

Biblum: A Python-Based Framework for Advanced Bibliometric Analysis

Authors:

Lan Umek, Dejan Ravšelj

Presenter:

Lan Umek

lan.umek@fu.uni-lj.si

Where is Slovenia?



Some data about Slovenia

- > 2.1 M inhabitants
- Largest city: Ljubljana (282,994 inhabitants)
- Area: 20,273 km²
- Slovenian language (slovenščina)
- Dual form (between singular and plural)



Landscape of Slovenia



University of Ljubljana

- established in 1919
- 40,000 students
- 300 study programs
- 23 faculties, 3 arts academies
- 6,000 higher education teachers, researchers, assistants and administrative staff
- listed amongst the top 500 universities in many lists



Faculty of Public Administration

- one of the youngest members of the University of Ljubljana as a fully faculty
- established in 1956 as School of Public Administration
- since 2003 the Faculty of (Public) Administration
- 3 undergraduate and 2 master study programs, 1 PhD programme
- app. 1,000 students
- 70 higher education teachers, researchers, assistants and administrative staff
- EAPAA accreditation - European Association for Public Administration Accreditation

Motivation

- build a comprehensive Python library for bibliometric analysis that reflects and extends the capabilities of R library bibliometrix
- develop and implement methods for (sub)group analysis in bibliometric research
- link prediction models (classification, regression) to bibliometric analysis

About biblum

- Available at (from June 1st)

<https://github.com/lan-umek/biblum>

- Main classes:

- BiblioStats, BiblioAnalysis
- BiblioGroup, BiblioGroupAnalysis
- BiblioGroupClassifier

- Special classes:

- ScienceGroupAnalysis, AreaGroupAnalysis, FiledGroupAnalysis
- SDGGroupAnalysisGoals, SDGGroupAnalysisPerspective, SDGGroupAnalysisDimension

Python libraries used

■ Standard Library

- os
- re
- datetime
- math
- collections
- itertools
- textwrap

Python libraries used

■ Third-Party

- pandas
- numpy
- openpyxl
- scipy
- sklearn
- nltk
- sentence_transformers
- networkx
- matplotlib
- seaborn
- wordcloud
- squarify
- venn
- upsetplot
- holoviews
- pptx
- docx

Bibliometric Analysis

■ Initialization parameters

```
In [1]: import biblium as bb
...: ba = bb.BiblioAnalysis(f_name="dataset.csv", db="scopus", res_folder="results",
...:                        preprocess_level=0, exclude_list_kw=None,
...:                        synonyms_kw=None, lemmatize_kw=False, default_keywords="author",
...:                        asjc_map_df=None, lang_of_docs="en", fancy_output=False, label_docs=True,
...:                        dpi=600, cmap="viridis", default_color="lightblue")
```

■ Simpler version

```
In [1]: import biblium as bb
...: ba = bb.BiblioAnalysis(f_name="dataset.csv", db="scopus")
```

Dataset for demonstration

- Top 200 documents (in terms of global citations) from Scopus related to the query:

TITLE-ABS-KEY(bibliometric*)

Main info

■ Calling

```
In [2]: ba.get_main_info(include=["descriptives", "performance", "time series"], performance_mode="full", stopwords=None,
excluded_sources_references=None)
```

■ Result in xlsx format

Variable	Indicator	Value
	Number of documents	200
	Total citations	126384
	H-index	200
	Average year	2011,15
	G-index	200
	C5	200
	C10	200
	C20	200
	C50	200
	C100	200
	First year	1976
	Q1 year	2007
	Median year	2012
	Q3 year	2017,25
	Last year	2022
	Number of cited documents	200
	A-index	631,92
	R-index	355,5052742
	H(2)-index	29
	W-index	57
	T-index	68,19759931
	D-index	29202

Preprocessing of variables

■ Processing of keywords

```
In [3]: ba.process_keywords(exclude_list=["research", "impact"], synonyms={"bibliometric analysis": "bibliometrics"},  
lemmatize=False)
```

