

Datasheet for ‘Children’s Books’ dataset*

Peter Fan

30 November 2024

First sentence. Second sentence. Third sentence. Fourth sentence.

Extract of the questions from @gebru2021datasheets.

Motivation

1. *For what purpose was the dataset created? Was there a specific task in mind? Was there a specific gap that needed to be filled? Please provide a description.*
 - The author of the dataset Alex Cookson created the dataset of children’s books with given ratings by the readers for empirical Bayes estimation of reader’s ratings. The dataset includes all details as per the specific ISBN of the books, whereas the genre of the book and the age of the target audience are not given in the dataset, which is a gap in the dataset. However, no assumptions are made in this study to fill the gap.
2. *Who created the dataset (for example, which team, research group) and on behalf of which entity (for example, company, institution, organization)?*
 - Alex Cookson has created the dataset as a part of his data repository for data-science-related data.
3. *Who funded the creation of the dataset? If there is an associated grant, please provide the name of the grantor and the grant name and number.*
 - N/A
4. *Any other comments?*
 - N/A

Composition

*Code and data are available at: [LINK](#).

1. *What do the instances that comprise the dataset represent (for example, documents, photos, people, countries)? Are there multiple types of instances (for example, movies, users, and ratings; people and interactions between them; nodes and edges)? Please provide a description.*
 - The dataset represents ratings of the readers of children’s books, whereas each instance or record represents a single book with a specific ISBN. .
2. *How many instances are there in total (of each type, if appropriate)?*
 - Total 9240 instances
3. *Does the dataset contain all possible instances or is it a sample (not necessarily random) of instances from a larger set? If the dataset is a sample, then what is the larger set? Is the sample representative of the larger set (for example, geographic coverage)? If so, please describe how this representativeness was validated/verified. If it is not representative of the larger set, please describe why not (for example, to cover a more diverse range of instances, because instances were withheld or unavailable).*
 - The sampling procedure and source of the data have not been mentioned clearly.
4. *What data does each instance consist of? “Raw” data (for example, unprocessed text or images) or features? In either case, please provide a description.*
 - Each record contains character and numerical type data to represent the book title, publisher name, ratings, published year, page number and book cover type
5. *Is there a label or target associated with each instance? If so, please provide a description.*
 - Each instance is label with ISBN it works as primary key of the dataset.
6. *Is any information missing from individual instances? If so, please provide a description, explaining why this information is missing (for example, because it was unavailable). This does not include intentionally removed information, but might include, for example, redacted text.*
 - In the process of data cleaning it was found that in 92 records publish year is missing, 32 ratings are missing, in 206 records publisher name is missing in 293 records number of pages is missing, and in 262 records the type of cover is missing.
7. *Are relationships between individual instances made explicit (for example, users’ movie ratings, social network links)? If so, please describe how these relationships are made explicit.*
 - Each instances are explicit and they are not associated with each other, because each instances represent unique children’s book.
8. *Are there recommended data splits (for example, training, development/validation, testing)? If so, please provide a description of these splits, explaining the rationale behind them.*

- There was no recommendations about data split.
9. *Are there any errors, sources of noise, or redundancies in the dataset? If so, please provide a description.*
- There was no errors, noise or redundancies.
10. *Is the dataset self-contained, or does it link to or otherwise rely on external resources (for example, websites, tweets, other datasets)? If it links to or relies on external resources, a) are there guarantees that they will exist, and remain constant, over time; b) are there official archival versions of the complete dataset (that is, including the external resources as they existed at the time the dataset was created); c) are there any restrictions (for example, licenses, fees) associated with any of the external resources that might apply to a dataset consumer? Please provide descriptions of all external resources and any restrictions associated with them, as well as links or other access points, as appropriate.*
- The dataset is self-contained.
11. *Does the dataset contain data that might be considered confidential (for example, data that is protected by legal privilege or by doctor-patient confidentiality, data that includes the content of individuals' non-public communications)? If so, please provide a description.*
- The dataset contains no confidential information.
12. *Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening, or might otherwise cause anxiety? If so, please describe why.*
- No, the dataset does not contain any such information.
13. *Does the dataset identify any sub-populations (for example, by age, gender)? If so, please describe how these subpopulations are identified and provide a description of their respective distributions within the dataset.*
- No specific age group is mentioned. However, considering the fact the dataset about children's books, it mostly contains ratings of young readers.
14. *Is it possible to identify individuals (that is, one or more natural persons), either directly or indirectly (that is, in combination with other data) from the dataset? If so, please describe how.*
- No
15. *Does the dataset contain data that might be considered sensitive in any way (for example, data that reveals race or ethnic origins, sexual orientations, religious beliefs, political opinions or union memberships, or locations; financial or health data; biometric or genetic data; forms of government identification, such as social security numbers; criminal history)? If so, please provide a description.*
- No

