



ISSN 2278 – 0211 (Online)

Towards a Framework for Bias Prevention to Ensure Open Data Quality

Nansukusa Yudaya

Assistant Lecturer, Department of Software Engineering, Computer and Information Sciences,
Muni University, Uganda

Kalyankolo Umaru

Assistant Lecturer, Department of Electrical Engineering, Muni University, Uganda

Abstract:

The rapid growth of open data initiatives has emphasized their potential to enhance transparency, foster innovation, and support equitable decision-making across sectors. However, the quality and reliability of open data remain compromised by biases that alter outcomes and spread inequalities. This paper critically examines the systemic sources of bias, including sampling, annotator, and algorithmic biases, that undermine data integrity and decision-making processes. It proposes a comprehensive framework to mitigate these biases through standardized data management protocols, inclusive data collection practices, robust data stewardship, and cross-sector collaboration. The study also highlights the ethical imperatives and practical challenges of bias prevention, emphasizing the need to balance inclusivity with privacy and resource constraints. By prioritizing fairness, inclusivity, and dependability, the proposed interventions aim to enhance the credibility and societal impact of open data, reaffirming its role as a catalyst for equitable innovation and policy development. The findings underscore that addressing biases in open data is not only a technical necessity but also a moral imperative essential for sustaining its transformative potential.

Keywords: Bias prevention, algorithmic bias, data stewardship, data integrity, equity

1. Introduction

The concept of open data has emerged as a focal element in contemporary discussions surrounding transparency, accountability, and innovation across various sectors (Okunleye, 2024). Open data refers to digital data that is freely accessible for anyone to use and distribute, promoting transparency, accountability, and citizen engagement in public governance (Salubi, 2024). It is also referred to as publicly accessible information from the public sector, provided in reusable formats without cost or restrictions, aimed at fostering transparency, citizen engagement, and innovation-driven public value (Gao et al., 2023). Numerous studies confirm that the open data initiative seeks to foster transparency, accountability, and innovation through the provision of access to a diverse array of datasets, encompassing governmental statistics, scientific inquiries, and environmental data (Lnenicka et al., 2024; Okunleye, 2024; Park & Gil-Garcia, 2022).

Open data holds unmatched potential to shape policy, foster innovation, and enhance governance (Lee, 2024). However, the reliability and fairness of decisions based on open data depend on its quality. According to (Baumann et al., 2023), bias refers to systematic deviations in data or algorithms that lead to unfair, inaccurate, or prejudiced outcomes in machine learning models, often resulting from historical inequalities, measurement inaccuracies, underrepresentation of certain groups, or the omission of relevant variables, thereby affecting the fairness and reliability of decision-making processes. It is also described as a systematic error in data collection, analysis, or interpretation that can distort findings, often stemming from selective data inclusion or poorly managed model risks, ultimately compromising research validity (Borissova, n.d.). Bias, whether arising during data collection, processing, or dissemination, undermines the quality of open data, leading to inaccurate insights and reinforcing existing inequalities (Wesson et al., 2022). This paper asserts that preventing biases in open data is essential for maintaining its quality and proposes a framework for systematically addressing bias.

2. Sources of Bias in Open Data

Biases in open data arise from multiple sources, including:

- **Sampling Bias:** Biases arise from decisions about which data should be labeled and the demographic composition of the annotator pool (Hettiachchi et al., 2021). Open datasets often lack diversity in representation, leading to skewed results that favor specific groups (Obermeyer et al., 2019). Representation bias is the systematic underrepresentation or misrepresentation of certain groups within a dataset, which can lead to skewed outcomes and inaccurate analyses (Shahbazi et al., 2023). In AI, for example, data biases in open datasets are often biased

towards majority groups, undermining AI fairness and generalization, necessitating equitable representation and robust fairness corrections for trustworthy models (Langbridge et al., 2024).

- **Annotator Bias:** Individuals with strong opinions or specific demographics (e.g., age, gender) can produce biased annotations, often influenced by the personal attributes and perspectives of the annotators (Hettiachchi et al., 2021).
- Because of selection biases or inadequate coverage, data biases are frequently present in open datasets (Fabbrizzi et al., 2022) in numerous domains such as Artificial Intelligence (AI), health, transport etc.
- **Algorithmic Bias:** The use of biased algorithms in processing open data exacerbates existing inequalities (Noble, 2018). According to (Shimron et al., 2022), algorithmic bias arises when machine learning models trained on public data perform unrealistically well due to hidden data processing pipelines, such as zero-padding or compression, which alter the data in ways that make algorithmic evaluation misleadingly optimistic.

