

Transitioning from benchmarks to a real-world case of information-seeking in Scientific Publications - DATASHEET

Chyrine Tahri ♣◇ Aurore Bochnakian ◇ Patrick Haouat ◇ Xavier Tannier ♣
♣ Sorbonne Université, Inserm, Université Sorbonne Paris-Nord, LIMICS, Paris, France
◇ ERDYN, Paris, France

This document is based on *Datasheets for Datasets* by Gebru *et al.* [1].

There are 1811 instances.

MOTIVATION

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 data was created to illustrate a real-world case of information-seeking in scientific papers on a biomedical theme. It was created intentionally to be compared with benchmarks of the same domain.

Who created this dataset (e.g., which team, research group) and on behalf of which entity (e.g., company, institution, organization)?

It was created by Dr. Aurore Bochnakian, an expert in immunology and human health.

What support was needed to make this dataset? (e.g. who funded the creation of the dataset? If there is an associated grant, provide the name of the grantor and the grant name and number, or if it was supported by a company or government agency, give those details.)

The creation of the dataset was supported by ERDYN.

Any other comments?

None.

COMPOSITION

What do the instances that comprise the dataset represent (e.g., documents, photos, people, countries)? Are there multiple types of instances (e.g., movies, users, and ratings; people and interactions between them; nodes and edges)? Please provide a description.

The instances are scientific publications details: titles, abstracts (if available), and PMID, together with a binary relevance judgement (in-scope, out-of-scope) corresponding to the task described in the paper: studying the impact of vitamin B on human health.

How many instances are there in total (of each type, if appropriate)?

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 (e.g., 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 (e.g., to cover a more diverse range of instances, because instances were withheld or unavailable).

All instances are present.

What data does each instance consist of? “Raw” data (e.g., unprocessed text or images) or features? In either case, please provide a description.

All instances consist of raw text.

Is there a label or target associated with each instance? If so, please provide a description.

The label is in-scope/out-of-scope relevance judgement.

Is any information missing from individual instances? If so, please provide a description, explaining why this information is missing (e.g., because it was unavailable). This does not include intentionally removed information, but might include, e.g., redacted text.

Some instances are missing abstracts, these were not available.

Are relationships between individual instances made explicit (e.g., users’ movie ratings, social network links)? If so, please describe how these relationships are made explicit.

None.

Are there recommended data splits (e.g., training, development/validation, testing)? If so, please provide a description of these splits, explaining the rationale behind them.

There are no recommended data splits as all instances represent one illustration case of zero-shot text search.

Are there any errors, sources of noise, or redundancies in the dataset? If so, please provide a description.

None.

Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g., 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 (i.e., including the external resources as they existed at the time the dataset was created); c) are there any restrictions (e.g., licenses, fees) associated with any of the external resources that might apply to a future user? Please provide descriptions of all external resources and any restrictions associated with them, as well as links or other access points, as appropriate.

The data is self-contained for the purpose it was created (relevance search on titles and abstracts). The PMID however can be searched on the PubMed database to retrieve the corresponding paper. We provide no guarantee that these will remain constant over time.

Does the dataset contain data that might be considered confidential (e.g., 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.

None.

Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening, or might otherwise cause anxiety? If so, please describe why.

None.

Does the dataset relate to people? If not, you may skip the remaining questions in this section.

None.

Does the dataset identify any subpopulations (e.g., by age, gender)? If so, please describe how these subpopulations are identified and provide a description of their respective distributions within the dataset.

N/A.

Is it possible to identify individuals (i.e., one or more natural persons), either directly or indirectly (i.e., in combination with other data) from the dataset? If so, please describe how.

N/A.

Does the dataset contain data that might be considered sensitive in any way (e.g., data that reveals racial 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.

N/A.

Any other comments?

None.

COLLECTION

How was the data associated with each instance acquired? Was the data directly observable (e.g., raw text, movie ratings), reported by subjects (e.g., survey responses), or indirectly inferred/derived from other data (e.g., part-of-speech tags, model-based guesses for age or language)? If 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 entirely observable as raw text, except for relevance judgement that was annotated by the expert. The data was collected from PubMed using the following query: ("vitamin B"[Title/Abstract]) AND (health[Title/Abstract] OR growth[Title/Abstract]).

Over what timeframe was the data collected? Does this timeframe match the creation timeframe of the data associated with the instances (e.g., recent crawl of old news articles)? If not, please describe the timeframe in which the data associated with the instances was created. Finally, list when the dataset was first published.