■ Processing of title and abstract

```
In [4]: ba.process_text_vars(stopwords_file=None, lang="en", remove_numbers=True, remove_two_letter_words=True)
```

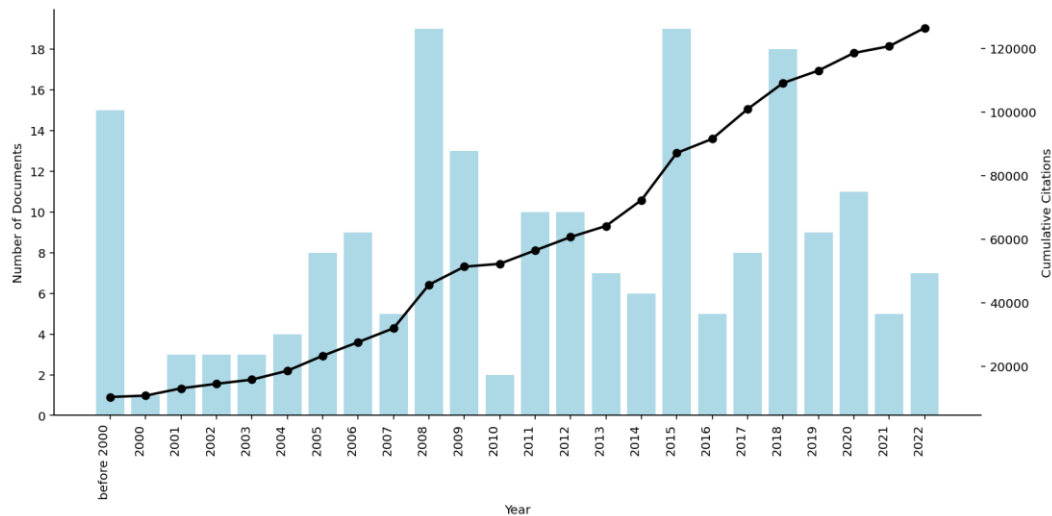
Scientific production

- Computation of production (number of documents, total number of citations)

```
In [2]: ba.get_production(relative_counts=True, cumulative=True, predict_last_year=False)
```

- Plotting

```
In [4]: ba.plot_scientific_production(cut_year=2000)
```



Counting occurrences

■ Initialization parameters

```
In [5]: ba.count_sources() # single type
...: ba.count_author_keywords() # list type
...: ba.count_ngrams_abstract(ngram_range=(1, 2)) # text type
```

■ Result

	Word - Phrase	...	Number of occurrences
15408	research	...	352
805	analysis	...	218
17650	study	...	266
15729	result	...	114
9750	journal	...	286
...
12	1960s term	...	1
11	1960s approximately	...	1
9	1940s schlesinger	...	1
8	1940s	...	1
7	1930s evolved	...	1

[20309 rows x 5 columns]

Performance indicators

■ Initialization parameters

```
In [6]: ba.get_sources_stats(items_of_interest=None, exclude_items=None, top_n=20,  
...: counts_df=None, count_method=None,  
...: regex_include=None, regex_exclude=None,  
...: value_type="string", indicators=False,  
...: missing_as_zero=False, mode="full")
```

■ Result (automatically saved to xlsx format)

```
In [7]: ba.sources_stats_df.columns  
Out[7]:  
Index(['Source', 'Number of documents', 'Total citations', 'H-index',  
       'Average year', 'G-index', 'C5', 'C10', 'C20', 'C50', 'C100',  
       'First year', 'Q1 year', 'Median year', 'Q3 year', 'Last year',  
       'Number of cited documents', 'A-index', 'R-index', 'H(2)-index',  
       'W-index', 'T-index', 'Pi-index', 'Gini index', 'HG-index', 'Chi-index',  
       'Tapered H-index'],  
      dtype='object')
```

