



Max Planck Institute
for Innovation and Competition

Max Planck Institute for Innovation and Competition Research Paper No. 19-03

Michael E. Rose and John R. Kitchen

pybliometrics: Scriptable bibliometrics using a Python interface to Scopus

Max Planck Institute for Innovation and Competition Research Paper Series

pybliometrics: Scriptable bibliometrics using a Python interface to Scopus*

Michael E. Rose^{†1} and John R. Kitchin²

¹Max Planck Institute for Innovation and Competition

²Carnegie Mellon University, Department of Chemical Engineering

Abstract

We present a wrapper for the Scopus RESTful API written for Python 3. The wrapper allows users to access the Scopus database via user-friendly interfaces and can be used without prior knowledge of RESTful APIs. The package provides classes to interact with different Scopus APIs to retrieve information as diverse as citation counts, author information or document abstracts. Files are cached to speed up subsequent analysis. The package addresses all users of Scopus data, such as researchers working in Science of Science or evaluators. It facilitates reproducibility of research projects and enhances data integrity for researchers using Scopus data.

Keywords: Scopus · software · Python · bibliometrics · scientometrics

*An earlier version was entitled "scopus: Scriptable bibliometrics using a Python interface to Scopus".

[†]Corresponding author: Michael.Rose@ip.mpg.de.

1 Motivation and significance

Scopus is one of the citations, bibliometrics and abstracts databases that have become standard in the field of scientometrics and bibliometrics (Falagas et al., 2007; Harzing and Alakangas, 2016; Mongeon and Paul-Hus, 2016; Zeng et al., 2017). A search on Google Scholar in January 2019 for "Scopus" returned about 1,990,000 results, about 24,200 of which date from 2018 alone. For example, a number of recent studies in the field of Economics of Science exemplify the importance of Scopus (Andrikopoulos et al., 2016; Thursby et al., 2018; Catalini et al., 2018; Baruffaldi et al., 2016; Sauermann and Haeussler, 2017; Gush et al., 2018; Heckman and Moktan, 2018).

To ease the use of the Scopus data, we develop `pybliometrics`, a Python package to access the RESTful APIs that Scopus provides. The alternative to our solution is either to download information from Scopus manually, or to use one of the few solutions in languages less prominent than Python. The latter includes solutions in PHP and bash (Montoya et al., 2018), DSpace (Atmire and DSpace developers, 2017) or JavaScript on Apache (Berry and Eprints Scopus Screen plugin developers, 2012). The functionality in Python however puts `pybliometrics` in the hands of scientists who already use Python (e.g. for data analysis, visualization, etc.). The manual download on the other hand has its own limitations. For example, only 2000 records can be downloaded from a search. Manual work is also error prone while an automated process avoids idiosyncratic errors, is more efficient, allows reproducibility and acts as "authoritative definition" of the final dataset (Gentzkow and Shapiro, 2014).

`pybliometrics` comes with a consistent, simple interface, does not require a server, and can easily be integrated with everything in the data science ecosystem of Python, including visualization tools and Machine Learning.

2 Software Description

`pybliometrics` builds on the requests package developed by Reitz and Requests Devel-

opers (2019) and is available from the Python Project Inventory (PyPI). It can be used cross-platform inside a Python interpreter or in script mode.

2.1 Software Architecture

`pybliometrics` implements eight RESTful APIs of the Scopus database (Table 1). These cover the three search APIs (Abstract Search, Author Search, Affiliation Search), the three retrieval APIs (Abstract Retrieval, Affiliation Retrieval, Author Retrieval) and two out of five metadata APIs (Citation Overview API, Serial Title API).

Each class is designed in a consistent way. The required parameter is a search query string (for search classes) or the identifier for Scopus entities (for retrieval and metadata classes) such as authors, documents or affiliations. Only the Citation Overview class requires an additional parameter, namely the start year. Query and retrieval results are cached for subsequent analysis.

2.2 `pybliometrics.scopus` exceptions

`pybliometrics` defines a number of exceptions depending on the type of error received. Though they all inherit from exceptions provided by the `requests` package, they allow error-specific handling.

