

Projet Python

Putting order to the Scientific Literature

Cyril Verluise

May 30, 2018

- 1 Statement of purpose
- 2 Project implementation
 - Data scraping
 - Data scraping
 - Name Disambiguation and Data Validation
 - Ranking(s)
 - Searching
- 3 Perspectives
- 4 Wrap-up

1 Statement of purpose

2 Project implementation

- Data scraping
- Data scraping
- Name Disambiguation and Data Validation
- Ranking(s)
- Searching

3 Perspectives

4 Wrap-up

Statement of purpose

The World Wide Web creates many new challenges for information retrieval. It is very large and heterogeneous (...) However unlike at document collections the World Wide Web is hypertext and provides considerable auxiliary information on top of the text of the web pages such as link structure (...) we take advantage of the link structure of the Web to produce a global "importance" ranking of every web page
L. Page, S. Brin, R. Motwani and T. Winograd (1999)

Statement of purpose

- Scientific literature is growing fast
- **But:** existing sorting and ranking methods are somehow disappointing.
 - 1 widely-spread scientific indexes such as the *h-index*, *i-index*, etc are mainly based on the number of citations without taking into account the field, or even the quality of the citing articles
 - 2 search engines such as Google Scholar are often disappointing and mostly work as a big articles' indexes

→ mining the scientific literature is "more an art than a science"

Statement of purpose

There are risks:

- Harm the efficiency of scientists unintentionally engaged in research that already existed
- Harm research spillovers of research since non-scientists interested in research such as engineers, operation researchers, students, etc might be unable to find their way to the right existing knowledge

... and hope !

Like the World Wide Web, scientific literature is endowed with a directed link structure (citations) together with meta-information → attempt to use algorithms developed for the WWW to scientific articles and/or authors

1 Statement of purpose

2 Project implementation

- Data scraping
- Data scraping
- Name Disambiguation and Data Validation
- Ranking(s)
- Searching

3 Perspectives

4 Wrap-up

- IDEAS: the largest bibliographic database dedicated to Economics and available freely on the Internet. It indexes over 2 500 000 items of research
- Based on RePEc a large volunteer effort to enhance the free dissemination of research in Economics which includes bibliographic metadata and citations from over 1 900 participating archives, including all the major publishers and research outlets
- Each article is accompanied with a set of useful features: title, author(s), date, journal, keywords, JEL, **AND** citations and references

Structure	Content
<i>root</i>	Homepage
+ article	List of editors' repositories
+ editor	List of journals' repositories
+ journal	List of articles
+ id	Article's pages

Table: Ideas-Repec architecture

Ex: <https://ideas.repec.org/a/oup/qjecon/v1y1886i1p1-27..html>

Focused on Top-30 journals (all time) articles → 82 000 articles. Each article's page was crawled and parsed to collect graph components (title, authors, citations and references) and attributes (date, keywords, JEL-Code)

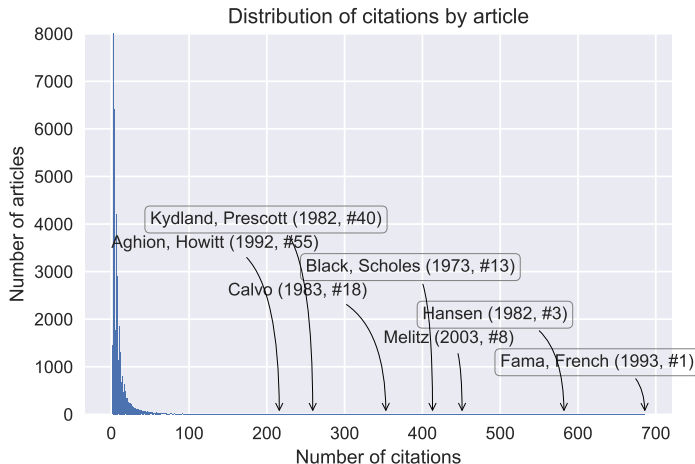
Name Disambiguation

- **Issue:** authors' name have different formats. Ex: Joseph Stiglitz appeared under no less than twenty aliases, sometimes including typos, such as J. Stiglitz, J. E. Stiglitz, Joseph Stiglitz ...
- Author names' disambiguation function, takes two authors as input and declares them similar or not based on how likely they are to be the same
 - 1 Levensthein distance with a threshold to match last names
 - 2 match the first names and the middle names the same way¹