The query was carried in December 2022.

What mechanisms or procedures were used to collect the data (e.g., hardware apparatus or sensor, manual human curation, software program, software API)? How were these mechanisms or procedures validated?

The PubMed database <https://pubmed.ncbi.nlm.nih.gov/> was queried to construct the collection of instances.

What was the resource cost of collecting the data? (e.g. what were the required computational resources, and the associated financial costs, and energy consumption - estimate the carbon footprint. See Strubell *et al.*[2] for approaches in this area.)

Unknown to the authors of the datasheet.

If the dataset is a sample from a larger set, what was the sampling strategy (e.g., deterministic, probabilistic with specific sampling probabilities)?

N/A.

Who was involved in the data collection process (e.g., students, crowdworkers, contractors) and how were they compensated (e.g., how much were crowdworkers paid)?

Authors of the paper were the only individuals involved in the data collection process.

Were any ethical review processes conducted (e.g.,

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.
None.

Does the dataset relate to people? If not, you may skip the remainder of the questions in this section.
None.

Did you collect the data from the individuals in question directly, or obtain it via third parties or other sources (e.g., websites)?
N/A.

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.
N/A.

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.

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.

Has an analysis of the potential impact of the dataset and its use on data subjects (e.g., 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.
None.

Any other comments?
None.

PREPROCESSING / CLEANING / LABELING

Was any preprocessing/cleaning/labeling of the data done (e.g., 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 remainder of the questions in this section.

No, the data was left as raw text as it was collected.

Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g., to support unanticipated future uses)? If so, please provide a link or other access point to the “raw” data.
N/A.

Is the software used to preprocess/clean/label the instances available? If so, please provide a link or other access point.
N/A.

Any other comments?
None.

USES

Has the dataset been used for any tasks already? If so, please provide a description.
Not outside of the paper experimental setup.

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.
None.

What (other) tasks could the dataset be used for?
Other types of retrieval models can be applied to this dataset.

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 future user might need to know to avoid uses that could result in unfair treatment of individuals or groups (e.g., stereotyping, quality of service issues) or other undesirable harms (e.g., financial harms, legal risks) If so, please provide a description. Is there anything a future user could do to mitigate these undesirable harms?

The instances represent the collection retrieved from PubMed in December 2022. Future users shall not refer to these instances exclusively if seeking a more complete information about the theme of vitamin B impact on human health, as new publications are subject to appear meanwhile.

Are there tasks for which the dataset should not be used? If so, please provide a description.
N/A.

Any other comments?
None.

DISTRIBUTION

Will the dataset be distributed to third parties outside

of the entity (e.g., company, institution, organization) on behalf of which the dataset was created? If so, please provide a description.

Yes, the dataset is publicly available and can be downloaded from https://www.erdyn.com/wp-content/uploads/2023/05/relevance_data_vitaminB_health.xlsx.

How will the dataset will be distributed (e.g., tarball on website, API, GitHub)? Does the dataset have a digital object identifier (DOI)?

The dataset is distributed on the creating entity's webpage. It does not have a DOI.

When will the dataset be distributed?

The dataset is released in May 2023.

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.
None.

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.
None.

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.
None.

Any other comments?

None.

MAINTENANCE

Who is supporting/hosting/maintaining the dataset?

The dataset is supported and hosted by ERDYN.

How can the owner/curator/manager of the dataset be contacted (e.g., email address)?

They can be contacted through the following email address: farah.hadj-idris@erdyn.fr

Is there an erratum? If so, please provide a link or other access point.
None.

Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete instances)? If so, please describe how often, by whom, and how updates will be communicated to users (e.g., mailing list, GitHub)?

Unknown to the authors at this time.

If the dataset relates to people, are there applicable limits on the retention of the data associated with the instances (e.g., were 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.

N/A.

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 users.

Unknown to the authors at this time.

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 other users? If so, please provide a description.

Unknown to the authors at this time.

Any other comments?

None.

REFERENCES

- [1] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Dauméé III, and Kate Crawford. Datasheets for Datasets. *arXiv:1803.09010 [cs]*, January 2020.
- [2] Emma Strubell, Ananya Ganesh, and Andrew McCallum. Energy and Policy Considerations for Deep Learning in NLP. *arXiv:1906.02243 [cs]*, June 2019.