it.

D’Ignazio and Klein (2020) extend this argument through the lens of intersectional feminism, emphasizing how data practices can reinforce or challenge existing power dynamics. Their critique of data science as predominantly white and male-dominated not only questions whose interests are served by data-driven decisions but also which voices and experiences are marginalized or excluded. This perspective calls for a more critical engagement with data, one that acknowledges and seeks to rectify these imbalances.

Randy Au (2020), meanwhile, focuses on the process of data cleaning as an analytical act that imposes interpretation on the dataset. He argues that no data are truly raw; they are always already shaped by human actions and decisions. The choices made during data cleaning—what to exclude, normalize, or correct—are not merely technical adjustments but analytical judgments that significantly affect the outcome of data analysis. This view challenges the notion of data neutrality and underscores the responsibility of data scientists to engage deeply with their datasets, understanding their limitations and biases.

* Available at: <https://github.com/xyalicemiao/week-9-tut/>

The Role of the Data Scientist

The works of Jordan, D'Ignazio and Klein, and Au collectively argue for a more active role of the data scientist in interpreting and presenting data. Far from being passive conduits for data to speak, data scientists are interpreters and storytellers who mediate between data and their implications for the world. This role requires a deep understanding of the data's origins, a critical awareness of one's own biases, and a commitment to ethical principles that prioritize fairness, transparency, and inclusivity.

Engaging with data in this way also involves recognizing the limits of what data can tell us. Data are not self-explanatory; they require context, theory, and human insight to yield meaningful conclusions. The insistence on letting data speak for themselves risks oversimplifying complex phenomena and overlooking the interpretive work necessary to understand the underlying patterns and causes.

Collaborative and Multidisciplinary Approaches to Data Science

The articles collectively hint at the necessity of collaborative and multidisciplinary approaches to data science. Jordan's discussion about the need for a new engineering discipline that integrates social sciences and humanities with computational methods suggests that tackling complex data science challenges requires expertise beyond traditional STEM fields. Similarly, D'Ignazio and Klein's application of feminist theory to data science underscores the value of diverse perspectives in uncovering and addressing biases in data practices. By promoting interdisciplinary collaborations, data science can benefit from a wider range of insights and methodologies, leading to more innovative, inclusive, and ethically sound approaches to understanding and acting on data. This aspect invites discussion on how different disciplines can contribute to the development of data science and how collaborative efforts can be fostered within academic, industry, and policy-making contexts.

Towards a More Responsible Data Practice

The critique of the notion that data can speak for themselves leads to a broader call for responsible data practice. This involves not only technical proficiency but also ethical engagement with the societal implications of data analysis. Data scientists must consider the potential impact of their work on different communities, especially those historically marginalized or vulnerable to harm. This includes designing studies and algorithms that respect privacy, avoid perpetuating biases, and strive for equitable outcomes.

Moreover, responsible data practice requires transparency about the limitations and uncertainties of data analysis. Acknowledging these limitations can foster a more nuanced understanding of what data can and cannot tell us, encouraging a more humble and questioning approach to data science.

The Importance of Data Provenance and Transparency

Drawing from Jordan's critique of the oversimplification of AI and data science, as well as Au's emphasis on the intricacies of data cleaning, another critical aspect to discuss is the importance

of data provenance and transparency. This involves understanding and documenting where data come from, how they are collected, and any transformations they undergo before analysis. Discussing data provenance and transparency can lead to better awareness of the biases and limitations inherent in datasets, thereby fostering more responsible use of data in scientific and decision-making processes. It also underscores the need for clear communication about the origins and manipulations of data, which is essential for reproducibility and trust in data-driven findings.

Ethical Considerations in Data Science Practices

D'Ignazio and Klein's focus on intersectional feminism as a lens for examining data science invites a broader discussion on the ethical considerations that must guide data practices. This includes reflecting on who benefits from data collection and analysis, who might be harmed, and how the harms can be minimized. Ethics in data science also encompasses issues of consent, privacy, and data security, especially in contexts where the collection and use of data could lead to surveillance, discrimination, or other forms of harm to individuals or communities. Discussing ethical considerations demands that data scientists not only follow legal requirements but also engage with deeper moral questions about justice, equity, and respect for persons.

Conclusion

In conclusion, the works of Jordan (2019), D'Ignazio and Klein (2020), and Au (2020) provide compelling arguments against the notion of letting data speak for themselves. They highlight the deeply human aspects of data science, from the collection and preparation of data to their analysis and interpretation. These insights call for a more engaged, critical, and ethical approach to data, one that recognizes the power of data to shape our understanding of the world and our impact on it. Far from being a limitation, this approach can enrich data science, making it more robust, equitable, and responsive to the complexities of human life.

References

- Au, Randy. 2020. “Data Cleaning Is Analysis, Not Grunt Work.” *Data Cleaning IS Analysis, Not Grunt Work - by Randy Au*. Counting Stuff. <https://counting.substack.com/p/data-cleaning-is-analysis-not-grunt>.
- D’Ignazio, Catherine, and Lauren F. Klein. 2020. “Data Feminism.” *Data Feminism*. The MIT Press. <https://data-feminism.mitpress.mit.edu/>.
- Jordan, Michael I. 2019. “Artificial Intelligence-the Revolution Hasn’t Happened Yet.” *Harvard Data Science Review*. PubPub. <https://hdsr.mitpress.mit.edu/pub/wot7mkc1/release/10>.