LLM description

■ Initialization parameters

```
In [7]: ba.llm_describe_df(ba.sources_counts_df.head(5), provider="openai", model="gpt-3.5-turbo", prompt="Describe the dataframe in one paragraph (with numbers and top sources).")
```

■ Result

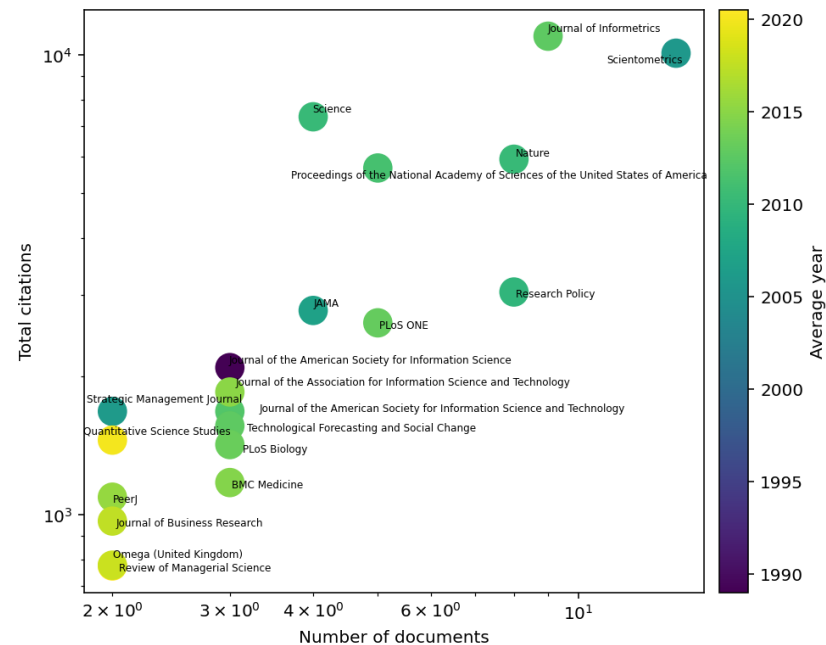
The collection is dominated by a handful of key outlets: Scientometrics contributes the most documents with 14 (7 %), followed by the Journal of Informetrics with 9 (4.5 %). Both Nature and Research Policy each account for 8 documents (4 % apiece), while the Proceedings of the National Academy of Sciences rounds out the top five with 5 documents (2.5 %). Together, these five sources represent nearly a quarter of the total document set.

Plotting top items

■ Example of scatter plot

```
In [8]: ba.plot_top_items_multi(["sources"], kind="scatter",
...: y="Total citations", color_col="Average year",
...: label_col="Source") # many **kwargs possible
...: # different possibilities for different kinds
...: # (scatter, barh, lollipop)
```

■ Result



Co-occurrences

■ Construction of network

```
In [12]: ba.get_author_keyword_cooccurrence(items_of_interest=None, top_n=20, normalization=False, network=True,
...: top_items_df=None, top_items_col=None, count_func=None, count_attr=None,
...: output_attr_prefix=None, value_type="List", separator="; ")
```

■ Setting partitions

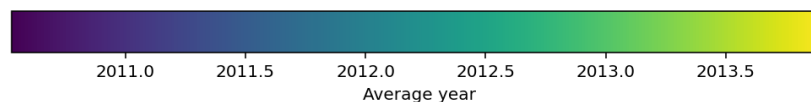
```
In [5]: ba.add_partitions_to_network(louvain_kwargs=None, greedy_kwargs=None, label_kwargs=None,
...: girvan_kwargs=None, k_clique_kwargs=None, kernighan_kwargs=None, walktrap_kwargs=None,
...: edge_betweenness_kwargs=None, infomap_kwargs=None, leading_kwargs=None, leiden_kwargs=None,
...: spinglass_kwargs=None)
```

■ Setting vectors

```
In [6]: ba.add_vectors_from_dataframe(["Average year", "H-index"])
```


Wordcloud

```
In [47]: ba.visualize_text("words from abstract",color_by="Average year")
```



