

text2sdg: An open-source solution to monitoring sustainable development goals from text

Dominik S. Meier¹, Rui Mata^{1, 2}, and Dirk U. Wulff^{1, 2}

¹ University of Basel ² Max Planck Institute for Human Development

DOI: [10.21105/joss.03988](https://doi.org/10.21105/joss.03988)

Software

- [Review](#) ↗
- [Repository](#) ↗
- [Archive](#) ↗

Editor: [Pending Editor](#) ↗

Submitted: 10 December 2021

Published: 11 December 2021

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](https://creativecommons.org/licenses/by/4.0/)).

Summary

Monitoring progress on the United Nations Sustainable Development Goals (SDGs) is important for both academic and non-academic organizations. Existing approaches to monitoring SDGs have focused on specific data types, namely, publications listed in proprietary research databases. We present the text2sdg R package, a user-friendly, open-source package that detects SDGs in any kind of text data using several scientifically-developed query systems. The text2sdg package facilitates the monitoring of SDGs for a wide array of text sources and provides a much-needed basis for validating and improving extant methods to detect SDGs in text.

Statement of need

The United Nations' Sustainable Development Goals (SDGs) have become an important guideline for both governmental and non-governmental organizations to monitor and plan their contributions to social, economic, and environmental transformations. As the latest UN report (UN, 2021) attests, progress is still needed and the availability of high-quality data will be critical to identify areas requiring most attention moving forward. One promising way to monitor progress on SDGs is to screen the increasing amount of digitally available text using automatized, natural language processing methods. This approach has taken hold, for example, in scientometric efforts that monitor the SDGs in academic publications (e.g., Jayabalasingham et al., 2021). These efforts have so far been spearheaded by for-profit organizations with methodologies that are only partly publicly available and cannot be easily applied beyond academic publishing databases. In what follows, we describe some of these approaches and discuss their shortcomings before introducing our open-source solution, the text2sdg R package, which is designed to help monitor work on the sustainable development goals from any text source.

There are currently five leading approaches to monitoring SDGs from text. The most influential of these was developed by the Elsevier SDG Research Mapping Initiative and uses Lucene-like queries to map several features available from a proprietary database of scientific publications (i.e., Scopus; <https://www.scopus.com>), including most importantly the abstracts of the publications, to SDGs (Jayabalasingham et al., 2021). The Elsevier system has been found to detect millions of SDG-related publications (Agnew et al., 2020) and is used by the Times Higher Education Impact Rankings to rank over 1000 universities worldwide according to their SDG-related research output. Elsevier recently partnered with Aurora, a network of universities that had independently developed a query system to detect SDGs in publications (Vanderfeesten et al., 2020). Other extant query systems include those of OSDG (Bautista, 2019), Siris (Duran-Silva et al., 2019), and the Sustainable Development Solutions Network (Sustainable Development Solutions Network, 2021). Despite the effort put into developing

41 these systems, they are not without shortcomings. First, they are not directly applicable to
42 text sources other than academic citation databases (e.g., Scopus). Second, the systems
43 lack user-friendly and transparent ways to communicate which features are matched to which
44 documents or how they compare between a choice of query systems.

45 We alleviate these shortcomings by providing an open-source solution, `text2sdg`, that lets
46 users detect SDGs in any kind of text using any of the above-mentioned systems or, even,
47 customized, user-made query systems. The package provides a common framework to im-
48 plement the different extant approaches and makes it easy to quantitatively compare and
49 visualize their results.

50 Features

51 We showcase the potential of `text2sdg` with an analysis of a publicly available database
52 of research projects funded by the Swiss National Science Foundation (<https://p3.snf.ch>).
53 To do this, we first applied the packages' main function, `detect_sdg()`, to 26,811 English
54 abstracts written by the research project authors, while setting the `system` argument to
55 `c("Aurora", "SIRIS", "Elsevier", "SDSN", "OSDG")` in order to recruit all five query
56 systems. We then used the packages' `plot_sdg()` function to visualize the frequency of SDG
57 hits and the `crosstab_sdg()` function to analyze the correspondence between query systems
58 and SDGs.

59 Figure 1 shows the results. We highlight three main findings. First and foremost, the re-
60 sults show it is possible to systematically map SDGs to text from sources other than citation
61 databases. Second, as can be seen in panels B-D, the results suggest important and sizable dif-
62 ferences between query systems in both the number of hits and the profiles of most researched
63 SDGs, suggesting it is important to question the results from any single approach. Third, the
64 results suggest that the SDGs vary considerably in their overlap (panel E), emphasizing the
65 promise and challenges of tackling different SDGs simultaneously.