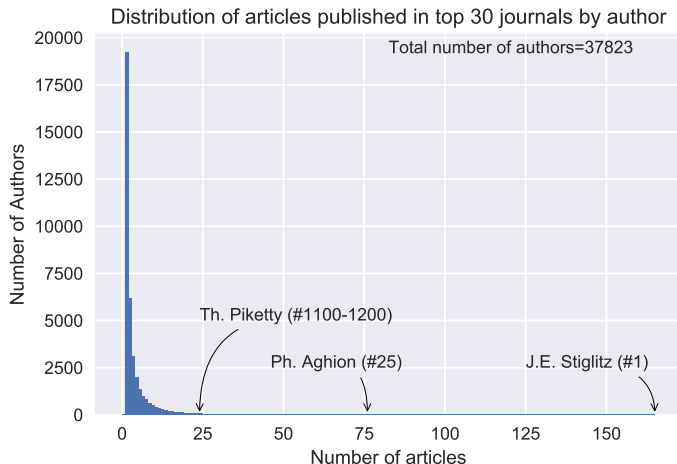
→ from 52 000 to 37 000 different names

¹special treatment of cases when the first name and/or the middle name are only one character or either one does not exist ...

Data validation

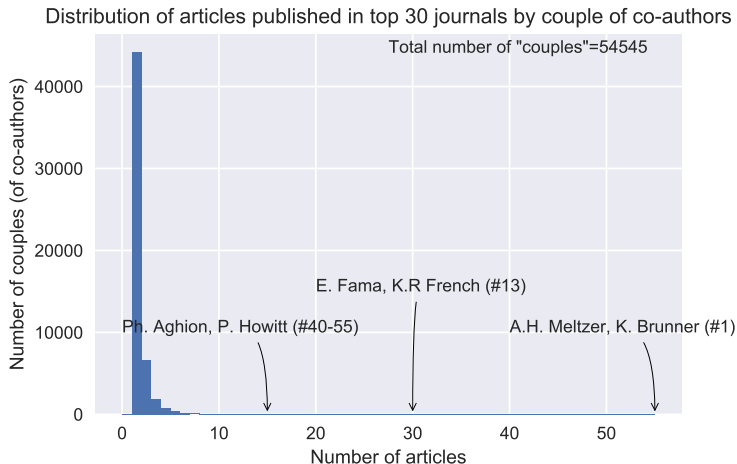


Data validation



Nb: 50% of authors that published at least one article in the top-30 journals published only once in these journals

Data validation



Some words on Rankings

Two "competing" algorithms

- PageRank (Page and Brin, 1998): ranking (\approx random surfer) on whole graph + query filtering/rearranging
 - Graph is potentially huge and G mat is not sparse²
 - But matrix*vector is easy to distribute
 - Once and for all
- HITS (Kleinberg, 1999): query \rightarrow iterative subgraph building \rightarrow Hubs and Authorities
 - Smaller graph
 - Sparse matrix
 - One query, One full computation process

²But can be written efficiently

3 different rankings³

- Number of citations (benchmark)
- PageRank
- Pagerank with time discount⁴
- HITS⁵

³ \approx 50 000 articles have at least 1 citation, ie are ranked

⁴ Accounts for scientific citations specific time structure

⁵ full graph

2 approaches:

- ① "Word filtering" on global ranking
- ② Query specific ranking
 - Natural implementation of HITS (topic and article's similarity)
 - Personalized PaegRank (article's similarity)

- 1 Statement of purpose
- 2 Project implementation
 - Data scraping
 - Data scraping
 - Name Disambiguation and Data Validation
 - Ranking(s)
 - Searching
- 3 Perspectives
- 4 Wrap-up

Avenues for further work

- **Search engine !** (tf-idf, word2vec, ...)
- App deployment (Django, Flask, ...)
- Extend dataset (arXiv, read pdf, ...)
- ...

- Database could prove useful for many other research projects
- CitNet module as well⁶

⁶We provide full documentation and standard installation tools in this goal

Outline

- 1 Statement of purpose
- 2 Project implementation
 - Data scraping
 - Data scraping
 - Name Disambiguation and Data Validation
 - Ranking(s)
 - Searching
- 3 Perspectives
- 4 Wrap-up

Technology used:

- Scraping
- Disambiguation
- Ranking + Network analysis
- Search
- overall Python embedding