Advanced features

- Topic modelling
- Factorial approach
- Sentiment analysis
- Bibliometric laws: Lotka, Bradford, Zipff
- K-field plot (extension of Sakey diagram 3-field plot)
- Citation network of documents with main path analysis

Bibliometric group analysis

■ Initialization parameters

```
In [7]: bg.BiblioGroupAnalysis(f_name=None, db="", df=None, group_desc=None,  
...: res_folder="results-groups", preprocess_level=0, exclude_list_kw=None,  
...: synonyms_kw=None, lemmatize_kw=False, default_keywords="author",  
...: lang_of_docs="en", fancy_output=False, label_docs=True,  
...: group_colors=None, **kwargs)
```

■ Main difference: group_desc parameter

- a categorical variable
- „list-like“ variable
- binary dataframe
- Year (with n_perdiodes or cut_points)

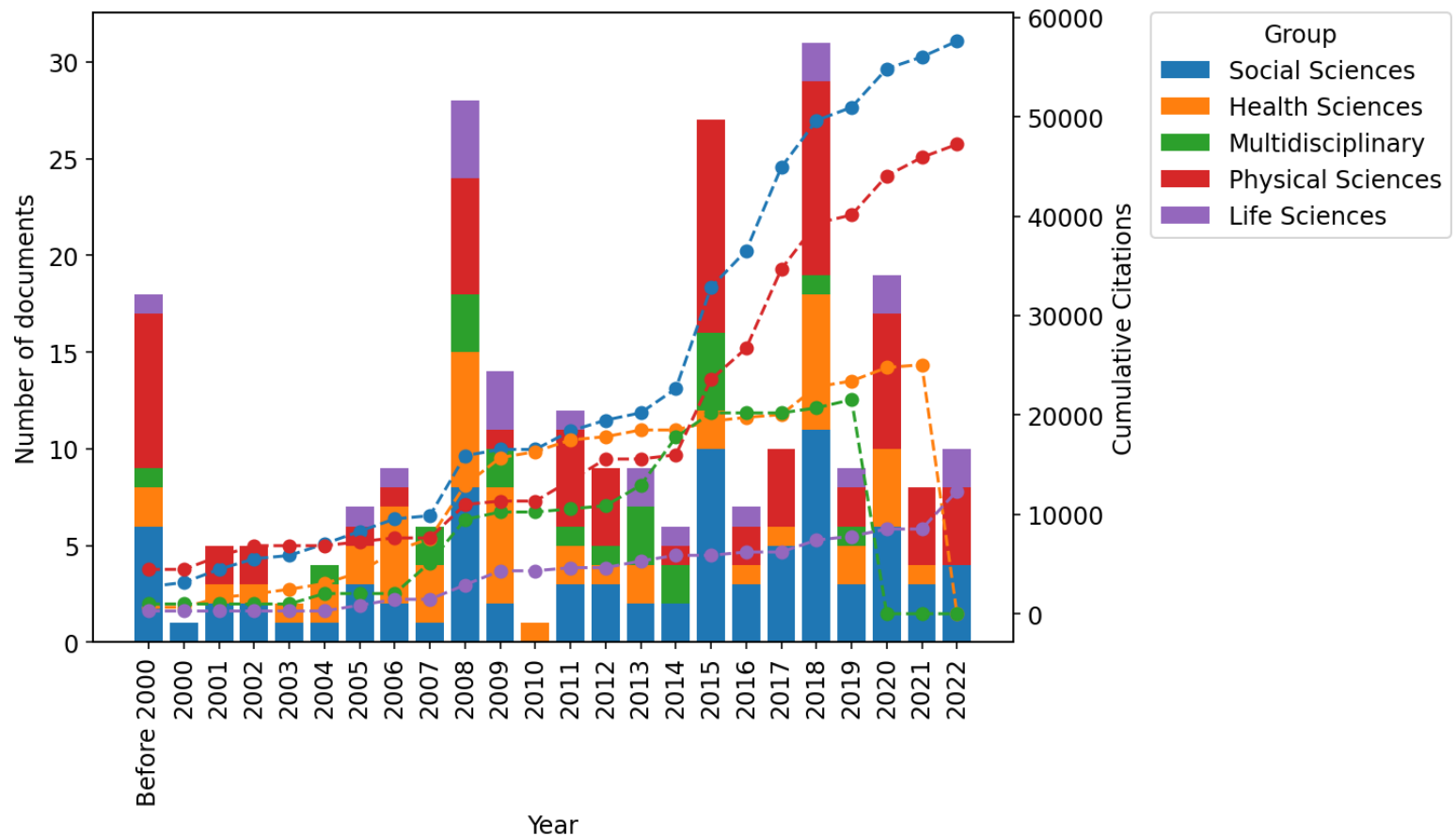
Bibliometric group analysis

- Example of initialization

```
In [16]: ba.add_sciences_scopus()  
...: bg = bb.BiblioGroupAnalysis(db="scopus", df=ba.df, group_desc="Science")
```

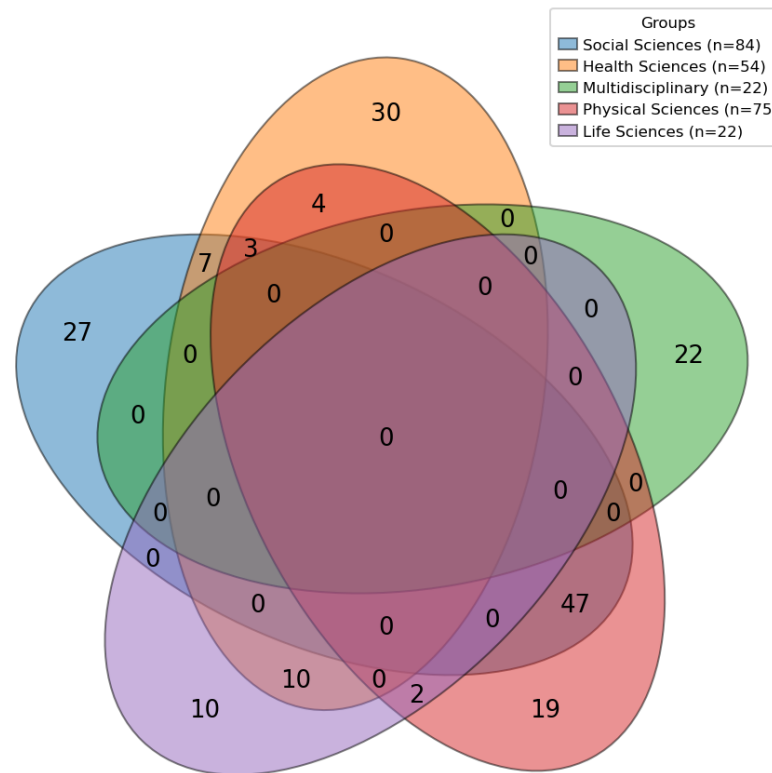
Scientific production

```
In [53]: bg.plot_stacked_production_by_group(cut_year=2000)
```