16. *Any other comments?*

- N/A

Collection process

1. *How was the data associated with each instance acquired? Was the data directly observable (for example, raw text, movie ratings), reported by subjects (for example, survey responses), or indirectly inferred/derived from other data (for example, part-of-speech tags, model-based guesses for age or language)? If the data was reported by subjects or indirectly inferred/derived from other data, was the data validated/verified? If so, please describe how.*

- The data was indirectly reported from online book reviews. No verification was done by the source of the data.

2. *What mechanisms or procedures were used to collect the data (for example, hardware apparatuses or sensors, manual human curation, software programs, software APIs)? How were these mechanisms or procedures validated?*

- The data collection process of the author of the dataset is not mentioned. However, this study collected the data from the GitHub data repository.

3. *If the dataset is a sample from a larger set, what was the sampling strategy (for example, deterministic, probabilistic with specific sampling probabilities)?*

- All data was included from the source dataset from the GitHub.

4. *Who was involved in the data collection process (for example, students, crowdworkers, contractors) and how were they compensated (for example, how much were crowdworkers paid)?*

- The author of this study

5. *Over what timeframe was the data collected? Does this timeframe match the creation timeframe of the data associated with the instances (for example, recent crawl of old news articles)? If not, please describe the timeframe in which the data associated with the instances was created.*

- 20th November 2024 to 21th November 2024

6. *Were any ethical review processes conducted (for example, by an institutional review board)? If so, please provide a description of these review processes, including the outcomes, as well as a link or other access point to any supporting documentation.*

- N/A

7. *Did you collect the data from the individuals in question directly, or obtain it via third parties or other sources (for example, websites)?*

- All data was included from the source dataset from the GitHub data repository.
8. *Were the individuals in question notified about the data collection? If so, please describe (or show with screenshots or other information) how notice was provided, and provide a link or other access point to, or otherwise reproduce, the exact language of the notification itself.*
- No
9. *Did the individuals in question consent to the collection and use of their data? If so, please describe (or show with screenshots or other information) how consent was requested and provided, and provide a link or other access point to, or otherwise reproduce, the exact language to which the individuals consented.*
- N/A
10. *If consent was obtained, were the consenting individuals provided with a mechanism to revoke their consent in the future or for certain uses? If so, please provide a description, as well as a link or other access point to the mechanism (if appropriate).*
- N/A
11. *Has an analysis of the potential impact of the dataset and its use on data subjects (for example, a data protection impact analysis) been conducted? If so, please provide a description of this analysis, including the outcomes, as well as a link or other access point to any supporting documentation.*
- No
12. *Any other comments?*
- N/A

Preprocessing/cleaning/labeling

1. *Was any preprocessing/cleaning/labeling of the data done (for example, discretization or bucketing, tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing of missing values)? If so, please provide a description. If not, you may skip the remaining questions in this section.*
- After collection, the dataset was cleaned by removing the missing value and then selecting records with more than the valid number of ratings. Then columns within the variables of interest were selected for final analysis.
2. *Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (for example, to support unanticipated future uses)? If so, please provide a link or other access point to the “raw” data.*
- The raw data was saved in data/01-raw data and the cleaned data was saved in data/02-analysis_data.