3. Preventing Bias in Open Data

- **Implementing Standardized Data Management Protocols:** Establishing standardized protocols for data collection, storage, and dissemination can significantly improve data quality and minimize variability (Dong et al., 2024; Johnson & Rostain, 2020). Standardized data protocols are crucial for maintaining high data quality, especially in large-scale, multi-center, and international projects, facilitating consistent data integration, enhancing semantic interoperability, and streamlining the merging of data from various studies (Rinaldi et al., 2022).
- **Promoting Data Stewardship and Maintenance:** Data stewardship involves the practical aspects of managing research data to maintain its quality as an asset while also ensuring that it remains accessible to the relevant community and retains a high standard of quality (Arend et al., 2022). In the health sector, data stewardship involves data collection, analysis, storage, ensuring data quality and integrity, and sharing, as well as protecting the privacy of study subjects, who are usually patients (Aksenova et al., 2024). Data maintenance involves regular updates, error corrections, and oversight to ensure data accuracy over time (Velayutham et al., 2021). Open data platforms should implement version control systems to track changes and facilitate user feedback, allowing users to report inconsistencies or missing information (Rygge, 2024). Also, ensure diverse representation in datasets by employing stratified sampling techniques and reaching underserved populations (Neves et al., 2020).
- **Creating Guidelines for Inclusive Data Collection:** To address representation biases, it is essential to establish guidelines that encourage inclusive data collection that incorporates gender, diversity, inclusion, and protection perspectives (Hajjaj PhD et al., n.d.). Encouraging data providers to evaluate and document potential biases in their datasets would also enhance transparency and accountability (Hajjaj PhD et al., n.d.; Kam et al., 2024). In the same regard, it is advised to implement ethical guidelines for open data governance, emphasizing fairness and accountability (Floridi et al., 2021).
- **Encouraging Cross-Sector Collaboration and Knowledge Sharing:** Collaborations between data providers, users, and standard-setting bodies can lead to the development of best practices that address the unique challenges of open data (Cabral, 2024). For example, collaborations between government agencies, academic institutions, and non-profit organizations can foster knowledge exchange and create a shared understanding of quality standards towards preventing biases (Laiti, 2024). For example, the Data Across Sectors for Health (DASH) initiative promotes cross-sector collaboration and knowledge sharing by connecting stakeholders across healthcare, social services, and public health to share data on social determinants of health (SDOH), providing resources through the "All In" network to enhance readiness for data sharing, and fostering an integrated, collaborative environment that strengthens networks and empowers communities to address health equity with more informed interventions (O'Neil et al., 2020).
- In summary, preventing biases in open data is not merely a technical challenge but an ethical imperative. By addressing biases at the systemic level, stakeholders can improve the credibility, inclusivity, and impact of open data. Bias prevention contributes directly to the core dimensions of data quality: accuracy, reliability, and usability (Vetro' et al., 2016; Zainuddin & Akhir, 2024). Without this focus, the promise of open data to foster equitable innovation and policy-making remains unfulfilled.

4. Challenges in Bias Prevention

While bias prevention in open data is essential, it is not without challenges. These include:

- **Identifying bias sources:** According to (Hettiachchi et al., 2021), bias can take many forms, such as algorithmic, annotator, sampling, temporal, behavioral, and task design biases, which are difficult to fully identify and understand in complex datasets. More to this, different biases can interact, making it challenging to isolate and address individual sources. Public datasets created for one purpose are frequently used for another, introducing mismatches between the data's original context and its new application, resulting in misleading algorithm performance (Shimron et al., 2022).
- **Quantifying Bias:** Developing reliable metrics and tools to measure bias accurately in annotated data is a significant challenge (Hettiachchi et al., 2021). Existing evaluation metrics often compare processed data (e.g., compressed or interpolated images) without accounting for alterations, masking the true quality of model performance (Shimron et al., 2022).
- **Economic and Practical Constraints:** Implementing bias prevention measures requires investments in training and technology, which usually come at high costs (Johnson & Rostain, 2020). Implementing strategies to prevent bias, e.g., balanced sampling and dynamic task assignment, can increase costs and require more resources (Hettiachchi

et al., 2021). Detecting and mitigating bias requires extensive computational resources for analyzing and validating algorithms across multiple scenarios, which may be prohibitive for some researchers (Shimron et al., 2022).

