



This appendix explores the creation of the Data Ethics Decision Aid (DEDA), a framework designed to evaluate government data projects (Franzke et al., 2021). It takes into account their social implications, embedded values, and the responsibilities of governments during the era of data-driven public management (Franzke et al., 2021). For this particular analysis, which explores the engagement metrics of music videos on YouTube, the following questions shall be addressed and elaborated on: 12, 14, 21, 35/36 and 40.

12. In what ways have you checked the quality of the data?

The following data quality checks have been performed:

- Data validation: after obtaining the datasets, they were qualitatively inspected. For instance, five videos were indexed from the YouTube data to confirm whether they were music videos and if their transcripts matched.
- Duplicate removal: the dataset was scanned for and eliminated any duplicate entries to ensure each observation was unique, especially when merging datasets.
- Handling missing values: missing values were addressed in each column by appropriate methods, such as dropping them.
- Outlier detection: each column was examined for extreme values that could affect the analysis, utilizing techniques like box plots and scatter plots to identify and determine the best approach for handling outliers.
- Ensuring data type consistency: each column has been verified for the correct data type, converting numerical data to float or integer types and categorical data to the appropriate string or category types.
- Intercoder reliability: to examine the data and the models used, 50 randomly selected videos were used to manually code the complexity of language.

Since this test was of an exploratory nature, the process involves trial-and-error. Therefore, ``print()`` and ``head`` were frequently used to verify the code's behavior to ensure that it was inline with the objective. It is essential to acknowledge that while these checks were performed, the dataset could not definitively be examined for biases, since the algorithm of YouTube is considered to be a 'black box' (Arthurs et al., 2018). Therefore, we cannot be certain whether the data is representable for what is portrayed on the platform, as we only know the input and output of the data. In all probability the data retrieved was influenced by the YouTube page-rank algorithm, meaning the recommender system determined what content was relevant to the query.

14. Should the data be made anonymous or pseudo-anonymous? Or generalized?

To safeguard data privacy and confidentiality, the dataset was reduced by eliminating all variables not directly related to the research question. Consequently, the data was anonymized by substituting channel IDs with pseudonymous values. These steps were implemented to avoid the risk of data spillover or breaches, which could lead to unauthorized access to sensitive information (Tucker, 2019).

21. Are there any obligations to (not) publish the data? If you were to provide open access to (parts of) the data, what opportunities and risks might arise?

If data contains sensitive or confidential information, there may be obligations to refrain from publishing it. Legal or ethical considerations might mandate keeping the data confidential to safeguard the privacy of individuals or organizations. In this instance, the dataset consists of raw transcripts of the lyrics extracted from videos, which could be misunderstood if separated from the video context. Therefore, it is crucial to weigh the potential risks and benefits of making the data accessible to the public. While open access to the data could enable other researchers to conduct additional analyses and validate or challenge the original study's findings, it also poses risks such as data breaches or unauthorized use.

35 / 36. Does this project use personal data? If not, continue with 'Bias'. Do the data provide insight into the personal lives of citizens?

This project exclusively utilizes public data obtained from the YouTube API concerning music, ensuring the exclusion of personal data. The collected data encompasses channel IDs, video IDs, and engagement metrics as likes, views and comment counts, with no extraction of comments or user identifications. Additionally, channel IDs have been anonymized to safeguard publisher identities. Consequently, the project does not delve into the personal lives of individuals, and the gathered data remains unrelated to any personal information.

40. Is there a risk that the project could contribute to discrimination against certain people or groups?

While no personal data, like comments or user names, even seemingly impartial algorithms can amplify existing biases and discrimination rooted in factors like race, gender, and socioeconomic status. In this project, the emphasis lies on researching engagement metrics of music videos, which may lead to preferences for particular types of music or creators on the platform. To alleviate this risk, it is crucial to ensure that the analysis and recommendations stem solely from objective data, free from any biases or prejudices.

Bibliography

- Arthurs, J., Drakopoulou, S., & Gandini, A. (2018). Researching YouTube. *Convergence*, 24(1), 3–15. <https://doi.org/10.1177/1354856517737222>
- Franzke, A. S., Muis, I., & Schäfer, M. T. (2021). Data Ethics Decision Aid (DEDA): A dialogical framework for ethical inquiry of AI and data projects in the Netherlands. *Ethics and Information Technology*, 23(3), 551–567. <https://doi.org/10.1007/s10676-020-09577-5>
- Tucker, C. (2018). Privacy, Algorithms, and Artificial Intelligence. In *The Economics of Artificial Intelligence: An Agenda* (pp. 423–437). University of Chicago Press.
<https://www.nber.org/books-and-chapters/economics-artificial-intelligence-agenda/privacy-algorithms-and-artificial-intelligence>