Automated Forward Citation Snowballing using Google Scholar and Machine Learning

Hung-Yuan (Vincent) Cheng ^{1, 2}

¹ Department of Population Health Sciences, Bristol Medical School, University of Bristol, UK ² NIHR Bristol Biomedical Research Centre, Bristol, UK

Motivation

- A systematic review summarises the results of available studies and provides a high level of evidence on the effectiveness of healthcare interventions to inform recommendations for healthcare.
- In a systematic review, searching for studies is one of the most crucial steps. Forward citation snowballing (refers to identify new studies based on those papers citing the study being examined) is an effective method to look for new studies and doublecheck.
- Google Scholar is a comprehensive database for forward citation snowballing.

Problem

Google Scholar's searching algorithm and display are not reproducible and transparent. The search results can come up with many irrelevant studies, leading to labour- and time-consuming screening.

Implementation

Machine learning model

- Existing screening results: excluded (n=1076) vs included references (n=45)
- Four separate machine learning models trained by individual inputs: 1) authors; 2) title; 3) journal and 4) abstract from each reference
- Text data processing and cleaning [nltk]
- Data exploration [scattertext]
- Machine learning (ML) model training via Tokenization, Tf-IDF, and RandomForestClassifier [sklearn]
- Accuracy ~ 86% in four models
- Four ML models form a majority rule system to exclude irrelevant references

Google Scholar scraper







NHS

Google Scholar does not provide an easy interface for downloading a large amount of references.

Aims

- **1.** To observe Google Scholar searching algorithm and citation pattern for forward citation snowballing
- To build an automated Google Scholar scraping system
- 3. To filter out irrelevant studies by machine learning

System overview

An automated forward citation snowballing system using Google Scholar search and machine learning



Scrape key data from Google Scholar search results, including Authors, Journal, Title, Abstract, Publication Year, URL, URL of cited, Number of cited [selenium/request-html/urllib]

Analysing data and undertaking meta-analyses JJ Deeks, JPT Higgins... - Cochrane handbook for ..., 2008 - Wiley Online Library iter contains sections titled: Introduction Types of data and effect measures Study designs and identifying the unit of analysis Summarizing effects across studies Heterogeneity Investigating heterogeneity Sensitivity analyses Chapter information ☆ 99 Cited by 2556 Related articles All 3 versions ≫

Automated forward citation snowballing system

- Started with key words, "Lassa fever ribavirin"
- Retrieved up to 10 references from "cited by" link in each potential reference
- Ten iterations:

10

50

120

270

470

668

909

1434

2437

4675

10

47

114

254

414

555

699

944

1134

1384

IT

0

2

3

6

7

8

IT

0

1

2

3

4

5

6

Glossary

IT: iteration

N: Number

Metric

Raw RR =

True RR =

PO: potential

CM: cumulative **DUP: duplicate**



Reflection

- Duplicates occurred more frequently in later iterations, reflecting decreased true retrieval rate and inefficient retrievals.
- After manually screened the retrieved references, useful references appeared more frequently in the first 4 iterations.
- Machine learning model still has room to improve its performance.

Future work

- Automated network plots and other visualisation features
- Tests on a variety of topics and parameters to determine optimal performance
- Other searching databases, such as Microsoft Academy, CrossRef, and Semantic Scholar

Acknowledgement

Thanks to the Jean Golding Institute, University of Bristol for their generous support. Also, thanks to authors behind open source code for contribution and Auto-synthesis group, University of Bristol for inspiration.