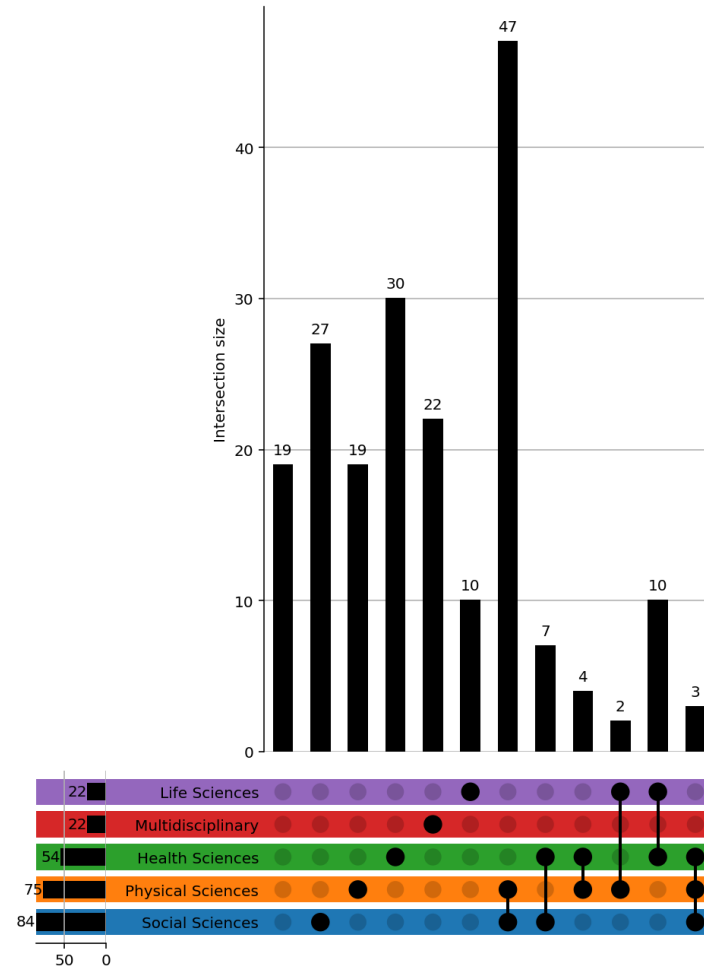
Overlap visualization - Venn

```
In [20]: bg.plot_group_venn()
```



Overlap visualization – upset plot

```
In [23]: bg.plot_group_upset()
```



Counting and performance analysis (in development)

- same idea as for the global analysis
- user can select from long and wide output

Saving the results

```
In [1]: ba.save_reports(formats=["docx", "xlsx", "pptx", "tex"], f_name="bibliometric report")
```

The screenshot shows an Excel spreadsheet with a table of contents. The table is structured as follows:

Level	Sheet Name	Link
Level 1	Main Information	Main Information
Level 1	Global Performances	Global Performances
Level 1	Scientific Production	Scientific Production
Level 2	Sources Counts	Sources Counts
Level 2	Document Types Counts	Document Types Counts
Level 2	Country Counts	Country Counts
Level 2	Authors Counts	Authors Counts
Level 2	Author Keywords Counts	Author Keywords Counts
Level 2	Index Keywords Counts	Index Keywords Counts
Level 2	Affiliation Counts	Affiliation Counts
Level 2	References Counts	References Counts
Level 2	Words Abstract Counts	Words Abstract Counts
Level 2	Words Title Counts	Words Title Counts
Level 3	Sources Stats	Sources Stats
Level 3	Country Stats	Country Stats
Level 3	Authors Stats	Authors Stats
Level 3	Author Keywords Stats	Author Keywords Stats
Level 3	Index Keywords Stats	Index Keywords Stats
Level 3	Words Abstract Stats	Words Abstract Stats
Level 3	Words Title Stats	Words Title Stats

The Excel interface shows the 'Table of Contents' sheet is active, and the bottom navigation bar includes tabs for 'Table of Contents', 'Main Information', 'Global Performances', 'Scientific Production', and 'Sources C...'. The status bar at the bottom indicates 'Ready' and 'Accessibility: Good to go'.

Association analysis (in development)

- analysis of two binary dataframes
 - group_matrix: description of groups
 - df_items: binary dataframe of items of interest
- matrix multiplication provides 2x2 tables
 - compute associations from these tables
 - Jaccard, Yule's Q, Sokal-Michener
 - statistical tests for association
 - Fisher's exact test, chi-squared test

Future work

- develop biblium as pip package
- Tkinter app
- Integration to open source data mining software Orange
- <https://orangedatamining.com/>

Future work

- presentation at ISSI conference (June 2025 in Yerevan, Armenia)
- publication in Scientometric journal
- implementation of subgroup discovery algorithms to biblium