- Ethical Considerations: Collecting detailed demographic or behavioral information to prevent biases can raise privacy and ethical issues (Hettiachchi et al., 2021). Also, the integration of biased algorithms into decision-making processes, especially in sensitive areas like healthcare or criminal justice, can perpetuate systemic inequalities (Shimron et al., 2022). Balancing inclusivity with privacy concerns is not easy, as it poses significant challenges (Floridi et al., 2021).

This study highlights that even when addressing biases in open data comes with challenges, the long-term value of bias prevention outweighs the costs and hurdles involved.

5. Conclusion

The mitigation of bias is fundamentally essential for preserving the integrity and efficacy of open data. Open data is inherently designed to democratize information access, stimulate innovation, and enable evidence-informed decision-making across various sectors. Nevertheless, the emergence of biases within datasets, whether arising from distorted sampling, inconsistent methodologies, or biased algorithms, diminishes trust, perpetuates disparities and undermines the principles of transparency and inclusivity that open data endeavors to uphold.

This paper emphasizes that the mitigation of biases is a pivotal factor influencing the integrity of open data quality, which extends beyond mere technical precision, integrating dimensions of fairness, inclusivity, and dependability. Through the implementation of inclusive methodologies that prioritize diverse representation, the standardization of data collection and processing techniques to ensure uniformity, and the integration of ethical principles within data governance, stakeholders can proactively address and alleviate bias. Such initiatives not only enhance the quality of data but also improve the credibility and applicability of open data in fostering equitable decision-making. Therefore, preventing biases is not just about ensuring the technical soundness of open data but about affirming its role as a motivation for positive societal transformation. It is a critical step toward achieving data-driven decision-making that reflects the diverse realities of all populations and secures the trust of future generations in the power of open data.

6. References

- i. Aksenova, A., Johny, A., Adams, T., Gibbon, P., Jacobs, M., & Hofmann-Apitius, M. (2024). Current state of data stewardship tools in life science. *Frontiers in Big Data*, 7, 1428568.
- ii. Arend, D., Psaroudakis, D., Memon, J. A., Rey-Mazo'n, E., Schu'ler, D., Szymanski, J. J., Scholz, U., Junker, A., & Lange, M. (2022). From data to knowledge-big, data needs stewardship, a plant phenomics perspective. *The Plant Journal*, 111(2), 335–347.
- iii. Baumann, J., Castelnovo, A., Crupi, R., Inverardi, N., & Regoli, D. (2023). Bias on demand: A modelling framework that generates synthetic data with bias. *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, 1002–1013.
- iv. Borissova, G. (n.d.). Balancing open data and data protection.
- v. Cabral, S. (2024). Cross-sector and public-public collaborations. In *Strategy for public and non-profit organizations: An applied perspective* (pp. 251–284). Springer.
- vi. Dong, S., Sahri, S., & Palpanas, T. (2024). Data quality awareness: A journey from traditional data management to data science systems. *arXiv preprint arXiv: 2411.03007*.
- vii. Fabbri, S., Papadopoulos, S., Ntoutsis, E., & Kompatsiaris, I. (2022). A survey on bias in visual datasets. *Computer Vision and Image Understanding*, 223, 103552.
- viii. Floridi, L., et al. (2021). *Ethics, governance, and policies in artificial intelligence*. Springer.
- ix. Gao, Y., Janssen, M., & Zhang, C. (2023). Understanding the evolution of open government data research: Towards open data sustainability and smartness. *International Review of Administrative Sciences*, 89(1), 59–75.
- x. Hajjaj PhD, M., Burrows, L., Rogers, T., Elias, S. K., Seng, S., et al. (n.d.). Inclusive data management: Reporting, storing, and sharing of information on beneficiaries in the mine action sector. *The Journal of Conventional Weapons Destruction*, 28(1), 6.
- xi. Hettiachchi, D., Sanderson, M., Goncalves, J., Hosio, S., Kazai, G., Lease, M., Schaekermann, M., & Yilmaz, E. (2021). Investigating and mitigating biases in crowdsourced data. *Companion Publication of the 2021 Conference on Computer Supported Cooperative Work and Social Computing*, 331–334.
- xii. Johnson, R. A., & Rostain, T. (2020). Tool for surveillance or spotlight on inequality? Big data and the law. *Annual Review of Law and Social Science*, 16(1), 453–472.
- xiii. Kam, J. W., Badhwar, A., Borghesani, V., Lee, K., Noble, S., Raamana, P. R., Ratnanather, J. T., Tan, D. G., Oestreich, L. K., Lee, H. W., et al. (2024). Creating diverse and inclusive scientific practices for research datasets and dissemination. *Imaging Neuroscience*, 2, 1–14.
- xiv. Laiti, A. (2024). Knowledge sharing in non-governmental organizations: Social and financial education in Tunisia.
- xv. Langbridge, A., Quinn, A., & Shorten, R. (2024). Overcoming representation bias in fairness-aware data repair using optimal transport. *arXiv preprint arXiv: 2410.02840*.
- xvi. Lee, P.-C. (2024). towards a more resilient ecosystem: A case study of open government data in Taiwan. *Journal of Library & Information Studies*, 22(1).
- xvii. Lnenicka, M., Nikiforova, A., Luterek, M., Milic, P., Rudmark, D., Neumaier, S., Santoro, C., Flores, C. C., Janssen, M., & Bol'ivar, M. P. R. (2024). Identifying patterns and recommendations of and for sustainable open data initiatives: A benchmarking-driven analysis of open government data initiatives among European countries. *Government Information Quarterly*, 41(1), 101898.