2.3 Configuration

A configuration file stores two important pieces of information: the authentication credentials and the paths to cache files. Upon first usage, and if the configuration file does not exist, `pybliometrics` guides through the process of creating it. The file is stored in a hidden folder named `".scopus"` in the user's home directory (`~/`).

Table 1: Overview of classes in **pybliometrics**

API	Type	Rest.	Class in pybliometrics	Purpose
Affiliation Search	Search		AffiliationSearch	Search for affiliation entities
Author Search	Search		AuthorSearch	Search for authors
Scopus Search	Search		ScopusSearch	Search for documents
Abstract Retrieval	Retrieval		AbstractRetrieval	Retrieve information on documents
Affiliation Retrieval	Retrieval		ContentAffiliation-Retrieval	Retrieve information on affiliations
Author Retrieval	Retrieval		AuthorRetrieval	Retrieve information on authors
Citations Count Metadata	Metadata	x		Retrieve document citation counts
Citations Overview	Metadata	x	CitationsOverview	Retrieve document citation counts by year
PlumX Metrics	Metadata			PlumX metrics (social media mentions) for documents
Serial Title Metadata	Metadata	x	SerialTitle	Retrieve information on serial titles
Subject Classifications	Metadata			Search for and retrieve Subject Classifications
Author Feedback	Other	x		Provide author feedback (i.e. corrections)

Notes: Table lists currently (January 2019) available RESTful APIs of the Scopus database. "Rest." indicates that the subscriber either needs additional permission from Elsevier for this API (i.e. the used key needs more privileges) or all views of the API are restricted.

2.4 Software functionality

In general, information the user is interested in is stored in properties. Most properties are of object type string or list of namedtuples (e.g. in case multiple results with multiple fields of information are returned). The names of properties are aligned to the field names

in the Scopus documentation.

3 Illustrative Examples

3.1 Automated download of information

Often one is only interested in getting a list of publications and the information Scopus provides on them. The following code provides the time necessary to download all publications in the journal Science in 2010.

Source Code 1: Example code to time the download of retrieving an article list.

```
1 from datetime import datetime
2
3 import pandas as pd
4 from pybliometrics.scopus import ScopusSearch
5
6 # Download
7 start = datetime.now().replace(microsecond=0)
8 s = ScopusSearch("ISSN(0036-8075) AND PUBYEAR IS 2010")
9 end = datetime.now().replace(microsecond=0)
10 print(end-start)
11 # 0:04:29
12 print(len(s.results)) # Number of papers
13 # 2260
14
15 # Reusing
16 start = datetime.now().replace(microsecond=0)
17 s = ScopusSearch("ISSN(0036-8075) AND PUBYEAR IS 2010")
18 end = datetime.now().replace(microsecond=0)
19 print(end-start)
20 # 0:00:00
```

The total number of results (i.e. the number of publications) equals 2260, which is already more than a manual download allows. Depending on download speed, the time to download the information takes less than 5 minutes. However, since `pybliometrics` caches the information, reusing it takes less than second only.

3.2 Creating a collaboration network

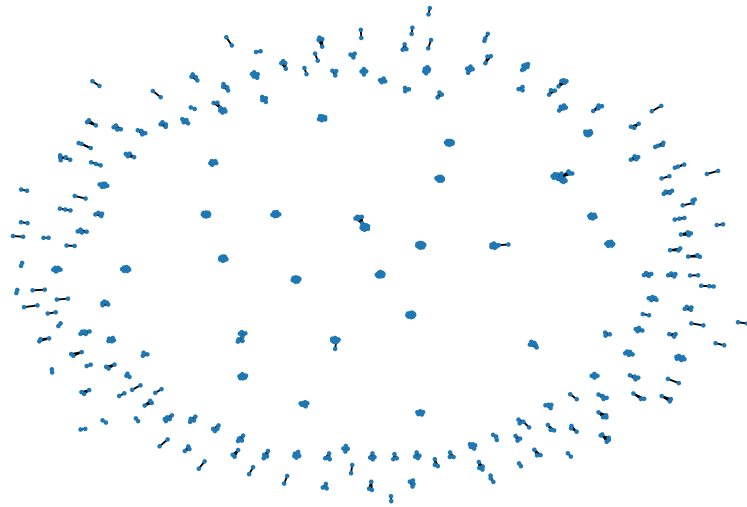
A common task in bibliometric analysis is to create a social network of co-authors. The following example uses NetworkX (Hagberg et al., 2004) (Version 2.3) to create the co-author network for the SoftwareX journal:

Source Code 2: Example code to produce a co-author network using scopus.

```
1 from itertools import combinations
2
3 import networkx as nx # Version 2.3
4 import matplotlib.pyplot as plt # Version 3.0.3
5 from pybliometrics.scopus import ScopusSearch
6
7 # Obtain authors
8 s = ScopusSearch('SOURCE-ID(21100422153)')
9 print(len(s.results)) # Number of papers
10 # 228
11 authors = [i.author_ids.split(';') for i in s.results]
12 # Create list of pairwise combinations
13 combs = [list(combinations(i, 2)) for i in authors]
14 edges = [i for j in combs for i in j]
15 # Create network
16 G = nx.Graph()
17 G.add_edges_from(edges)
18 print(nx.info(G))
19 # Name:
20 # Type: Graph
21 # Number of nodes: 734
22 # Number of edges: 1636
23 # Average degree: 4.4578
24 # Draw and save network
25 nx.draw(G, node_size=2)
26 plt.savefig('network.pdf', bbox_inches='tight', figsize=(50, 50))
```

The resulting graph is depicted in figure 1. The network consists of 734 unique authors that appear on 228 papers. The nodes in the graph are connected through 1636 edges. Because we use author profile IDs provided by Scopus, we do not need to disambiguate author names.

Figure 1: Collaboration network for the SoftwareX journal



Notes: Figure shows a co-author network using publications in the journal SoftwareX until 2019. Each blue dot represents an author. Two authors are connected by an undirected and unweighted black link. Graph layout according to an adopted Fruchterman-Reingold algorithm using NetworkX.

3.3 Analyzing citation distributions

Citations to papers are the foundation of many evaluation exercises. The class `CitationOverview()` provides access to the (restricted) Citation Overview API. The following exercise gathers yearly citation counts for three publications until early 2018 cited in the Nobel decision of Economist Richard H. Thaler. We use pandas (McKinney, 2010) (Version 0.24.2) and seaborn (Waskom et al., 2018) (Version 0.9.0). Additionally, we perform a two-sided Kolmogorov-Smirnov test for the similarity of citation distributions using scipy (Jones et al., 2001) (Version 1.2.1).

Source Code 3: Example code to analyze citation distributions using scopus.

```
1 import matplotlib.pyplot as plt # Version 3.0.3
2 import pandas as pd # Version 0.24.2
3 import scipy # Version 1.2.1
4 import seaborn as sns # Version 0.9.0
5 from pybliometrics.scopus import AbstractRetrieval, CitationOverview
6
7 # Save yearly citation counts in DataFrame
8 data = pd.DataFrame()
```