3. *Is the software that was used to preprocess/clean/label the data available? If so, please provide a link or other access point.*
 - Only R and packages of R software such as tidyverse, janitor, arrow etc. were used for analysis.
4. *Any other comments?* -N/A

Uses

1. *Has the dataset been used for any tasks already? If so, please provide a description.*
 - The author Alex Cookson used the dataset of children's books with given ratings by the readers for empirical Bayes estimation of reader's ratings.
2. *Is there a repository that links to any or all papers or systems that use the dataset? If so, please provide a link or other access point.*
 - <https://github.com/tacookson/data?tab=readme-ov-file>
3. *What (other) tasks could the dataset be used for?*
 - The dataset could be used for prediction and analysis of children's book ratings using the physical attributes, publishing details of the book.
4. *Is there anything about the composition of the dataset or the way it was collected and preprocessed/cleaned/labeled that might impact future uses? For example, is there anything that a dataset consumer might need to know to avoid uses that could result in unfair treatment of individuals or groups (for example, stereotyping, quality of service issues) or other risks or harms (for example, legal risks, financial harms)? If so, please provide a description. Is there anything a dataset consumer could do to mitigate these risks or harms?*
 - No
5. *Are there tasks for which the dataset should not be used? If so, please provide a description.*
 - No
6. *Any other comments?*
 - N/A

Distribution

1. *Will the dataset be distributed to third parties outside of the entity (for example, company, institution, organization) on behalf of which the dataset was created? If so, please provide a description.*

- No
2. *How will the dataset be distributed (for example, tarball on website, API, GitHub)? Does the dataset have a digital object identifier (DOI)?*
 - N/A
 3. *When will the dataset be distributed?*
 - N/A
 4. *Will the dataset be distributed under a copyright or other intellectual property (IP) license, and/or under applicable terms of use (ToU)? If so, please describe this license and/ or ToU, and provide a link or other access point to, or otherwise reproduce, any relevant licensing terms or ToU, as well as any fees associated with these restrictions.*
 - No, Dataset Link: https://github.com/PeterFan/Predicting_Popularity_of_Childrens_Books/data
 5. *Have any third parties imposed IP-based or other restrictions on the data associated with the instances? If so, please describe these restrictions, and provide a link or other access point to, or otherwise reproduce, any relevant licensing terms, as well as any fees associated with these restrictions.*
 - N/A
 6. *Do any export controls or other regulatory restrictions apply to the dataset or to individual instances? If so, please describe these restrictions, and provide a link or other access point to, or otherwise reproduce, any supporting documentation.*
 - No
 7. *Any other comments?*
 - N/A

Maintenance

1. *Who will be supporting/hosting/maintaining the dataset?*
 - The author of the research
2. *How can the owner/curator/manager of the dataset be contacted (for example, email address)?*
 - Author: Peter Fan, Contact: peteryc.fan@mail.utoronto.ca
3. *Is there an erratum? If so, please provide a link or other access point.*
 - Link: https://github.com/PeterFan/Predicting_Popularity_of_Childrens_Books

4. *Will the dataset be updated (for example, to correct labeling errors, add new instances, delete instances)? If so, please describe how often, by whom, and how updates will be communicated to dataset consumers (for example, mailing list, GitHub)?*
 - No
5. *If the dataset relates to people, are there applicable limits on the retention of the data associated with the instances (for example, were the individuals in question told that their data would be retained for a fixed period of time and then deleted)? If so, please describe these limits and explain how they will be enforced.*
 - No
6. *Will older versions of the dataset continue to be supported/hosted/maintained? If so, please describe how. If not, please describe how its obsolescence will be communicated to dataset consumers.*
 - No
7. *If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for them to do so? If so, please provide a description. Will these contributions be validated/verified? If so, please describe how. If not, why not? Is there a process for communicating/distributing these contributions to dataset consumers? If so, please provide a description.*
 - No
8. *Any other comments?*
 - N/A

1 References