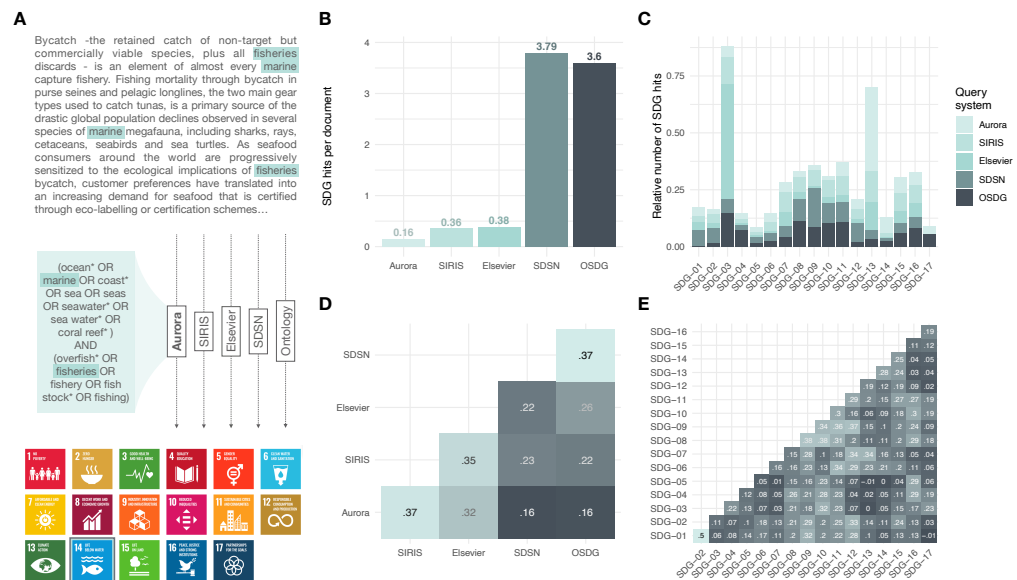


Figure 1: Analysis of 26,811 research projects funded by the Swiss National Science Foundation using the five query systems currently available in `text2sdg`. Panel A illustrates the identification of SDGs based on a query search: The example shows an excerpt of one abstract with some terms highlighted that match a query of the Aurora query system indicating SDG-14 (i.e., conserve and sustainably use the oceans, seas and marine resources for sustainable development). Panel B shows the number of hits per document for the different systems made available in `text2sdg`, which reveals striking differences in numbers of hits. Panel C shows the relative number of hits per SDG cumulatively across systems. Panel D shows the correlations between query systems, which reveal overall small to medium levels of correspondence between them. Panel E shows the correlation between detected SDGs over all query systems.

To facilitate further development of SDG query systems, `text2sdg` additionally includes the `detect_any` function, which permits users to run self-specified queries. We hope this function can be instrumental in future efforts to advance work that validates existing and future query systems and, therefore, can contribute to advancing the sustainable development goals.

Learn more about the background and functionality of the package at <https://dwulff.github.io/text2sdg>. See the package vignette (<https://dwulff.github.io/text2sdg/articles/text2sdg.html>), which includes a reproducible analysis pipeline involving all functions and the package's projects dataset consisting of a random subset of research project abstracts.

References

- Agnew, K., Francescon, D., Martin, R., Rhannam, M., & Schemm, Y. (2020). *The Power of Data to Advance the SDGs. Mapping research for the Sustainable Development Goals*. Elsevier. https://www.elsevier.com/__data/assets/pdf_file/0004/1058179/Elsevier-SDG-Report-2020.pdf
- Bautista, N. (2019). *SDG ontology* [Data set]. <https://doi.org/10.6084/m9.figshare.11106113.v1>
- Duran-Silva, N., Fuster, E., Massucci, F. A., & Quinquilla, A. (2019). *A controlled vocabulary defining the semantic perimeter of Sustainable Development Goals* (Version 1.2) [Data set]. Zenodo. <https://doi.org/10.5281/zenodo.3567769>

- 84 Jayabalasingham, B., Boverhof, R., Agnew, K., & Klein, L. (2021). *Identifying research*
85 *supporting the United Nations sustainable development goals* [Data set]. [https://doi.org/](https://doi.org/10.17632/87txkw7khs.1)
86 [10.17632/87txkw7khs.1](https://doi.org/10.17632/87txkw7khs.1)
- 87 Sustainable Development Solutions Network. (2021). *Compiled list of SDG keywords* [Data
88 set]. <https://ap-unsdsn.org/regional-initiatives/universities-sdgs/>
- 89 UN. (2021). *The Sustainable Development Goals Report 2021*. United Nations. [https:](https://doi.org/10.18356/9789210056083)
90 [//doi.org/10.18356/9789210056083](https://doi.org/10.18356/9789210056083)
- 91 Vanderfeesten, M., Otten, R., & Spielberg, E. (2020). *Search Queries for "Mapping Research*
92 *Output to the Sustainable Development Goals (SDGs)" v5.0* (Version 5.0) [Data set].
93 Zenodo. <https://doi.org/10.5281/zenodo.3817445>

DRAFT