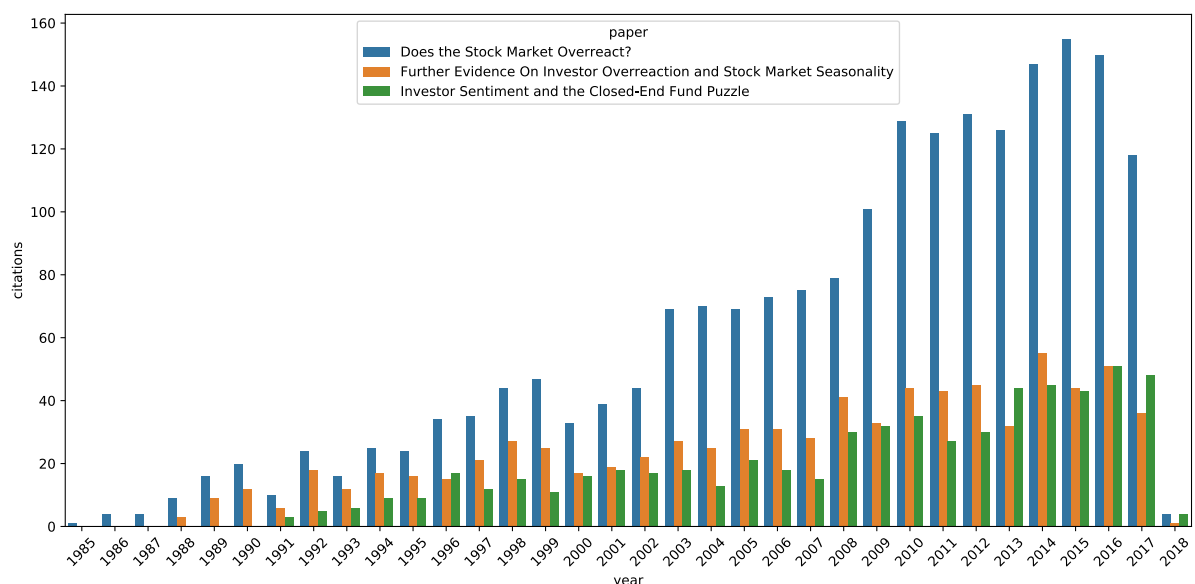
```

9 for eid in ("2-s2.0-84900013243", "2-s2.0-84977703147", "2-s2.0-84977736029"):
10     ab = AbstractRetrieval(eid)
11     year = ab.coverDate[:4] # Year from format YYYY-MM-DD
12     co = CitationOverview(eid, start=year)
13     data = data.append(pd.Series(dict(co.cc), name=ab.title))
14 # Prepare data for plotting
15 data = data.T.astype(float)
16 data.index.name = "year"
17 data = data.reset_index()
18 # Plot and save
19 fig, ax = plt.subplots(figsize=(15, 7))
20 melted = data.melt(id_vars="year", var_name="paper", value_name="citations")
21 sns.barplot(data=melted, x="year", y="citations", hue="paper", ax=ax)
22 plt.savefig("citations.pdf", bbox_inches="tight")
23 # KL-similarity of distributions
24 scipy.stats.ks_2samp(data.iloc[:,0], data.iloc[:,1])
25 # Ks_2sampResult(statistic=1.0, pvalue=2.922351375337816e-16)

```

The resulting citation distribution for the three papers is shown in figure 2. With a p -value $< 2.9 \times 10^{-16}$, the test rejects the null hypothesis that the citation distributions for the first and the second paper are the same.

Figure 2: Citation distribution for three selected papers



Notes: Figure shows the yearly citation count of three selected papers until early 2018.

3.4 Topic modelling using abstracts

A final exercise connects to the emerging field of machine learning and natural language processing. Researchers interested in understanding research topics can use `pybliometrics`, for instance, to obtain large amounts of scientific abstracts quickly and analyze them with specialized software. In source code 4 we find topics among a corpus of scientific articles using Latent Dirichlet Allocation. We use `nltk` (Bird et al., 2009) (Version 3.4.3) and `sklearn` (Pedregosa et al., 2011) (Version 0.20.2). We use abstracts of 226 papers published in SoftwareX until June 2019 and assume 3 topics.