- xviii. Neves, F. T., de Castro Neto, M., & Aparicio, M. (2020). The impacts of open data initiatives on smart cities: A framework for evaluation and monitoring. *Cities*, 106, 102860.
- xix. Noble, S. U. (2018). *Algorithms of oppression: How search engines reinforce racism*. New York University Press.
- xx. Obermeyer, Z., Powers, B., Vogeli, C., & Mullainathan, S. (2019). Dissecting racial bias in an algorithm used to manage the health of populations. *Science*, 366(6464), 447–453.
- xxi. Okunleye, O. J. (2024). The role of open data in driving sectoral innovation and global economic development. *Journal of Engineering Research and Reports*, 26(7), 222–243.
- xxii. O’Neil, S., Hoe, E., Ward, E., Goyal, R., & Staatz, C. (2020). Data across sectors for health initiative: Promoting a culture of health through cross-sector data networks (tech. rep.). *Mathematica Policy Research*.
- xxiii. Park, S., & Gil-Garcia, J. R. (2022). Open data innovation: Visualizations and process redesign as a way to bridge the transparency-accountability gap. *Government Information Quarterly*, 39(1), 101456.
- xxiv. Rinaldi, E., Stellmach, C., Rajkumar, N. M. R., Carocchia, N., Dellacasa, C., Giannella, M., Guedes, M., Mirandola, M., Scipione, G., Tacconelli, E., et al. (2022). Harmonization and standardization of data for a pan-European cohort on the SARS-CoV-2 pandemic. *NPJ Digital Medicine*, 5(1), 75.
- xxv. Rygge, M. K. (2024). Towards common data environments; measures to improve project management with digital management platforms [Master’s thesis, NTNU].
- xxvi. Salubi, O. (2024). Open data.
- xxvii. Shahbazi, N., Lin, Y., Asudeh, A., & Jagadish, H. (2023). Representation bias in data: A survey on identification and resolution techniques. *ACM Computing Surveys*, 55(13s), 1–39.
- xxviii. Shimron, E., Tamir, J. I., Wang, K., & Lustig, M. (2022). Implicit data crimes: Machine learning bias arising from misuse of public data. *Proceedings of the National Academy of Sciences*, 119(13), e2117203119.
- xxix. Velayutham, A., et al. (2021). Overcoming technical challenges and implementing best practices in large-scale data center storage migration: Minimizing downtime, ensuring data integrity, and optimizing resource allocation. *International Journal of Applied Machine Learning and Computational Intelligence*, 11(12), 21–55.
- xxx. Vetro, A., Canova, L., Torchiano, M., Minotas, C. O., Iemma, R., & Morando, F. (2016). Open data quality measurement framework: Definition and application to open government data. *Government Information Quarterly*, 33(2), 325–337.
- xxxi. Wesson, P., Hswen, Y., Valdes, G., Stojanovski, K., & Handley, M. A. (2022). Risks and opportunities to ensure equity in the application of big data research in public health. *Annual Review of Public Health*, 43(1), 59–78.
- xxxii. Zainuddin, Z., & Akhir, E. A. P. (2024). Systematic literature review of data quality in open government data: Trend, methods, and applications. *IEEE Access*.