Source Code 4: Example code for topic modelling with LDA where scopus provides abstracts of published papers.

```
1 from string import digits, punctuation
2
3 import nltk # Version 3.4.3
4 import pandas as pd # Version 0.24.2
5 from pybliometrics.scopus import ScopusSearch
6 from sklearn.feature_extraction.text import CountVectorizer as CV
7 from sklearn.decomposition import LatentDirichletAllocation as LDA
8
9 # Obtain abstracts
10 s = ScopusSearch('SOURCE-ID(21100422153)') # Retrieve and cache query results
11 abstracts = [i.description for i in s.results if i.description]
12 print(len(abstracts)) # Number of abstracts
13 # 226
14 # Stem words
15 stemmer = nltk.snowball.SnowballStemmer('english')
16 def tokenize_and_stem(text):
17     text = text.translate(str.maketrans({p: "" for p in punctuation + digits + "@"}))
18     return [stemmer.stem(t) for t in nltk.word_tokenize(text.lower())]
19 # Vectorize abstracts
20 params = {"stop_words": "english", "tokenizer": tokenize_and_stem, "ngram_range": (1, 2),
21          "max_df": 0.6}
22 vectorizer = CV(**params)
23 matrix = vectorizer.fit_transform(abstracts)
24 terms = vectorizer.get_feature_names()
25 print(len(terms)) # Number of terms
26 # 17536
27 # Perform LDA
28 lda = LDA(n_components=3).fit(matrix)
29 # Print ten most important terms per topic
```

```

30 for topic in lda.components_:
31     text = ", ".join(terms[i] for i in topic.argsort()[::-10:-1])
32     print(text)
33 # model, provid, function, simul, calcul, comput, librari, implement, user
34 # data, softwar, model, tool, develop, applic, method, user, provid
35 # softwar, data, develop, model, provid, analysi, simul, process, code

```

The first topic covers technical terms used to describe the presented software, the second topic includes terms that present the software while the third topic seems to be about terms authors use to describe what the software can be used for.

4 Impact

`pybliometrics` has several features that make it optimal for scientists interested in using Scopus data. It is not necessary to understand RESTful APIs or to know how to parse XML or json. The `pybliometrics` code is open-source. The only requirement is to have some understanding of Python. Due to its availability on PyPI, `pybliometrics` is widely available.

Using `pybliometrics`, scientists working in the field of Science of Science (Fortunato et al., 2018) (alternative names include Economics of Science (Stephan, 2010) and Sociology of Science) can increase both the reproducibility and exactness of their projects. Reproducibility increases because users obtain data in the same way from the same source and it becomes transparent where the data originates from and how it was obtained (Sandve et al., 2013).¹ Exactness increases because it becomes very easy to integrate updated data into the analysis.

¹It should be noted though that using data from a database where additions, deletions and corrections are frequent, limit reproducibility per se.

5 Conclusions

The capability of `pybliometrics` to help researchers in bibliometric studies is based on the consistent, simple interface automating all the data retrieval, json serizalization, caching and error handling. It allows researchers to speed up their analyses in reproducible ways and science evaluators to automate the retrieval of evaluation data.

Acknowledgments

We thank the various contributors listed at <https://github.com/scopus-api/scopus-graphs/contributors>. For help preparing this paper we thank Stefano Baruffaldi, Felix Pöge and Nurzhan Sapargali.

References

- Andrikopoulos, A., Samitas, A. and Kostaris, K. (2016), ‘Four decades of the Journal of Econometrics: Coauthorship patterns and networks’, *Journal of Econometrics* **195**(1), 23–32.
- Atmire and DSpace developers (2017), DSpace Elsevier patch, mimeo.
URL: <https://atmire.github.io/Elsevier/>
- Baruffaldi, S. H., Visentin, F. and Conti, A. (2016), ‘The productivity of science & engineering PhD students hired from supervisors’ networks’, *Research Policy* **45**(4), 785–796.
- Berry, R. and Eprints Scopus Screen plugin developers (2012), Eprints Scopus Screen plugin, mimeo.
- Bird, S., Loper, E. and Klein, E. (2009), *Natural Language Processing with Python*, O’Reilly Media Inc.
- Catalini, C., Fons-Fosen, C. and Gaulé, P. (2018), ‘How Do Travel Costs Shape Collaboration?’, *NBER Working Paper Series* **24780**.

- Falagas, M. E., Pitsouni, E. I., Malietzis, G. A. and Pappas, G. (2007), ‘Comparison of PubMed, Scopus, Web of Science, and Google Scholar: strengths and weaknesses’, *The FASEB Journal* **22**(2), 338–342.
- Fortunato, S., Bergstrom, C. T., Börner, K., Evans, J. A., Helbing, D., Milojević, S., Petersen, A. M., Radicchi, F., Sinatra, R., Uzzi, B., Vespignani, A., Waltman, L., Wang, D. and Barabási, A.-L. (2018), ‘Science of science’, *Science* **359**(6379), eaao0185.
- Gentzkow, M. and Shapiro, J. M. (2014), Code and Data for the Social Sciences: A Practitioner’s Guide, mimeo March.
URL: <https://web.stanford.edu/~gentzkow/research/CodeAndData.pdf>
- Gush, J., Jaffe, A., Larsen, V. and Laws, A. (2018), ‘The effect of public funding on research output: the New Zealand Marsden Fund’, *New Zealand Economic Papers* **52**(2), 227–248.
- Hagberg, A., Schult, D., Swart, P. and Networkx developers (2004), NetworkX, mimeo.
URL: <https://networkx.github.io/documentation/stable/index.html>
- Harzing, A.-W. and Alakangas, S. (2016), ‘Google Scholar, Scopus and the Web of Science: a longitudinal and cross-disciplinary comparison’, *Scientometrics* **106**(2), 787–804.
- Heckman, J. and Moktan, S. (2018), ‘Publishing and Promotion in Economics: The Tyranny of the Top Five’, *NBER Working Paper Series* **25093**.
- Jones, E., Oliphant, T., Peterson, P. and SciPy developers (2001), SciPy: Open Source Scientific Tools for Python, mimeo.
URL: <http://www.scipy.org/>
- McKinney, W. (2010), Data Structures for Statistical Computing in Python, in S. van der Walt and J. Millman, eds, ‘Proceedings of the 9th Python in Science Conference’, pp. 51–56.
- Mongeon, P. and Paul-Hus, A. (2016), ‘The journal coverage of Web of Science and Scopus: a comparative analysis’, *Scientometrics* **106**(1), 213–228.
- Montoya, F. G., Alcaide, A., Baños, R. and Manzano-Agugliaro, F. (2018), ‘A fast method for identifying worldwide scientific collaborations using the Scopus database’, *Telematics and Informatics* **35**(1), 168–185.

- Pedregosa, F., Michel, V., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Vanderplas, J., Cournapeau, D., Varoquaux, G., Gramfort, A., Thirion, B., Grisel, O., Dubourg, V., Passos, A., Brucher, M., Perrot and Duchesnay, E. (2011), ‘Scikit-learn: Machine Learning in Python’, *Journal of Machine Learning Research* **12**, 2825–2830.
- Reitz, K. and Requests Developers (2019), Requests: HTTP for Humans, mimeo, Cary, NC.
- Sandve, G. K., Nekrutenko, A., Taylor, J. and Hovig, E. (2013), ‘Ten Simple Rules for Reproducible Computational Research’, *PLoS Computational Biology* **9**(10), e1003285.
- Sauermann, H. and Haeussler, C. (2017), ‘Authorship and contribution disclosures’, *Science Advances* **3**(11), e1700404.
- Stephan, P. E. (2010), The Economics of Science, in ‘Journal of Economic Literature’, Vol. 34, pp. 217–273.
- Thursby, J. G., Haeussler, C., Thursby, M. C. and Jiang, L. (2018), ‘Prepublication disclosure of scientific results: Norms, competition, and commercial orientation’, *Science Advances* **4**(5), eaar2133.
- Waskom, M., Botvinnik, O., O’Kane, D., Hobson, P., Ostblom, J., Lukauskas, S., Gempert, D. C., Augspurger, T., Halchenko, Y., Cole, J. B., Warmenhoven, J., Ruiter, J. d., Pye, C., Hoyer, S., Vanderplas, J., Villalba, S., Kunter, G., Quintero, E., Bachant, P., Martin, M., Meyer, K., Miles, A., Ram, Y., Brunner, T., Yarkoni, T., Williams, M. L., Evans, C., Fitzgerald, C., Brian and Qalieh, A. (2018), ‘mwaskom/seaborn: v0.9.0 (July 2018)’.
- Zeng, A., Shen, Z., Zhou, J., Wu, J., Fan, Y., Wang, Y. and Stanley, H. E. (2017), ‘The science of science: From the perspective of complex systems’, *Physics Reports* **714–715**, 1–73.